

United States of America

Food and Drug Administration

2014 Next Generation Standards Conference

Bethesda, Maryland

Thursday, September 25, 2014

1 PARTICIPANTS (CONT'D):

2 Next Generation Sequencing Devices and Clinical
3 Applications:

4 ZIVANA TEZAK, PhD, Session Chair

5 JUSTIN ZOOK, PhD

6 HEIKE SICHTIG, PhD

7 JONAS KORLACH, PhD

8 JAY SHAW, PhD

9 KAREN GUTEKUNST, PhD

10 KENDAL DINSMORE

11 IRA LUBIN, PhD
12 FACMG, Team Lead, Genetics, OPHSS Centers for
Disease Control and Prevention

13 JUSTIN ZOOK, PhD
14 Researcher, Genome in a Bottle National Institute
of Standards and Technologies

15 MYA THOMAE
16 Vice President, Regulatory Affairs Illumina

17 HEIKE SICHTIG, PhD
18 Principle Investigator, Regulatory Scientist, CDRH
Food and Drug Administration

19 Food Safety and Pathogen Detection:

20 ERROL STRAIN, PhD, Session Chair

21 HEIKE SICHTIG, PhD

22

1 PARTICIPANTS (CONT'D):

2 MARC ALLARD, PhD
3 Research Microbiologist, CFSAN
4 Food and Drug Administration

5 BILL KLIMKE, PhD
6 Staff Scientist, NCBI
7 National Institutes of Health

8 KRISTIN G. HOLT, DVM, MPH
9 FSIS Liaison to CDC Food Safety and Inspection
10 Service, OPHS
11 U.S. Department of Agriculture

12

13

14 * * * * *

15

16

17

18

19

20

21

22

23

24

25

26

1 P R O C E E D I N G S

2 MS. VOSKANIAN-KORDI: Hello everyone.

3 If everyone could get seated, we're going to go
4 ahead and start.

5 MS. KHAN: Hi. Good morning everyone.

6 Welcome to Day 2 of the conference. I hope
7 everyone was able to stay dry, but you have a lot
8 of, you know, coffee around, across the hall as
9 well as in another café.

10 So basically, today our session on
11 biologics product evaluation is going to cover
12 both the use of NGOs for novel virus discoveries,
13 as well as some development of bioinformatics
14 tools, as well as optimization of bioinformatics
15 pipelines.

16 I'm the first speaker and I'll start the
17 first session. The format of this session,
18 actually, will be a little bit different than the
19 ones you had yesterday. We're going to have four
20 talks, and then after each talk we'll have a
21 couple of minutes for questions that are
22 immediately relevant to the talks.

1 Then we are planning to have a
2 discussion session for the last 30 minutes. That
3 will give opportunity for more audience
4 participation as well as -- I've created some
5 bullet topic points for discussion of the audience
6 can't come up with some hot topics on their own.

7 So anyway, for my talk, before I start
8 I'm going to present the same disclaimer as you've
9 seen for the other FDA speakers that my comments
10 are an informal communication, and these comments
11 do not bind or obligate the FDA.

12 This is a very general outline of my
13 talk. I'm going to start with presenting the need
14 for considering new virus detection methods in
15 CBER, and then present FDA ongoing internal and
16 external efforts for advancing new technologies
17 for evaluation of biologicals. Then try to
18 present the current status of the technologies and
19 further challenges.

20 So the question is why is CBER
21 interested in using NGS, and the main concern in
22 biologicals, safety concern, is still advantageous

1 viruses. This class of agents includes a variety
2 of different types of viruses such as exogenously
3 acquired viruses which can be replicating or
4 latent, genetically inherited viruses such as
5 indigenous retroviruses. Viruses that can be de
6 novo-generated such as some novel recombinant
7 viruses.

8 Virus like particles which may contain
9 co- packaged, unwanted RNA or DNA. As well as, in
10 some cases, the presence of a reverse
11 transcriptase activity that is known to be
12 produced from avian as well as some insect cells.
13 The concern there would be whether this activity
14 is associated with retro viral particles or retro
15 transposons, and if those particles are
16 infectious.

17 So there are strategies that need to be
18 used to mitigate the risk of advantageous viruses,
19 and my office, the Office of Vaccines, uses a
20 four-pronged strategy indicated here. This is to
21 initially identify potential safety concerns do
22 enable development of a comprehensive testing plan

1 and risk mitigation strategy. This includes
2 characterization and history of the cell
3 substrate, and cell banking, and use of qualified
4 raw materials.

5 As you notice that these points are, you
6 know, relevant to cell based products, but some of
7 them can also be relevant to other biological
8 products. For cell banking, it's important to
9 develop well-characterized cell banks, and then to
10 use certified or tested animal drive biological
11 materials such as serum, trypsin, antibodies. I
12 should also point out that in cell-based products
13 these two are the main sources of potential
14 contamination.

15 Additionally, the incorporating of steps
16 during manufacture for viral clearance and purity
17 can be done in the case of some of the inactivated
18 products. This includes process validation for
19 clearance of potential viral contaminants in final
20 product, as well as reduction of residue cellular
21 materials such as RNA, DNA, and proteins.

22 Testing is very important and critical,

1 especially for live viral vaccines where you
2 cannot incorporate inactivation steps or viral
3 clearance steps because the product is a virus.
4 In such cases includes extensive testing for known
5 and unknown agents in the starting materials, in
6 advantageous agent testing, and different steps in
7 manufacturing process, and at steps with the
8 greatest potential for contamination, and using
9 various sensitive and broad detection assays.

10 This slide indicates the conventional
11 assays that are used for evaluating advantageous
12 viruses in vaccine cell substrates as well as in
13 certain other biological products. The general
14 assays include in vivo and in vitro assays,
15 transmission like turn microscopy, as well as a
16 highly sensitive PTR based reverse transcriptase
17 assay for detection of retro viruses.

18 Additionally, specific assays can be
19 used for specie specific viruses based upon the
20 raw materials that you're using, the cell line.
21 As well as sources of potential introduction of
22 viruses, and virus specific assays such as PCR and

1 infectivity assays. These are described in detail
2 in the guidance for industry, the 2010 Cell
3 Substrate guidance that is indicated at the bottom
4 as a reference.

5 In certain cases, there are additional
6 assays, which is recommended on a case by case
7 basis. These, actually, are more relevant to some
8 of the new cell substrates that have been
9 introduced for vaccines. These include
10 oncogenicity assays if a cell substrate is
11 tumorigenic, and to look for tumor inducing
12 viruses, and chemical induction assays, if there
13 are concerns for the presence of latent viruses,
14 such as indigenous retrovirus and latent DNA
15 viruses since such viruses may become activated
16 due to the manufacturing conditions.

17 However, the list of novel cell
18 substrates for vaccines continues to length. Here
19 there are some examples of cell lines from
20 different species that are currently being used.
21 As you can see, there are a variety of species.
22 Some of them, actually, are quite challenging in

1 terms of -- for regulatory review, and especially
2 these cell lines A549 hela CM, which are human
3 tumor derived cell lines.

4 So therefore, the use of such cell
5 lines, and the variety of different cell types and
6 species poses a advantageous agent challenge for
7 the use of novel cell substrates that includes
8 detection of known unknown as well as unexpected
9 viruses. These can be tumor inducing viruses, as
10 I mentioned, the case of tumor cell lines, latent
11 viruses, occult viruses, and novel viruses. So
12 the bar is quite high now, especially for viral
13 vaccines in terms of using technologies that can
14 detect all possible agents.

15 Therefore, the need for considering new
16 technologies for broad novel virus detection.
17 I've listed three here, by no means are these the
18 only ones, but these are the ones that are
19 established and, I guess, available for use
20 currently.

21 Next generation sequencing technologies
22 have demonstrated success in detection of novel

1 viruses, and can be used for both known and novel
2 viruses. There are many different platforms, and
3 the challenges of using these are the topic of
4 this meeting. Additionally, there are micro
5 arrays that are useful for the detection of known
6 and related virus sequences. There's also broad
7 range PCL with mass spec that actually has not
8 been mentioned here, but this uses long PCR
9 primers that are specific for virus families. All
10 of these have demonstrated success in novel virus
11 discoveries as indicated in the last bullet for
12 each one.

13 So some of the potential applications of
14 these technologies for biologicals are for cell
15 substrate investigations and characterization. As
16 you know, it could be a whole genome sequence, if
17 you wish. Transcriptome analysis, as well as
18 analysis of virome on the supernatant.

19 Additionally, these can be potentially
20 used for evaluation of animal derived raw
21 materials to characterize the starting material,
22 as well as for product stability or purity, and

1 other possibilities in the future.

2 So, this slide lists various efforts
3 that have been done in the past, and more recently
4 also, in addressing advantageous virus concerns
5 for vaccines. As you can see, there were PDA
6 meetings from 2010 as well as more recently that
7 addressed advantageous viruses and mitigation
8 strategies earlier. Now more recently, for
9 example, from 2011 and 2013. These have also
10 included discussions on the use of emerging
11 technologies, and the challenges posed by these
12 for use for regulatory applications.

13 Additionally, the IABS is active in
14 discussions on new technologies for advantageous
15 agents, and also, in 2013 with the WHO study group
16 held in Beijing that also raised cell substrate
17 issues and discussions on using DNA based
18 technologies.

19 As you can see, all of these do have
20 proceedings published as well as -- or those that
21 are in progress. In our meeting in 2013 the
22 proceedings are expected to be published in

1 November/December issue this year.

2 There are other group efforts. These
3 include the formation of an advanced virus
4 detection technology users group that was in
5 October 2012. That was a coordinated effort by
6 FDA industry and PDA. The mission of this was to
7 advance the tools for next generation of viral
8 risk evaluation by providing an informal
9 scientific forum for discussions and scientific
10 collaborations.

11 This group does comprise of scientists
12 from various aspects, from industry, from
13 regulatory and other government agencies,
14 academia, and technology service developers, as
15 well as providers. Most recently this year this
16 group has now been converted to a PDA interest
17 group, and the participation, actually, has even
18 broadened. We have about 100 members now. We do
19 have very active discussions and sharing of
20 knowledge, as well as data. So if anyone's
21 interested please feel free to contact me.

22 Additionally, it was mentioned yesterday

1 in Dr. Wilson's talk that in 2013 the FDA
2 genomics working group was formed with a mission
3 to prepare the FDA to address IT and scientific
4 challenges to facilitate its readiness for HTS
5 data review. This is being done through
6 establishment of robust infrastructure as well as
7 through collaborations.

8 So the result of the various meetings
9 was to identify some of the challenges for --
10 common challenges, I should say, shared by
11 regulators as well as by industry for regulatory
12 applications of emerging technologies. This
13 included assay sterilization including assay
14 specificity, and sensitivity using appropriate
15 model viruses as well as sample handling and
16 processing.

17 It included bioinformatics, that was
18 much discussed yesterday, in terms of data
19 analysis, as well as data storage and transfer.
20 In the terms of data analysis, it was recognized
21 that there was a need for guidelines for
22 acceptable quality of Reed's parameters for short

1 reed assembly, approaches to identify a novel
2 virus that has minimal nucleic acid sequencing
3 model to known viruses.

4 As well as a need for a publicly
5 available complete, curated reference virus
6 database, and we've heard about databases
7 yesterday as well. Of course, as was discussed,
8 the need for deciding the format and security
9 measures for data storage and transfer.
10 Additionally, one thing that was not mentioned so
11 far is the need for a follow-up strategy to
12 evaluate a positive signal to determine its
13 biological relevance and significance of the
14 signal.

15 Some of these challenges have been
16 undertaken by the advanced virus detection
17 technologies group. A pilot spiking study for
18 assay generalization was initiated in 2013. The
19 viruses were selected based upon their commercial
20 availability and relevant viruses were selected to
21 encompass different physical and chemical
22 properties, and a human cell line was chosen for

1 background nucleic acid.

2 Three labs are participating in this
3 study, two industry, and my lab, and the FDA, and
4 using at least two different NGS platforms. There
5 are different sample types were created
6 independently based upon shared information
7 regarding virus titer particle count and viral
8 genome copies.

9 There were different types of samples,
10 such as viral nucleic acid spiked into background
11 of cell nucleic acid, and viral particles spiked
12 into cell lysates. The results are currently
13 being obtained and will be shared after internal
14 review and analysis, and publications are planned
15 for broader sharing of these results.

16 The results from these studies, we hope,
17 will provide important information for further
18 development of virus reference materials for NGS
19 technologies. This was also an important point in
20 the questions that were raised yesterday.

21 Another effort that was initiated by the
22 group that is ongoing in my lab is creating a

1 complete virus reference database. The approach
2 that we have undertaken is to download all the
3 Genbank entries, except phage bacteria and
4 contigs. Virus related entries were selected
5 based upon positive key words, and then tagged
6 with non-viral negative words for manual review
7 and decision making.

8 We had a Version 2 database that was
9 internally tested and revised to improve format
10 and data presentation. There's a Version 2.5 that
11 was distributed through the NCBI for initial
12 evaluation by the Advanced Virus Technologies
13 Group. Comments were received, and additional
14 revision is ongoing. Further efforts are directed
15 towards clustering to reduce redundancy, to add
16 viral related sequences from other databases, and
17 planning a web portal for broader distribution and
18 for public use.

19 However, community efforts will be
20 needed for accurate annotation, curation, and
21 continued maintenance of the virus database.
22 That's where we seek all of your help.

1 One thing I just wanted to point out
2 that would be noted is that there's some viruses
3 that are misannotated in the Genbank has host
4 because one of the questions that you might pose
5 is why do we need another database? We have, you
6 know, the NT database, and we have our own private
7 databases.

8 Well, we are trying -- the effort is to
9 have a complete virus database. One of the main,
10 I guess, efforts is to include an accurate and
11 annotate indigenous retroviruses. Based on my
12 expertise in this subject matter I'm able to
13 contribute toward that effort.

14 So some examples are indicated here that
15 are currently taxonomically misannotated. It's
16 highlighted. Like indigenous retroviral DNA from
17 cat is indicated as the host. So if one were to
18 look for viruses under taxonomy, based on taxonomy
19 you will not pick this up.

20 Similarly, for the RK indigenous
21 retrovirus LTR you would not pick that up or a
22 mouse memory tumor virus or a retro transposing

1 like sequence from the ZMAs. So these are just
2 some examples as to why we need to continue
3 efforts to have a much more accurately annotated
4 database for viruses.

5 Additionally, we had posed a follow-up
6 strategy of a positive signal to determine its
7 biological relevance. I won't go through the
8 points here, but basically you can read them. I
9 will just mention we have been successfully
10 applying this proposed strategy in our own lab
11 when we are evaluating cell substrates.

12 I would just like to show this
13 application of this strategy successfully in our
14 recent discovery of a novel rhabdovirus in assay
15 lines cell lines. This was based upon MPS efforts
16 or data from the transcriptome, from nucleic acids
17 from the filtered supernatant, from total nucleic
18 acids that were concentrated from the filtered
19 supernatant. As well as from sized fractionation
20 filtered supernatants.

21 This was done on the SF9 cell line which
22 is from (inaudible) insect, and the fall army

1 worm. This cell line is broadly used for Bactra
2 virus produced products, vaccines as well as
3 therapeutics.

4 The strategy that we use, we were trying
5 to encompass detection of all possible, even
6 distantly related sequences, so the data generated
7 from the MPS was analyzed using blast initially,
8 by using blast N. We used it against the NRNT
9 database. We used blast X against the protein
10 database, and then we found that we have to use T
11 blast X, actually, I'll indicate shortly why,
12 against the NRNT as well as the virus database.

13 Then further then the results evaluated
14 by further binomephomatic analysis. We used the
15 CLC genomics workbench. Then we were able to
16 successfully identify the novel rhabdovirus.

17 This just gives you details of the data
18 obtained from the different samples indicated
19 here. This is from the transcriptome, from the
20 soup, and from the concentrated TFF sample. The
21 contigs from each resulted in a large fragment
22 indicated here, 5 to 6 KBs.

1 However, we were only able to detect
2 this particular fragment. Then had to manually,
3 actually, walk up the genome to assemble the other
4 sequences here which resulted in the assembly of a
5 full length, almost full length, rhabdovirus.
6 That's indicated here.

7 This indicates the struggle in trying to
8 assemble the sequence. Because, basically, the
9 extent of homology or the extent of identity of
10 this whole 13 KB sequence to any known sequence
11 was only within a 295 base pure region in the L
12 protein gene.

13 Even there the identity was only 66
14 percent. The most identity 66 percent in the
15 nucleotide level was found to this task strip
16 virus which is a new member of the
17 mononegrovirus. The 4 percent coverage
18 indicates the 295 base pure region that showed
19 this identity, so the rest of the genome did not
20 have any identity to any known virus.

21 But using Blast P, indicated in this
22 column here, we were able to get 53 percent

1 coverage, but 27 percent identity to the lettuce
2 yellow model virus. You said well, what happened
3 to the taster. Well, it was not in this database
4 because that virus did not have a complete L
5 protein available, so it was not included in the
6 protein database.

7 Our T Blast X results were the most
8 useful because they were able to identify the
9 tastra virus from the NT database with more
10 coverage. The identity was 31 to 47 percent.
11 Therefore, we did recognize that there is a need
12 to use all possible tools if you're looking for
13 novel virus detection.

14 This is just a phylogenetic tree
15 indicating the novel viruses up here. With the
16 tastra virus it is distinct. It's more closely
17 related to plant rhabdoviruses, than to animal
18 rhabdoviruses. This result has been published
19 this year.

20 So this just comes back to our original
21 scheme for virus detection. As I mentioned, T
22 Blast X was useful. However, this was extremely

1 tedious. Literally it takes about two weeks to
2 get the data. This is where we think that having
3 a reference viral database would be very useful
4 because it would weed out all of the cellular
5 background that's currently in the NT database.
6 So we hope that in the future this will make using
7 T Blast X or any other format engines much faster
8 and more specific for virus detection.

9 I just wanted to show a picture by
10 cryoTN that yes, we did see a virus particle, as
11 it is one of the checkpoints for our strategy.
12 Infectivity analysis was done under conditions of
13 mammalian cell culture as well as insect cell
14 culture. We were not able to show replication in
15 an mammalian cell lines. However, there was
16 limited infection in the insect cell lines.

17 These are some additional examples of
18 efforts in the Office of Vaccines. As I
19 mentioned, the effort on the rhabdovirus has been
20 published. We also have looked at other
21 technologies. This is also published on the broad
22 range PCR electro spray ionization mass

1 spectrometry system as well as virus arrays for
2 virus detection.

3 Additionally, there's been efforts in
4 Phil Krause's lab, who's my co-chair for this
5 session, he'll come up after me to introduce the
6 other speakers. They have recently published
7 optimization of virus detection in cells using
8 MPS.

9 Also, in Costa Chumacoff's lab they're
10 using deep sequencing approaches for genetic
11 stability evaluation. They've published on
12 influenza A viruses, and I think Dr. Wilson's
13 talk also referred to their earlier work, and
14 looking at polio virus vaccine stability. In 2010
15 they published a PNA paper on that.

16 Okay. So in conclusion I just want to
17 mention that the new technologies may be useful
18 for supplementing the conventional assays for
19 novel virus detection. I should say these are not
20 my personal conclusions. There are, sort of,
21 conclusions of the field, and the field meaning
22 from the various meetings we've had. These are

1 the generally conclusions.

2 Then the results from ongoing studies
3 can provide scientific direction for additional
4 work that may be needed for advancing new
5 technologies for regulatory application,
6 particularly the development virus reference
7 materials, and that is my own addition here. Due
8 to the complex nature of the new methods and
9 associated bioinformatics, collaborative efforts
10 among industry and regulators should continue to
11 identify and address the technical challenges and
12 efforts toward harmonization. These should be
13 initiated among regulator health agencies sooner
14 than later.

15 I think that this meeting is one of the
16 formats, you know, towards that as well. So I
17 should just acknowledge the people in my lab,
18 Helen Mont, Terry Galvon, and Saya Chiziabod who
19 were involved in the rhabdovirus study. Helen
20 actually spearheaded that study. Aisha Aljahana,
21 he has led the efforts on the virus database
22 development, and Laryn Talifaria on the spiking

1 studies that I described.

2 Rodney Bristor from NCBI has really been
3 extremely helpful in having ongoing discussions,
4 and whatever needs to be done in terms of
5 addressing any comments or concerns about the
6 database. Yesterday he mentioned he'll try to do
7 it as soon as possible if he can't do it the next
8 day. I should say he generally gets it done the
9 same day, so please, you know, if you have any
10 comments do share it with him.

11 Also, the members of the Advanced Virus
12 Detection Technologies group. There are
13 sub-leaders for each of the different topics that
14 we're addressing. Siemon Ng Sanofi leads the
15 efforts on the sample preparation and processing.
16 John Pol Casart on the virus selections and
17 standardization. Kavetha Bekari on the virus
18 database, and Adam Palermo on the optimization of
19 the pipelines.

20 Again, these are just some of the --
21 these are just the sub-group leaders, and, of
22 course, each of them have a lot of other excellent

1 experts involved, and some of them are
2 participants in this conference. Thank you.

3 I'm going to turn the session over to my
4 co-chair Phil Krause now. Oh, questions. We do
5 have time for a couple of questions or maybe not,
6 but we can take them.

7 QUESTIONER: Very nice presentation.
8 Thank you very much for that. I leave with
9 uncertain when you getting to the slides and the
10 conclusion. Basically, in the last slides you
11 mentioned the new technology will be helpful to
12 generate a hypothesis and to suggest new assays,
13 but it's not ready for the regulatory
14 applications. Do you want to elaborate a little
15 bit more on the last slides?

16 MS. KHAN: I think we were trying to do
17 this toward regulatory applications. Is there a
18 mistake on the slide? These efforts are ongoing
19 so that we can, you know, move ahead with these
20 technologies towards regulatory applications since
21 there doesn't seem that they can be, you know,
22 very useful in certain situations. We can

1 definitely discuss it more in the discussion
2 session if I haven't been clear about that.

3 QUESTIONER: So the question is if they
4 are building a gold standard reference database
5 for biogenomes is Dr. Khan's group working closely
6 with the virus pathogen database and analysis
7 resource group? Seems like there's a lot of
8 potential for redundancy if not. That's cited as
9 the IPRBRC.org.

10 MS. KHAN: Basically, I should mention
11 that in terms of calling it a gold standard, I
12 think it will always be work in progress. I think
13 that's for any database. We hope that it will be
14 useful, you know, for the applications that this
15 is intended, which is for biologicals.

16 So basically I am aware that there are
17 many database efforts, and more are being brought
18 to our attention as we discuss them publicly. But
19 I think what I've also heard is that the different
20 databases are maybe useful for specific
21 applications and being created for specific needs.
22 Therefore, it may not be as global as what we

1 want, which is to look for any virus, whether
2 it's, you know, regardless whether it's a pathogen
3 or not. You know, because we want to know if it's
4 infectious, in the end.

5 So even if it's a sequence of a virus
6 we're still interested in that because then the
7 follow-up will be whether if there's a related
8 virus can that virus be infectious. So I think
9 ours is a little bit more global. But definitely
10 we want to coordinate efforts and reduce
11 redundancy if we can. But I'm hoping that once
12 the different database efforts, you know, come
13 nearer to finish then we can be at a point of
14 sharing, you know, the databases, and maybe come
15 up with a more global.

16 Definitely the goal of our database is
17 to work with NCBI, at least in their current --
18 you know, to incorporate our database and theirs
19 together towards, you know, a virus reference
20 database. So hopefully, you know, as long as
21 we're aware and we keep talking about what the
22 databases are and what the goal are, hopefully,

1 when these techniques are being applied, at least
2 initially, we're finding more of these potential
3 positive results that request a huge amount of
4 follow-up, and, perhaps, not as many of the kinds
5 of results that would necessarily be concerning
6 from a safety perspective.

7 So the other speakers in the session are
8 going to address this issue from their own
9 perspectives. I'm very pleased that the next
10 speaker is Charles Chiu from University of
11 California at San Francisco, who's going to talk
12 about next generation sequencing assays for
13 pathogen discovery in the infectious disease
14 diagnostics.

15 He, in particular, has experience using
16 these techniques in the context of a CLIA
17 certified lab where the stakes are very high as
18 well, and so his perspective is going to be very
19 valuable here. Thanks, Charles.

20 MR. CHIU: I just want to start by
21 thanking the organizers, especially Arifa and
22 Phillip for inviting me to speak with you today.

1 I'm going to talk about, continue, sort of, what
2 was brought up by Dr. Khan, and discuss, sort of,
3 our experience with using next generation
4 sequencing.

5 My laboratory spans a wide variety of
6 different, kind of, areas of focus. Part of it is
7 in development of NGS or next generation
8 sequencing based diagnostics. We also do pathogen
9 discovery. As well as, looking for, kind of,
10 advantageous agents. I'll actually give you
11 examples of all of those during this talk.

12 I want to start by just disclosing that
13 I do get research funding from Abbott Diagnostics,
14 from Amazon.com, to establish NGS based pipelines,
15 and from Rubicon Genomics.

16 So the general idea is that we want to
17 use NGS as a way to cast a wide net to detect
18 pathogens. This is opposed to traditional methods
19 like culture and PCR where typically you're only
20 targeting one agent or a limited number of agents.
21 We're just using NGS to really pick up the full
22 spectrum of agents that can cause disease, whether

1 they're viruses, bacteria, fungi, or parasites.

2 You might think that there probably are
3 many different applications or potential use
4 applications by which NGS would be very useful.
5 In clinical diagnostics, and pathogen discovery,
6 and in advantageous agent detection. Our approach
7 is to use, what I'm referring to, as unbiased
8 metagenomic next generation sequencing.

9 The general idea is that this is a
10 method that's -- we're, essentially, just
11 sequencing everything in the clinical sample.
12 We're not using specific primers and probes. The
13 idea is that by sequencing everything in the
14 sample, in theory, you can reconstruct the
15 sequence and assemble the sequence of really any
16 and all pathogens, including viruses, bacteria,
17 fungi, and parasites.

18 Now, part of what's made this
19 method/approach attractive is the fact that now
20 that there are these portable instruments that can
21 be deployed into clinical laboratories, and are
22 actually being deployed. As well as the fact that

1 now you have enormous sequencing capacity such as,
2 you know, 30 to 40 million reads on the miseek,
3 aluminum miseek or more than a billion reads in
4 the hiseek as well as other technologies.

5 In addition, there are not methods to
6 bar code primers so you can multiplex clinical
7 samples to reduce cost per run. However, a big
8 challenge with this whole field is that it's
9 really required the development of both customize,
10 specialized protocols as well as novel
11 bioinformatics analysis tools. I'll be focusing
12 primarily on the bioinformatics analysis here.

13 Our general approach is once you have
14 your raw sequencing data, as you might expect, the
15 vast majority of reads correspond to the hosts.
16 So if you're looking at human samples, clinical
17 samples, probably 99.9 percent of your reads are
18 going to be from human.

19 This is really a needle in a haystack
20 problem. We're not interested in the human
21 fraction. Although, you know, geneticists and
22 oncologists will probably be interested in doing

1 that. We're interested in the fraction that
2 corresponds to potential pathogenic
3 microorganisms.

4 So the general approach is called
5 computational subtraction. We're identifying for
6 reads that corresponds to human. Then we aligned
7 those remaining reads to all sequences in --
8 pathogen specific databases are all sequences in
9 reference databases. This tends to be a huge
10 computational bottleneck. As, kind of Arifa just
11 spoke, after the last talk, this is something that
12 can take, easily, days to weeks using traditional
13 blast based approaches.

14 In addition, the magnitude of the
15 bottleneck actually tends not be highly dependent
16 on the clinical sample type. So if you -- it's a
17 little bit ironic that if you think of tissue,
18 like, say, human tissue where the vast majority of
19 our sequences are human. It's actually quite
20 efficient to just subtract all of the human reads
21 from it. Then you might end up with a small
22 fraction that, hopefully, correspond to the

1 pathogens of interest that you're looking for.

2 On the other hand, take the other
3 extreme, which is a stool sample, which is,
4 essentially, an environment metagenomic sample.
5 Where even if, and this is actual data, where even
6 if you eliminate all human reads, all viral reads,
7 and all bacterial reads you still have more than
8 50 percent of your reads that are presumably
9 something else.

10 They could be unidentified, unsequenced
11 bacteria. They might be environmental sequences.
12 They might correspond to insects or animals, or
13 some components of someone's diet. So this is a
14 big problem because the magnitude of the
15 bottleneck really depends on the clinical sample
16 you're interested in looking at.

17 So to give you some idea of the
18 bioinformatics challenge, this is actually a graph
19 showing, kind of, how long or how expensive it is
20 to analyze 100 million sequences. So if you take,
21 for instance, a single machine, like a, you know,
22 a 16 core server, it takes about three months

1 which is not very conducive to, kind of, short
2 turnaround times for diagnostics.

3 Now, you can, kind of, move up the chain
4 where you can go from a network to a cluster, and
5 then finally to a cloud. On the Amazon cloud, for
6 instance, you can actually recruit up to 5,000
7 cores of data. You can actually get that amount
8 of sequence analyzed by using a blast-based
9 pipeline in about six hours.

10 But the problem with using a cloud then
11 is that the cost rapidly mount out. So this gives
12 you an example. It actually costs almost \$400
13 just to do the analysis in 6 hours by recruiting
14 5,000 cores.

15 So we quickly realized that we needed to
16 find a better way to do this. So we developed a
17 pipeline called SURPI, sequence based ultra-rapid
18 pathogen identification. What it does is actually
19 harnesses two algorithms. One of which was
20 developed from a former post-hoc in my lab who
21 collaborated with UC Berkeley and Microsoft to
22 develop an algorithm called SNAP, which is

1 essentially a rapid or ultra-rapid nucleotide
2 aligner.

3 We also do de novo assembly. We use,
4 kind of a translated nucleotide or protein search
5 to look for, kind of, you know, highly divergent
6 viruses for pathogen discover.

7 I want to talk a little bit about SURPI.
8 Specifically, Dr. Brister actually gave a talk
9 yesterday that was a very good talk about the
10 construction of viral reference database. We
11 actually, for SUPRI, we actually use all of
12 Genbank NT. We're not using all of Genbank NT
13 because we're able to do it. They're actually
14 some good reasons why we're choosing to use
15 Genbank NT instead of, say, a curated viral
16 reference datable.

17 Part of it is that by using that for
18 pathogen discovery it gives us immediate access to
19 all of the newly described partial sequences or
20 gene-only viral sequences that are added to the
21 database. The second thing is that this,
22 actually, provides more efficient viral discovery.

1 What do I mean by that?

2 What I mean by that is by aligning your
3 sequences to everything in Genbank NT you can then
4 quickly identify all sequences that are non-viral.
5 Whether they're human, bacterial, fungal, insect,
6 algae, animals. The idea behind doing that is
7 that then you can take the remaining sequences and
8 very efficiently use a small subset of sequences
9 to do de novo assembly of viruses, and potentially
10 to do faster, kind of, nucleotide or translate
11 nucleotide or protein searches which usually is,
12 kind of, the computational taxing part of the
13 algorithm.

14 Another big issue is that with curated
15 databases such as Refseek there's a loss of
16 specificity. The reason that exists is because
17 viruses are incredibly diverse. The problem with
18 doing that is in order to capture the diversity of
19 viruses with using, I'll say a viral reference
20 database, you need to use a very low stringency
21 cut off.

22 By doing so you then end up with a

1 larger number of false positive. So this is why
2 our solution was just, simply, to use everything
3 in Genbank NT. That has its own problems, and
4 I'll go into that in a little bit.

5 But the advantage of doing that is that
6 we are able to do it because we have an algorithm
7 called SURPI which is incredibly fast. So this
8 gives you an idea. This is SURPI which uses this
9 algorithm called SNAP. Aligning nucleotide
10 alignment to the human database, which is about 3
11 gigabases in size.

12 You can see here it takes SNAP about two
13 minutes to align 25 million reads, and about 7
14 minutes and 30 seconds to align 125 million reads.
15 The comparable times for blast would be 16 and 80
16 days. For a billion sequences it's actually
17 fairly reasonable. It can align a billion
18 sequences in about an hour and a half, and it
19 would take 2.2 years on exact same machine for
20 blast to do that alignment.

21 Because it's so fast we're actually able
22 to also do alignment to NT. I can tell you about

1 10 million reads by aligning to NT which is -- at
2 that time it was about 42 gigabases. It's
3 actually now 70 gigabases, NT. You can actually
4 align 10 million reads to all of NT in about 10
5 minutes. Basically with, kind of, the use of
6 SNAP.

7 So we felt that we've dealt with, kind
8 of, the issue of -- at least partly dealt with the
9 issue of this computational challenge of this
10 bottleneck by implementing something like the
11 SURPI pipeline where we can analyze data sets from
12 10 minutes to 15 hours.

13 Now, 15 hours was kind of a worst-case
14 scenario where we took 100 million reads from a
15 highly complex stool sample. Basically, we had to
16 kind of go through everything. Couldn't really do
17 a computational subtraction with it. You can see
18 here that the costs are now reasonable. So I can
19 get this on the cloud, for instance, I can
20 actually get the same data analyzed, about \$4
21 versus more than nearly \$700 previously.

22 I also want to mention that we've

1 actually published SURPI. We published SURPI
2 about four months ago. But in the meantime, we've
3 added some bioinformatics enhancements to SURPI.
4 One is that, you know, one of the challenges of
5 using all of NT is there are a lot of misannotated
6 sequences in NT.

7 As you might expect, there are sequences
8 that have been misannotated as virtual. So we
9 actually have a rapid filtering algorithm that can
10 take the output, the raw output, and then rapidly,
11 kind of, re annotate, you know, sequences that had
12 been misannotated.

13 We've also included a method that you
14 can do, kind of meta data tagging. You can
15 basically provide additional information in the
16 SAM file that will give you some metadata
17 information. Such as, whether or not this is a
18 clinical pathogen, where is this from, the host,
19 etcetera.

20 You know, I think that Dr. Brister
21 talked a little bit about that as well, about
22 intents to incorporate that with the databases.

1 We've developed, also, a front-end visualization
2 tool for SURPI, a web-based interface. Then we've
3 also used SNAP in a different application. We've
4 actually used it not only for aligning to
5 reference databases, but also for doing taxonomic
6 classification.

7 That actually becomes really, really
8 important because, as you might imagine, viruses
9 have conserve regions and divergent regions. The
10 problem with the conserve regions is that conserve
11 regions basically can represent a wide number of
12 viruses at various taxonomic levels.

13 So, for instance, you can have sapid
14 species, more than 100 different species of
15 viruses. They may share the exact same read or
16 short region, and by doing accurate taxonomic
17 classification you can then really be able to
18 identify exactly what virus you actually have in
19 your sample.

20 In addition, we've also been able to do
21 some enhancements to the memory usage of SNAP, so
22 that now you can, for instance, put all of Genbank

1 NT in memory. We're able to do that. You do need
2 a 1 terabyte memory server, but that's -- actually
3 now, actually it's not -- I just purchase a
4 server. It's about \$8,000. It's a reasonable
5 price, actually, to get a 1 terabyte memory
6 server. That enables you, for instance, to align
7 -- 10,000 reads, everything in NT in about 5
8 minutes versus 40 minutes before in the published
9 paper.

10 This gives you some idea of the
11 front-end visualization tool. So this is actually
12 automated output coming out of the pipeline. So
13 I'm not doing any, sort of, manual, kind of
14 curating, making it look nice. So what it'll do
15 is it will not only identify, kind of, the closest
16 hit. This is actually a clinical sample from a
17 patient with Ebola virus infection. Not from the
18 West African outbreak, from a previous outbreak in
19 the Democratic Republic of the Congo.

20 You can see here that the it basically
21 found that the closest hit was Ebola Virus Zaire.
22 It actually gave you the actual GI corresponding

1 to it. It gives you the pure wise identity. It
2 will take the reads that were generated and plot
3 them across the genome of the virus. It will also
4 give you, what I call a pure wise identity plot.

5 It will tell you, kind of, at a glance,
6 kind of, how difference is this virus from the
7 viral genome in the database. That's a quick way
8 to identify recombination events, etcetera, and
9 also map that against the genome of it's --
10 genomic structure of the virus as from Genbank.

11 The nice thing about this is it will do
12 that for all of the viruses that were detected in
13 your raw data, and it will do it in automated
14 fashion. We're hoping to develop some additional
15 tools that will allow you to edit, kind of, this
16 output that comes out of the pipeline.

17 In addition, we've had the pipeline --
18 also, it can be used as input into what we call
19 Kroma plots. These are taxonomic plots that,
20 essentially, show you, kind of, the distribution
21 of the viruses that were detected at different
22 taxonomic levels.

1 This gives you an example. This is
2 actually a stool sample where you can see that,
3 you know, about 90 percent of the reads were from
4 norovirus. This is a norovirus infection as well
5 as a par echovirus which is kind of a secondary
6 infection in the same sample.

7 In addition, we've been really very
8 interested in developing what I call a clinical
9 work flow in a CLIA laboratory. So we have,
10 actually, a mechanism for actually developing a
11 process for actually generating this data and
12 going from sample to sequence in about 12 hours to
13 2 days. We can take, pretty much any sort of
14 sample type, serum, whole blood, respiratory
15 secretions, cervical, spinal fluid, etcetera, and
16 be able to, kind of, put it through, kind of, a
17 rapid automated pipeline.

18 We've currently been working on some
19 spiking experiments where we're trying to
20 calculate using probative analysis what the limits
21 of detection are. You can see here that for
22 limited detection of HIV as an example of an RNA

1 virus we're getting about less than 100 copies per
2 mil form. This is spiked into plasma.

3 For CMV, which is a DNA virus, it's
4 about 1,000. Although, we are improving that.
5 Actually, the latest data suggests that we are
6 also able to get, kind of, a similar level of
7 sensitivity.

8 This gives you an idea of how we're
9 using the approach to actually diagnose a large
10 variety febrile illnesses. So these are
11 respiratory secretions and serum from patients in
12 the hospital with various diseases. You can see
13 here that we can detect plasmodium falciparum,
14 haemophilus, influenza, rhinovirus, menuavirus,
15 adenovirus, adenovirus senna type 7, salmonella,
16 typhoe, the cause of typhoid fever. These are all
17 clinical samples that we did metagenomic next
18 generation sequencing from.

19 I also want to mention that there also
20 was a part about sample prep. We've been working,
21 also, on ways to actually improve the sensitivity
22 of metagenomic next gen sequencing. We're doing

1 this by combining random application with what we
2 call probase enrichment.

3 On the left you actually see a serum
4 sample from a patient with Hantavirus pulmonary
5 infection from the recent Yosemite outbreak. You
6 can see here that this was actually a case that
7 was PCR negative from the serum. But, if you
8 actually do the probase enrichment you can now
9 detect sequences, specific sequences to
10 Hantavirus.

11 Similarity, on the right is basically a
12 sample from a patient with Zaire Ebola virus
13 infection. You can see that it was unenriched.
14 It's at the limits of detection of PCR with a
15 cycle threshold of 35. You detect very few reads,
16 but then once you enrich it you can actually get,
17 kind of, full genome coverage.

18 So we think this is a viable technique
19 that could be used similar to what's being used in
20 cancer now with cancer gene panels, you know, NGS
21 coupled with pro base enrichment. We think that
22 this is a viable technique to use for this viruses

1 as well.

2 I just want to briefly go over -- end up
3 with like three examples of how we've actually
4 used this technique for various applications. So
5 the first is, basically, identification of
6 contaminants or advantageous agents. We basically
7 reported last year that we found a virus that we
8 called powerful like hybrid virus.

9 We were screening hepatitis serum
10 samples by NGS, and this happens to be a very
11 divergent virus. It only has about, you know, 31
12 percent identity on the amino acid level. We then
13 looked at three different data sets, and we were
14 able to, basically, assemble and identify this
15 virus from three different data sets. If you
16 actually look at -- one of the data sets is
17 actually water, okay?

18 That was concerning because, you know,
19 we were able to assemble the genome of a virus
20 from water. That's a problem. So we quickly
21 looked at that. We actually looked at the
22 phylogenetic analysis of this virus. It happens

1 to be some virus that happens to be between the
2 powerful virada and the circle virada. It happens
3 to be a powerful virus, circle virus hybrid.

4 If you actually look at the data we then
5 went through all of our NGS data, and we were
6 finding this virus everywhere. We found it in,
7 you know, MRSA cultures. We found it in human
8 samples. We found it in animal samples. We found
9 it in stool. We found it in blood. So this was
10 very suspicious.

11 Basically, it was around the same time
12 that it was actually reported that they had found
13 a hybrid DNA virus by a group at the NIH. They
14 actually also reported this in PNAS. They had
15 found a hybrid DNA virus with seronet from
16 patients also with seronegative hepatitis.

17 We did further work and we actually
18 traced that this virus actually could be traced to
19 extraction spin column from a certain
20 manufacturer. We actually found this. I one
21 batch we found the virus in every single spin
22 column from that manufacturer. In fact, we're

1 still actually seeing this virus now in spin
2 columns, so it hasn't been a problem that's been.
3 But if you're doing par viral discovery you
4 probably should take a look for this virus.

5 This actually happens to be the same
6 virus, actually, as the one that was reported in
7 PNES by the NIH which was called NIH-CQV which we
8 now think, if you actually screen environment
9 metagenomics, you screen through this. So I ended
10 up screening through 78 public datasets, about 213
11 million reads. Out of all of those reads, the
12 only reads that had hematology to those viruses
13 were found, actually, in coastal sea water off of
14 California and off of the Pacific Coast.

15 Actually, we have a mechanism, we
16 actually think that this actually happens to be a
17 diatom related virus. It happens to be the silica
18 is contaminated. We think the silica that are
19 typically used as components of most spin columns
20 by manufacturers, okay? So we think that this is,
21 ultimately, the source of this virus which we
22 believe to be an environmental contaminant.

1 I want to go over another quick example.
2 We identified a new hemorrhagic fever virus called
3 the Congo virus. It appeared to be the cause of a
4 very small outbreak cluster of hemorrhagic fever
5 in Africa. I also want to point out, you know,
6 similar to Arifa who had discovered a very diverse
7 rhabdovirus, this was a very diverse rhabdovirus.
8 In fact, it only had about 30 percent amino acid
9 identity.

10 But if you have enough coverage, so
11 basically we generate about 140 million alumina
12 reads. You can actually de novo assemble the
13 whole virus. Even though there is no reference
14 sequence, or at the time, there was no reference
15 sequence in the database.

16 So we were able to assemble about 98
17 percent of the viral genome solely by looking for
18 overlaps and doing de novo assembly. This gives
19 you an idea of the difficulty. So what you see
20 here is in grey, this is the map of the virus, the
21 coverage map of the virus. In grey are parts of
22 the virus that have absolutely no homology to

1 anything else in the database at the time.

2 If you look here in pink, pink were
3 basically regions of the virus that were only
4 identified on the basis of translate nucleotide or
5 protein alignment. It just shows you, kind of,
6 the challenge with viruses because they have so
7 much diversity. You can't rely, necessarily, on
8 nucleotide alignment. You actually have to go to
9 translated nucleotide alignment in order to
10 identify these viruses.

11 I want to end by just one last case
12 which is a case that's actually in press now at
13 CID. This was a 55 year old male who presented --
14 who had a bone marrow transplant that presented
15 with a rapidly developing encephalitis that
16 initially began as hearing loss. He had, kind of,
17 a negative work up, negative cervical spinal
18 fluid, negative MRI scan. He had developed
19 nausea, fatigue over several weeks, and had a
20 negative MRI, but then repeat MRI several months
21 later was positive.

22 He had developed at that time,

1 basically, abnormal signal in his thalamus and his
2 mid-brain, as you can see here by the white
3 arrows. So for this we actually took his brain
4 biopsy, generated 381 million sequence reads. I'm
5 not sure if I overdid it, but I was really worried
6 because it was tissue. Then we ran it through the
7 SURPI pipeline.

8 To give you an idea of, like, SURPI it
9 took about hours to analyze the 381 million reads.
10 But what we found, basically, we were able to
11 assemble directly from that the full sequence of a
12 novel astrovirus. Now, it actually turned out
13 that this was not really a novel astrovirus. It
14 had actually been reported by two different groups
15 before. It actually is an astrovirus which is
16 called astrovirus VA-1 or HMOC. For those of you
17 that are in astroviruses, found in children with
18 diarrhea.

19 We were also able to show that this is
20 by subsequent RNA, DNA insitu hybridization. You
21 can see that, shown there in the red, basically
22 that this is using a probe for the astrovirus. It

1 appeared to directly infect, kind of, the neurons.
2 So this patient had an invasive astrovirus
3 encephalitis.

4 It's actually a little bit concerning
5 because it had been reported before by Ann
6 Lipkin's group that there was a boy who was
7 immunocompromised just like our patient who has
8 astrovirus encephalitis, and he actually died from
9 it. Our patient is the first adult who is
10 actually characterized with this neurovasive
11 astrovirus infection.

12 He was started on ribavirin and
13 intravenous immunoglobulin, unfortunately, we
14 don't know where he got the virus. Presumably it
15 was community acquired. Unfortunately, despite
16 treatment he did deteriorate, but he passed away
17 about four months after the diagnosis. But at
18 least the fact that we were able to identify it,
19 an agent, suggests that we did have four months to
20 intervene if there had been a therapy.
21 Unfortunately, there is no established therapy for
22 astroviruses.

1 So I just want to end by talking -- by
2 acknowledging members of my laboratory. Erik
3 Samoya is the CLIA technician who's developing the
4 NGS assay. Samina Cash is a biologist who, and
5 microbiologist, who worked on the SURPI pipeline.
6 Scott is a bioinformaticist who development the
7 SURPI pipeline. We collaborate with Eric Delword
8 on the NIHCQV study.

9 I'd like to also acknowledge funding
10 from the NIH and Abbott file discovery ward, and
11 amazon.com grant, and a UC discovery grant
12 encephalitis. Thank you very much.

13 QUESTIONER: Thank you for the great
14 talk. I have two questions. First is, there are
15 some cases where viral load or titer will drop and
16 you'll still have symptoms, so there must be some
17 consideration in when you take the samples when
18 you're doing the sequencing for anything that's
19 present too?

20 Then the second question has to do with
21 dark matter. Are you seeing anything that looks
22 like virus that you can't identify at all? I know

1 from previous studies we've seen long stretches of
2 nucleotides that encode proteins that look viral
3 in nature, but no identity in the databases.

4 Thanks.

5 MR. CHIU: So the answer to your second
6 question for us, we definitely see it. We do see
7 these long contigs that have been formed that have
8 no resemblance to anything else in the database,
9 and we don't know whether they're virus or
10 something else. For those we've been archiving
11 them and, you know, in the hope that eventually we
12 -- you know, somebody or -- if we would not be
13 able to definitively identify what they are. So
14 we are seeing dark matter.

15 The first question is a really good one
16 in that you're right that it's really dependent on
17 when you get a clinical sample. Because this is a
18 direct detection method just like PCR. It's
19 really dependent on there being the virus or the
20 pathogen in your clinical sample.

21 Unfortunately, we've been focused mainly
22 on what I call Hail Mary samples. Basically, you

1 know, samples of last resort where they've been
2 extensively tested for other agents. As you might
3 expect, we're getting a very low percentage of
4 positives.

5 But I think that where this assay really
6 belongs is a, kind of, early second line assay.
7 So you would do, kind of, the frontline tests
8 upfront, like bacterial culture, respiratory virus
9 testing, and then as -- if it's still negative and
10 you're considering sending out tests, more
11 esoteric testing then you probably would want to
12 have a test like this in place, a second line
13 testing. But not as, kind of, a test of last
14 resort.

15 MR. KRAUSE: Just to make a quick
16 comment. We're running a little over time, so I'd
17 like the questions to be quick and the answers to
18 be quick. Maybe we take two more questions.

19 MR. CHIU: Okay.

20 QUESTIONER: While you're generally
21 looking for viruses, you're also, basically,
22 containing the full human genome of these

1 patients. How do you protect patient privacy, and
2 how do you store your data so that it would be
3 protected?

4 MR. CHIU: Yeah, that's a great
5 question. So essentially, I mean, these are all
6 protected under, kind of a VPN. That's for the
7 hospital VPN, so it's considered clinical data.
8 It's considered, kind of, HIPPA protected data.
9 That's how we're keeping it.

10 Now, with respect to actually providing
11 this information publicly, we've been removing the
12 human sequences and then depositing it into SRA.
13 Otherwise it has to go into something like DD gap.

14 MR. KRAUSE: We'll go to the last
15 question over here.

16 QUESTIONER: Great talk.

17 MR. CHIU: Thank you.

18 QUESTIONER: How do you see this
19 changing with advances in long read technology,
20 and will that have an impact on some of the things
21 you're doing?

22 MR. CHIU: I think long read technology

1 is very promising, especially for doing genomic
2 assembly, and for identifying bacterial genomes.
3 There is still an issue with having short reads
4 because, as you might expect, you might have
5 chimeric assemblies and things that may make it
6 difficult to interpret the data. So I think
7 having long read technology is going to be really,
8 really important in the future, especially for
9 larger genomes.

10 For viruses, probably less useful.
11 Because, I think, the reads of about 100 bases are
12 probably adequate for viral identification. But
13 certainly for genome assembly. Thank you.

14 MR. KRAUSE: Thanks, Charles. That was
15 really excellent. I hope those who wanted to ask
16 questions will catch up with you during the break.
17 Having heard from government, and academia in a
18 clinical setting, we're now fortunate to have two
19 speakers from industry who will be talking about
20 using these kinds of techniques in production
21 flow.

22 The first of these speakers is Robert

1 Charlebois, who's a senior scientist Sanofi
2 Pasteur. He's going to be talking about
3 standardizing a polymorphic next generation
4 sequencing pipeline. So, Robert.

5
6 DR. CHARLEBOIS: The right arrow to
7 move? Okay. Thank you for the opportunity to
8 speak to you today and also, thanks to previous
9 speakers for providing a lot of the motivation and
10 the background for this kind of work.

11 So Arifa, this morning, talked about the
12 importance of exploring NextGen sequencing assay
13 for example, for advantageous agent detection but
14 also other applications in industry to get, in our
15 case, vaccine products out the door. But it's not
16 easy.

17 The main problem for industry is that
18 it's a moving target. So the technology is
19 evolving which is very good and it has a lot of
20 promise which is very exciting but in terms of
21 developing assays, we have to worry about not stag
22 -- like not getting stuck in a particular version

1 of the technology when next year it's going to be
2 so much better. And we also don't completely
3 understand all of its strengths and weaknesses.
4 It is still new.

5 It's not the same way of looking at the
6 problems as we are used to looking. And so, we
7 need to learn about it and as we learn and as it
8 evolves, we need to still control and standardize
9 this technology so that we can actually apply it
10 in a quality environment.

11 Now, some of the assays that we do are
12 research grade. You can have a lot more freedom.
13 Some of them are, you know, characterization
14 studies, some of them are release tests. And so,
15 we have different levels of quality that we need
16 to match. And we don't want to not be able to
17 take advantage of the technology.

18 And so, yesterday we heard a couple of
19 talks that talked about object-oriented
20 programming. At heart, I'm a C++, a hardcore C++
21 programmer and so, I understand these concepts
22 pretty well. And so, the way that we're

1 approaching the standardization of an assay is
2 essentially to make each step of that assay a
3 standalone part of the overall pipeline.

4 And so, we're essentially making an
5 abstract functional module where you have some --
6 where that particular step doesn't care what came
7 before and doesn't care about what comes after.
8 And so, by doing that, by let's say validating the
9 interfaces, the input and the output, then we can
10 have the independence of the individual steps.
11 And so, that makes it, I think a lot easier to
12 validate the individual steps.

13 So eliminating dependency but we still
14 want flexibility. We have a lot of different
15 kinds of assays we want to run. We want to be
16 able to reuse the steps for different assays. We
17 want to be able to improve the steps and so, we'd
18 have different instances of each of these
19 particular steps. I'll provide more concrete
20 examples in a few minutes.

21 And so, this is essentially what a
22 pipeline looks like to us. The bottom part has

1 several streams because some of the, I mean, it
2 isn't -- some of those steps are not in common.
3 Like, if you're looking for genetic stability,
4 it's quite different than looking for adventitious
5 agent detection or sequencing your cell line for
6 characterization.

7 But some of the steps at the top are in
8 common to all of them. And so, you would have,
9 you know, your nucleic acid extraction and your
10 conversion to a double-stranded DNA. And I'll
11 talk about these steps and the challenges of each
12 of them shortly. And the sequencing, and then,
13 you have your domain specific sequence analysis.
14 And so, those are the big automated bioinformatics
15 pipelines, the things that we've been talking
16 about as big data.

17 And then, the next step after that, the
18 domain specific data analysis, we're no longer in
19 big data mode. And so, there we're in the more
20 traditional kinds of analyses where we have, you
21 know, lists or trees or variant profiles. We've
22 collapsed the huge data into something that we can

1 use our ordinary tools for like molecular modeling
2 or what have you.

3 And then, there can certainly very well
4 be some domain specific follow-up work depending
5 on what we find. Most of the time, of course, in
6 industry the assays end up being pretty boring. I
7 mean, you get what you expect. But occasionally,
8 you have this surprise and then, we mobilize and
9 do an investigation and try to figure out what
10 might have gone wrong.

11 It's not quite the Pandora's box as
12 people feared initially. The false-positives that
13 we see, for example, in add agents detection are
14 very limited in number and very manageable, at
15 least in our experience. And so, it's not so
16 scary. Initially, people were terrified but it's
17 not so scary today after we've learned a bit.

18 And so, I'm going to talk about the
19 interfaces too. So step one sounds sort of
20 straightforward and for some of the applications
21 like RNA seq or looking at bacterial DNA, it's
22 pretty straightforward. But some of the other

1 applications like nucleic acid extraction for add
2 agents detection where you have single-stranded
3 RNA and DNA in packaged in viral capsids or in
4 bacterial cells or fungal spores or what have you,
5 it's not so obvious how to do this without bias.

6 We don't want to throw the baby out with
7 the bath water. Through enrichment techniques,
8 for example, and there have been some talks in
9 earlier conferences where people have tried like
10 nuclease and they've seen signals disappearing.
11 And so, we have to really worry about that and so,
12 that's one of those boxes that is evolving. We
13 have a working version that we like that is
14 working well for us but we know that we can
15 improve, that we can improve the sensitivity
16 because most of the background in our case won't
17 be human but it'll be the cell line that we use
18 for manufacturing like a monkey or a hamster or
19 what have you.

20 But from the point of view of the
21 pipeline, the only thing that the next step needs
22 to worry about is that it got nucleic acid. And

1 so, it might need to what kind of nucleic acid it
2 is. And so, your test requisition would say, you
3 know, you're looking at RNA. But the point is
4 that it's all standardized to nucleic acid.

5 And then step number two, there's not
6 much to talk about. That is very straightforward.
7 DNA, if you have double-stranded DNA to start,
8 then that's a no-op. If you have RNA to start
9 you'll turn that into double-stranded DNA and so
10 on. If you have a complete unknown, then you go
11 through all the steps of reverse transcription,
12 second strand synthesis, maybe even some
13 amplification.

14 Step three is, so it's shown here as a
15 black box, library prep and sequencing but there's
16 opportunities for modularization with that too.
17 So if you have a particular instrument in-house,
18 there might be different kits for library
19 preparation. And so, you can sort of separate out
20 library preparation from the sequencing.

21 If you have different instruments, you
22 might have different kits and different

1 procedures, different SOPs and so, there's
2 opportunity to modularize with that. And I'm not
3 showing that here because then, the talk would be
4 too long. But in any case, we can decide, and so,
5 this -- we may not decide as a large community but
6 we can decide in-house that our output of this
7 step will always be FASTQ, for example.

8 And then, given that then it doesn't --
9 this step doesn't care about what comes after.
10 And so, we produce FASTQ and the next step has to
11 worry about dealing with FASTQ files.

12 Step four, so we've heard and we're
13 going to hear some more about pipelines. And
14 again, this can be modularized. So I'm not
15 presenting our bioinformatics pipeline which is
16 sort of the focus of this session. It's been
17 presented at last year's PDA Conference and will
18 be published in the PDA journal, fingers crossed.
19 But here again, you can modularize.

20 So for example, our particular pipeline
21 has two steps. We have as our first step that
22 takes FASTQ data from whatever kind of instrument.

1 It will do the quality filtering. It will do the
2 adapter trimming. It will do assembly. We do de
3 novo assembly on the primary data. It works for
4 us. And then, we do some annotation of the
5 sequences.

6 And then, that produces a FASTA file.
7 And then, the second part of our bioinformatics
8 pipeline, there's also an automated self-contained
9 module that takes a FASTA file that's properly
10 annotated and produces lists of potential
11 contaminants. It actually identifies everything
12 that's in the sample. You know, whether it's
13 hamster or bacterial or fungal or what have you.

14 And right now, it just produces lists
15 but our goal for standardization will be to
16 produce something like.xml or at least, like.xml
17 as a technology but the goal is to have like an
18 ontological or defined vocabulary kind of output
19 so that we can standardize. The reason I put.xml
20 there as a proposal is because 50 years from now,
21 if our product's still on the market, anybody can
22 still read an.xml file because it's human

1 readable. They can write a parser for it.

2 If we produce some kind of binary or
3 proprietary output, good luck reading that 50
4 years from now. I mean, we'd have to preserve the
5 software and the hardware and that was discussed
6 yesterday. So that can be a challenge. We can
7 actually, you know, print out on paper an.xml
8 because I'm thinking, you know, 50 years ago, 1964
9 there were punch cards with machine language
10 instructions for hardware that is sort of
11 interesting but not really relevant today.

12 So let's try to keep this because, you
13 know, as a manufacturer, we have to store these
14 records for a long time because we may need to
15 revisit them. And so, something that it is, I
16 mean, it is a heavy format at it was discussed
17 yesterday. But I think it's something that we
18 could preserve for the long, long term. And then,
19 the language ontologies tend to use, you know,
20 well-defined vocabulary so it's easy to understand
21 in the future.

22 The data analysis, so this can be a lot

1 of different kinds of things and you probably will
2 need adapters or you will need adapters from our
3 defined vocabulary to feed into these new, these
4 various tools like molecular modeling or alignment
5 and so on. So, sorry, in the previous step what
6 we're specifically not doing is things like
7 aligning to the reference genomes, looking at
8 coverage, look at percent identity and all that.
9 Because that can be -- we don't want to lock that
10 into this particular module because if we want to
11 improve that part, we don't want to have to
12 revalidate that entire step.

13 So if all we do is say this sequence is
14 from this taxon, then the next step, then we can
15 go further and do the additional bioinformatics
16 analysis on it. So we still do that work but we
17 just don't do it in the same step. And so,
18 there's lots of things we can do here. We'll need
19 adapters coming in and adapters going out. Again,
20 producing some kind of defined vocabulary output.

21 So I don't know if we can all agree on
22 such kind of vocabulary or ontologies but within a

1 particular company, we certainly could. And that
2 would, at least, make it easy for let's say a
3 regulatory agency so they'd know, you know,
4 information coming from Sanofi always looks like
5 this and information coming from Merck always
6 looks like that or what have you.

7 The follow-up is a lot more open-ended
8 because it really depends on what you find. You
9 may be able to use some assays or some techniques
10 that you've already developed before. You may
11 have to develop some new techniques. It's
12 essentially problem-solving and investigation and
13 a lot of the time, there's just nothing to do.
14 And so, if there's nothing to do you just simply
15 write your report, your C of A. If there's
16 something to do, then you have to follow it up.
17 And so, we will hope to make those kinds of
18 reports, semi-automated and standardized for at
19 least, you know, for our own safe if not for our
20 customer's sake.

21 So now I have three simple enough
22 examples of just showing what this would look

1 like. And so, suppose we want to look -- to
2 characterize the genetic stability of a
3 single-stranded RNA virus. So we produce those
4 products and so, we have -- we would have our
5 standard validated RNA extraction procedure, our
6 TPCR. We would sequence, let's say with an
7 Illumina or some other instrument that would give
8 us the sufficient resolution for looking at
9 genetic stability.

10 We could have our statistical variant
11 calling pipeline that compares various samples
12 against each other. So there could be a database
13 in there of what we found before. Looking for,
14 let's say, lot to lot variability or any kind of
15 trends. Whatever we find, then, could feed into
16 -- could be mapped onto a 3D structure so that the
17 biologist in the risk assessment can look at it
18 and say, there is a problem or there's no problem
19 and go from there and write the report.

20 Example number two, adventitious agent
21 detection, so let's say we're working with
22 supernatant from a master cell bank. Our current

1 approach is to do total nucleic acid extraction
2 with no enrichment. It is the needle in a
3 haystack problem. We have a lot of host nucleic
4 acid. And we just sequence it at very high depth
5 so that we get the resolution that we need.

6 We're still working, as Arifa mentioned,
7 there's some spiking studies ongoing to see what
8 kind of sensitivity we get, the LOD that we can
9 get from no enrichment and see if we have to
10 enrich to get -- to match the existing methods.
11 So we have our taxonomic classifier sequences. We
12 get candidates. We then characterize them,
13 alignments and so on and then, follow-up work.

14 And then finally, and just to show you
15 some reuse, suppose we would want to look at the
16 genetic identify of bacterial endocriduction
17 cells. So we have our E. coli seed and we're
18 producing our protein product and we want to make
19 sure after the fermenters have grown up that it's
20 all still all the same thing. And so, here it's
21 very simple. We have a DNA extraction. This --
22 it's already DNA and we're reusing this myseek

1 module for example. We don't have a myseek but
2 I'm just using that as an example because it would
3 be appropriate technology here.

4 Again, we have all of the standard
5 interfaces. We could use our classifier again.
6 We don't have to reinvent something. Because
7 we're trying to find out what this bacterium is.
8 And so, in our database we have this bacterium and
9 we have other strains of this bacterium. We can
10 resolve them with our classifier. We decide that
11 it's okay and we go on.

12 So this sounds very repetitious. So my
13 talk so far has been very repetitious and that's
14 the point is the reuse and the standardization in
15 doing things the same way over and over again so
16 that we don't have to reinvent the wheel every
17 time we want to do something new. And so, just to
18 summarize in two slides, the module is responsible
19 for accepting the format that it gets. And it's
20 responsible for producing the format that the next
21 step expects.

22 And so, if you can prove invalidation

1 that it deals with input A and that it produces
2 output B, then you don't have to so much worry
3 about all of the combinations of all of the
4 modules. You will, of course, in a validation
5 have to an end-to-end pipeline run to make sure
6 that the whole system works well. But then,
7 that's not so onerous as mixing and matching all
8 of the different components and all of their
9 different parameters. So that's our proposal and
10 the advantage is it provides a lot of flexibility
11 in testing. And so, your test requisition may be
12 asking for lots of different versions of things
13 and we can provide that with this kind of an
14 approach assuming that the different modules are
15 validated at the right level of quality.

16 Then we can run the assay that we're
17 asked to do. It allows reuse and so, that
18 accelerates new tests. And it allows evolution so
19 we can improve an individual test without
20 impacting the rest.

21 And then, finally, I want to acknowledge
22 people that I work with so Lucy and Lauren and

1 Simon and I are on the advanced virus detection
2 technologies user interest group that was
3 mentioned earlier. Greg and Sam have recently
4 joined and they're working on the bioinformatics
5 validation and sample extraction optimization
6 respectively. Thank you.

7 DR. KRAUSE: Questions?

8 SPEAKER: (inaudible) or rather a
9 comment. Thank you exactly this -- this is
10 exactly the kind of input which we need to proceed
11 with our strategization in bioinformatics,
12 harmonization efforts. So please do join the
13 group and we need --

14 DR. CHARLEBOIS: I signed up.

15 SPEAKER: Exactly, perfect. That's the
16 comment. This is the exactly the reverberance we
17 needed.

18 SPEAKER: So we heard yesterday that for
19 variant detection that they often use three
20 different methods to validate or to come up with
21 some reasonable assessment of true variation.
22 What do you intend to do with variant detection?

1 I noticed that part of your pipeline. How will
2 you validate that? Which method will you select?
3 Those kinds of things, how are you dealing with
4 those issues?

5 DR. CHARLEBOIS: So the variant
6 detection, so we're primarily have the problem of
7 viral quasi-species and so, the viral seed that we
8 start with is already quite variable. And we --
9 what we want to do is to ensure that
10 manufacturing, that it's consistent and so, lot to
11 lot, variability needs to be minimized or if there
12 is any, we need to understand what the impact of
13 that is.

14 And so, it's more of a like we do
15 Chi-square tests, for example. That's what we're
16 exploring to essentially look at the profiles for
17 each nucleotide in the genome to see if there's
18 any actual statistically significant difference
19 among them. So it's a bit different than, I
20 think, the problems that were described yesterday
21 and in terms of validating that were we've
22 traditionally been using Sanger sequencing

1 approaches which is -- which has been accepted in
2 the community.

3 And so, we've done, you know,
4 head-to-head comparisons of that. For other kinds
5 of variants, I think it was mentioned on the very
6 first day the neurovirulence, monkey
7 neurovirulence in polio, and so, that requires a
8 very sensitive measurement of the variability at a
9 very particular nucleotide. And so, there's a
10 specific test called MAPREC that does that.

11 And so, we're essentially doing our
12 comparisons head-to-head against that assay which
13 is the way that we've been doing it and
14 submitting. To make sure that they correlate very
15 well, that they give exactly the same result. So
16 we're always comparing our existing methods with
17 what we've been doing.

18 In some cases like with this add agents
19 detection, we can do things that we couldn't do
20 before. But we have -- for the things that we
21 could do before, we have to do at least as well
22 and it has to be consistent. So that's pretty

1 much all I can say to that.

2 DR. KRAUSE: So our final talk of the
3 session is also an industry talk from John
4 Thompson at Merck who will be talking about viral
5 metagenomics analysis, BLAT-based approach and
6 adventitious virus detection interest group
7 activities. So he'll be providing us both an
8 industry perspective and an idea of what the
9 interest group has been doing and doubtless will
10 encourage you all to join at the end.

11 So John? And so, depending on the
12 timing after John's talk, we may have a short
13 panel discussion which may give you an opportunity
14 to revisit some of these issues and ask questions
15 that may be interesting to hear all of the
16 speakers in this group discuss together. But
17 we'll see what the clock looks like.

18 DR. THOMPSON: Thank you. Thanks to
19 Arifa and Philip for inviting me to this session.
20 So I'm going to describe a pipeline that we've set
21 up and been working with in a research mode
22 locally. In the process, we curated the public

1 data and developed an automatic way to curate the
2 public data that I'll describe.

3 And we've used an older algorithm,
4 BLAT-based and I'll go into the reason we've done
5 that. And some of the methods we use to sort of
6 calibrate the sensitivity of that and shake it
7 down. And then, I'll spend about a third or the
8 last half of my talk talking about some of the
9 activities of the adventitious virus detection
10 group that Arifa introduced.

11 So like other people who have spoken
12 today, we've centralized on the NCBI virus
13 resources and the virus genome databases is the
14 RefSeq collection of genomes. It's a
15 representative sequence approach generally with
16 one record per virus. And then, we supplement
17 that with data from the NCBI non-redundant
18 nucleotide, the NT collection which collects data
19 from all the public resources.

20 And it includes many additional virus
21 genomes that represent additional virus diversity
22 that we want to pick up. And it also contains a

1 large number of partial viral sequences, fragments
2 of viruses that have been deposited in the
3 resource. And so, we've developed a system for
4 combining record from both sources but in prior
5 work in this area, we found that the greatest
6 source of false positives are the partial viral
7 sequences that are in the NT collection.

8 They're often deposited with vector and
9 host sequences attached to them. And we found
10 that many times we're getting matches to the
11 vector and host sequences and not the viral
12 portion of the sequence and they represent a huge
13 false positive burden that the analysts need to
14 wade through. So we developed this system to sort
15 of capture the diversity represented in the
16 partial sequences but avoid as much as possible
17 the false positive issues associated with those.

18 And so, at the top right we've got the
19 RefSeq collection that's going -- of whole genomes
20 that are going into our virus database and then we
21 take the -- on the left, take the GenBank NT
22 collection and we parse it into whole viral

1 genomes, partial viral sequences and then, we keep
2 the non-viral NTs separate for use in our pipeline
3 as well.

4 And so, we take the whole viral genomes,
5 remove exact duplicates and combine them with the
6 RefSeq collection of whole genomes. Then what we
7 do is take the partial sequences and align them to
8 the whole genomes and we discard any partial
9 sequences that align to the whole genomes as being
10 already represented by a whole genome. And in
11 that way, we're only keeping about three percent
12 of the partial sequences. And so, we've
13 eliminated most of the source of our false
14 positives but we're keeping those partial
15 sequences that add diversity.

16 And so, in our algorithm choice and this
17 is now almost two years ago, which is a long time
18 in NextGen sequencing technology, we initially
19 rejected the short read aligners, BWA, Bowtie in
20 this and some 70 similar algorithms are out there
21 now that are all optimized for matching with near
22 identical sequences to a reference sequence. And

1 we decided that those were too stringent to
2 capture the viral diversity for our purposes and
3 we fell back on a decision between BLAST and BLAT.
4 I think almost everybody who is in the room should
5 have known, know about and probably used BLAT --
6 BLAST in their past sometime.

7 BLAT is a BLAST-like alignment tool that
8 if you've ever aligned a sequence to the genome
9 using the UCSC genome browser, BLAT is the
10 underlying algorithm. It's designed to be faster
11 for high homology alignments. But it is tunable
12 into the range that we needed for lower identity
13 matching as well. And so, that's what we
14 optimized or what we centered on for a number of
15 technical reasons that I won't completely cover
16 today.

17 But this is showing how we started to
18 define the sensitivity and specificity parameters
19 that we needed to use with BLAT and in this
20 example, we're using retroviral sequences, HIV.
21 We took the two HIV records from the reference --
22 from the RefSeq collection and one for HIV-1 and

1 one for HIV-2 as our database in this trial. And
2 then, we picked 10,000 random HIV sequences from
3 the NT collection to have additional diversity to
4 see what kind of identity threshold we needed to
5 adopt to pick up the diverse sequences as
6 represented in NT.

7 And we selected that 10,000 sequences to
8 be less than 90 percent identical to the reference
9 sequences to make it a challenging test. And we
10 also picked up a false positive dataset of 10,000
11 random non-viral sequences from the non-viral
12 section of the NT collection. And we found that
13 at 80 percent identity, we were picking up about
14 82 percent of those 10,000 reads. We picked up
15 another hundred reads by dropping to 70 percent
16 and another 10 reads by dropping to 60 percent.
17 So we decided that 70 percent was sort of our
18 threshold of diminishing returns and adopted that
19 for further work.

20 Then some additional sensitivity tuning,
21 we've done using the whole RefSeq collection of
22 viral sequences as a database and using the entire

1 NT viral collection of sequences as a query test
2 to see what proportion of the diversity
3 represented in the NT collection would be detected
4 by the reference database. And at 90 percent
5 threshold that was only 61 percent of the NT
6 collection viral sequences were detected but at 70
7 percent we were picking up over 90 percent of
8 those sequences. So we think the sensitivity of
9 BLAT at 70 percent is pretty appropriate.

10 And this describes the pipeline we
11 established using BLAT. We took a two-step
12 approach rather than compare -- because it's
13 engineeringly (sic) difficult to compare to a 50
14 gigabyte or larger database. So we take the
15 clean, raw reads and align to the smaller virus
16 database first. And then, we take the smaller set
17 of virus-like reads and counter screen them
18 against the NT viral database, both using the BLAT
19 algorithm.

20 And out of that we get a set of reads
21 that match best to the viral sequences including
22 we keep ties as well. And then, we go through a

1 process that some similar to what others have
2 described of consolidating those reads and
3 reporting on a target basis what the depth and
4 breadth of coverage is to help interpret the data
5 for the biologist downstream.

6 So as I said, even running a large
7 Illumina dataset against just the viral database
8 in one fell swoop is difficult. And we end up
9 splitting both of those sequence sets into
10 multiple smaller pieces and doing all the pairwise
11 combinations and then, on the back end of that
12 merging and sorting the data to get the best hits.
13 And so, there's a great deal of splitting and
14 merging that goes on in that process and we have a
15 whole job control system that manages that and
16 resubmits jobs that fail for hardware reasons and
17 so forth.

18 The processing times we've been
19 achieving are not as good as what Charles is
20 getting with the SURPI pipeline and that's part of
21 the -- part of why this pipeline is really showing
22 some signs of age at this point, I think, and may

1 need some further development but with Roche,
2 we're processing with Roche datasets of about
3 600,000 reads. We were getting things through in
4 three hours. With Illumina reads at 26 million
5 reads, it's now taking three to five days to
6 process and that depends largely on the amount of
7 host reads that need to be aligned or that might
8 be in the database as well. The automated viral
9 database curation process that I described runs in
10 basically an overnight run.

11 And one last test I'll describe is an
12 in-silico spike-in analysis that we did to sort of
13 model a RNA seq experiment. We simulated 230
14 megabytes of human transcriptome reads and then,
15 spiked-in a simulated set of 1X to 100X coverage
16 of the three viruses listed in the table here.
17 And we were primarily concerned that the counter
18 screening against the NT collection could, in
19 theory, remove some real virus reads if the NT
20 collection was contaminated with some poorly
21 annotated or mis-annotated virus sequences and so
22 forth.

1 But what we found was that with the
2 versions of the databases we're using right now,
3 we didn't lose any of our virus reads going
4 through this. We recovered 100 percent of our
5 reads at every level. And so, we think the
6 bioinformatics is actually working fairly well in
7 this process.

8 The other thing that comes out of that
9 spike-in analysis is a question of how do you
10 interpret when you have a whole virus. And the
11 important question is what depth of sequence do
12 you need before you expect to see full coverage.
13 And from the simulated data, that answer is about
14 10 percent. That's probably overly optimistic
15 compared to real data that will have less even
16 coverage.

17 But the point is that around -- with
18 this dataset, around 10X coverage gets you to 100
19 percent coverage. So if you're looking at a
20 fragment of a virus, alignments to a fragment of a
21 virus that are very deep, the hundreds or
22 thousands of reads deep and you don't see coverage

1 across the entire virus, you can logically
2 conclude that that virus is not really present.

3 So some of the future considerations for
4 this pipeline are we initially avoided host
5 filtering because of the potential to remove
6 virus-like sequences based on endogenous
7 retroviruses in the host genomes. But I think
8 scalability issues are forcing that -- forcing us
9 to reconsider that option. The suggested
10 constraints that we would apply for host filtering
11 would be to use really high-stringency alignment
12 methods and possibly consider masking the
13 endogenous, the known endogenous, retroviral
14 elements so that those are -- types of sequences
15 are not being removed by the host filtering.

16 And we've also considered going to de
17 novo assembly. Early on we avoided that because
18 we can align individual reads but when you do a de
19 novo assembly you need a certain depth of sequence
20 to form a contig. So we thought there's a
21 potential for some loss of sensitivity there. And
22 so, we didn't go that route initially but I think

1 that's really attractive because the longer contig
2 sizes greatly improve your ability to evaluate the
3 presence of about a -- coverage of a whole virus
4 as well as identify potentially novel viruses.

5 So just to summarize it that at this
6 point then, we've presented what's a practical
7 automated viral database curation approach that
8 relies on the existing sequence sources and
9 balances capturing virus diversity with the need
10 to avoid a source of large -- large source of
11 false positives. And we've shown that a
12 BLAT-based viral detection pipeline can work
13 fairly well but there are definitely scalability
14 issues as the size of datasets continue to
15 increase on us.

16 And then, I've presented an approach to
17 sort of calibrating and testing the sensitivity of
18 that -- the pipeline. So for the rest of the talk
19 I'm going to summarize some of the activities,
20 some of the bioinformatic activities of the virus
21 detection interest group that Arifa introduced.
22 And so, I can pretty much skip this slide. Arifa

1 introduced this already.

2 This group was formed in October 2010
3 and I think it's a great example of several
4 companies, multiple companies, in the industry
5 recognizing that this is sort of a precompetitive
6 space. We're all going to have to go through this
7 and come to some sort of agreement on a
8 standardized process eventually. And we can
9 accelerate that process by sharing information
10 that these companies, I think, have logically
11 considered a precompetitive space.

12 So the virus detection group has been
13 broken into four subgroups and I'm going to focus
14 on the C and the D group. But this, as Arifa
15 mentioned, the sample A group is involved in
16 evaluating sample preparation methodology and the
17 B group has been setting up a spike-in, excuse me,
18 a spike-in experiment. And those two groups
19 together are sort of setting up the lab side of
20 the process of how do we -- how do you optimize
21 isolating viral nucleic acid.

22 And then, the C and the D groups have

1 been working on the bioinformatics aspect. C has
2 been focused on standardizing viral sequence
3 databases. And D has been focused on comparing
4 and contrasting the bioinformatics pipelines.

5 So the focus on subgroup C which is
6 chaired by Kavitha Bekkari of Merck, the group did
7 an initial survey of data sources and used and
8 there's a great commonality, I think. Most people
9 are using the NCBI RefSeq collection and NTB
10 sources. And I've described those a little bit
11 already.

12 The NT collection, again, this slide is
13 old. Charles told us this is up to 70 gigabytes
14 now. And I described the distribution of whole
15 genomes and partial sequences so I'm not going to
16 dwell on that here.

17 I want to speak for a minute about the
18 genomes neighbors that Rodney Brister set up
19 because this is what we use to identify the whole,
20 the punitive whole viral genomes from the NT
21 collection. The problem is that the sequences are
22 coming in so fast that the curation tends to lag

1 behind the sequence collection and that there are
2 many potential, full-length viral genomes in the
3 NT collection that do not yet have that complete
4 genome tag on them.

5 So the approach that Rodney and his
6 group took were to compare coding regions between
7 the RefSeq entries and the NT collections. And if
8 all the coding regions seem to be present, then
9 they tag it as a genome neighbor or as an
10 otherwise known or otherwise thought of as a
11 punitive whole genome. And we capture -- we use
12 that genome neighbor approach to capture the whole
13 genomes from the NT collection. And it's quite a
14 bit better in terms of the numbers of genomes
15 we're picking up than just relying on the complete
16 genome tag.

17 So Arifa's lab de -- Arifa described
18 this so I don't have to spend a lot of time on
19 this either but this is a slide from Arifa's group
20 and where she and Aisha Aljanahi have been at the
21 PDA have been working on a keyword approach with a
22 series of positive and negative keyword selections

1 so that you, up until now or in our pipeline,
2 we've been relying on the taxonomy tags of GenBank
3 to pick out viral sequences. And so, this is an
4 attempt to, again, avoid relying -- overly relying
5 on the annotation and pick up viral sequences more
6 comprehensively and that's being evaluated in the
7 context of the interest group.

8 So subgroup D is chaired by Adam Palermo
9 at Genzyme Sanofi and has been setting up and
10 well, initially surveying the different methods
11 that people use and working toward setting up a
12 comparison of different methods. And so, Adam put
13 together this slide after trying to sort of
14 generalize the input we got from different
15 companies and different groups that are working on
16 this that pretty much everybody is starting with
17 reads, doing some kind of quality filtering.

18 Some people filter against the host.
19 Some people skip that step and then, some people
20 take the high quality reads and preassemble them
21 with de novo or other methods. Some people take
22 the raw reads into the database alignment step.

1 Some people are aligning not just to a viral
2 database but to a more comprehensive pathogen
3 database.

4 And then, you have a set of punitive
5 contaminating, potentially contaminating reads,
6 that are indicative of some kind of contamination
7 that are being put through some kind of a
8 post-processing analysis to determine how reads
9 are mapping to each target and what extent of
10 coverage of those targets is taking place so that
11 the biologists can begin to sort through and
12 interpret those results.

13 But as you can see, there's considerable
14 variation and choices made at each step by
15 different companies and it's really sort of
16 impossible to a priori define what the best
17 pipeline strategy is. And that's sort of the
18 basis for the interest group in the first place.

19 So this group has turned focus on
20 setting up a pipeline evaluation and of between
21 spike-in simulations and public datasets, we've
22 actually settled on some public datasets as a

1 first trial and several groups are participating
2 in this comparison. So we've picked several sub
3 -- public datasets. One is virus particles from a
4 human fecal sample that represents a complex
5 metagenomic sample and should be likely to have
6 many things that may not be in the databases yet.
7 So it should be a very interesting dataset to
8 analyze in these pipelines.

9 The second one is an HIV infection
10 experiment that includes a human cell line, plus
11 and minus HIV infections. So with the HIV
12 infection, there'll be a massive predominance of
13 HIV sequences in that dataset that hopefully
14 everybody will pick up. And the uninfected
15 control serves as sort of a negative control that
16 hopefully doesn't contain much in it and will be a
17 test of how well each pipeline deals with false
18 positives.

19 The third dataset is also an uninfected
20 human cell line but it's a particularly large
21 dataset. This may be 5 to 10 times larger than
22 what most people are currently using. This is

1 being collected on the new Illumina X Ten machine,
2 120 gigabytes of data from this cell line.

3 So there's several companies and
4 organizations that are currently participating on
5 that. We've come to some agreement on how to
6 summarize the data at the hit or target level.
7 And then, in particular cases where there are
8 discrepancies between pipelines, we're likely to
9 delve down into the read level on selected samples
10 or targets.

11 So the different groups have downloaded
12 these different datasets and are in the process of
13 running those now. In the next few weeks or so,
14 we'll start comparing results.

15 So the last slide I'm going to present
16 is kind of, I think, a consensus from discussions
17 in the interest group of some rules or some
18 thoughts about interpreting the data and that's --
19 I posed it as a series of three questions. And
20 the first question is is the complete virus
21 present? And for that, we, I think, generally
22 agree that evaluating the depth and breadth of

1 coverage is critical for determining whether a
2 complete virus is present but it's hard to
3 establish precise thresholds and to a large
4 degree, that's going to depend on the different
5 library preps that are being used and how even
6 their coverage is.

7 And so, that's going to be something
8 that evolves with the methods that are being
9 applied. The second question is is the viral
10 titer increasing during the culture and I think
11 this is even more important than the depth and
12 breadth of coverage. Because we've seen in
13 Arifa's nice example that you can get partial
14 coverage from a known virus from a novel virus.
15 So even if we had partial coverage of a virus but
16 we say that titer increasing during a culture run,
17 during a time course from a culture run for a
18 biologic or a vaccine preparation, that would be a
19 great cause for concern and for -- generate
20 follow-up activity certainly.

21 And the second -- the third question is
22 really is the detected virus infectious to humans?

1 So if we found a full length virus or we found a
2 partial virus that was increasing during a
3 culture, we would be relying on the information
4 from the public databases and collective knowledge
5 about where the human infectivity risk lies.

6 So that's it. I want to thank the
7 collaborators, my collaborators at Merck, Kavitha
8 Bekkari did a lot of the bioinformatics for this
9 pipeline and Paul Duncan provided a lot of
10 biological guidance. Joe gave us engineering
11 support for a lot of this splitting and merging
12 we're doing. And I've also listed the interest
13 group organizers and subgroup chairs and
14 certainly, many other people in the interest group
15 deserve credit as well. Thank you.

16 SPEAKER: So I actually have two
17 questions from the online audience. The first is
18 a little broad. It may be difficult to answer
19 concisely but what is your opinion on the pros and
20 cons of k-mer based viral detection?

21 DR. THOMPSON: So k-mer is something
22 that's been -- I've been introduced to but we

1 haven't had time to play with yet. It's -- from
2 what little I understand of it so far, I think
3 it'll be a greatly different way of looking at the
4 data but it's potentially much more scalable. So
5 I'm anxious to see somebody, if not us, go that
6 direction.

7 DR. KHAN: Okay, and the second is what
8 is the FDA's perspective? Also not directly to
9 you, I guess, on applying NGS first sequence
10 variant detection in cell banks and production
11 cells used for the production of other biologics,
12 i.e. antibodies?

13 DR. THOMPSON: Well, I'm not from the
14 FDA so I don't know that I can speak to that.

15 DR. KRAUSE: Neither -- do I need this
16 microphone? I'll speak into your tie. So neither
17 Arifa nor I are from the part of the FDA that
18 deals with that issue as well but I think we're
19 all faced with the same question of how it is that
20 these methods can be standardized in a way that
21 they provide reliable results. And I do know that
22 there's great interest, in particular, in the

1 monoclonal antibody production field in favor of
2 using these kinds of techniques to rapidly detect
3 potential fermenter contamination events. And so,
4 I think the discussion we have here today is
5 directly relevant to that issue.

6 MS. VAN 'T VEER: I'm Laura van 't Veer.
7 I had a que -- my background is more in cancer,
8 molecular testing. And I was wondering, in the
9 cancer next generation sequencing there are a
10 number of these dream challenges organized and one
11 is ongoing actually on pipelines for next
12 generation sequencing REM files to FASTQ files
13 which I think they're probably very specific
14 elements for viral testing but some might be
15 actually very similar to what you're doing.

16 And since, Merck is also working in the
17 space of oncology and actually in next generation
18 sequencing, as I'm aware, so do you actually
19 exchange part of the pipeline would actually would
20 fit nicely with the previous speaker who actually
21 showed that there are modules you could sort of
22 validate across maybe different organisms,

1 different testing capabilities?

2 DR. THOMPSON: Well, I think that's kind
3 of what we're doing within the interest group. We
4 have this trial going of different bioinformatic
5 pipelines. We have the spike-in test going where
6 different labs are applying different methods for
7 isolating nucleic acid. So I think that's sort of
8 essentially what we're doing within the interest
9 group.

10 Ah, I see, look for commonality between
11 the cancer work and the virus work. That's a good
12 idea.

13 DR. KRAUSE: Well, great, I'd like to
14 thank all the speakers but I'd also like to ask
15 them to come up to the front. We started a little
16 bit late and there were two hours assigned to the
17 session and so, we have about 10 minutes for a
18 brief panel discussion.

19 And so, the two questions that I'd like
20 to put to the panel, one of them is, which the
21 speakers can be thinking about as they get up, one
22 of them is where do they see the most immediate

1 promise for these techniques in biological product
2 evaluation? Where do they see the most immediate
3 applications where these things are going to be
4 used first based on everything that they've seen?

5 And the other question is a more
6 technical one which is that we heard from several
7 different speakers about different methods of
8 analyzing next generation sequencing data. We
9 heard about the value of casting a very wide net
10 using programs like Tblastx and we heard from, I
11 think, Charles was the winner in terms of his SNAP
12 program which was able to analyze data very, very
13 quickly. And of course, we all know that the
14 faster you analyze your data, there may be some
15 tradeoff in what it is you're looking at versus
16 not.

17 And so, I guess the other question that
18 is -- and then, we had the BLAT discussion which
19 is somewhere in-between, I think. And so, the
20 other question that I would like to sort of hear
21 some broader discussion on is where is the right
22 location in this intensity of analysis tradeoff

1 for this particular application, what are the
2 parts of the analysis that we can perhaps do
3 without that may be present in some of the more
4 complicated and time intensive programs? And how
5 does one arrive at the optimal algorithm for
6 analyzing these kinds of data in a virus detection
7 setting?

8 So maybe we can start with that second
9 question since it's fresh in everybody's mind and
10 then, as we close then people can give their
11 vision for the future.

12 DR. THOMPSON: Okay, all right.

13 DR. KHAN: Your microphone is not on.

14 DR. THOMPSON: There we go, ah hah. So
15 the question was?

16 DR. KRAUSE: The question is is that
17 there are many different programs that do the --
18 and algorithms that do the sequence analysis. How
19 do you decide among them in terms of the time
20 versus intensity of sequence analysis tradeoff and
21 arrive at the right degree of intensity of the
22 algorithm for looking at the data?

1 And what tradeoffs, in particular, can
2 easily be made with this kind of data that maybe
3 you couldn't make with other kinds of sequence
4 analysis?

5 DR. THOMPSON: So we felt that working
6 at the read level was going to give -- offer us
7 the most sensitivity and that's where we started.
8 And then we had, you know, long- term plans to
9 move into the de novo space. And I think, you
10 know, some people may be beating us to that and I
11 think we're going to be looking -- we're going to
12 be agnostic about the algorithms and look at what
13 other people are doing with de novo and see if the
14 advantages really outweigh the potential
15 detriments in that.

16 So, you know, I think we do have to pit
17 these different approaches against each other and
18 figure out what the pros and cons of each are.

19 DR. CHIU: Yeah, I mean, I think that
20 what's really needed with respect to centerization
21 is actually more precise benchmarking of
22 algorithms. And this is a big problem because

1 everyone is using their own algorithm. They're
2 even cobbling the algorithms together into
3 pipelines and yet, there's very little published
4 data on actually benchmarking, you know,
5 algorithms even ones that are commonly used. And
6 I'm actually quite surprised by how little I
7 actually see about comparisons.

8 And this is important because it makes
9 it impossible to determine which algorithm is
10 better and which one's worse or for different
11 characteristics. And the truth is that there's no
12 one algorithm which is the best algorithm. It
13 really depends on the application.

14 And so, the second point I want to make
15 really quickly is that I think it's really
16 important that this is not a yes or no answer.
17 That there's one algorithm that's the best for
18 every possible application and this is precisely
19 why I think it's really important to discuss how
20 we're going to construct pipelines and how we're
21 actually going to analyze the data.

22 Because I actually think that a tiered

1 approach would probably be the best at least for
2 adventitious agents where you can use a very
3 precise and sensitive algorithm to get rid of,
4 say, host and then use kind of a more broad
5 encompassing algorithm maybe translate nucleotide
6 alignment to identify novel agents.

7 DR. KRAUSE: So can you, Charles, just
8 add a little bit also about -- because I was
9 struck by how fast your SNAP algorithm is. And of
10 course, that may have to do with your amazing
11 computational power. But it sounds like probably
12 you're making some tradeoffs in that algorithm.
13 Can you say a little bit about how that works?
14 What is it doing differently from more
15 time-intensive algorithms?

16 DR. CHIU: I mean that's a good
17 question. Basically, it's -- you would -- the
18 best way actually would be to speak with Bill
19 Bolosky from Microsoft who actually developed the
20 algorithm. He actually gave a talk where he
21 talked about the details and if anyone's
22 interested in specific details, there are some

1 changes to the algorithm beyond what's typically
2 done with other kind of seed-based algorithms.

3 It is a seed-based algorithm. So it is
4 a k-mer based algorithm and I can talk about that
5 in more detail. But I think the tradeoffs are
6 that it is still a nucleotide aligner meaning that
7 you still need to have very good alignment. It
8 can tolerate some degree of mismatch but not a
9 huge amount of mismatch.

10 So you can set that as the editus. And
11 so, it's important to realize that that's not the
12 single algorithm that I think is going to work in
13 all cases. I mean, clearly it doesn't do
14 translated nucleotide alignment, for instance.
15 Which is -- or, like, it doesn't do anything,
16 like, similar to BLASTX or Tblastx and there needs
17 to be. And that's why we really use it to
18 identify known sequences upfront and then, we use
19 a different aligner such as BLASTX for novel
20 sequences.

21 DR. KRAUSE: Okay, Robert?

22 DR. CHARLEBOIS: So we were very

1 concerned about using anything, any aligner that
2 was too stringent like BLAT or anything else. And
3 so, we wanted to be very conservative initially.
4 We can change our minds later. And so, we're
5 using BLAST with a very low, like, $1E$ minus 5
6 threshold.

7 And we compensate for that later on by
8 doing some statistical analysis on the hits to
9 clean it up and we get quite good results. But
10 the thing is, like, for example, I had done an
11 alignment of all of the members, all of the
12 current members of a viral family. And the
13 longest contiguous stretch that they had in common
14 was two nucleotides.

15 And that is alarming if you're thinking
16 about potentially finding a new member of that
17 particular viral family. So I mean, individuals
18 will have, you know, stretches that are more
19 common. So we wanted to pick something very
20 relaxed to start and then, clean up the signal
21 after in a subsequent step.

22 So we do have very different approaches

1 to things and they may all end up at the same
2 place in the end. And that's going to be very
3 interesting to find out. So our different
4 approaches -- so for clinical samples, time is
5 really critical. For a manufacturing thing where
6 you want to release your product in three months
7 from now, it might not be so critical. Of course,
8 we still want to do it faster and we're going to
9 explore faster ways but we wanted to start
10 ourselves with a base that was non-controversial
11 using a, you know, really standard tools like
12 BLASTN and BLASTX that nobody would argue with.

13 They might say, oh, it's really slow and
14 we'd say, yes, it is. But then, and then,
15 learning about these new techniques like I read
16 the paper on serpa and it's quite exciting. But
17 I'm a little scared because I want to see it head
18 to head, you know, to see if it misses things.

19 DR. KRAUSE: Is there something you'd
20 like to add, Arifa? You can say no.

21 DR. KHAN: I agree with everything
22 that's been said. I think it is challenging. I'm

1 really excited, you know, to hear, you know, the
2 different, you know, people sharing their
3 information. I think this is really what's
4 important is to openly discuss, you know, the
5 different approaches people have.

6 So as Robert said, you know, so at the
7 end of the day we can hopefully make a
8 science-based decision whether all or any one is
9 or which one may be applicable for certain
10 situations. That's all.

11 DR. KRAUSE: Well, great, so I guess
12 we're almost at the end of our time here.

13 DR. KHAN: I think there's a question
14 from the --

15 DR. KRAUSE: Oh, there is a couple of
16 questions from the audience though, so why don't
17 we just take those. We'll start there and then
18 we'll go there and then, we will finish up.

19 SPEAKER: I'll be loud, it's okay.
20 Okay, thank you. So NT is doubling very, very
21 quickly almost. So the BLAT and the BLAST-based
22 approaches, I mean, yes, you know, they have been

1 the standard for many, many years. So do you see
2 within the next year, for example, having huge
3 challenges to continue? So what is the timeline
4 by when you think things needs to be resolved that
5 which pipeline is maybe one of the better ones and
6 is it scalable?

7 DR. KRAUSE: Answers? I think obviously
8 we're faced with a field that is moving very
9 rapidly and what we say today almost surely will
10 not be true even a year from now. And it's also
11 the case that there are many sequencing projects
12 that have important sequencing that aren't even
13 present in NT.

14 And so, if you cast a wide enough net
15 even now things get very complicated. And I think
16 all we can do is as people in the field is work
17 together to figure out how this can be managed.
18 And that's the importance of meetings like this.
19 To the degree that viral databases can be culled
20 from the genomic databases upfront and so, work
21 can be done that doesn't need to be repeated each
22 time, that obviously can be helpful also.

1 But I think we all recognize that and as
2 with all of big data, as the data continues to
3 increase potentially exponentially, the complexity
4 of these issues keeps increasing.

5 DR. KHAN: I also wanted to mention, I
6 mean, that's a very interesting -- it's a very
7 important point but I think it does indicate the
8 need for others to get involved. And so, these
9 challenges are not only for us but I think perhaps
10 we'll get others to develop further tools that
11 would help us facilitate the analysis and meet the
12 challenges. So I'm hopeful that the discussions
13 here will maybe go out broader to others who are
14 non-biologists maybe more, you know, computer
15 experts or bioinformaticians or who can help
16 develop tools that might address these current and
17 future challenges.

18 DR. KRAUSE: So Robert had a comment
19 also but we need to pass this down to him.

20 DR. CHARLEBOIS: So we can never play
21 catchup with the sequence data and we may need, in
22 the future, to adopt a heuristic approach. And

1 so, Rodney Brister yesterday talked about, you
2 know, the -- and also, others have talked and even
3 yourself have talked about having a representative
4 genome. And so, if we, I mean you could
5 potentially miss something but if you have a
6 representative genome based on sequence space
7 definitions, you could first see, you know, can it
8 potentially match something in that area. And
9 then, you focus in on the next level to get more
10 detail.

11 And so, heuristic approaches always have
12 costs and tradeoffs but it might be a way to just
13 organize the biological sequence space so that we
14 can deal with this explosion of data in the
15 future. I think there's -- we're going to have to
16 pay that price, I think.

17 DR. KRAUSE: So that's a very thoughtful
18 answer. And I see there's one more question. I
19 also see somebody standing by the stage with a
20 hook. So we'll go with this question and
21 hopefully just one or two very quick answers to
22 it.

1 SPEAKER: Do you see any platform-based
2 biases that might be introduced particularly in
3 relation to the choices of parameters that you use
4 for each of those steps in your processing?

5 DR. KRAUSE: Short answer? Yes.
6 Anybody else?

7 DR. THOMPSON: Well, I think that's
8 exactly what we're trying to sort out by beginning
9 to compare different pipelines and different
10 choices that people are making.

11 DR. KRAUSE: Okay, well, thank you very
12 much to the audience. Thank you very much to all
13 of the speakers. This has been a great session
14 and it's part of a great meeting. So thank you.

15 DR. KHAN: So our next group is the
16 clinical biomarkers and personalized medicine that
17 is charged by Dr. Eric Donaldson and co-chaired
18 by Dr. Raja Mazumder. Our first speaker will be
19 Eric Donaldson, so the stage is yours.

20 DR. DONALDSON: All right. Good
21 morning, and welcome to the Clinical Biomarkers
22 and Personalized Medicine session.

1 I'm going to give a brief introduction
2 to the session and give a little bit of background
3 on two new drug applications that were submitted
4 last year, reviewed, and approved that use
5 next-generation sequencing for their resistance
6 analysis. These came through the Division of
7 Antiviral Products and it's sort of a launching
8 point for our discussion.

9 First, a disclaimer that myself and the
10 members of the panel do not represent official FDA
11 policy. This is an ongoing discussion and we're
12 here to learn. The outline for the talk this
13 morning is I'll give about a 10-minute
14 introduction. I've asked each of our
15 distinguished speakers to give about 15 minutes,
16 and then the idea being that if there's time at
17 the end we'll have 20 minutes for discussion. I
18 know that we're already running late in time so if
19 it turns out that we don't have time, then I have
20 a question that I'd like to pose to all speakers
21 and everybody in the audience at the end.

22 But I wanted to start this session with

1 just a few definitions so that we're all on the
2 same page. Now, the official FDA definitions are
3 listed in the guidance, which is bullet point
4 number four on this slide. But just for the sake
5 of conversation this morning, personalized
6 medicine is customization of health care based
7 upon analyses of individual biomarkers to identify
8 the most relevant therapies for an individual. We
9 heard this described yesterday as precision
10 medicine, which is kind of along the same lines.
11 The biomarker is a measurable indicator of a
12 biological state or conditioned response, and a
13 companion diagnostic is a medical device or assay
14 that identifies patients for a specific treatment
15 and/or dose based on the discovery of those
16 biomarkers that are predicted for the drug
17 response, and companion diagnostics are paired
18 with a specific drug and often approved together
19 in collaboration with CDRH and CDER.

20 In clinical practice, personalized
21 medicine is generally used to select subgroups of
22 patients that are most likely to respond to a

1 specific treatment, and there are some examples of
2 these, particularly in the field of oncology,
3 where a number of important biomarkers determine
4 which drug is used for a specific type of cancer.
5 For example, EGFR mutations to treat lung cancer.
6 Some of our speakers will go into details about
7 these particular ones in more detail in their
8 talks so I won't spend a lot of time there.

9 Another example is for cystic fibrosis
10 with (inaudible), which is indicated for gating
11 alleles -- one or more of the following gating
12 alleles that are listed there.

13 Now, I've included two antivirals at the
14 bottom. These are not specifically examples of
15 personalized medicine in a sense that they don't
16 have a companion diagnostic to go along with them.
17 However, both of these have viral determinants
18 that determine which drug or which patients or
19 which patients that are infected with which
20 viruses should receive those drugs. For example,
21 Maraviroc is for treating HIV-1 infected patients
22 with the CCR5 tropic virus only, and Simeprevir,

1 which was approved last year, is for treating HCV
2 patients, Hepatitis C virus patients with genotype
3 1, but with the genotype 1A, subjects with Q80K
4 present at baseline in the protease do not respond
5 to Simeprevir, so it's actually recommended in the
6 label that if you're infected with genotype 1A,
7 you be screened prior to receiving Simeprevir to
8 see if that polymorphism is present.

9 These are examples of viral determinants
10 that are biomarkers, but viral resistance is also
11 a biomarker that we're very much interested in in
12 the Division of Antiviral Products. And there are
13 several viruses -- I've just listed a few here --
14 that are important. I mentioned the Q80K
15 polymorphism, but in addition, R155K is a common
16 resistance- associated substitution to the
17 Hepatitis C virus protease, and if a subject has
18 that or a patient has that at baseline, they've
19 actually lost the entire class of HCV proteases
20 that are approved. So that's an important
21 biomarker for understanding whether or not a
22 patient will respond to a particular therapy.

1 With Sofosbuvir, which is a polymerase inhibitor
2 that was recently approved last year, L159F and
3 C316N may be important markers present at baseline
4 that predict failure with Sofosbuvir. That's
5 still yet to be determined.

6 Human immunodeficiency virus, baseline
7 resistance profiles are often taken for treatment,
8 experienced patients who have failed heart, to
9 establish which regimens they will respond to
10 currently. And then of course, with influenza A
11 virus, the H275Y substitution and the H1N1
12 neuraminidase is predictive of resistance to also
13 Tamivir.

14 So identifying resistance-associated
15 substitutions and baseline resistance
16 polymorphisms is essential for determining the
17 safety and efficacy of a drug and for identifying
18 subgroups of patients that will benefit the most
19 from using that drug.

20 So in our division, we take that very
21 seriously, and in general, the sponsor will
22 provide all of the information on resistance

1 analysis from the cross clinical trials in the
2 format of the summary. Our division also requests
3 that they send the raw sequence data, which to
4 this point has been Sanger data, Sanger population
5 sequencing, and generally what we ask for is a
6 baseline sample and a time of failure or multiple
7 times of failure after failing a treatment or a
8 regimen in clinical trial so that we can compare.
9 And we're asking for sequences of the target. So
10 for a polymerase inhibitor, we're asking for the
11 polymerase gene.

12 And then our viewers, the virology
13 viewers, actually do an independent analysis of
14 the resistance data because we feel it's extremely
15 important to understand exactly what's going on
16 and we want an independent assessment of
17 resistance so that we can describe that accurately
18 in the label.

19 So as you can imagine, next-generation
20 sequencing adds a whole level of complexity to
21 this resistance analysis. And we're seeing within
22 our division, it started with these two

1 submissions. We now know of at least 25
2 submissions that are coming in the next year or so
3 that will include next-generation sequencing,
4 driven much by cost as Carolyn Wilson pointed out
5 yesterday.

6 But what we love about next-generation
7 sequencing is it allows you to take a much deeper
8 look into the data. So with population
9 sequencing, you get sort of an average consensus
10 of what's in the sample. With next-generation
11 sequencing, you get an opportunity to look at the
12 minor viral variance in the population, and so you
13 might be able to identify baseline factors that
14 are present before treatment that will emerge and
15 later become associated with resistance.

16 But, of course, with next-generation
17 sequencing comes the technology challenge, which
18 really a lot of times comes down to
19 interpretation. And so with Sanger you have a
20 nice trace that is arguably one interpretation for
21 one sample. With next-generation sequencing, you
22 have millions of viral sequences in the tube,

1 hundreds of alignment algorithms, variable
2 filtering criteria and assemblies, no standard
3 analysis pipeline, which means many different
4 possible interpretations.

5 What the crux is for us is with variant
6 detection, and we heard yesterday that somebody is
7 using three different variant detection methods.
8 And the idea is that we want to try to identify
9 real variants versus those that may be sequence
10 error or introduced by PCR or any number of
11 different ways. The problem is that without a
12 standardized analysis pipeline, different sponsors
13 are using different variant calling techniques,
14 and these all were designed for different purposes
15 in some cases. There's very little
16 standardization of cross methods for calling these
17 variants with little overlap among the different
18 callers. In fact, Ara Attal did a study recently
19 comparing five different variant detection systems
20 and they were sequenced by alumina 15XMs from four
21 different families and then compared the five
22 different methods and found that there was 57.4

1 percent congruence across those concordance. It
2 was even worse when it came to insertion
3 solutions; it dropped down into the mid 20 percent
4 concordance. And each one of these methods
5 actually identified 0.5 to 5.1 percent of unique
6 variants of their own.

7 And so, as you can imagine, if you're
8 comparing the results from one detection system to
9 the other, there may be major discrepancies. And
10 when it comes to identifying those key variants
11 that are important for resistance, we want to make
12 sure we know what we have and what is the truth.
13 And Oscar Wilde puts it best, "The truth is rarely
14 pure and never simple," particularly when it comes
15 to next-generation sequencing.

16 So when these submissions were coming,
17 we had several questions that we asked, and we
18 have attempted to find ways to work through these
19 solutions. And for experimental design, we wanted
20 to know what exactly is being sequenced. Is the
21 target amplified? Is it multiplexed? Are there
22 minor variants preserved? Are there experimental

1 artifacts? Is there something in the way that
2 they're preparing the samples that might leave
3 behind important information? Or is there
4 something in the sample that shouldn't be there?
5 And in fact, with those two submissions, we saw
6 examples of PCR contamination across contamination
7 of multiplex samples, both of which could have an
8 impact on whether or not you call a minor variant
9 real or not.

10 In addition, for the data analysis, how
11 is the data being analyzed? Is it reproducible?
12 What programs or parameters are being used? What
13 we learned from this first early experience is
14 that both sponsors used entirely different
15 pipelines and many of the components of the
16 pipeline were proprietary so that we did not have
17 access to those.

18 And then when it comes down to
19 regulatory review, how do we make a decision based
20 on the data? What data are needed? Are the
21 summary data from one analysis pipeline sufficient
22 to do this? And we felt early on that this was

1 not. Knowing of the issues that I've described
2 earlier, we felt that it would be important to
3 have something to compare it to. And so our
4 strategy was to come up with an internal analysis
5 pipeline that we could then ask the sponsor for
6 the fast queue sequences, do a quick analysis of
7 the data to see if our results match those of the
8 sponsors, and to identify any discrepancies in
9 that way.

10 That then introduced another problem --
11 data transfer and storage. These are large files,
12 a terabyte or more. So how are we going to store
13 these as an agency? How are these data stored
14 long term? Is encryption required? How is the
15 data integrity managed? And we have to think
16 about this in the context of regulatory data which
17 is stored for 30 years in some cases, so how do
18 you maintain data integrity over a 30-year period?
19 And then, of course, data confidentiality. How
20 are we protecting the proprietary and patient
21 information associated with this? And so this was
22 really a major collaborative work to put together

1 this analysis pipeline. And I'm just going to
2 quickly show that it was a collaboration between
3 CDER, CBER, CDRH, with support from the Office of
4 the Chief Scientist.

5 So now that was our experience as a
6 division with two applications that came in last
7 year. We know that this field is exploding and
8 expanding, and so we've gathered a distinguished
9 panel of speakers who we've asked to discuss from
10 their perspectives what the different issues are
11 with next-generation sequencing and clinical data
12 or drugs.

13 And so I want to just introduce the
14 panel now. Dr. Andrea Ferreira-Gonzalez is
15 Professor of Pathology and Director of the
16 Molecular Diagnostic Laboratory at the Virginia
17 Commonwealth University Medical Center. We have
18 Dr. Andrew Grupe, Senior Scientific Director from
19 Quest Diagnostics. Dr. Charles Sawyers will be
20 joining us, hopefully, remotely. His airplane was
21 canceled. His flights out of New York were
22 canceled, and so his slides are here. He's going

1 to try to call in and present the presentation.
2 But he's the Chair of Human Oncology and
3 Pathogenesis Program, Marie-Josée and Henry R.
4 Kravis chair of the Memorial Sloan-Kettering
5 Cancer Center. And then Dr. Laura van't Veer,
6 leader of the Breast Oncology Program and
7 Associate Director of Applied Genomics from UCSF
8 Helen Diller Family Comprehensive Cancer Center.

9 So we're very happy to have this panel
10 of distinguished guests. Our first topic will be
11 presented by Andrea Ferreira-Gonzalez, and her
12 topic is going to be "Next-Generation Issues
13 Related to the Clinical Lab."

14 MS. VOSKANIAN-KORDI: Just a small note,
15 due to logistics issues since Dr. Sawyers is doing
16 it remotely, we're going to do him last.

17 DR. FERREIRA-GONZALEZ: Thank you very
18 much. I would like to thank the organizing
19 committee, but also the Association of Molecular
20 Pathology for inviting me here today to share some
21 of the issues that we currently have in the
22 clinical laboratory as we apply next-generation

1 sequencing technology.

2 My laboratory is currently using this
3 technology for the identification of mutations or
4 genetic changes, specifically for
5 cancer-actionable genes. So that's how I'm going
6 to do the context of my presentation.

7 I can pass the slides. So let me
8 introduce first why we have started incorporating
9 next-generation sequencing technology in the
10 clinical laboratory. And again, specific for the
11 cancer application. Mutation screening has been
12 an integral part of the analysis of oncology
13 patients for many, many years, and in the past,
14 the molecular diagnostics laboratory has had some
15 interventions in helping the diagnosis,
16 identifying the different markers or prognostics,
17 or even detections or determining what is the best
18 therapy for patients, but it was very limited.
19 And the reason it was very limited is because
20 cancer is a very complex disorder with many, many
21 different genetic alterations and many different
22 genetic genes. And it was very difficult for the

1 laboratory to have a comprehensive way to look at
2 all these different changes at one time.

3 With the advent of next-generation
4 sequencing, now we have a tool that we can apply
5 to really look comprehensively at what is the
6 genetic status of the particular cancer. And what
7 we're seeing is that more and more every day we
8 have at least a new gene that has been associated
9 with different cancers. And in the past, we have
10 actually associated specific gene mutations to
11 specific cancers, and the cancer is in a specific
12 location. And as we do more and more
13 next-generation sequencing, what we're starting to
14 identify is that actually there is a huge amount
15 of commonality between some of the genetic changes
16 that occur in a different multiplicity of cancers
17 and we are starting to see a movement towards the
18 new description of cancer in a more molecular way.
19 Very similar to what is happening in the
20 hematologic oncology processing where the WHO is
21 starting to classify them using these molecular
22 targets.

1 So I put here a list of the different
2 mutated genes that have clinical utility, and I'm
3 not just talking about clinical validity, but
4 actually clinical utility. So finding the status
5 of that particular gene in a particular cancer or
6 condition is very important to not only helping
7 the diagnosis, prognosis, but also in the
8 therapeutic determination of where these patients
9 are. And many of these markers of each tumor has
10 been identified, that play a very specific role in
11 very specific cancers, but now we're starting to
12 see that there actually might be a role for some
13 of these other different cancers.

14 As we continue to the application of
15 these technologies, we are also faced with the
16 decrease in the invasiveness of the procedures
17 that we are putting our patients through. Fifty
18 percent of the specimens can detect cancer or
19 laboratory today for doing cancer genotyping
20 analyses or fine needle aspirations, and due to
21 the fact that we use fine needle aspirations, the
22 tissue that we're getting to the laboratory are

1 extremely small, are extremely limited, that then
2 we have much DNA that we can extract, allowing us
3 to do a large number of different tests with
4 different platforms. We have to start having to
5 consolidate and get the most information that we
6 got.

7 But with that, it's also very important
8 that when we have single gene testing today, we
9 have a large number of different platforms that
10 can be used that require a large amount of DNA, in
11 some instances we are not available in the
12 laboratory. And we are seeing an increase in the
13 test volume for these different clinical
14 applications due to the fact that we have this
15 amount of information now with clinical utility.

16 So let me show you what has happened in
17 my laboratory in the last year. In the past, we
18 used to have three different tests for the
19 detection of KRAS, EGFR, or BRAF, and these are
20 laboratory-developed procedures, the first one and
21 the second one, and this is a cryogen procedure
22 for the detection of EGFR mutations that require

1 if you needed to ascertain these three particular
2 genes in one patient sample, a large amount of
3 DNA. And in some instances, we didn't have enough
4 to be able to do the three tests, and we had to
5 choose which one we would do first, and that was
6 the only information we can get from patients.

7 Earlier this year we developed,
8 implemented, validated a test that is commercially
9 available using the Ion PGM Cancer Panel v2 that
10 allows us to ascertain the mutational status of 50
11 different genes that could be applied for
12 different types of cancers. As we implemented and
13 validated these types of assets in the laboratory,
14 we went through a very rigorous procedure. To
15 date, even though we ascertained 50 different
16 genes, we only have reported seven for specific
17 cancers being colon, lung, or melanoma due to the
18 clinical utility of these genes has been
19 demonstrated for a limited number of these genes.

20 As we brought the clinical testing into
21 the laboratory, we followed their rigorous
22 procedure for not only test development but also

1 test validation. And as we went through the test
2 validation, we were able to obtain a significant
3 amount of information that was necessary to
4 develop the quality management and that was then
5 implemented to monitor the quality of the testing
6 as we move forward, but also maintaining the
7 accuracy of the testing that we do. During the
8 test development -- we were able to do the test
9 development, and where you do test validation,
10 sometimes simultaneous or not, but we had to
11 validate not only the wet part of the laboratory
12 testing but also the dry part of the laboratory
13 testing which I consider the pipeline.

14 So we did a validation of the pipeline
15 by itself and then we did validation of the wet
16 part of the laboratory and incorporated the
17 validation of the pipeline and then validated the
18 entire process together from the wet part and the
19 pipeline together, and then we had to do
20 modifications as we went along until we got to the
21 right type of conditions necessary to meet the
22 specifications that we require from the design of

1 the testing that we're bringing to the clinical
2 laboratory.

3 We were lucky to start to see a merge in
4 number of different guidelines for different areas
5 of testing, starting with inherited disorders, but
6 we're starting to see some evolving for cancer
7 applications. But these are very important. The
8 CDC started working early on, and we're going to
9 hear more from Dr. Ira Lubin, who is present here
10 today. The Division of Laboratory Sciences and
11 Standards got involved very early in developing an
12 extraordinary successful program called the Gate
13 RM Program, but also convene a group of experts to
14 start looking at some of the issues necessary for
15 the validation of these next-generation sequencing
16 applications, and what are some of the
17 recommendations or some of the practices that are
18 starting to emerge?

19 MM9 is a great document from the
20 Clinical Laboratory Standards Institute that
21 allows to provide some information; also, guidance
22 on how to bring next-generation sequencing. And

1 there are other ones, too. The Association for
2 Molecular Pathology, which I'm a member of the
3 Whole Genome Analysis Workgroup, has been working
4 very diligently in trying to deal with these
5 issues, and we have produced several white papers
6 and are working right now on another white paper
7 on how to validate or best move forward with this
8 type of testing in the clinical laboratory.

9 But as we go through the testing, I'm
10 going to show one of the cases that we did in the
11 laboratory just to show you why it is so needed
12 this type of testing. And as you can see here, we
13 got a final aspiration specimen from a 79-year-old
14 Caucasian female, who had a history of right lower
15 lobe nodule, and she was actually followed for a
16 year and a half. She had a negative PET scan that
17 then continued to evolve and change in size, and
18 after the final aspiration procedure, there was a
19 diagnosis of invasive lung carcinoma which was
20 modularly differentiated, but the point to put
21 here is the size of the actual cells that are
22 actually needed to be analyzed within the final

1 aspiration specimen. Just to give you an idea,
2 this is a glass line and this is the marked area
3 by an anatomic pathologist letting us know that
4 this is the area that we need to analyze within
5 the tissue, that then we need to extract the DNA
6 and apply it to the next-generation sequencing.

7 And this is important because with that
8 amount of tissue, we could do EGFR testing,
9 determine if these individuals could be eligible
10 for some of the molecular type drugs that have
11 occurred in living proof. But there is a large
12 number of other genes that have the mutation, and
13 some of them are exclusively -- one you have a
14 mutation and one you don't it in the other one.
15 But if we have to apply this test for the EGFR
16 that we had last year, from this specimen we
17 couldn't get enough DNA to actually be able to
18 analyze it. Some people recommend using KRAS
19 first as a triage and then go to the EGFR, but I
20 think it's better just to see if you actually have
21 EGFR.

22 So we did apply the next-generation

1 sequencing assay to these particular patients, and
2 we actually found a number of different mutations.
3 And as you can see here, even though we do test
4 for 50 different genes, we only report seven of
5 those. And you can see that the patient had an
6 AKT1 KRAS mutation. There were pathogenic, and it
7 had a BRAF mutation that was likely pathogenic.

8 So this individual still could receive
9 some of the molecular-targeted drugs for EGFR, but
10 at the same time, if this didn't work, we already
11 know that there are other drugs that are going
12 through clinical trials that might have relevance
13 for the AKT1 mutations and so forth. And it gives
14 you a complete picture of the status of that
15 particular cancer sample.

16 As we go through the data analysis, and
17 this was kindly provided by Dr. Lubin, again, he
18 presented these at the annual meeting of the
19 Association of Molecular Pathology last year, we
20 go through a different number of steps, as we do
21 the data analyses, and through the validation of
22 these particular tests, we have to validate each

1 of the different steps of this pipeline alone, and
2 then we have it with the wet all together. We
3 actually multiplex or barcode our specimens. We
4 can run up to seven specimens per sequencing
5 reaction which allows a higher throughput, but it
6 gives us enough depth in the reading of the
7 electronic file that gives us a sensitivity of 4
8 percent. We have to have 4 percent of the DNA
9 that we analyze mutated to be able to be detected.

10 So you have to really determine or
11 validate that your algorithms to be able to
12 demultiplex the data -- see, all the data coming
13 out, you have to demultiplex what reads correspond
14 to each of the different patients and make sure
15 there is no crosstalk. That needs to be
16 validated. And then the secondary alignment,
17 variant calling, and then from there, the tertiary
18 analysis and where you have the variant alleles
19 and you have to annotate these and prioritize
20 them. And then the final clinical interpretation
21 which is the most labor-intensive, to be honest
22 with you.

1 And then again, I'm not going to go into
2 detail on this, but you can see here the different
3 steps that are needed to go through and all the
4 information of the genes that have to be annotated
5 to remove those genes not associated with the
6 presentation at that time. You have to have
7 patient information, not only the phenotype but
8 also some of the variants, and you start having to
9 annotate. That allows you to remove variants that
10 are unlikely to alter in function, and then you
11 have to combine all this information to develop at
12 least the prioritized list, and then compare it to
13 the clinical grade assessments and clinical result
14 reporting. It takes a longer time to do the data
15 analysis sometimes actually than generating the
16 wet laboratory data.

17 There are annotation databases that we
18 use continuously. And as you can see here, I put
19 a very limited number of different databases that
20 we can use to start trying to get information.
21 So, the ultimate information, is this change
22 clinically relevant for the current presentation

1 of my patient?

2 We use a large number of different
3 databases -- dbSNP, Clean Bar, HGMD, Cause Big, My
4 Cancer Genome, but we have to use all of them. We
5 cannot use a single one because all of them have
6 caveats. Some of them are better to ascertain
7 different changes and the information might not be
8 updated and so forth. But it takes a large number
9 of these annotation databases to be able to build
10 confidence in some of the data that we're trying
11 to produce.

12 So we do have some interpretation issues
13 when we get to that point. In some instances,
14 there is very limited evidence on the clinical
15 utility of specific mutation for that particular
16 cancer that we're seeing it. There are few
17 mutations that are listed in consensus management
18 guidelines. Institutional and national efforts to
19 gather and cure it are just starting. For
20 example, some of the efforts that were mentioned
21 earlier today at the NIH label, My Cancer Genome.
22 Sometimes we don't know what the clinical

1 significance or well study mutations in different
2 tumor types. And as we start doing genotyping
3 with a certain number of genes, we're starting to
4 see these mutations coming up new, and we don't
5 know what the significance is in that particular
6 organ. Actually, the clinical significance in
7 other mutations in targetable genes is very
8 important. For example, unreported mutations in
9 EGFR, they want to confer susceptibility to the
10 different molecular targeted drugs that we see.
11 So this is the interpretation of the information
12 that we need to start trying to get a handle on.

13 As we move clinically, there are a
14 number of different missing pieces in the
15 different pipelines that we use and even the
16 laboratory information systems. We don't have
17 good tools to actually very rapidly and easily
18 determine the adequacy of the amplicons. How the
19 amplicons are working and actually how your PCR
20 reactions, for example, in our case, are working
21 without having to really start digging into the
22 info informatically. There's not an easy way to

1 ascertain these. There's no easy way to visualize
2 the results as we go through the data analysis.

3 Conversion of variants into standardized
4 nomenclature is highly -- I think is highly
5 desirable. We are supposed to be reporting all
6 the gene changes used in standardized
7 nomenclature, so I think it's very important that
8 as we look at these pipelines, that we start using
9 the common language that we have to use
10 clinically, which is standardized nomenclature.
11 We have to be able to monitor frequency of
12 specific mutations in our own samples. We have to
13 develop our own databases. There is not a
14 commonplace where we can go start depositing some
15 of our information to monitor over time what
16 happens with these changes and then how they
17 compare with other sites. Their link into genomic
18 databases is also very important.

19 Inclusions or exclusions of variants for
20 the reporting interface with the laboratory
21 information system. Today, we have no way to
22 represent this information within the laboratory

1 information system or the electronic medical
2 record, so we are actually putting PDF or paper as
3 an image into the electronic medical record that
4 cannot be leveraged. We cannot develop clinical
5 decision support tools and so forth from this
6 information.

7 And then generation of clinical reports
8 replaces in different forms of these reports and
9 there is not clear guidance or good guidance to be
10 able to determine what should we put into the
11 different clinical reports.

12 Additional problems, as we see that as
13 we come up with the current technologies, it's
14 very difficult to start aligning with confidence
15 some of these reads. So the long range actually
16 is going to help as we move forward. But it's
17 very difficult to identify variants, so we don't
18 have very good tools to be able to identify
19 variants.

20 The different data formats is a
21 limitation today. How do you describe the quality
22 of data generated, changes from the different

1 platforms from the different vendors, so sometimes
2 it's hard to know when you're talking to
3 individuals using different platforms, when we
4 talk about the quality, just something basic as
5 the quality of your data, but how do you actually
6 talk apples to apples?

7 I think some of the opportunities that
8 we have is actually looking at interoperability of
9 prodigals conference, where we can call get
10 together and look at the different prodigals that
11 we're using for the different steps and then come
12 up with maybe common prodigals or recommendations
13 to which we can actually use and then start
14 comparing.

15 I think also the development of
16 benchmark tools for variant call might be really
17 relevant or important to move forward, where we
18 can actually all work together to develop this
19 tool. I know the NIST is having a significant
20 effort in trying to work towards that goal in
21 developing this type of benchmark tool that we can
22 use.

1 As you can see here in the little graph
2 on the left, you can see the different software
3 that you can actually use to create or generate a
4 bond file. And as you can see from previous
5 speakers, there are even more than this. And they
6 produce different formats, different information
7 that it's very difficult to actually then start
8 comparing between the laboratories.

9 So one of the challenges I think will be
10 standardization of all these different topics that
11 I just mentioned. But I think also it will be
12 important to look at difference reference
13 materials, because as we develop the tools and we
14 try to come up with a common language, we have to
15 be able to prove that these common languages or
16 tools can talk to each other or actually produce
17 accuracy. So developing of reference materials,
18 and not just only materials as a form of DNA or
19 form of a specimen, but also electronic files that
20 can be used as reference material to calibrate all
21 our different aspects. So the publication of the
22 Genome in a Bottle that NIST is working very

1 diligently that I think now you're going to have a
2 speaker that is going to talk specifically to this
3 will be very well received. These would allow us
4 to calibrate the pipelines in the electronic file,
5 but also calibrate the test because we can use
6 materials with the electronic file that allow us
7 to (inaudible) through the development,
8 validation, quality control, and so forth.

9 The reference materials for it I think
10 will be important that we can all get together and
11 talk about what are our needs and what materials
12 can be developed that can be used for validation
13 and compliance for quality control, and I think
14 maybe doing something similar to what I heard
15 earlier in the other group, having different
16 laboratories working together with exchange
17 specimens, exchange electronic files at different
18 levels and see what you obtain at the end, is
19 starting to come up with common censuses or some
20 commonalities that we can all come up with a tool
21 that can be used for the different areas of the
22 testing.

1 And with that I would like to end my
2 presentation. Thank you.

3 (Applause)

4 SPEAKER: I have a quick question in
5 terms of the bioinformatics. How do you version
6 control the annotation databases, like CFNC and
7 dbSNP that change every few months?

8 DR. FERREIRA-GONZALEZ: Yes, it's very
9 difficult to keep that. Not only the version of
10 the software that the manufacture is giving us for
11 the instrumentation, it's a huge issue. The
12 version of the human genome reference. You know,
13 we were 19 to 20 and we had to remap all our
14 algorithms to that. So it's very hard to maintain
15 that and it's very challenging. So we are trying
16 to keep some of the reference on our own local
17 server, but it, again, becomes very challenging.
18 So having to remap to the current version all the
19 algorithms that need to be revalidated, so there
20 is a huge issue there on that.

21 SPEAKER: Do you have a criterion for
22 when you upgrade to a new version of dbSNP, for

1 instance, or do you just automatically take the
2 latest one?

3 DR. FERREIRA-GONZALEZ: We automatically
4 take -- that's what we are doing all the time.

5 MS. VOSKANIAN-KORDI: I have a few
6 questions from the online audience. The first is,
7 how are you evaluating the abundance of
8 false-positive semantic variant calls during test
9 validation?

10 DR. FERREIRA-GONZALEZ: So what we do is
11 we have specimens that we've already
12 characterized. We did for colon, for lung, and
13 for melanoma, so specimens that had already been
14 characterized for very specific hotspots. We knew
15 the sequences around those areas. So as we went
16 through the validation, we knew what we needed to
17 obtain, and then as we started getting information
18 or data out of the instrumentation we started to
19 -- that they were called by the software, we
20 started identifying artifacts. But then we
21 started veiling up masks to block or just ignore
22 to be honest with you as we go through that

1 analysis.

2 MS. VOSKANIAN-KORDI: And, sorry, let me
3 find the second question.

4 What is the decision-making process on a
5 negative finding from a clinical sample? Are
6 validations done or required for a negative
7 sample?

8 DR. FERREIRA-GONZALEZ: So we have
9 samples that are already being characterized for
10 specific genes that were negative. So as we run
11 them through the instrumentation and they come in
12 negative, they're fine. So those samples that we
13 knew were positive that might come up negative,
14 they will have to be confirmed with a third method
15 or another method to see the discrepancy between
16 the two different values.

17 SPEAKER: One of the IT challenges or
18 bits that need to be fixed or filled out that you
19 mentioned was converting a variant into standard
20 nomenclature. It seems a lot of people are using
21 standard nomenclatures, like HgBS or even the
22 variant call format, but the problem isn't that

1 standards don't exist but when you create such a
2 thought, you can create a different nomenclature
3 from the same variant so that two labs producing
4 the same variant can produce different HgbS or
5 different variant call format, and those things
6 are difficult to compare. So the challenge to me
7 isn't really a standard nomenclature but a
8 directly comparable nomenclature that I can say
9 that these two things are identical or not.

10 Can you talk about that a little bit?

11 DR. FERREIRA-GONZALEZ: I think you're
12 bringing up an important issue. One of the things
13 that we see that sometimes will have discrepancy
14 (inaudible) because sometimes you can go by the
15 difference of a base where it has to do with which
16 reference you're using of the sequence. So maybe
17 the nomenclature needs to go back to the reference
18 of the sequence and from the reference version
19 that you have, then how do you compare to one
20 versus the other one. But you have to provide the
21 reference, the specific reference version that
22 you're using.

1 That brings also the point from the
2 previous individual that asked, version control.

3 SPEAKER: Thank you for the talk. In
4 terms of reference materials, what would you
5 consider the optimal reference material? Is it a
6 plasma-based cellular disease mimic or is it
7 something else that hasn't emerged yet?

8 DR. FERREIRA-GONZALEZ: Well, you know,
9 ideally with a patient specimen. Plasma has a
10 certain advantage. You can manage what you
11 actually have specifically in different locations,
12 but it doesn't have all the structure, secondary
13 tertiary structure that you might have from the
14 whole human genome. I think the closest you can
15 actually get to a real patient specimen is the
16 best way to come out to some of this reference
17 material. Sometimes it's not going to be
18 possible, but the closest that we can get to that
19 will be very useful.

20 Using cell lines, again, it's not the
21 best because isolated cell lines are not in the
22 context of the patient specimen, but at least you

1 have the whole genome structure there that you
2 have to deal with on the wet part and even the
3 electronic part of the sequence that you're
4 obtaining.

5 DR. DONALDSON: Due to time, we're going
6 to move on. So thank you.

7 (Applause)

8 DR. DONALDSON: Our next speaker for
9 this panel is Andrew Grupe. He'll be talking
10 about "Sequencing and Truth: Utility for Common
11 Reference Standards."

12 DR. GRUPE: Thanks, Eric, for the
13 introduction. And thanks to the organizers for
14 the invitation.

15 I wanted to take the opportunity today
16 and combine actually two different things. One,
17 talk about our experience at Quest Diagnostics,
18 but also I've been chairing the Reference Material
19 Selection and Design Working Group of the Genome
20 in a Bottle Consortium, so I wanted to combine
21 those two aspects and really talk about reference
22 materials and the utility of reference materials

1 in this section.

2 And maybe a show of hands here. Before
3 Andrea's presentation or before you came to the
4 meeting, who has heard about Genome in a Bottle?

5 All right. So that's maybe 20 percent.
6 And who of you have attended a work session of
7 Genome in a Bottle in the past?

8 So that's maybe five or six of you. All
9 right. Very good.

10 So we all know that sequence information
11 is becoming more and more important in order to
12 make clinical diagnoses or also to make treatment
13 decisions. And in the past, just like Andrea
14 already iterated, that was kind of a linear
15 hierarchical process that was being used where
16 initially, sanger sequencing was used or
17 genotyping, and you had your sample and you had
18 your most likely gene that you wanted to analyze,
19 so you looked at that gene and you worked your way
20 down, looking at gene number two, number 3, number
21 4. And that could entail a long time until you
22 found a result, so the patient had to wait a long

1 time before they were informed.

2 Do I need to do something different?

3 MS. VOSKANIAN-KORDI: It's this one
4 here. We're having an issue here.

5 Is this better? I'm also closer.

6 DR. GRUPE: I can try to get closer to
7 the mic. That's a situation where I'd like to be
8 shorter.

9 So next-generation sequencing technology
10 avoids some of that conflict so you don't deplete
11 the sample, you have a much faster turnaround time
12 to get information back to the clinician.

13 We've seen that the times to find
14 relevant clinical markers and translate them into
15 practice has become shorter and shorter over time.
16 In the beginning, we were talking about decades.
17 Now we're talking about years in order to
18 translate that information. So therefore, I think
19 it behooves us to think about the whole process in
20 terms of analytical validation, clinical
21 validation, and also the regulatory review
22 process, how we can address that timeframe and

1 provide still optimal and timely care for
2 patients.

3 So next-generation sequencing, as you
4 all know, we look at a really large universe. We
5 can generate so much more sequence compared to
6 Sanger sequencing. So validation of that
7 technology becomes much more complex than for a
8 sanger-type of assay or a simple genotyping assay.
9 And so therefore, I think it is really critical to
10 have well characterized reference materials
11 available for this process, and those reference
12 materials, they probably will differ depending on
13 the application that you would like to use. So
14 there will be different reference materials for
15 inherited diseases or oncology if you think about
16 noninvasive prenatal testing or infectious
17 diseases.

18 So I want to spend some time on the
19 Genome in a Bottle Consortium and the reference
20 materials that are being discussed there and being
21 developed. I want to talk about synthetic
22 constructs and some ideas about their utility.

1 And then also simulated data.

2 When you look at the National
3 Comprehensive Cancer Network Guidelines for
4 recommendations on how to treat cancer patients,
5 you see an evolution there in terms of the types
6 of mutations that are listed. When you go back
7 and look into the recent past, you might see a
8 single gene listed for certain cancer types, or
9 you might actually see specific code-ons or
10 specific mutations listed that are relevant. If
11 you look today, that information is really
12 increasing in terms of the number of genes, but
13 now we're not only looking at specific mutations
14 that are in the guidelines but we're looking at
15 whole axons or whole genes, so that adds another
16 complexity level also in terms of the analysis and
17 finding those mutations because you're on longer
18 just limited to having a specific sample with a
19 specific mutation, but you need to analyze the
20 sequence as a whole.

21 So the increase in terms of
22 clinically-relevant genes and mutations, that is

1 also expressed in the number of tests that are
2 being offered. Andrea gave an example. Here is
3 an example of what we do at Quest Diagnostics.
4 Traditionally, and still we're offering specific
5 panels that look at specific cancer types and have
6 specific genes for when we do Sanger sequencing.
7 And then recently, we've also released a panel
8 that looks at 34 genes and analyzes those genes
9 using Next Gen sequencing technology. We all know
10 that Next Gen sequencing technology comes along
11 also with potential errors. We've seen one slide
12 from Eric. Here's another example where groups
13 have looked at the error rate using different
14 platforms and actually also using different
15 genomes to analyze that. So again, stressing the
16 point that if you do sequencing on such a large
17 scale, you really need well very well
18 characterized, high confidence reference samples
19 that are available.

20 Now, when you think about errors within
21 the process, we're always talking about the wet
22 lab process and we're talking about the dry lab or

1 bioinformatics pipeline process. Where do these
2 errors really occur, or what processes can
3 contribute to these errors? In the case of
4 oncology, the tissue is usually formalin fixed and
5 embedded in paraffin, so that can lead to
6 deamination of the sequence, so there is a
7 potential error that gets introduced or changed in
8 the sequence. When you share your DNA, there have
9 been reports that that can introduce oxidative
10 damage.

11 We all know about PCR-based errors. If
12 your sequence, you're variant that you're trying
13 to find is close to a homopolymer or in a high
14 GC-rich region, that's going to cause problems.
15 The error rate will be associated with the REDPs
16 that you have, with the quality of the alignment,
17 the trimming that you do, and then actually,
18 pseudo genes, they also will mess up your data and
19 make your interpretation more difficult. So when
20 you have reference samples, you can strategically
21 place those reference samples and analyze your
22 process by looking at the reference samples and

1 see how they perform. And also, looking
2 downstream in the process when you actually don't
3 have samples for specific mutations, you can think
4 about using a synthetic dataset for your analysis
5 to test your bioinformatics pipeline.

6 So we've been talking about reference
7 samples, but what might be the different
8 requirements for reference samples that you might
9 want to think about? So when I think about DNA
10 sequencing, I want to look at single nucleotide
11 variants. I want to look at small INDELs,
12 probably copy numbers and rearrangements.

13 So when I look at inherited disorders,
14 there is a specific -- well, there are fixed
15 variables that you identify. So if you look at
16 SNPs or single nucleotide variants, they are most
17 likely to be homozygous or heterozygous only. If
18 you look at copy numbers, it's going to be a
19 homozygous deletion or heterozygous deletion, or
20 you're going to have two, three, four copies of
21 that chromosome or chromosomal region that you
22 analyze.

1 And so in the Genome in a Bottle
2 Consortium, which was initiated by NIST and
3 stakeholders from academia, from the government
4 and industry, we're discussing, and we are
5 developing reference materials, reference methods,
6 mostly reference data to really use that
7 information and assess the confidence in human
8 genome variant calls and to enable the performance
9 assessment of clinical genome sequencing.

10 I'm sure Justin, who is here, who will
11 talk in the afternoon, will go into more detail,
12 but the characterization of those reference
13 materials is based on multiple platforms, on
14 Illumina, on Life Technology Platforms, PacBio,
15 Complete Genomics, and others. So there is really
16 a very comprehensive dataset that is being
17 assembled and analyzed to come up with those very
18 high confidence genotype calls.

19 Here is a list of the reference
20 materials that we have selected so far. NA12878,
21 that's probably something that is familiar to at
22 least some of the old-timers who were involved in

1 the Human Genome Sequencing project. There has
2 been a lot of data been generated using Sanger
3 sequencing, but there has been new data been
4 generated on the next- generation sequencing
5 platforms, and there are high- confidence genotype
6 calls available that Justin actually has put up on
7 the server and that are available to be used right
8 now. The sample itself is available from Coreal
9 right now, but it will also be released by NIST
10 shortly, probably towards the end of the year.

11 And then there are additional samples
12 that are in different stages of the process. For
13 the European Ashkenazi trio samples, there will be
14 initial sequences available later this year, and
15 currently, these samples are also available
16 through Coreal, but in the future they will also
17 be available as reference materials through NIST.

18 The reason why we've selected a number
19 of different reference materials, different
20 ethnicities, was that we didn't want to focus on
21 only a single sample because you can over-optimize
22 your algorithms when you use a single sample, and

1 we wanted to have a variety of genetic information
2 context but which we can stress the system and
3 analyze the systems.

4 Coming back to this table, and if you
5 think about not only inherited diseases, but if
6 you consider cancer, somatic mutations,
7 noninvasive prenatal tests or infectious diseases,
8 now you're no longer left with homozygous or
9 heterozygous genotype calls but you're left with
10 fractions, allele fractions where you want to
11 report on allele frequencies that are maybe 4
12 percent, 5 percent, depending on where the cutoff
13 is for your test. Copy numbers are also going to
14 take the whole range, between zero and any number,
15 not just whole numbers. So what your test needs
16 to do is your test really has to have a much
17 higher performance, has to have higher sensitivity
18 and has to maintain the specificity as a test for
19 the inherited diseases.

20 So that obviously also has implications
21 for the reference materials that you might want to
22 consider. I've walked you, on a previous slide,

1 through some of the errors that can get introduced
2 when processing your samples. For inherited
3 diseases, many of those errors are probably not
4 that drastic, that they will not have a big impact
5 on your end result, but you have to keep in mind
6 they will be additive also. So when you look at
7 lower allele frequencies, your test has a much
8 higher sensitivity, so small mistakes that are
9 being made throughout the process, they accumulate
10 and you will see them in the end, and they will
11 impact your results.

12 So reference materials that we discussed
13 at the Genome in a Bottle workshop were to simply
14 mix some of the reference material that have been
15 selected to get a range of allele frequencies that
16 you can analyze, although those allele frequency
17 differences would be based on polymorphisms that
18 are present in those reference materials.

19 In the case of cancer, we thought it
20 would be very useful to have specific mutations
21 because some specific mutations play such an
22 important role to have these available, and there

1 is utility for cell lines that contain
2 human-relevant mutations. There are some vendors
3 who offer these. And then we've been starting to
4 discuss synthetic constructs where mutations have
5 been engineered in plasmids, and these can be
6 used to test your process. I'll come back to
7 those in a minute.

8 Here are some examples of how we've used
9 some of this at Quest Diagnostics early on in the
10 development of this next-generation sequencing
11 test, OncoVantage. So on the right hand, the
12 table shows you that we've used those cell line
13 references with engineered mutations to determine
14 can our test detect those mutations, and then how
15 well do we match up the frequencies for these
16 mutations that have been in the reference material
17 established using droplet PCR.

18 The other way we've used these cell
19 lines was to analyze and to optimize our
20 bioinformatics pipeline. If you think back a year
21 or two years, there has been a real need to be
22 able to identify short deletions or insertions

1 using the different alignment processes, and some
2 platforms simply were not able to do that. And so
3 we've looked into these platforms and tried to
4 optimize the parameters so that we were able to
5 detect those deletions because we thought those
6 were really critical for a test like this, and we
7 would not be able to make a test available without
8 being able to detect those deletions.

9 FFPE treatment, I've already mentioned
10 it. It impacts your DNA sequence. So within the
11 working group of Genome in a Bottle, we've
12 discussed experiments and how to test this.
13 Actually, ABRF, headed or represented by Don
14 Baldwin at our workshop meetings -- ABRF stands
15 for Association of Biomolecular Resource
16 Facilities-- is planning to take cell lines and
17 compare the results of FFPE- treated cell lines
18 and untreated cell lines to determine what is
19 really the impact of FFPE treatment as the end
20 result.

21 In terms of synthetic controls, in my
22 mind, and I think in a lot of people's mind, there

1 is the question, well, are synthetic controls --
2 how can we use them? What are the limitations?
3 And are they really useful in terms of
4 substitution or supplementing patient samples
5 within the process?

6 So one hypothesis that we want to test
7 is are synthetic controls good surrogates for DNA
8 isolated from a clinical sample? And in order to
9 do that, I think we have to acquire evidence to
10 either support or refute that hypothesis. And so
11 what we want to do is we want to come up with a
12 study design that will enable us to answer that
13 question, come up with answers to that question.
14 And we're actually looking also for input from
15 multiple stakeholders, because if there is
16 information like that, we would hope that a lot of
17 labs and companies would be using that information
18 and translating that information into practice.
19 And we would anticipate that a study like this
20 would not only be conducted by individual groups
21 but by multiple different groups in order to put
22 that on solid footing.

1 Simulated sequences, the value really is
2 that very often we don't have samples available,
3 patient samples, that carry specific mutations, or
4 we have specific questions towards our
5 bioinformatics pipeline. So we can use a dataset
6 and manipulate it in a way that we engineer
7 mutations into that dataset. And that can be
8 specific to a platform. The platform that you're
9 using in your lab can be specific for the variant
10 type that you want to analyze, and it can be
11 specific to the context of the sequence. So you
12 can put it next to homopolymer sequences, et
13 cetera. And you can determine whether your
14 pipeline will detect that specific variant, what
15 are the frequencies that are attainable with these
16 types of variants, and which adjustments do you
17 have to make to your pipeline.

18 Here is one of the examples for two
19 mutations that we've generated such simulated data
20 for. What we wanted to find out -- well, what is
21 really the average frequency that we will detect
22 within our pipeline, but also what is the range of

1 frequency that we will detect versus the expected
2 frequency to give us a good idea about the
3 sensitivity of the pipeline.

4 In summary, one point that I haven't
5 made but I think is obvious to everybody is when
6 you look at sequence and compare research
7 applications and clinical applications, there's an
8 important difference between research and clinical
9 applications where you have to have a much higher
10 performance level on the clinical side versus the
11 research side. That information is being used to
12 make critical decisions for patients, so you have
13 to be sure about what you're reporting.
14 Next-generation sequencing will really benefit
15 from a well-characterized group of reference
16 standards, common reference standards that are
17 suitable for analytical assay validation.

18 One point that I want to make about the
19 Genome in a Bottle reference materials is that one
20 of our selection criteria also was that there is a
21 fairly broad-based consent available for these
22 samples. So if somebody -- another group wanted

1 to come in and produce derivatives of that
2 material, put other mutations in there, or you can
3 think about different things, that the consent
4 would allow for that. Synthetic constructs and
5 simulated sequence data, they all can be very
6 useful to assess your pipeline and assess
7 different parts within your pipeline, and I think
8 the next steps really are to discuss and determine
9 what the limitations or what the use of synthetic
10 standards can be specifically also in the oncology
11 area, specifically also in the oncology area for
12 analytical validation, but also to run them every
13 time when you do a sequencing test to include that
14 within your test.

15 And then for simulated datasets, I think
16 there also had to be agreement. How do you
17 assemble the simulated datasets? Which
18 information can you draw from these simulated
19 datasets and come to that acceptable conclusion
20 within the community?

21 With that I want to conclude, and thank
22 you for your attention.

1 (Applause)

2 SPEAKER: Hi, Andrew, that was a
3 fantastic overview.

4 On your last point, in your conclusion,
5 the simulated data, I mean, obviously simulated
6 data is fantastically powerful. You can generate
7 thousands of fake samples, millions of variants,
8 and use them for validating pipeline updates.
9 We've, in our diagnostics lab, done that
10 extensively and it's very useful. The challenge,
11 of course, is it's difficult to simulate all the
12 things that are going on biochemically in the raw
13 data coming in. So, for example, if base
14 modifications, because of the FFP treatment, are
15 local sequence-context specific, or pull-down
16 efficiency, which of course varies by mutation
17 type, have you looked at all of that issue, and do
18 you have any thoughts about how to generate
19 simulated data that might better reflect the
20 underlying biochemistry?

21 DR. GRUPE: Steve, I think you make a
22 very good point. There is a lot of complexity

1 that goes into the whole process and so it
2 requires really consideration of many of those
3 points. I think it will be difficult to have a
4 standard simulated dataset for that purpose,
5 because that simulated dataset probably has to be
6 consistent with the workflow that you really want
7 to test. So that's probably something that you
8 have to take into consideration when generating
9 those simulated datasets.

10 SPEAKER: So we've been trying to work
11 on such a thing. As you point out, it's a
12 challenging problem, but we'd love to see if a
13 consortium around that might be of interest.

14 SPEAKER: Another consideration for
15 cancer samples, you can have at the same time
16 gross genomic changes, a mixture of (inaudible).
17 And what I've seen in the clinic is that you
18 should -- you usually did cytogenetics and then
19 the molecular oncology part. Now, everything is
20 done at the same time, so that might affect the
21 results, too, and something that has to be taken
22 into account, I think.

1 DR. GRUPE: Yeah. I agree with you.
2 Specific cell lines probably have endogenous
3 rearrangements. Specifically, if you take cancer
4 cell lines, they are a mess, but if you look at
5 specific mutations within those, I think there is
6 value for those cancer cell lines. The reference
7 materials that come from normal donors, they might
8 be different because they will be EBB-transformed
9 B cells most of the time and they probably will
10 have less of those limitations that you refer to.

11 MS. VOSKANIAN-KORDI: So a few more from
12 the online audience.

13 One was, you spoke to this a little bit,
14 but a question on whether or not it is actually
15 suitable or appropriate to use cell lines versus
16 patient lines as a reference, and the comment from
17 the audience was that it is not clear what the
18 different institutions or government agencies'
19 position is on this or whether there is policy on
20 this yet, and my understanding is that you're
21 actually working to study this to kind of help
22 form guidance for the policy. So the policies are

1 not yet in place, but if you could just clarify an
2 opinion on that, that would be --

3 DR. GRUPE: Yeah. I think I can't make
4 any FDA policy.

5 MS. VOSKANIAN-KORDI: Sure.

6 DR. GRUPE: I'm the wrong person to
7 answer that question. And I think we just need
8 some studies in order to substantiate our
9 decisions one way or the other.

10 MS. VOSKANIAN-KORDI: Well, and then the
11 other question is what is Quest's position on how
12 to handle potential false positives in patient
13 samples that show deleterious impact?

14 DR. GRUPE: That's a good question.
15 That's a piece of discussion also. For example,
16 for the OncoVantage test, right now we're only
17 analyzing tumor tissue and not blood samples, so
18 we don't necessarily get direct information on the
19 germ line DNA other than through the tumor. It's
20 also limited to those 34 genes, so you're not
21 looking at a whole genome scan or whole genome
22 sequencing for these samples, so therefore,

1 incidental findings would also be limited.

2 (Applause)

3 DR. DONALDSON: Our next speaker will be
4 Dr. Laura van't Veer, and her talk will be "Impact
5 of Next-Generation Sequencing on Drug and
6 Diagnostic Co-development."

7 MS. VAN'T VEER: Thank you very much for
8 the opportunity to have me speak here on this
9 interesting topic and I will particularly focus on
10 the diagnostic code developments of next
11 generation sequencing and other molecular tests
12 because in that aspect they're all the same. I
13 think we are in an exciting era, the era of here
14 are my genes and what is my drug. And people are
15 pouncing on our doors and want to know the answer
16 and sooner than later is actually what is
17 required. So how do we get there?

18 So companion diagnostics and in
19 particular companion therapeutics there is an
20 exciting time ongoing in oncology because there
21 are over 1000 targeted drugs in development.
22 We're approaching 50 that approach full approval

1 in oncology. There are numerous biomarker
2 technologies for screening of diagnostics that
3 would apply to those targeted drugs. Some are
4 really specific, others are more broad. But we're
5 discussing today of course standards, are these
6 tests all reliable and how do we make them
7 reliable? And the clinical trials that are
8 currently ongoing some immediately test drug
9 biomarker combination but sometimes that's
10 actually done after the fact. And then there is
11 the issue of efficacy of off label drug use. So
12 once you find a mutation in one tumor type where
13 you may have a PMA approved diagnostics for a
14 targeted drug you may find that mutation in
15 another cancer type and then what do you do. And
16 so off label drug use of course is always possible
17 but it may I think particularly in this setting of
18 next generation sequencing we may want to have
19 some further guidance.

20 So how do we manage all this so we make
21 fast progress to actually deliver on the promise
22 of precision medicine? The biomarkers for agent

1 selection has been extensively documented in a
2 report from the Institute of Medicine under
3 leadership of Gil Omenn which actually nicely,
4 besides only this figure, lays out a landscape for
5 diagnostic tests and how do they come from being
6 identified to actually really clinical utility.
7 And also in the end as you can see here in the
8 bottom that they end up in practice guidelines and
9 reimbursements. And although we're focusing here
10 at this meeting and what are the requirements for
11 FDA clearance and FDA approval there is actually a
12 whole other part that we're not addressing today
13 which is in the bottom to get this in these
14 practice guidelines and being reimbursed. And
15 having experience myself that actually the
16 requirements for each of these steps are not the
17 same which is confusing. So once we have actually
18 a biomarker that can be used for agent assignment
19 there are two ways how that can be tested in
20 prospective clinic trials. One option where the
21 test does not direct patient treatment and you can
22 just test it in the background, or when you really

1 test this to direct patient management, a so
2 called FDA investigational device exempt status is
3 needed. And so those are being tested for some of
4 the targeted drugs currently in practice. And
5 sometimes maybe a particular test has already been
6 cleared and you can use it as part of the IDE as
7 well for a new indication. So the targeted drugs,
8 they're over a 1000 that under development or in
9 use or actually targeting what we call the
10 hallmarks of cancer which has been nicely
11 described over the years and consists of processes
12 like sustaining proliferation, evading the immune
13 surveillance, and many new drugs have seen the
14 lights that are very effective in model systems
15 and some start to show their efficacy when used in
16 cancer patients.

17 So accompanying those tests are a number
18 of prognostics and today's topic is predictive
19 tests that are either marketed as LDTs, 510Ks, or
20 PMAs and that's for all of these different
21 indications. And so you can see here for which
22 cancer types are actually already tests being

1 used. But there is sort of a flip side to this
2 because many of the companion diagnostics are
3 based on maybe a single clinical trial or
4 retrospective prospective assessment. Then there
5 comes the issue of me-too test because if you can
6 identify a mutation with one technology and it's
7 easy to do with another then sometimes this is
8 being used as a laboratory developed test which is
9 of interest to the FDA with their latest
10 announcements to actually also oversee these type
11 of LTDs. Sometimes one trial per diagnostic test
12 is not sustainable for the future, and one test
13 per approved drug might also not be sustainable.
14 The regulatory oversight is variable so there is
15 lots to improve to actually make us all go to a
16 more structured and efficient process.

17 I'd like to take you to one example in
18 breast cancer where we tried to develop such a
19 structured process also in communication with the
20 FDA on several steps. In breast cancer there are
21 also a number of targeted drugs in development
22 that are guided for their use by the biology of

1 the tumor. And for this particular example we've
2 been working with the FDA guidance to use
3 pathological complete remission as an early end
4 point for accelerated approval. So oncology that
5 means you can give a treatment to a patient when
6 the tumor is still in the body and you're
7 observing in a six month period usually where at
8 surgical resection at the end of the periods
9 there's still tumor left in the patients or
10 whether that's all gone, being identified as
11 pathological complete remission. There's an end
12 point that is only accepted for these high risk
13 breast cancer patients.

14 High risk breast cancer patients, there
15 are number of new strategies that either target
16 the immune environment, target the (inaudible)
17 gene, target stems cells, the PI 3- kinase
18 pathway, are several ways to target the
19 (inaudible) to a receptor. The trial that all
20 describe where we are using this is the I SPY 2
21 trial, investigation of serial studies to
22 predictor therapeutic response with imaging and

1 molecular analysis where the trial PIs are Laura
2 Esserman and Dawn Berry. I'm chairing the
3 biomarker work of this particular trial.
4 Investigational arms, you can see here the
5 different processes in the middle column that were
6 also shown in the hallmarks of cancer that are all
7 options for high risk breast cancer patients to be
8 tested in this Phase II trial. Under several drug
9 arms, investigational arms in this trial in
10 parallel up to five in comparison to the standard
11 of care. The trial design is novel because it's
12 uses as I said PCR as an endpoint, an early
13 endpoint. Final assessment will be by long term
14 survival, but the first endpoint is PCR to assess
15 the response rate of investigational agents in
16 combination with standard therapy as compared to
17 standard therapy alone. Randomization is done
18 adaptively and I will not show the details of
19 that, that can be a whole different talk, but
20 suffice it to say that it uses a pre-specified and
21 automated algorithm where the randomization
22 probability to each of the arms of these different

1 targeted drugs as the study proceeds is based on
2 the molecular signature of the tumor, and actually
3 based on MRI ongoing through the trial, and
4 finally the assessment of the pathologist in the
5 surgical specimen of PCR. This same algorithm
6 triggers the decision to graduate when
7 approximately 60 to 120 patients are enrolled
8 because if at that time the adaptive randomization
9 doesn't show superiority over standard of care for
10 this particular trial it is decided that
11 investigation was being dropped because the
12 incremental increase of success is too low for it
13 to go forward. It doesn't mean there's no effect,
14 but not in this trial process as chosen how to use
15 it.

16 But in addition, and that's more
17 important for today's talk that we implemented the
18 framework to qualify promising biomarkers for
19 companion diagnostics to make it as efficient as
20 possible. The trial design in a simplified
21 version here, there are two investigational arms,
22 the middle one and the bottom one as compared to

1 the top which is standard of care for breast
2 cancer patients before surgery, but we have up to
3 five -- and actually currently five
4 investigational arms in addition to the standard
5 of care. So here you can see that at different
6 time points MRI assessment is done as shown in the
7 bottom as well as biopsies. And for eligibility
8 the biopsy before randomization is used to
9 actually define eligibility and is used within the
10 adaptive randomization. There are tests that are
11 being used for eligibility and randomization are
12 either FDA cleared or are available as standard of
13 care like some immune is to chemistry tests. And
14 if the patient shows to be high risk by either one
15 of these features as indicated here the patient is
16 on the study and this one in particular
17 combination of these factors is used in the
18 adaptive randomization leading to defining eight
19 subtypes that include ones like triple negative
20 for breast cancer, hereto two positive,
21 (inaudible) very high in those. The evaluation of
22 the investigation for the standard is being done.

1 So there's a long story in between these
2 two slides which is to show you that last year the
3 first drug connection, this year the second drug
4 were reported at meetings to have graduated
5 successfully from the trial arm. Here I'm showing
6 you a PARP inhibitor Veliparib which was given in
7 combination with Carboplatin where you can see
8 that three different signatures all were to
9 negative patients, hormone receptor positive for
10 two negatives and hormone or triple negative were
11 being evaluated ongoing through the trial and
12 going the trial it was identified that the triple
13 negative group had a more favorable response than
14 either the full group or the hormone receptor
15 positive for two negative. So there was an
16 enrichment ongoing for this particular subtype
17 which meant that the drug graduated only after 115
18 patients had been randomized to this treatment arm
19 or to their concurrent controls. And so that's a
20 real efficiency increase because the number of
21 patients needed to come to this conclusion of
22 graduation because of the enrichment has been

1 decreased enormously. In the right two columns
2 you can see the probability that in a phase three
3 trial this drug will be better than standard of
4 care.

5 What I would like to point out to you
6 because it's important for the qualifying
7 biomarkers that although we were successful in
8 identifying the tripe negative group as a
9 graduating signature where 52 percent of the
10 patients achieved a pathological complete
11 remission as compared to 26, that still means that
12 only half of the patients benefit from the drug so
13 there is room for improvement. So one of the
14 programs and very active programs we're running in
15 parallel is how can we use the biology to identify
16 companion diagnostics and we have a process to
17 qualify biomarkers to be tested on a qualifying
18 level based on earlier evidence. And it's all to
19 see if we can increase this 52 percent in an
20 investigational way to a higher proportion by
21 applying those biomarkers. So the way how we set
22 up the system is that we're using established

1 biomarkers or ones under an IDE so it's an
2 investigational device exempt status for
3 stratification and randomization as I explained.
4 But the qualifying biomarkers need to be run under
5 CLIA and for some of them we actually have an IDE
6 platform that we're using to qualify these
7 biomarkers. So these are hypothesis testing not
8 yet used for treatment assignments, but they could
9 be subsequently in a Phase III trial be validated
10 and then become companion diagnostics.

11 Of course in the background there is a
12 lot of exploratory analysis which is more
13 hypothesis generating and once the biomarkers are
14 used as an IDE this facilitates as part of the IND
15 the process towards the FDA PMA approvals. So the
16 companion diagnostics that we tested in the I Spy
17 trial for the PARP inhibitor Veliparib in
18 combination with carboplatin are listed here and
19 you can see that the biology speaks because PARP
20 inhibition is effective in the DNA repair or
21 blocks DNA repair processes and is considered to
22 be very effective if that is deficient. So here

1 are a number of tests that we have been qualifying
2 and I cannot show you the results yet. They will
3 be presented this year in December, but just for
4 the purpose of this talk assume that some of these
5 qualified what do you actually need to do next.
6 So this is like an overall scheme as was published
7 in clinical cancer research earlier this year
8 showing that if you go from an adequate test
9 performance in the laboratory setting and if you
10 then validate it in a clinical validation you need
11 to go all the way to clinical utility. And on top
12 of that the regulatory oversight of course can be
13 several. In Europe it's mainly ISO certification,
14 in the U.S. CLIA. There can be a European CE
15 marking or here in the United States it's a 510K
16 or a PMA process. For treatment assignment it is
17 the PMA. And of course as I said before
18 guidelines and reimbursements in the end will
19 actually define the clinical utility.

20 So where are we with I Spy 2? I think
21 for those that I just showed you, although not all
22 of them I can say, but some of them actually made

1 it into the final phase of clinical validation.
2 So that means those that have been successfully
3 qualified were run under CLIA and some actually
4 were from an IDE platform, they can now be used in
5 a Phase III trial to really show the clinical
6 utility of the tests in a predefined way even
7 though you can still choose if you want to have
8 treatment assignment based on the diagnostics or
9 not. If that is successful then of course as I
10 started my talk and it's very much of interest and
11 for our next generation sequencing there can be a
12 lot of off label use which could eventually lead
13 to additional indications if sufficient evidence
14 has been accumulated. But if you use a particular
15 test and particular technology and that's
16 particular for next generation sequencing, once
17 you have either 510K or PMA in case of treatment
18 assignments, if there are improvements in
19 technology one always needs to keep in mind that
20 you need to update the PMA if that's the
21 technology you're using. So that's a separate
22 discussion.

1 Just to end to show you one slide which
2 was inspired because of the whole database session
3 yesterday because besides that you can do all of
4 these analyses in isolation and here you can see
5 the I Spy clinical trial data and the forced
6 database we're using as well as molecular systems
7 and hospital systems. We also actually are
8 thinking, and this part is actually functional,
9 we're also thinking how to actually be better in
10 analyzing data by using semantics. And this is a
11 slide that I got from Sue Dubman who's actually
12 overseeing the I Spy IT systems. So I think
13 clinic trials very much also needs to be much more
14 documented in an automated way so that analysis
15 can be done more efficiently.

16 I'd like to end with to actually
17 acknowledge the I Spy 2 trial study team which has
18 a large group of working group tiers but also many
19 individual molecular oncologists are active as
20 agent chaperones in the trial, the project
21 management office, and the oversight, the sponsor
22 and the funding and the oversights where you can

1 see from the FDA, I see there Janet Woodcock and
2 Richard Pastor, but I'd also like to acknowledge
3 actually CDRH previously Steven Gutman, but
4 currently Bob Becker, and also the work of Liz
5 Mansfield in guiding how to use next generation
6 sequencing which is currently being employed for
7 one of the drugs in the trial. And I would like
8 to end to say that also I think lots of the
9 thinking how to implement all of these diagnostics
10 is based on the work with Jon Retzlaff and Rasika
11 Kalamegham from AACR where we have a very active
12 working group between academia and industry,
13 particularly pharma and diagnostic companies. So
14 when you think about what are the efficiencies and
15 what are the requirements that we need to actually
16 come to companion diagnostics.

17 So I hope I didn't leave you too much in
18 a sort of twist here. I'm happy to take any
19 questions. (Applause)

20 MS. VOSKANIAN-KORDI: Due to time
21 limitations we're only going to take one or two
22 questions so if anyone has any.

1 Okay. Our next speaker is going to tune
2 in remotely so I'm going to set him up if --
3 introduce him maybe.

4 MR. DONALDSON: Okay. Thank you very
5 much, Laura.

6 (Applause) Our next speaker will
7 be Dr. Charles Sawyers and he'll be
8 sharing in the age -- the topic of
9 his talk will be Sharing in the Age
10 of Clinical Genetic Data.

11 MS. VOSKANIAN-KORDI: I want to see if
12 he can hear us first. Dr. Sawyers? Is the volume
13 on?

14 DR. SAWYERS: This is Charles Sawyers in
15 New York. I apologize for not being able to be
16 there in person. I don't see my slides appearing
17 on the screen yet. And I cannot hear my voice
18 over the background so I don't know if I'm live.
19 Okay, my slides are up. Thank you. I apologize
20 for not being able to appear in person. I've been
21 listening to the session, enjoying the
22 conversation. I'm going to cover this topic from

1 a slightly different perspective and that is
2 leveraging existing clinical data to help us
3 discover new biomarkers in a validated way.

4 So I'd like to begin by reminding you
5 that many of the companion diagnostics that we've
6 heard about, particularly in lung cancer, came
7 about from clinical observations that set sets of
8 patients were responding to targeted agents and
9 that led to investigations of the molecular
10 characteristics of those responders versus non
11 responders which elucidated the predictive
12 biomarkers. As many of you are well aware that
13 prior to clearance of information sequencing
14 instruments many cancer centers have invested
15 heavily in sequencing large numbers of patients
16 either with very focused panels of genes or of
17 larger, you know, exomes. I can tell you that at
18 Memorial Sloan Kettering here in New York we've
19 committed to sequencing 10,000 patients over the
20 next 12 months just at our center alone and I know
21 of many similar efforts at other large cancer
22 centers. So I mention this because it presents I

1 think an unprecedented opportunity to try to
2 leverage this existing data to figure or discern
3 more clinical correlations between different
4 mutation patterns.

5 So on the next slide which I don't know
6 that I can control -- can someone advance the
7 slide for me please? Got it. I just want to
8 refer you to this paper that was published earlier
9 this year in Nature by Gaddy Getz and colleagues
10 to introduce you to this, you know, concept called
11 saturation analysis. So while this was done for a
12 different purpose it led to a thought experiment
13 that I'd like to share with you what I've done
14 recently with Gaddy. So what's showing here on
15 this slide is from his paper, an attempt to
16 leverage the 5000 tumors on which whole exomes had
17 been completed and are in publicly accessible data
18 bases estimate how many tumors would you need to
19 sequence to discover, you know, it includes all
20 cancer mutations. And of course it depends on the
21 frequency of the mutation and the tumor type.
22 What you can see on the red line here is that

1 mutations present in 20 percent of patients with
2 that particular tumor type can -- all those
3 mutations can be discovered, you know, with
4 several hundred patients. But most mutations are
5 not present at 20 percent so we have a significant
6 number of tumors left to go to define the universe
7 of mutations. And what this paper shows you is
8 that it detects mutations that are present at a
9 two percent frequency and you need to do about
10 2000 samples of each tumor type. If you assume
11 there are 50 different types of tumors it's only
12 100,000 tumors. Well 100,000 is a large number,
13 but we all know with the cost of sequencing today
14 that's an attainable goal and one that many
15 consortia I think will achieve over the next
16 several years. But the real reason I show you
17 think is to ask the question can we apply this
18 thinking to biomarker discovery. How many
19 patients would we need to discover a biomarker?
20 We're watching through a well known example in
21 oncology which is drugs that target mutant BRAF
22 present here in less than 50 percent of patients

1 with melanoma. So this is a famous waterfall plot
2 from a New England Journal paper that described
3 the response to the BRAF mutant patients treated
4 with Vemurafenib. And you can see that most --
5 you know, a huge number of patients are responding
6 which is why this drug was such a game changer
7 when it was first studied. You can also see by
8 the red line that only about half of the patients
9 who the response biomarker have tumor shrinkage
10 greater than 40 percent which would meet the
11 criteria for objective response that's established
12 among clinical trialists. Now we all recognize
13 that there's a smaller fraction of patients which
14 I've depicted here as 20 percent which are having
15 the best responses. And the question we'd all
16 like to know is there a second biomarker, what I
17 would presume biologically is a modifier of
18 response that might predict for those patients who
19 have the best response. And if you knew the
20 answer to that maybe you'd treat them with
21 Vemurafenib alone and you'd treat the other 80
22 percent with some kind of combination. So how

1 many patients would it take to discover such a
2 biomarker? That's the question. So using
3 saturation analysis Gaddy ran the same
4 mathematical algorithms and what you can see here
5 is that -- of course it depends on the frequency
6 of the biomarker in the population, but if it's
7 present in one out of five patients remarkably you
8 need to only 200 patients to discover that
9 biomarker. If it's present in 5 percent of the
10 population you need 1200 patients and if you scale
11 this all the way you hit a saturation at 1 percent
12 of biomarker frequency, you need 6000 patients.
13 Now depending on your point of view that's a large
14 number, but I can guarantee you there's 6000
15 patients who'll think (inaudible) lower number
16 receiving Vemurafenib in the U.S. today. So if we
17 can just leverage the existing clinical practice
18 that's going on we're in a very powerful position
19 to refine ability to give these drugs more
20 precision.

21 So there's a medicine area that I'd also
22 like to walk you through which is the phenomenon

1 that you're well aware of is that most of these
2 drugs are effective for short periods of time and
3 the patients develop resistance. So again using
4 the same drug in the same disease the patients
5 tend to respond for about a year before relapse.
6 So in another paper published earlier this year by
7 Levi Garraway and colleagues in Boston there
8 reported the results of whole exome sequencing
9 from a small cohort of patients who had a biopsy
10 done, pretreatment, and then one had relapsed.
11 Now we know that melanomas have a large number of
12 mutations at baseline due to UV exposure most of
13 which are passenger mutations. So the biological
14 evolution of this disease is depicted here with
15 the heterogeneous nature of the tumor at
16 presentation, each of these shades representing
17 different sub-clones which have been produced
18 through different driver mutations that accumulate
19 over time in a background of passenger mutation.
20 So what Gaddy and I did was look at this published
21 data and noted that the background mutation rate
22 in tumors at diagnosis was roughly 12 mutations

1 per megabase, but now post treatment patients, you
2 know, roughly in a year now have 13 mutations per
3 megabase which across the genome the enormous
4 number of mutations to attribute causality for
5 resistance. However, you know, it's fairly
6 straightforward to realize that what is most
7 important here are the new ones or the additional
8 mutations that have appeared since treatment
9 started and that's a very small number and
10 therefore allows you to make a calculation as to
11 how many patients would it take to discover causes
12 of resistance and it's only fair in existing
13 patients to discover all of the resistance genes
14 for Vemurafenib that would be present at at least
15 a five percent frequency in patients with a power
16 of greater than 90 percent.

17 The last though I'd like to leave you
18 with is another clinical scenario which gathered a
19 lot of momentum about two years ago when two of my
20 colleagues at Sloan Kettering, Barry Taylor and
21 David Solit, reported a whole exome sequencing
22 analysis of a single patient who was treated on a

1 clinical trial of a drug called everolimus which
2 inhibits mTOR in a trial of patients with bladder
3 cancer. This was the only patient on the trial
4 that had a complete remission, that had a
5 remarkable complete remission. And what they
6 found when they sequenced that patient were
7 mutations in two genes, TSC 1 and NF 2 for which
8 biology, you know basic cancer biology conducted
9 in laboratories had clearly predicted dependence
10 upon t he mTOR pathway. So this is the classic
11 (inaudible) story, an anecdote, but by knowing
12 this and going back and looking at the remainder
13 of the patients in the trial they were able to
14 identify four other patients who had TSC 1
15 mutations who had partial responses. So this had
16 led to now a prospective study of everolimus in
17 patients who have TSC 1 or TSC 2 or NF 2
18 mutations. Since that report many additional
19 similar stories have been published or are on
20 their way to publication with different drugs and
21 the National Cancer Institute is launching an
22 extraordinary responder initiative which has

1 already led to queries from several hundred
2 investigators who had several hundred patients who
3 would fit the criteria for such a response.

4 So, you know, the summation of all of
5 this is that there are remarkable clinical
6 experiments being done. The data is sort of out
7 there on the table, but it's been difficult to
8 harness. And we are very involved as a community
9 in efforts to build the infrastructure to allow
10 such data to be shared and leveraged more easily.
11 So I want to call your attention to an
12 organization called the Global Alliance for
13 Genomics and Health on which I serve the executive
14 committee which was initiated just in the past
15 year based on a series of meetings held about 18
16 months ago amongst leaders not just in cancer but
17 in mendelian diseases, infectious diseases and
18 other communities. You can go to the website for
19 details. The first in person meeting was held
20 earlier this year in London. The next one will be
21 held in October at the Human Genetics meeting in
22 San Diego. This slide is a bit out of date but

1 there are well over 150 partner organizations now
2 spanning 20 countries and 6 continents which have
3 a number of working groups which you can query on
4 the web which are coming up with data sharing
5 standards, consent standards, and infrastructure
6 for data sharing. I'd also like to let you know
7 that very soon the American Association of Cancer
8 Research for which I just completed my term as
9 President is structuring a trial project amongst
10 several cancer centers in the U.S. and Europe to
11 share sequencing data of the type that I've
12 described that's already happening at our centers
13 to try to leverage the power of this approach.

14 I'm going to stop there and you can tell
15 me what to do next. And if there are questions
16 I'll try to take them. Thank you. (Applause)

17 MR. SIMONYAN: Questions? Maybe if
18 there are questions we can ask them to read them
19 for him because he cannot hear.

20 MS. VOSKANIAN-KORDI: That microphone is
21 not on.

22 SPEAKER: Very different in that my talk

1 was sort of edifying we need to have structure
2 at trials to actually validate these biomarkers
3 whereas his presentation shows us that we need
4 sort of the world to actually understand what are
5 the right biological markers to understand
6 response and resistance. So my question to him
7 would be how does he think about validating the
8 findings from the global alliance? So do we need
9 trials or do we just need our common sense to
10 actually start using the -- it would be like an
11 off label testing that you start to promote. The
12 sort of complexity of the mutations is so huge
13 that to make sense out of it you need to combine
14 the whole testing from the whole world. And so
15 how do we validate it?

16 (Audio Interruption)

17 DR. SAWYERS: -- from the extraordinary
18 responders program would make their way into the
19 clinic. I don't know if there is a formal
20 position, but the approach that, you know, in the
21 academic community we're thinking is the
22 extraordinary responder insight is an insight that

1 leads them to a more formal, you know, test of the
2 drug in the companion diagnostic mode for
3 approval. That at least would be the way the
4 everolimus example would be proceeding in which a
5 set of biomarkers for the TOR pathway, TSC 1 and 2
6 in particular in a clinical trial would move
7 forward.

8 (Audio Interruption)

9 MR. DONALDSON: Okay. In the interest
10 of time we're going to not have the panel but I
11 would invite you all to continue this discussion
12 with the working groups that will be forming and
13 if you have any ideas over lunch please let me
14 know.

15 Thank you very much to the panel and
16 thank you for your patience and for part of your
17 lunch hour. (Applause)

18 MS. VOSKANIAN-KORDI: We will try to
19 start again at 1:00 o'clock. I know it's a little
20 late; I understand some of the speakers had a
21 little more to say than the time allotted. If you
22 guys could gather back in here at 1:00 it would be

1 really appreciated. If a lot of people aren't in
2 here we'll try to hold it back a little bit but
3 just would like to say on the time frame. Thank
4 you.

5 (Recess)

6 DR. TEZAK: Okay. We're going to start
7 the session. I want to thank everybody for having
8 a really short lunch and coming back. We're just
9 a little bit late. So next session is Next
10 Generation Sequencing Devices and Clinical
11 Applications. And actually the session that just
12 ended was like a really, really intro to our
13 session so I want to thank Eric and all the
14 speakers from the previous session for giving
15 really good intro both to the NIST standards and
16 Genome in a Bottle Consortium efforts and also to
17 CDC led efforts for a standardization of next gen
18 sequencing, data formats, and testing.

19 So I'm Zivana Tezak; I'm at FDA. I'm in
20 CDRH Office of In Vitro Diagnostics and my
21 co-chairs are Justin Zook, who's going to speak a
22 little later from NIST, and Heike Sichtig who is

1 from FDA and she's going to also speak a little
2 later. What we're going to do is -- so this is
3 the agenda and probably we're going to have to
4 shorten maybe Heike's presentation and the panel
5 maybe to 30 minutes so that we're still done
6 around 3:00. So the first talk will be -- I'm
7 going to introduce all the speakers first, the
8 same as what Eric did. So the first talk will be
9 from Dr. Ira Lubin from CDC, and he's been
10 actually as previous speakers mentioned very
11 involved in leading and forming various consortia
12 on standardizing the human genome sequence
13 representation and the formats. Then the next
14 talk will be Dr. Justin Zook from NIST who has
15 been leading Genome in a Bottle Consortium. So we
16 have the data format standardization, the first
17 talk and the second is actually about reference
18 materials. And then the third talk is going to be
19 from Mya Thomae who is VP for Regulatory Affairs
20 in Illumina and she's going describe the first
21 four clearances for next gen sequencing devices
22 that Illumina paved the way for others. And then

1 the last talk will be Dr. Heike Sichtig from
2 microbiology division and she's going to talk
3 about the infectious diseases efforts and thinking
4 that FDA has in some research efforts. And that's
5 going to be followed by the panel discussion. And
6 we are very happy that we have four current
7 micrograde developers, all the representatives on
8 the panel. And Justin agreed to moderate the
9 panel and our panel is going to be Dr. Korlach
10 from PacBio, Dr. Shaw from Ion Torrent, Mya Thomae
11 again from Illumina, and Kendal Dinsmore from
12 Complete Genomics.

13 So without further ado we're going to
14 start the session.

15 DR. LUBIN: Okay. Back from lunch and
16 reenergized. So I want to thank the organizers
17 for the opportunity to present on one of the
18 efforts that we have underway at CDC in which we
19 have established a national work group to explore
20 the development of a clinical grade variant file
21 and the data standards therein that will be useful
22 -- that will support interoperability for

1 applications relevant for clinical and public
2 health settings. So the intention is to take a
3 look at variant files and by variant files we're
4 thinking about the VCF and similar files that are
5 currently used during the course of next
6 generation sequencing and how we can evolve these
7 to really present unambiguous descriptions for the
8 sequence elements that are representing therein,
9 preferably with a haplotype background and also to
10 provide the capability to communicate these
11 sequence elements, genomic representations among
12 clinical laboratories and within healthcare
13 infrastructure that conforms to adoptive standards
14 such as HL7 that will support systems
15 interoperability.

16 So the applications that we are seeking
17 to address, to develop use cases today at best are
18 challenging to put into practice and often require
19 custom solutions for those in which we can
20 actually address today. And this includes
21 exchange of sequence information among
22 laboratories for quality assurance. Proficiency

1 testing is probably the best example that folks
2 are knowledgeable about here. To be able to take
3 sequence generated from the laboratory either at
4 the stage when a variant file is being generated
5 or those results that are found to be clinically
6 relevant and deposit them into medical data bases
7 or disease registries to outsource variant file
8 data in a uniform manner for downstream analysis
9 and interpretation and ultimately to be able to
10 message that data to patient records making it
11 amenable for clinical decision support systems or
12 reanalysis at a future time when there are new
13 testing indications or as new knowledge is
14 developed that would then trigger the need to
15 reanalyze that data.

16 Looking at the broader sphere of how our
17 health IT system is evolving we're now at a time
18 when we are seeing health information exchanges
19 come on line. In some settings more successful
20 than others but this is the direction we're
21 heading. The continued integration of electronic
22 health records, data warehouses where clinical

1 information is being deposited and available for
2 analysis to provide a more timely improvement of
3 patient care and assess quality metrics and how
4 that care is being delivered. Also we have
5 patient portals coming on line and the capability
6 of patients being able to see their data in more
7 real time and mechanism in which we translate
8 what's coming off of the test and being made
9 available to patients.

10 So a little over a year ago we
11 established what we're calling the Clinical Grade
12 Variant File Specification Work Group and this was
13 established in collaboration with other Federal
14 agencies based on previous work groups we had that
15 saw a real need in trying to develop a more
16 standardized way for genomic representation both
17 in the variant file and how it can be communicated
18 to other settings. We're also aligned with the
19 HL7 Clinical Genomics Work Group and we have a
20 number of others throughout the community,
21 clinical laboratory professionals, informaticians,
22 and folks who are active with a number of

1 professional organizations as well as
2 translational researchers who are involved with
3 this group. Now we also realize that the
4 outcomes, the opinions of any work group is the
5 consensus of that work group. So we just released
6 a website VCFclin.org in which we will begin
7 publishing the outcomes of our discussions on that
8 site eliciting feedback from the broader community
9 that can inform our final recommendations because
10 as all of us has probably learned the opinion of
11 even a small group of people is probably not
12 shared by the larger community, so we really need
13 this consensus building process. Our intent is to
14 take the principles and recommendations that are
15 developed and use those to inform use cases
16 designed around the applications I noted earlier.
17 And this will form the framework for pilot studies
18 in an effort to move the recommendations and
19 principles into practices. This is not something
20 that can happen quickly but in taking sort of this
21 structured approach to identifying what needs to
22 be done and then figuring out how to do it the

1 hope is that we can move forward in an orderly
2 fashion that will be least disruptive but useful
3 to the broader community.

4 So in thinking about the variant file as
5 it's currently integrated into next generation
6 sequencing we were thinking about this in three
7 stages, three parts of the process for performing
8 next generation sequencing. The first is to
9 consider the data that populates the variant file
10 and what needs to be considered in preparing that
11 data before the variant file consumes it. The
12 second is the variant file itself and to what
13 extent we need to further constrain data within
14 that variant file. And for some applications such
15 as laboratory exchange of data that variant file
16 may be the end point assuming that we can get the
17 variant file into a format that can be readily
18 understood by other laboratories for the
19 comparison of data. And third is to take the data
20 in the variant file and to translate it into
21 formats amenable for messaging through an IT
22 system using adopted standards such as HL7. So

1 for example we would want after downstream
2 analysis and interpretation on the laboratory side
3 we would want to certainly translate the outcomes
4 of that in HGVS which is recognized by HL7 along
5 with other established descriptors. And that
6 would be the format that those particular test
7 results would be messaged.

8 So this is a basic representation of the
9 work flow when performing next generation
10 sequencing and in thinking about what we need to
11 do standardize the way sequences are presented we
12 actually have to look back to the stage of
13 alignment and then work through this and see what
14 else we need to consider. So taking a look at the
15 first step which addresses the preparation of the
16 data for input into the variant file the work
17 group discussed primarily two issues, and that is
18 the assignment of genomic coordinates and variant
19 callers and how they need to be evaluated and set
20 up. And so here was recognition and broad
21 consensus within our work group that the alignment
22 or mapping of the sequence reads needs to be done

1 against a version reference assembly. And the
2 reality today is that a lot of labs are doing
3 panels even though there's a transition to exome
4 analysis today and those doing panels will often
5 align to a laboratory reference sequence. And
6 this sequence or may not be deposited in a
7 standard database such as RefSeq, HLA also has
8 their own database in which there are a uniform
9 set of methods used to cross-map that sequence
10 against the reference assembly. And so if the
11 sequence that you're aligning against is in RefSeq
12 for example you're pretty well set. You can go
13 back and NCBI has done that cross-mapping so you
14 can get your genomic coordinates. If it's a
15 sequence that is not deposited in these standard
16 databases then we would strongly recommend that it
17 get deposited and NCBI allowed to take that
18 through their methods to assign genomic
19 coordinates. Laboratories are discouraged from
20 doing that alignment themselves because aligners
21 vary, the way aligners are set up vary, so there
22 is the probability that if you were to do that

1 cross-mapping the genomic coordinates that a
2 laboratory would assign may have the potential for
3 higher error and those coordinates being wrong.

4 In terms of the variant callers, the
5 output of the variant callers are what populates
6 the variant files and the work group is
7 recommending that the data that be collected
8 include just not the variance but the reference
9 calls, the no calls, as well as phasing
10 information to allow the greatest possibility for
11 applying metrics so you can really understand the
12 quality of the sequence that is being generated,
13 as well as permit you the greatest opportunity for
14 interpreting that sequence and developing the data
15 set that you'll analyze to derive what is
16 clinically relevant. So one of the advantages of
17 taking the phasing information is when you have
18 variants close by to know whether they're on the
19 same allele or not in the absence of family
20 studies which you may or may not have, generally
21 not. You really need to use the phasing
22 information to show the allelic representations.

1 And in this figure we actually have under
2 alternate a single haplotype representation, but
3 you'll see that depending on how the variant
4 caller is set up you can have different
5 representations of this. In the absence of
6 phasing you would not know which output is
7 associated with another in terms of the allelic
8 association. This also addresses the question of
9 the size of the event that you should output when
10 you have control over that. And when possible
11 it's better to output large events where you can
12 see all the allelic associations in the same
13 segment, but when you don't have that opportunity
14 then you really do need phasing information that
15 you're putting out such as a base or just a few
16 bases at a time.

17 In thinking about the variant file now
18 the work group wanted to focus on content and not
19 format. And the rationale there is that there are
20 translators that will allow you to translate from
21 one format to another. But if the content is
22 sufficiently constrained then it should appear as

1 the same content from one file to another which
2 then it's amenable to being understood in
3 different settings and supports interoperability.
4 Today variant files in clinical settings are
5 sufficiently different that it's anecdotally it's
6 challenging to take a variant file from one
7 clinical lab and have it understood by another.
8 And this is the message that we have received
9 particularly from NCBI with respect to another
10 project in our group, the GeT-RM project that
11 collected data from various laboratories and
12 almost in all cases the laboratory needed to be
13 contacted to make sense of what was in the file
14 that was communicated. And it's a matter of the
15 implementation of existing specs, not so much the
16 specs themselves in that laboratories will adjust
17 the way they set up their variant file to
18 accommodate the platform they're using in their
19 downstream analytics, and because this is all
20 contained in the lab they can deal with the issues
21 that come up. When you talk about taking the data
22 within a variant file and messaging to an external

1 setting then you really need to think about the
2 content and have agreement on what it should it
3 look like and how it will be interpreted. The
4 good news though is that for several of the
5 features within the variant file, the designation
6 of the chromosome, the genomic coordinate, a
7 simple variance, these are fairly well established
8 and should be understood even in the current state
9 when there's files moved from lab to lab or lab to
10 another setting. There are limitations though.
11 In the new reference assembly I believe we have
12 about 170 or so alternate assemblies such as HLA
13 where there are multiple alleles that are highly
14 divergent compared to the primary assembly and
15 there has not yet been a consistent means to
16 represent these in a variant file, although
17 there's discussions under way with other groups
18 and how that might be undertaken. Also there are
19 issues with complex variants and haplotypes and
20 how they're represented. You can have different
21 representations and I showed you a basic one
22 earlier that they look different but it actually

1 turns out to be the same variant. So work has to
2 be done to assure that when dealing with these
3 complexities that there is some level of common
4 structure that these are represented and can be
5 communicated or otherwise tagged as something that
6 may be more difficult and needs to have a closer
7 look.

8 The other issue which we'll come back to
9 later deals with the metadata. And the metadata
10 supports the sequence, provides context to the
11 sequence. And the metadata that typically
12 currently exists in these files as text type data,
13 not structured data. So there is a lot of
14 interest, and some laboratories are doing this, in
15 trying to put that data in that XML format that
16 permits it to be more structured in mind and more
17 useful for computational processes.

18 One of the questions that comes up,
19 particularly from the clinical realm is how do you
20 assess the quality of this data. And this came up
21 earlier in Dr. Gonzalez's talk is that there are
22 no common metrics that transcend platforms,

1 pipelines or assays. So the work group could only
2 recommend that sufficient description be put into
3 these variant files or otherwise made available
4 that describes the filters or metrics in a uniform
5 manner. And what that uniformity is has yet to be
6 established to provide the ranges, the cut offs,
7 and what's being thought of as good versus poor
8 quality. And to put enough context in there that
9 someone who is reasonably knowledgeable, who's
10 outside the laboratory can look at the data and
11 make some determination on the quality of that
12 data.

13 So the variant file that -- a VCF for
14 example that constrains to the 4.1 specs and where
15 you're dealing with for the most part relatively
16 simple SNPs, INDELS, what have you, should be
17 exchangeable among laboratories. The challenge is
18 that there's enough variation at this point that
19 it's not. So there needs to be additional efforts
20 to encourage more uniformity to provide a means to
21 take whatever the lab has optimized for their
22 pipeline into a format that can be readily

1 messed to other laboratories or other settings.
2 In terms of messaging that through a healthcare
3 infrastructure it's possible to put an HL7 wrapper
4 around that and move it to where it needs to be.
5 For the results that are to be reported as
6 clinically actionable you do want to assign the
7 HGVS and in time assign LOINC codes and put this
8 in a message that is more amenable to the
9 standards that have already been established in an
10 implementation guide by HL7. And so I explain
11 that here and the sort of things that you want to
12 at least assign to your clinically relevant
13 variants are the ID numbers H, G, and C. Human
14 Genome Nomenclature Committee assigns ID numbers,
15 names, and symbols to each gene. The caveat is
16 the names and symbols sometimes change. So while
17 it's not standard practice it seems reasonable
18 that ID numbers should go along with all of those
19 with names or symbols perhaps. Very variable.
20 Version HGVS nomenclature, it's easy to say but
21 the reality is that HGVS nomenclature is not
22 always applied in a uniform manner and some of the

1 specs for HGVS are not well established for
2 certain types of variants, particularly the more
3 complex ones. So in order to have more uniformity
4 of HGVS this is another issue that needs to be
5 addressed in terms of thinking about how we can
6 promote use of that. And it is the adopted
7 standard for how variants are described. Other
8 identifiers, and we mentioned LOINC codes in time;
9 we're not quite ready for the full set yet I think
10 but as we move more into the electronic messaging
11 era we'll need to consider those. And I just
12 listed a number of documents there that contain a
13 number of the standards that have been currently
14 adopted by HL7. These are constantly getting re-
15 discussed and new documents are being developed.
16 And NAACCR has standards applicable to cancer
17 registries but they're weak in the genomics area
18 right now, but that should be changing within the
19 next two years as genomics becomes more a part of
20 the discussions of that and other groups.

21 The metadata as we mentioned is
22 important for providing context to the sequence

1 data that's included. When all of this is within
2 a single lab the metadata supports the analysis
3 that's done in the lab and typically there's
4 minimal metadata to support what needs to be
5 performed in terms of the downstream analytics.
6 When you're talking about taking the data in a
7 variant file and messaging it to an external
8 source where it can be analyzed and interpreted at
9 a later time by someone else the work group
10 considered four general categories that need to be
11 considered to provide more robust metadata within
12 that file and those are listed here. And I just
13 want to emphasize the idea that metadata together
14 with sequence data is what is needed.

15 And finally I'll just end by mentioning
16 the machine versus human readable data. At the
17 level of variant file software is deciding what is
18 deposited in that file. The data within the file
19 then undergoes further analytics with additional
20 software analysis afterwards. So that data needs
21 to be machine readable to a great extent. For
22 messaging, there you're talking about a different

1 format at least for what's clinically relevant in
2 terms of the chromosome, the HGVS nomenclature.
3 And when reporting to physicians you absolutely
4 need a format that physicians understand. So for
5 example when you're talking about cancer you need
6 to use the formal name of the biomarker so they
7 know what to base their decisions upon. Providing
8 genomic coordinates or refs or even RefSeq, HGVS
9 is not necessarily -- would be useful to the
10 physicians at all times, particularly those who
11 are less familiar with how all of this is
12 evolving. So in summary these kinds of data
13 standards are required to support systems
14 interoperability we have a work group that's
15 addressing this, a website to get community
16 feedback, and it's just not about the standards,
17 it's what you do in the laboratories as well. We
18 talked about the aligning and its importance in
19 the mapping in terms of identifying -- assigning
20 the right genomic coordinate, and the importance
21 of the metadata in being able to understand
22 sequence context.

1 So that's my world wind tour of what our
2 work group is doing and thank you for listening.

3 (Appause)

4 DR. TEZAK: Thank you, Ira, so much. So
5 I think this was a really great talk to show into
6 the weeds what kind of problems people are dealing
7 with here. Unfortunately we ran over so we have
8 no time for questions right now but Ira will be
9 around.

10 Next up Justin Zook from NIST.

11 DR. ZOOK: Okay. Thanks, Vahan, for
12 organizing this and I'm happy to tell you some of
13 what we've been doing in the Genome in a Bottle
14 Consortium. Andrew did a great job this morning
15 of giving an introduction to some of the
16 activities that we have and so I may skip quickly
17 over a couple part of the talk, but I'll talk
18 generally about Genome in a Bottle.

19 And so the question that we're trying to
20 address is that so you sequenced a genome, how
21 well did you. And these reference materials that
22 we're developing are designed to help answer that

1 question. So we formed this Genome in a Bottle
2 Consortium to develop the infrastructure for
3 performance assessment of next gen sequencing and
4 part of the purpose of this is to support science
5 based regulatory oversight. And so the FDA is
6 actually funding a fair amount of our work at NIST
7 in this area I know is interested in using the
8 reference materials that we're developing. When
9 we started this there weren't any widely accepted
10 set of metrics that you can use to know how well
11 you're doing in variant calling and so we're
12 working towards that but still aren't quite there
13 yet. And we're developing standards to address
14 this. So this includes human genomes that we're
15 making into reference materials, so these are
16 really well characterized whole human genome cell
17 line, extracted DNA from them. And these are
18 being characterized and will be disseminated by
19 NIST. And we're also developing tools and methods
20 that you need to use these reference materials.
21 So if you sequence this DNA and run it through
22 your pipeline how do you compare your variant

1 calls to our in a standard way.

2 A lot of the motivation for doing this
3 was because if you compare different sequencing
4 technologies or different bioinformatics methods you
5 end up with different answers for variant calls,
6 and so this is just a comparison of three
7 different platforms, bioinformatic pipelines, and
8 you end up with hundreds of thousands of variant
9 calls that are unique to an individual pipeline.
10 And similarly, this figure was shown this morning
11 also for exome sequencing pipelines, even with
12 those exact same sequencing data if you run
13 different bioinformatic software you can get
14 largely different answers. And so we wanted to
15 move beyond just doing these sort of venn diagrams
16 where you say oh, these disagree with each other,
17 to actually understand why they disagree with each
18 other and how do you know which one is right when
19 they disagree with each other.

20 So the reference materials that we're
21 characterizing like this are designed to
22 understand the performance for a particular part

1 of the measurement process. And that part is
2 depicted here in yellow on the right hand side.
3 So there are pre-analytical steps where you might
4 extract the DNA from the sample, and our samples
5 are already extracted DNA so they do not test that
6 part of the process at all. But they do test the
7 library preparation, the sequencing, alignment,
8 mapping steps, variant calling, and whatever
9 confidence estimates you associate with the
10 variant calls. And so it will test that sort of
11 black box there that's in yellow. If you get the
12 right answer that doesn't mean that you're perfect
13 it just means that you're doing well for this
14 sample in this particular situation and so you can
15 proceed with caution and so have this red light
16 over here that is sort of -- it gives you a yellow
17 light or a red light. So if you get the complete
18 wrong answer then you know that there's something
19 going wrong. These also are not designed to test
20 the downstream steps, like the clinical
21 interpretation or functional annotation at this
22 point. And these are just normal genomes that

1 we're characterizing right now.

2 So the reference materials that we have
3 in the pipeline right now are all 10 microgram
4 samples of DNA that are isolated from a large
5 batch of cells. So we have between 2000 and
6 10,000 tubes of DNA from each reference material
7 and each of them is from a single large batch.
8 And so this helps to control for any types of
9 mutations that might happen during the growth of
10 the cell lines because we've evenly distributed
11 those mutations over the entire batch. We have a
12 pilot genome, NA12878, that we're hoping to
13 release as a NIST reference material later this
14 year and we already have characterization
15 information that I'll be talking about in a bit.
16 Then we've also selected a couple of trios from
17 the Personal Genome Project for the reasons of
18 consent that Andrew talked about earlier. And so
19 one of these is of Ashkenazim Jewish origin and
20 the other is of Asian original. And these are
21 mother-father-son trios that we're developing
22 because you can do some error analysis by looking

1 at the trios and also do phasing of the variant
2 calls that you get.

3 So our goals for the data to accompany
4 these reference materials is to have approximately
5 zero false positive and false negative calls in
6 our high confidence regions. And we want to
7 include as much of the genome as possible in our
8 confident regions. So we don't want to just take
9 the inner section of those venn diagrams that I
10 showed earlier, we want to -- if they disagree
11 with each other we want to understand as much as
12 possible why they disagree and decide which one is
13 correct when we can. We also want to avoid bias
14 towards any particular platform and take
15 advantages of the strengths of each of the
16 platforms so we're using data from multiple
17 technologies in this process, and avoid bias
18 towards any particular bioinformatics algorithms
19 as much as possible.

20 So for our pilot genome we have 14 and
21 actually now even more data sets from 5 different
22 sequencing technologies. And so these are from

1 Illumina, Complete Genomics, SOLid, 454, and Ion
2 Torrent. And then we also will be incorporating
3 biodata as it's generated and I'll talk about that
4 a little bit more later. The process that we've
5 used to do this integration was described in a
6 recent paper in Nature Biotechnology and I'd
7 encourage you to look at that paper if you're
8 interested in learning more about the details of
9 the methods that we used, but in general we find a
10 list of candidate variant sites, so we go through
11 with a really low threshold and look for any site
12 that might possibly be a variant and we'll look
13 more closely at those. Then we find the
14 concordant variants and this sort of helps us to
15 see so if it's likely to be a true variant what
16 are the characteristics of these variants. And so
17 we use those characteristics of where if you have
18 a new sit that has sort of abnormal
19 characteristics where it doesn't look like the
20 true sites then we will trust that data set less
21 at that location than another data set that would
22 have sort of normal characteristics. So for

1 example if one of the data sets has strand bias at
2 a location then says it's a variant and another
3 set does not have strand bias and says it's not a
4 variant we would trust the one that does not have
5 strand bias. And then we assign a confidence
6 level and part of that process is also excluding
7 sites where all of the technologies have evidence
8 of bias. So if all the technologies have low
9 mapping quality then that would not fall in our
10 high confidence regions. Or if it's a segmental
11 duplication, or if it's in a potential structural
12 variant or in a repetitive region, we exclude
13 those regions. And so part of this last process
14 now more recently is also taking the work that the
15 Illumina Platinum Genomes Project is doing as well
16 as Real Time Genomics where they've analyzed the
17 pedigree that our pilot genome is part of so it's
18 a 17 member family pedigree, and they phased
19 everyone and looked to see which variants are
20 inherited properly as you'd expect from the
21 phasing, and that has helped to move some of our
22 less confident sites into the more confident sites

1 by ones that were consistent.

2 So after we have these calls questioned
3 a lot of people ask us how you compare it to these
4 Genome in a Bottle calls. And if you're not
5 asking that question you probably should be
6 because it's not actually a trivial process to
7 compare to the calls. So this is something that
8 the Genome in a Bottle Performance Metrics Working
9 Group has been working on and we've developed the
10 specification for what a tool that compares
11 variant call to ours should do in terms of what
12 the input should be, the outputs, what type of
13 licensing it should have, and the definitions of
14 the different performance metrics, like
15 sensitivity and specificity, and false positive
16 and false negative. And if you search on our blog
17 page you should be able to find this
18 specification.

19 More recently we've been working with
20 the Global Alliance for Genomic Health
21 Benchmarking Working Group which is a new working
22 group that was just formed in the past couple of

1 months that's sort of taking this work that the
2 performance metrics was doing and actually trying
3 to develop a tool that will do these comparisons.
4 And so we're further developing sort of the
5 standard definitions for these performance metrics
6 and developing some standard methods for
7 benchmarking that can be used so that if you
8 compare your calls to our you know it's done in
9 the same way as if someone compares their calls to
10 ours. We're also developing some standard data
11 sets that people can use if you just want to test
12 about informatics pipeline. And also looking at
13 how you stratify performance. So how you can
14 stratify in terms of how well you're doing for
15 SNPs versus INDELS versus complex variants and
16 then also for different regions of the genomes.
17 So exome versus genome, or repetitive regions
18 versus normal regions. And there are a variety
19 including some doing that.

20 So I just wanted to give a few examples
21 of some preliminary uses of our high confidence
22 genotypes. These are available from an FTP side

1 that's at NCBI and you can go to the
2 genomeinabottle.org website to find a link to this
3 FTP site. This FTP site includes both our high
4 confidence calls as well as a lot of data that
5 people have generated on these reference
6 materials. So if you just want to download FastQ
7 files and run them through a bioinformatics
8 pipeline you can do that. And we have exome data
9 and genome data and multiple technologies, and I
10 think potentially some targeted data also. And so
11 these data are already being used. And the data
12 that most people are comparing to are a set of
13 SNPs and INDELS that we have which cover about 77
14 percent of the genome. And so we have a BED file
15 that specifies the high confidence regions and you
16 should only compare your calls to ours in those
17 regions. There's a GCAT website that's been
18 developed where you can do interactive exome
19 comparisons and you can upload your data to that
20 site and compare you calls to ours and get a
21 variety of different types of performance metrics
22 from it and you can compare to other

1 bioinformatics pipelines that people have run as
2 well.

3 The short article, Perspective Peace,
4 about the first FDA -- of the authorization of
5 next gen sequencer by Francis Collins and Margaret
6 Hamburg mentioned that the reference materials
7 that we're developing played a role in this first
8 FDA authorization and we're hoping they'll
9 continue to play a role in the future. Mount
10 Sinai School of Medicine published a paper in BMC
11 Genomics earlier this year where they used our
12 calls to benchmark clinical exome sequencing.
13 Qiagen published a paper in BMC Genomics where
14 they looked at somatic mutation calling methods
15 and they mixed the reads from our reference
16 material in with reads from another genome to see
17 what levels of somatic mutations different
18 pipelines could detect.

19 So this looks like it got distorted a
20 little bit in the conversion but this lays out
21 what our development plans are for the reference
22 materials over the next year or so. So the first

1 line is related to our pilot reference material.
2 We plan to release the NIST reference material
3 which will be number 8398 by around the end of
4 this year. And we also are developing a set of
5 preliminary large deletions calls so we're
6 starting to move into looking at structural
7 variants and will be continuing to develop those
8 over the next year and the methods that we're
9 using to develop those so people can do
10 performance assessment with structural variance as
11 well. The next line is about the Ashkenazim
12 Jewish Trio that we're developing into reference
13 materials. We've already collected some data on
14 that that's on our FTP site and I'm actually about
15 to post a blog about the data that's there. We
16 have Illumina and Complete Genomics' data so far
17 but we're also collecting Ion data, bionano
18 genomics, optical mapping data, and solid data
19 that will all be publicly available. And we'll be
20 similarly generating SNPs and INDELS for those and
21 over the next few months we're excited to be
22 collaborating with Mount Sinai School of Medicine

1 to do 100x PacBio sequencing over this trio, so
2 this will be whole genome sequencing and we think
3 this will really help to get at some of the larger
4 types of structural variants as well as some of
5 the more difficult parts of the genome and help us
6 to characterize them. And that data will be made
7 publicly available in stages as it's produced so
8 that anyone who's interested in helping us to
9 analyze it can help. We're also getting Illumina
10 assembled long reads which is often called
11 molecular data and we will be continuing to
12 develop those and then we'll follow a similar
13 process for the Asian trio in terms of getting all
14 of this data together and characterizing them.
15 And hopefully we'll have those reference materials
16 released around the end of the year, end of next
17 year.

18 The future directions include for
19 germline mutations we want to look at additional
20 ancestry groups, African Americans and Hispanics
21 in particular we'd like to have represented as
22 well as a large family because we've found that

1 for NA12878 it's been really useful to have this
2 pedigree. And then we're also looking at somatic
3 mutations and Andrew talked about those earlier
4 this morning. I just wanted to highlight that
5 there are some commercial members of Genome in a
6 Bottle Consortium that I've already either already
7 developed FFPE cell lines based on the cell lines
8 that we're using as reference materials or they're
9 thinking about doing it. And potentially might
10 introduce mutations into them or modify them in
11 other ways.

12 So if you'd like to get involved there
13 are a number of ways you can do that. You can use
14 our integrated calls and give us feedback. For
15 NA12878 these calls are available on our FTP site
16 like I mentioned. You can help with sequencing
17 and analyzing the new Genome in a Bottle samples,
18 especially with structural variant calling or with
19 analyzing the data from the long read technologies
20 that we're getting. We also hold biannual
21 workshops in January on the Stanford University
22 Campus in California and August here on the NIST

1 campus in Gaithersburg, Maryland. And you can get
2 involved in the Global Alliance for Genomic Health
3 Working Group where we're developing methods to
4 get performance metrics. And you could also sign
5 up for genomic newsletters at genomeinabottle.org
6 and feel free to email me with any questions that
7 you have.

8 Thanks. (Applause)

9 DR. TEZAK: Thank you very much. So we
10 have time for a couple of questions if anybody has
11 any.

12 SPEAKER: No? Yes? There we go.
13 Sorry. A few questions again from the on line
14 audience. One, any recommendation of reference
15 materials for targeted mutation panels?

16 DR. ZOOK: So it depends sort of on how
17 targeted the panel is. So I think our reference
18 materials will likely have some mutations in
19 whatever targeted panel you're using, but I also
20 recognize that depending on how small the panel is
21 there might not be very many that are your panel.
22 And so I think for those panels there are some

1 spikens that Mickey Williams group that has
2 developed that Andrew had mentioned and I think
3 sometimes probably just using cell lines that have
4 mutations in those regions are what you need to
5 do. It ends up being a combination of different
6 methods of doing validation when you have a
7 targeted panel.

8 SPEAKER: Okay. Next, what interactions
9 occur between NIST and the FDA with regard to any
10 regulation of reference materials?

11 DR. ZOOK: We are non regulatory so we
12 don't set any of the guidelines at all for what's
13 done, but we hope that our reference materials are
14 useful in the regulatory process. So we
15 definitely talk with the FDA as part of the
16 process and they've funded us to do a lot of our
17 work, but we don't set the regulations.

18 SPEAKER: And if no one else has any
19 questions -- one more -- any plans to also
20 integrate human phasing data from Illumina?

21 DR. ZOOK: So our most recent set of
22 variant calls for NA12878 does include phasing

1 information though it's based on the pedigree
2 instead of using that. I'm assuming they're
3 referring to the molecular phasing data. One of
4 our next steps is to think about how we integrate
5 different forms of phasing information so you can
6 do phasing in a variety of different ways and so
7 we want to think about how to integrate that data
8 together but I've not done that yet.

9 DR. TEZAK: So our next speaker is Mya
10 Thomae from Illumina.

11 MS. THOMAE: Good afternoon, and thank
12 you so much for inviting me to speak today. What
13 I thought I'd do today is just go over the
14 clearances that were granted to Illumina at a
15 fairly high level. I haven't included a lot of
16 the data, but for folks that don't tune into
17 regulatory submissions all the time, there's
18 actual a fair amount of interesting work that's
19 been just done on the intended uses and the
20 various products that have come out.

21 As everybody knows, last November there
22 were four clearances that were released all in the

1 same day: Two for CF assays, one for an
2 instrument, and then one for a set of reagents.
3 FDA was kind enough to do a press release on this,
4 so we ended up in the FDA news. There was also,
5 as Justin said, an article in the New England
6 Journal of Medicine, and then Jeff Sharon also did
7 a blog on this. It was a pretty good day back in
8 November for Illumina.

9 The first thing that I thought I'd talk
10 about is the instrument platform, that's the
11 MiSeqDx. For the regulatory geeks in the crowd,
12 I've put in the product codes and the device
13 types. What's really interesting here is this was
14 done as a de novo 510(k). De novo 510(k)'s are
15 utilized when the risk of a device is similar to
16 other class II devices but there isn't necessarily
17 a specific, predicate device.

18 The instrument was cleared as a de novo,
19 and the very exciting thing that happened is FDA
20 gave it a class II instrument, which is similar to
21 what was done for previous multiplex instruments
22 like the arrays, but they gave it a status of

1 exempt from premarket notification. That means --
2 and we saw this this week with Life Technologies
3 -- is that folks that come in after Illumina and
4 have a similar intended use are able to get their
5 instruments on the regulated market without having
6 to do a 510(k).

7 This was an incredibly innovative
8 approach that FDA took to make sure that
9 sequencing instruments can get out there and a
10 fairly expeditious manner to be used in the
11 diagnostic marketplace, though the intended use of
12 the instrument reflects the data that was put into
13 FDA. It's very specifically for use with the
14 MiSeqDx universal reagent kit, which I'll cover in
15 just a minute, but it's also very specifically for
16 use with DNA and also with peripheral whole blood
17 samples. Similar to other intended uses, it is
18 very reflective of the data that was submitted to
19 FDA.

20 Along with this there's the universal
21 reagent kit 1.0. Again, a new regulation for
22 this: Reagents for molecular diagnostic test

1 systems. The reagents were given a class I status,
2 also exempt from premarket notification. Again,
3 this allows folks that have instruments, next-
4 generation sequencing instruments, and related
5 reagents to be able to get both the instruments
6 and the reagents registered with FDA in a very
7 least-burdensome manner. The intended use of the
8 universal reagent kit, pretty straightforward:
9 Library prep and sequencing.

10 Along with these, there were also two CF
11 assays that were cleared. The first CF assay is
12 pretty similar to other CF assays that have been
13 cleared in the past, and this one threw as a
14 traditional 510(k). For those that want more
15 data, you can find the CF assays in the
16 traditional 510(k) database on the FDA website,
17 but just so you know, I looked up the other day,
18 and I was like, 'Huh. There's only two Illuminas
19 there.' The de novos are actually now in a
20 separate database, so if you need to look up all
21 of the data that is in here you'll see it on the
22 right-hand side; there's one for 510(k)s and then

1 there's one for de novos. They're just in two
2 separate databases now.

3 The 139-Variant Assay went through as a
4 class II. This is how other CF assays have gone
5 through, but the interesting part about this assay
6 is that it is an extended panel. It's a much
7 larger panel than most of the other 510(k)s that
8 have been cleared, and that's reflective of the
9 next-generation sequencing technology. I won't
10 read that to you, but it's very similar to other
11 CF-intended uses.

12 Then there was also a second CF assay
13 that was cleared, this one for a slightly
14 different intended use, and I think that's
15 probably, again, class II, but I think the
16 intended use is actually the most interesting part
17 of this. Again, it's talking about specifically
18 one can be sequenced with this, but the intended
19 uses a bit different than the other CF panels in
20 that the test is intended to be used as an aid in
21 diagnosis of individuals with suspected cystic
22 fibrosis -- again, a little bit different from the

1 intended CF assay -- and the test is most
2 appropriate on the patient has an atypical or
3 nonclassical presentation of CF or when other
4 mutations have failed. This is a bit more of a
5 specific intended use that was put out there to
6 cover this additional content since it's a bit
7 different from how the normal panels are used.

8 That's actually all I have today. I'm
9 happy to answer any questions, but it was a really
10 big day. There was a lot of work that was done to
11 get four assays into FDA and get them all cleared,
12 and I know that was true on the FDA side as well.
13 I know we're trying to make up a little time, so
14 I'm happy to stop here or we can take some
15 questions or I can answer things as part of the
16 panel.

17 (Applause)

18 DR. TEZAK: We have time for several
19 questions.

20 SPEAKER: Mya, that was great. Can you
21 just comment a little more on the process by which
22 that particular intended use in your last slide

1 was decided and in particular the restriction to
2 patients with atypical presentation or for whom
3 the classic mutation panels were negative?

4 MS. THOMAE: I can't talk about
5 everything that happened, but my understanding is
6 -- and this happened actually before I was at
7 Illumina -- this could be done as an all-in-one
8 panel, but the entire content of the panel was
9 very difficult to cover in a single-intended use.
10 Zivana, I don't know if you want to comment
11 further on that, but it was going to be difficult
12 to find a single, intended use that could cover
13 both of them. I think that the discussion was
14 about making sure that each of the portions of the
15 panel was appropriate.

16 DR. TEZAK: Your question was why to CF.
17 The 139 actually has clinical utility because the
18 clinical utility comes from the CFTR2 database, so
19 it was really the regular regulation for CF as the
20 previous tests were cleared. The one which is for
21 the whole gene was a novel approach, and we didn't
22 look at every single variation that can come out

1 of it because it's the whole gene obviously. It
2 was the whole sequencing and not just the panel,
3 not just clinically significant mutations, and we
4 were looking at classes. It was differently
5 analytically validated, and the clinical utility
6 is not as clear.

7 SPEAKER: Just further clarification:
8 The one finds 139 mutations or variants, but
9 there's over 1000 CTFs that are known that might
10 be related, and the other one's sequence is a
11 whole other thing. Have we found any of them and
12 does the first thing only look at 139 or is it
13 probes?

14 DR. TEZAK: The first one does look at
15 only 139, that's the output. The tests are very
16 similar underneath, but the output is just 139 for
17 the first one.

18 SPEAKER: They're both sequencing?

19 DR. TEZAK: They're both sequencing.

20 SPEAKER: One just doesn't show
21 everything sort of?

22 DR. TEZAK: No, 139 only shows 139.

1 SPEAKER: All right.

2 SPEAKER: A few from the online
3 audience: Is the MiSeqDx approved for
4 tissue-derived DNA from tumors?

5 MS. THOMAE: (off mic)

6 SPEAKER: Okay, and what is the
7 regulatory implication MiSeqDx if someone wants to
8 get an RNA-Seq-based assay cleared?

9 MS. THOMAE: You need to come in and
10 talk to Illumina Business Development and Illumina
11 Regulatory, and we'll figure it out.

12 SPEAKER: Just as an added bonus, if you
13 could find the original press release, and we
14 could maybe make that -- you got it? Never mind,
15 sorry.

16 MS. THOMAE: Got it? Okay, great.
17 Thanks.

18 DR. TEZAK: Just the last question.

19 SPEAKER: Given the talk two talks ago
20 about the variability in terms of calling and
21 algorithms, etc., what would the extent of
22 computer systems validation done to the

1 instruments where the intermediate steps are all
2 validated or just the end result?

3 MS. THOMAE: There's a lot of different
4 kinds of validation that goes into a system like
5 this. There's definitely software validation that
6 occurred, that was part of the 510(k). There's
7 also analytical validation that's done, analytical
8 studies, and then there's also the clinical
9 studies. I think there's multiple levels of
10 validation that go into it, but I think at the end
11 the whole system is validated as part of the FDA
12 clearance process. (Applause)

13 DR. TEZAK: Our next speaker, Heike
14 Sichtig from FDA.

15 DR. SICHTIG: Hi, my name is Heike
16 Sichtig. I'm a principal investigator and
17 regulatory scientist with the FDA in the
18 microbiology division. I'm going to be talking
19 about how we are going to go about enabling
20 NGS-based technologies for clinical diagnostics.
21 This is strictly from the microbiology
22 perspective.

1 First I want to just show a quick
2 disclaimer where we know that the (inaudible) is
3 moving very quickly and we know that new policy
4 and regulatory issues are being brought up at the
5 moment. These thoughts presented in this
6 presentation are preliminary and do not represent
7 finalized FDA policy. I just wanted to make that
8 clear in the beginning. We are currently working
9 on a draft guidance, but until such time that it's
10 released, we really encourage everybody to come in
11 and use our pre-submission process for any
12 outstanding questions currently.

13 This presentation today I will outline
14 first because NGS is really a cutting edge tool
15 that is making it into the clinical market. I
16 want to outline possible approaches to validation
17 studies for next-generation sequencing systems and
18 really highlight the difference between
19 metagenomic versus targeted or what is also called
20 Custom Amplicon sequencing and then also talk
21 about the use of sequence outputs in combination
22 with databases to potentially evaluate

1 performance. Then I'm going to quickly talk about
2 the Interagency Working Group on NGS feasibility
3 that we are currently working on and then really
4 focus this presentation on the MicroDB reference
5 database that we are creating with funds from FDA,
6 and then also touch on these microbial reference
7 materials.

8 First, for the potential validation
9 strategies, this is a 10,000 foot view of how
10 we're going to go about that. For industry folks
11 in the audience, they already know that we really
12 look at everything from a systems perspective. We
13 really start from specimen collection and then go
14 all the way to the clinical (inaudible). If you
15 look at the green box, that would be something
16 that we look at that would come under regulatory
17 purview, at least from the microbiology
18 perspective.

19 At the FDA, for IBD devices, we really
20 do apply risk-based regulation. As you just
21 heard, the MiSeqDx clearance got classified as a
22 class II device, and for microbiology devices that

1 will most likely be the same case. There are some
2 class III devices where you would submit something
3 that is called a premarket approval, or PMA
4 application, and that would be for something like
5 the human papilloma virus or hepatitis viruses,
6 but most of the bacteria that we're looking out
7 would be class II.

8 Here is the FDA's general concept of how
9 we evaluate diagnostic devices. On the left-hand
10 side, if you focus there, there's a picture of the
11 samples which could be, for example, bacteria,
12 virus, or fungi. We would apply a simple
13 comparative performance, and on one side you have
14 the new device which could be an NGS device, and
15 on the other side you have the reference method
16 which could be a PC or amino acid or composite
17 thereof. Then you apply a 2 x 2 comparative
18 analysis where you want to show clinical
19 sensitivity and specificity. But really the
20 bottom (inaudible) currently is each possible
21 organism needs to be confirmed by this reference
22 method, and that becomes exceedingly burdensome

1 for metagenomic samples because showing every
2 possible organism in there is almost impossible.

3 Here again, as I mentioned previously,
4 we really look at evaluating the whole system from
5 sampling or through pre- analytical steps,
6 amplifications, sequence detection, assembly, and
7 all the way down to the clinical (inaudible). We
8 would also look for validation studies for each of
9 these steps, but these would be different for what
10 is called metagenomic sequencing versus targeted
11 sequencing.

12 For targeted sequencing, we already have
13 an approach called Multiplex, and there's a
14 guidance out if you want to look that up. That
15 would be a good first step to look at.

16 For these metagenomic sequencing
17 validation studies, we would most likely look at
18 the classical, two- armed approach that we applied
19 before where on one side we have in-house
20 analytical studies and on the other side we look
21 at clinical evaluations. For the in-house
22 analytical studies, we most likely will have to

1 take liberties -- we understand that -- for the
2 limit of detection and the inclusivity studies
3 because of the exceedingly burdensome showing of
4 metagenomic samples that they cannot detect all of
5 the organisms that are in the sample. We would
6 also look at interference studies, contamination
7 studies, and carryover studies.

8 On the other side, for the clinical
9 performance, we will stick with multisite clinical
10 users to look at the device, how it performs at
11 different sites, and this should really be
12 designed to test the data analysis pipeline all
13 the way down to the final core determination. The
14 performance could also be established possibly
15 through comparison to a database, but again, we
16 really encourage you to use the pre-submission
17 process to discuss this with us.

18 For precision and reproducibility, we
19 could suggest using a panel of representative
20 microorganisms. FDA is funding this to create a
21 microbial reference standard to be used as a tool
22 to show precision and reproducibility, but I will

1 touch on a later slide.

2 All of these studies were also discussed
3 at a workshop that we had on April 1st, and the
4 workshop agenda discussion paper and webcast is
5 available online. You can simply write down the
6 FR dock note or write me an email if you'd like to
7 get this information or just simply Google it.
8 There we also discussed how we are going to go
9 about working on this draft guidance, and this
10 will be coming out soon.

11 Now I'm going to change gears and just
12 briefly touch on the feasibility group that I'm a
13 part of where FDA, NIH, NIHID, DTRA, (inaudible),
14 Livermore, and CDC came together. We really
15 wanted to focus in on looking at NGS feasibility,
16 and then we looked at bacteria only. We want to
17 conduct a small pilot study -- which is still
18 ongoing -- to generate information to evaluate
19 quality of existing sequences in the public
20 domain. This is still in progress and we're still
21 churning through the data, but publication should
22 be coming out soon.

1 The overarching goal was to identify the
2 preexisting high quality deposits in the public
3 and then really use that information and build
4 from there. Then after we have established these
5 quality metrics, we would like to use that
6 information and set our sequence outputs for our
7 ongoing sequencing efforts that we're currently
8 doing that will be the reference database. Then
9 utilizing existing standards, if they're
10 available, for technical and (inaudible) of
11 metadata because why should be reinvent the wheel
12 when some working groups or other people have
13 already looked at these things. We really wanted
14 to pay attention on connecting antimicrobial
15 resistance phenotype to genomic deposits at the
16 clinical collection site.

17 A very preliminary view of what these
18 feasibility studies showed or what the suggestions
19 that we will be coming out with will look like are
20 that multiple levels of reference databases are
21 likely where you have, sort of, high-quality
22 genomes on one side where these can be used for

1 validation and clinical use, and then on the other
2 side you would add other available genomes for the
3 testing and development process only. Extensive
4 screening will also probably be required of human
5 and other hosts: Chimeras and artificial
6 constructs. The other suggestion would be to
7 really use separate bacterial, viral, and fungal
8 reference databases. This should all be publicly
9 available.

10 This really highlights the current need
11 for having a robust, standardized, and
12 high-quality microbial sequence database in the
13 public sector containing representative samples,
14 metadata, high-quality raw sequences, assemblies,
15 annotations all available for the public.

16 This is a slide given to us courtesy of
17 NCBI, and it really shows the exponential growth
18 of the number of genomes in GenBank, and at the
19 same time not the same growth for the number of
20 new species, but for diagnostic processes. The
21 representation of real species really needs to be
22 increased, and that really highlights our efforts

1 that we are spearheading with the targeted
2 sequence effort that is funded with limited funds
3 by FDA at the moment to identify and fill gaps.
4 All these (inaudible), assemblies, annotations,
5 and metadata that get generated are sent to NCBI
6 to be made accessible to the public. All these
7 results that are generated should be traceable so
8 they could be reevaluated as necessary.

9 We identified a very small set of 600
10 clinically relevant (inaudible) microorganisms
11 initially in collaboration with NCBI, and we are
12 applying a highly controlled and documented
13 approach, and we found a very, very good solution
14 with the Institute of Genome Sciences at the
15 University of Maryland where we are applying a
16 hybrid approach of using Pacific Biosciences
17 sequencing and Illumina Sequencing. We are also
18 using metadata and standards that were already
19 established by Sifsen and also a template that was
20 made by Dr. Lynn Bright for antimicrobial
21 resistance data that we also adapted for our data.
22 All this, again, is pushed out to NCBI. On the

1 left-hand side, you see a list of our
2 collaborations with clinical labs and repositories
3 where we currently receive our samples from.

4 Here is the sequencing strategy that we
5 are using for creating the reference data. One is
6 Illumina (inaudible) sequencing, and then on the
7 other hand we also perform Pacific Biosciences
8 long-insert sequencing for all these samples.
9 Then for assembly, we are doing various assembly
10 algorithms, and then there will be a
11 bioinformatics expert that will be sitting there
12 and deciding on which one gets pushed into NCBI.
13 But we will have all these assemblies available
14 later on. Also what they're doing is they're
15 applying (inaudible) and GCQ4QC, and then they're
16 using an automated pipeline for annotation, and
17 they also give us base modification detection data
18 that will be most likely uploaded as well to a
19 database.

20 Here are just some very preliminary
21 sequencing and assembly steps for the first batch.
22 Here on one side we had samples from Rockefeller

1 University that was, sort of, a uniform set of
2 (inaudible) samples. Here you can see that
3 already 23 out of 50 on single-contact status, but
4 they're still churning for the data and more might
5 qualify later on. Then on the other side you see
6 a more diverse sample set from Children's National
7 Hospital and Sifsen where you have now 18 genera
8 represented. Here currently we have 12 of the 41
9 in single-contact status, but again, they're still
10 working on it and as time goes on more might
11 qualify to be in single-contact status. That
12 would be our goal for these.

13 Here's a quick snapshot on how this data
14 will be represented at NCBI, thanks to Bud Klinke
15 for putting this data really together. The ID for
16 this is 231221 if you want to check it out at
17 NCBI. On the left-hand side, you have BioSample
18 database sort of, like, an example of how this
19 data will be represented. There are two samples
20 here that are shown: One is from the Children's
21 National Hospital and one is from Rockefeller
22 University. On the right-hand side you see the

1 sequencing archive database, and this is where
2 (inaudible) will be available. For BioSample, for
3 example, if you go into the Children's National
4 Hospital sample, there you see all the fields for
5 the metadata that would be required in order to
6 qualify for our database. Also, if we have the
7 antimicrobial resistance data available, they'll
8 be put in by a template that we adopted from Dr.
9 Lynn Bright.

10 Then all the assemblies that get
11 generated get put into the assembly database at
12 NCBI, and here's an example. For example, IGS
13 looked at one assembly, and they used the
14 asymmetry in GCQ which then suggested a
15 misassembly, and then various assemblies were
16 selected previously. This really showed that they
17 should use the HGAP 2 assembly to push forward to
18 NCBI.

19 Lastly, I just want to briefly touch on
20 the microbial reference databases. Justin had
21 talked about human reference samples, but they're
22 also, at the same time, working on microbial

1 reference materials that got funded by FDA. There
2 are four (inaudible) rated bacterial standards
3 that they're working on, and they're growing large
4 batches of these as you can see from the table.

5 Then on the right-hand side you see
6 where you can these extracted DNA for testing,
7 library preparation, sequencing, mapping,
8 (inaudible) calling, and confidence estimates.
9 Then these were also generated with multiple
10 platforms. The de novo was made with Pacific
11 Biosciences, and then also platforms like
12 Illumina, INTRON, and Ingen technologies were used
13 for validation.

14 So, just wanted to go very briefly
15 (inaudible) that in order to realize NGS as a
16 regulated device for microbial diagnostics. It is
17 necessary to ensure that the (inaudible)
18 information used for interpretation of results is
19 of suitable quality and that it should include
20 appropriate metadata. Efforts are currently
21 underway. There will be a difference in
22 validation strategy for targeted versus

1 metagenomic as I outlined previously, and we will
2 enable a streamlined approach for regulatory
3 evaluation as was highlighted by many, many
4 commercial developers and clinical end-users.

5 The Division of Microbiology devises
6 that FDA will continue with the current efforts to
7 augment existing sequence information in the
8 public domain. There are many people to thank for
9 it. As you can see on the slide, I want to
10 specially thank the FDA Micro team and the
11 collaborators and everybody involved. I'll take
12 questions.

13 DR. TEZAK: I think you have time for
14 maybe a couple of short questions. I want to ask
15 the panelists if they can come up here while we're
16 answering questions. If there are no questions we
17 can start the panel.

18 SPEAKER: If the red light is on that
19 means the microphone is working, and if you guys
20 could pass around the other phone so people can
21 hear us. Whenever you need the slides to be
22 changed let me know.

1 SPEAKER: We're excited to have a few
2 representatives from a variety of the sequencing
3 companies here to talk some about where they see
4 their technologies and sequencing technology in
5 general going in the future. I think I'll let
6 each of them introduce themselves, and maybe we'll
7 go from Mya down that way. First if you could
8 introduce yourself, where you're from, and then
9 just describe in a couple of minutes your
10 technology, your company, and where you see your
11 company going next with either the current
12 technology that you have or with new technologies
13 that you might develop.

14 MS. THOMAE: I'm Mya Thomae. I'm Vice
15 President of Regulatory for Illumina, also was
16 founder of MyRAQA which was purchased by Illumina
17 back in July. I think we've talked a lot about
18 Illumina Technology today, so think everybody
19 understands that. I think Illumina's been pretty
20 public that we're really going through a clinical
21 transformation process at Illumina and really
22 trying to bring NGS into the clinic.

1 I think we're just getting to start with
2 what we've done so far. The CF Assays are just
3 what we're starting with. I think there's massive
4 potential in all parts of the clinic, and we've
5 already announced some collaborations with PhRMA
6 to try and do some cancer-panel- type of work in
7 competing diagnostics as well.

8 I think the sky's the limit, which is
9 why I'm here, so think there's going to be a lot
10 of interesting things in the near future.

11 MR. ROSENFELD: Hi, I'm David Rosenfeld.
12 I work for Complete Genomics where I manage their
13 data analysis pipeline. Complete Genomics is
14 different; we're not an instrument vendor. We're
15 a genome-sequencing service, and we only do whole
16 human genomes. That allows us to focus our
17 efforts in various ways. We have our own
18 proprietary and, sort of, novel strategy for doing
19 that. We built our own instruments and we use our
20 own sequencing chemistry and our own assembly
21 algorithms.

22 About 2 years ago we were purchased by

1 BGI, so for the last 2 years we've actually been
2 developing technology for them. They have a lot
3 of our sequencing instruments in China now, and
4 they've actually used it to develop a noninvasive
5 (inaudible) assay that's been cleared by the CFDA
6 in China.

7 I can't speak to a lot of future
8 technology that we're working on. The one thing
9 that I can discuss is things that we've announced
10 already, and that involves our LFR technology.
11 That's something we call Long Fragment Read. It's
12 a way of basically separating long fragments of
13 DNA early in the process, tagging the shorter
14 fragments that come from that, and using that tag
15 information through the assembly process to get
16 much higher accuracy in our short variant calls.
17 We can basically count how many similar tags show
18 up in the assembly and filter using that. That
19 process also allows us to phase the genome, and it
20 allows us to assemble from a very, very small
21 number of cells, so it's a pretty exciting
22 technology. It has difficulties in

1 commercialization that we're working on, but
2 anyway that's what I have to say.

3 DR. SHAW: Thanks. I'm Jay Shaw. I'm a
4 Senior Director at Thermo Fisher. It was
5 discussed earlier; we just announced registering
6 the PGM Dx platform with the FDA. Similar to what
7 Illumina has done, we've gone with a genotyping
8 claim for isolating old genomic DNA from whole
9 blood. Plan to move into content. We've also
10 announced a collaboration with PhRMA for cancer
11 diagnostics, which will bring us into the arena of
12 detecting somatic mutation.

13 In terms of where we want to move in the
14 future, it's still a fairly complex technology.
15 We view automation as being really important;
16 simplifying the process for the end user and just
17 making it so that we can push it out further and
18 further into less-complex labs so that it can be
19 used there.

20 DR. KORLACH: Thank you. I'm Jonas
21 Korlach. I'm the Chief Scientific Officer at
22 Pacific Biosciences, and I was told that we would

1 be able to show a few slides, and I took advantage
2 of that offer.

3 At Pacific Biosciences we've taken a
4 little different approach in terms of the
5 sequencing performance in that, as you know, there
6 has been a fairly strong focus on the throughput
7 of the sequencing systems. At Pacific Biosciences,
8 we wanted to build a technology first and foremost
9 that provides the highest quality of data in all
10 the four areas that are relevant when evaluating
11 really any sequencing technology, and you see them
12 up there.

13 As you may know, we have, by far, the
14 longest reads and the least bias, which gives you
15 the opportunity to sequence the entire genome of
16 microorganisms or also the human genome. Then as
17 was mentioned in some other talks, we have the
18 unique capability of looking at the chemical
19 changes in the DNA, the epigenetic methylation,
20 and other changes as part of the sequencing
21 project.

22 In the next slide you see the

1 development by us and others that now make it
2 feasible in an automated fashion to generate
3 finished, microbial genomes and (inaudible)
4 outlined that with the University of Maryland for
5 comprehensively describing and characterizing
6 pathogens. This has now been used by the
7 community.

8 The next slide just shows a partial list
9 of publications that have been described that use
10 that capability to now obtain in a very efficient,
11 automated workflow and also cost efficient to get
12 comprehensive information about the pathogen of
13 interest. I just want to show one recent example
14 in the last slide which was a paper that came out
15 last week in Science Translational Medicine by
16 Julie Segre from across the road here at the NIH
17 Clinical Center using this technique to understand
18 the horizontal gene transfer, the plasma
19 trafficking, that happens in the
20 multidrug-resistant enterobacteriaceae like
21 klebsiella that caused these multi-drug resistant
22 hospital outbreaks.

1 This figure is from a commentary in
2 science showing that you really -- because these
3 events are complex and it's now realized that in
4 microbiology, hospital- acquired infections, food
5 safety, and in the human genome the variation is
6 not just limited to snips. There's large
7 structure variations. The acquisition of new
8 plasmids or rearrangements of plasmids play
9 crucial roles in understanding whether a
10 particular salmonella strain will become an
11 outbreak strain or not. With the technology, it's
12 now possible to understand that more
13 comprehensively, and that's being utilized not
14 only in terms of hospital-acquired infections.

15 I talked to Michael Art; they're using
16 their two Pacific Biosciences instruments to
17 generate finished genomes to have reference
18 strains for the major foodborne outbreaks, and
19 that has two advantages. First, you may actually
20 understand better what's going on and what makes
21 that particular strain an outbreak strain.
22 Secondly, having this complete genetic and

1 epigenetic information may give you more sensitive
2 and specific markers now to track the evolution of
3 the outbreak and to track the source in
4 conjunction with the other technologies.

5 With regards to where we're going, we
6 are for researchers only right now. As you may
7 know, in October of last year we announced a
8 strategic partnership with Roche for leveraging
9 the power of this technology into the diagnostics
10 market. Now, as you know, that take a few years,
11 but we are working with them very closely and
12 ultimately getting this technology into the
13 diagnostics base as well.

14 DR. ZOOK: I have one or two more
15 questions here that I'll just let whoever wants to
16 talk about the questions answer. If you in the
17 audience could be thinking about any questions
18 that you have also, we'll open it up for questions
19 for any of you also.

20 The next question that I had was what
21 types of sequencing you think will be done
22 clinically, say, 5 years from now or so? Do you

1 think it will be whole genome versus targeted
2 sequencing? Do you think it will be just human or
3 will there be a lot of other microbial sequencing?
4 Will it be single molecule or amplified or in
5 highly specialized labs or in more (inaudible)
6 pipelines or even maybe bedside sequencing? Just
7 interested in your thoughts as to where it might
8 be moving in the future.

9 DR. SHAW: I think from a technology
10 standpoint 5 years is a long way away, but 5 years
11 for a regulated device is a short period. As a
12 person in development, I always appreciate how
13 long it takes to do things. I don't think we'll
14 see this moving, in my opinion, that fast to
15 getting to the bedside in terms of 5 years.

16 The problem is you start the development
17 process, and you have to lock down. As the
18 company continues to make improvements, they're
19 always asking, 'Can you roll that in? Can you
20 roll that in?' The problem is you have to start
21 the process all over again. You have to make the
22 commitment early on, and it takes a while to get

1 this in. Within the next 5 years, my opinion is I
2 don't think we'll see bedside sequencing going on.

3 MR. ROSENFELD: That was, kind of, a
4 complicated question maybe. From our perspective,
5 we're all about scale and doing things at a big
6 scale. We certainly expect the clinical markets
7 to be much bigger than the research market was.
8 We think whole genome sequencing has a place
9 there. It's not the only thing that's going to be
10 there, and certainly not in 5 years. It won't be
11 the only thing that's going to be there, but we
12 think it's got a place there, especially around
13 cancer.

14 But a lot of it is going to be what
15 brilliant breakthroughs come through in the assays
16 that get developed that really drive a market
17 because that's where the investments going to be
18 put into on our part and our competitors' parts.
19 That's where you'll see the really rapid progress.

20 MS. THOMAE: I would agree with Jay that
21 I think bedside sequencing in the next 5 years --
22 we haven't even gotten PCR to the bedside yet.

1 But I don't know, Zivana, if you're going to cut
2 us a break on this one. I'm just kidding.

3 (Laughter) I think that's probably ambitious in
4 the next 5 years, but I do hope it's before I
5 retire.

6 I do think there's a lot of work in
7 human, but I think the microbiology applications
8 to this are very interesting as well. I really
9 appreciate the previous presentation on it, and I
10 think there's a lot to be learned.

11 DR. KORLACH: Yeah, I agree with what
12 was said, but I want to maybe add that I think
13 there will be a gap in the research arena in what
14 we will be able to do.

15 I'm going to make a prediction -- it may
16 not be true -- but I think I would expect that 5
17 years from now we will be able to generate de novo
18 genomes for not only microbes but also the host,
19 the human genomes. I think we'll eventually free
20 us from the reference and from a reference-based
21 research --I'm talking about the research setting
22 -- and would be able to de novo assemble human

1 genomes in an efficient and fast manner 5 years
2 from now.

3 That relates to then a challenge both
4 with respect to the regulatory environment -- how
5 do you deal with that? How do you deal with the
6 structural variation with everyone being a true
7 individual with respect to their genetic
8 blueprint? Then also bioinformatics challenge
9 because right now we don't really have the
10 infrastructure and the methods to interpret that
11 data and take that information from literally
12 hundreds of thousands of individuals for which
13 this information, both in terms of the host and
14 then the pathogen in a metagenomic setting, will
15 be created.

16 I think it's going to be very exciting.
17 There's a lot of work to do in all of those areas.
18 All the sequencing technologies are advancing
19 extremely rapidly, so I think we should be
20 thinking ahead and stimulating all the young
21 bioinformaticians to be aware of that and get
22 excited about that.

1 MR. ROSENFELD: If I could add to that,
2 that LFR process that we have actually does allow
3 for de novo sequencing of the human genome, and it
4 creates all kinds of interesting problems for data
5 formats and things like that where all these
6 resequencing standards that we have work a certain
7 way, but when you're outside of that and your
8 first step isn't mapped to the reference, all that
9 stuff doesn't make so much sense anymore.

10 DR. ZOOK: The last question I had was
11 much shorter. What challenges do you see with
12 getting the new technologies into the clinic in
13 the near future?

14 MS. THOMAE: I'm a regulatory person, so
15 I usually see the regulatory challenges first. I
16 think there's the issue of how do you find the
17 right validation strategies? Are you trying to
18 look at everything you might find or are you
19 trying to look at representative sorts of things?
20 What are your reference standards? This can
21 become a pretty complicated conversation.

22 Then I think there's a lot of

1 interesting challenges, and I know FDA is thinking
2 about this too. When you start to pour it over
3 content, for example in competing diagnostics,
4 that's already been approved for other PMAs, and
5 you want to maybe put that content on your cancer
6 panel, sort of, how do you do that? What's the
7 best way to do that? I think the regulatory
8 challenges are pretty interesting, but I think
9 there's a lot of good thinking being done about it
10 right now.

11 Then, of course, we can't always be
12 U.S.-centric. There's a lot of interesting
13 challenges around the world with the different
14 kinds of testing that's done. Some of it, like
15 the prenatal testing, has various different
16 cultural impacts. Not everybody maybe wants to
17 find out everything that Americans do. I think
18 there's a lot of very interesting questions to be
19 handled as this moves to the clinic.

20 MR. ROSENFELD: I think almost
21 everything is a challenge to get into the clinic.
22 It starts with quality, that's a challenge. I'm

1 not sure any of us are good enough yet. Then the
2 regulatory stuff, and then there's just a huge
3 challenge with acceptance by doctors and having
4 them be educated to understand what all this stuff
5 does.

6 We're really way ahead of a lot of
7 things in terms of technology and trying to push
8 it into the clinic when the clinic isn't asking
9 for it, so it's pretty hard sometimes.

10 DR. SHAW: I also think as we move
11 forward with these large, complex panels, it's
12 hard enough to develop a single analyte test. Now
13 we have to factor in for all of the analytes that
14 we have to deal with, so sample availability,
15 performance across all of these analytes is going
16 to be, I think, the real challenge in terms of
17 generating the body of data you need to prove the
18 performance of these systems.

19 DR. KORLACH: I fully agree. I'm a
20 method developer; I'm a little more removed, but
21 Eric Schadt, who is at Mount Sinai and just doing
22 clinical sequencing there with the different

1 technologies, he had a Mendelspod podcast this
2 week, and he was asked the same question. He
3 highlighted the need for better education and
4 openness of the clinicians. We need to train the
5 existing ones and bring up the young ones who are
6 having a much more open mind with respect to these
7 kinds of datasets and engaging in the opportunity
8 to use those for making clinical decisions.

9 I'm personally interested in maybe, for
10 this forum, two philosophical questions in how
11 this concept of personalized medicine and knowing
12 something about someone in a very detailed fashion
13 is going to work out with, on the other hand,
14 having lots of cohorts and trials and medicine
15 that works for lots of people and to save lots of
16 people and so forth. To me, that's an interesting
17 paradigm consideration on how we make that work
18 where eventually we will have the information that
19 we know for a certain person that drug works, but
20 how do you prove that? How do you validate that,
21 and how do you get that into an environment that
22 makes the public feel safe and feel good about the

1 process?

2 DR. ZOOK: Thanks, everyone. Are there
3 any questions from the audience? There are
4 microphones around.

5 SPEAKER: They're focused on the
6 bacteria, so are any of these methods applicable,
7 sort of, the soup of blood? You don't have to
8 isolate first? I want a follow-up question if
9 they're -- it wasn't clear because that would be
10 great, and there are some tests that now with DNA
11 testing.

12 MR. ROSENFELD: Like with our LFR test
13 that works from a small number of cells, you do
14 have to do a cell selection out of the blood, but
15 you're --

16 SPEAKER: You have to isolate the
17 bacteria?

18 MR. ROSENFELD: You're isolating either
19 a circulating tumor cell or (inaudible) DNA
20 (inaudible) way it works.

21 SPEAKER: I was talking about the
22 bacteria groups here and whether they can attack

1 just a raw sample and figure out what it is.

2 DR. KORLACH: Right now with Pacific
3 Biosciences I think typically it's done with
4 isolating the bacteria. I know of a few studies
5 where they have taken just the whole sample, and
6 then you sequence everything for identification
7 (inaudible) what's there, and I know that's also
8 been true for the other technologies. Joe Diresi
9 gave a very nice talk in February using the
10 Illumina system to sequence a patient, taking just
11 the blood and sequence everything and then find
12 that there was a certain pathogen in there.

13 I think that will likely become more and
14 more common that you don't need to do these
15 manipulations, and you can just take the DNA of the
16 entire mixed sample, sequence everything, and then
17 bioinformatically filter out what you're
18 interested in.

19 SPEAKER: The follow-up question is
20 there's a lot of tests currently for probes
21 mostly, probe-based or probe PCR, for identifying
22 any of a sequence of bacteria or whether they have

1 the various resistance genes. Will these methods
2 replace those? Are there advantages or
3 disadvantages assuming you isolate?

4 DR. KORLACH: Yeah, from my perspective,
5 my prediction is that these methods will be
6 replaced eventually as it becomes more standard
7 and cheaper and a, sort of, automated workflow
8 because in addition to what you think you already
9 know is there with the PCR method, you need to
10 know what those -- and the primary sites have to
11 be there and the genome and so forth. The
12 sequencing gives you de novo information and gives
13 you information about something that you didn't
14 know was there. It's a hypothesis-free approach,
15 which I think inherently is more powerful.

16 Right now, I believe there still is a
17 hybrid approach. When you have a de novo
18 reference you can design specific primers to a
19 certain outbreak that makes the tracking of that
20 more efficient. Ultimately, I think it'll just go
21 to sequencing.

22 SPEAKER: Do you see a role for, let's

1 say, paleobiology where you're looking at, for
2 instance, mummies or other human tissues for human
3 diseases or pathogens? You start out with small,
4 short sequences, so you're halfway there, but
5 they're from complex genomes.

6 DR. KORLACH: Maybe I'll start. We've
7 had a couple papers where that's been used. Now
8 obviously you don't need the long (inaudible) from
9 Pacific Biosciences, but the fact that you don't
10 have to amplify the DNA helps there. There have
11 been some studies there also looking uniquely at
12 the base modifications in these ancient samples.
13 From our point of view, that's been demonstrated
14 in the research setting and is feasible.

15 MR. ROSENFELD: I'm pretty sure no one's
16 used our technology for that. (Laughter)

17 DR. ZOOK: As far as you know.

18 MR. ROSENFELD: I think I would have
19 noticed. I do the QC on the genomes. I think I
20 would have noticed.

21 SPEAKER: Thanks to the efforts of all
22 of your companies, we have technology now that can

1 sequence willy- nilly. We all can talk about ways
2 to validate and do it medically responsibly.
3 Indeed, we can sequence far better than we can
4 clinically interpret. Most variants that aren't
5 obvious, common polymorphisms that you find in a
6 patient are VUSs.

7 Even worse, of the stuff we can
8 clinically interpret, we can diagnose far better
9 than we can treat. We can definitively say, we
10 know what's wrong with you. I'm so sorry, there's
11 nothing we know to do about it today.

12 My question is the fact that information
13 about clinical utility at variations that are
14 discovered during sequencing, both somatic and
15 constitutional, is changing very, very rapidly,
16 what's the best way to incorporate that into
17 clinical practice but to make sure it's done
18 responsibly? Let me start with Mya just because I
19 like to put you on the spot.

20 MS. THOMAE: Great question. I hadn't
21 prepared for that one. I think some of it's about
22 being willing to go back to the data that has been

1 generated for a patient as we understand more.
2 That gets tricky in terms of saving large
3 gigabytes of data, but I think you have to be
4 willing to figure out how to report things and
5 potentially be bucketing them appropriately at the
6 time the test is done but then be willing to go
7 back and relook at it.

8 For example, there's a number of
9 consortiums working on this in cancer and trying
10 to bucket things by known clinical relevance,
11 maybe clinical relevance, and we can find this,
12 but we don't know what it means. There is the
13 Actionable Genome Consortium, something Rick
14 Klausner's been working on, that's really trying
15 to do that for cancer, and I think that's probably
16 needed all over the place; in HLA and everything
17 else.

18 But I think it is a process of
19 revisiting it periodically. How that happens in
20 the clinic, how our payer system works to be able
21 to support something like that, how we make sure
22 physicians get this kind of information, and not

1 just at the large institutions but across the
2 board; those are questions that are pretty tricky.
3 They're definitely not FDA or regulatory driven.
4 I think they're more in how healthcare is
5 delivered.

6 I'm working on world peace, and I'm
7 working on fixing our healthcare system in the
8 United States at the same time. I've got a couple
9 PowerPoint presentations that I'd love if you
10 could work with me on those. It's pretty big
11 issues.

12 SPEAKER: Which is harder: World peace
13 or fixing healthcare? (Laughter)

14 MS. THOMAE: Fixing healthcare.

15 MR. ROSENFELD: I'm not an expert on
16 this, but it seems like different technologies are
17 going to move at different rates all the time. If
18 we want the understanding of what all the
19 polymorphisms mean, to speed up, it means you've
20 got to do bigger studies and you've got to, sort
21 of, free all that data that's already been
22 generated.

1 Right now it's all in little islands of
2 individual studies, and there's no way for people
3 to share it. That was discussed yesterday as part
4 of the cancer discussion and NCI discussion. It's
5 that way for almost every dataset that's out
6 there. I know we've produced almost 20,000
7 genomes for our customers, and there's not way
8 that can all be put together to do a study.
9 That's going to continue to be the case unless
10 somehow there's some kind of framework that's put
11 together to fix that. That would be, to my mind,
12 a big first step in trying to move that analysis
13 along.

14 DR. KORLACH: I fully agree, and just a
15 small note to add. What I'm impressed by and what
16 I take heart in is the work that's been done at
17 places like Mount Sinai, at the Brode where people
18 are looking at this in new ways, a very
19 integrative approach, lots and lots of data,
20 developing new math and new statistics to treat
21 those kinds of data. I think that's the first
22 step.

1 We don't really know about the genomic
2 variation at all-size scales in the human
3 population, and we need to find that out first by
4 generating a bunch of data. Then it's going to
5 take a couple of years to really be able to
6 understand and treat it appropriately to then make
7 informed decisions for the clinic.

8 DR. ZOOK: We could maybe take one last
9 question depending on how long it takes.

10 SPEAKER: This question is somewhat
11 related to the last one. Hopefully it's a simpler
12 answer. We've talked about personal genomes;
13 there's a 100,000 genome project in the U.K.
14 Given regulatory issues and reimbursement issues
15 and diagnostic issues, we can envision the future
16 where everyone is walking around with their
17 personal genome. Instead of doing a lab test, a
18 physician orders a bioinformatics test. I get
19 consent to look at a particular region and then
20 get the answer instead. Of doing all the
21 informatics up front, I do it as I need it.

22 In that future where everyone has their

1 own personal genome in their medical record, how
2 long do you think it'll take to get that? Will it
3 happen in my lifetime? Will I have it at some
4 point? Will it happen in my kid's lifetime? When
5 do we start planning our organizations around not
6 doing lab tests but doing bioinformatics tests on
7 a genome?

8 DR. KORLACH: My prediction is that
9 you're going to have that in your lifetime, and
10 you're going to have to do that --

11 SPEAKER: Ten years, five years?

12 DR. KORLACH: Something like that;
13 between 5 and years. I think you're going to have
14 to do it multiple times because when you develop a
15 tumor then you want to sequence that tumor too.
16 Then in terms of gene expression with the
17 different tissues, there's going to be changes
18 there. But 10 years, I'd be surprised if it'd
19 take longer than that.

20 It's going to be expensive and maybe not
21 reimbursed initially, but I know in Japan some of
22 our users there are already approaching a certain

1 customer base that are perfectly happy to spend a
2 few \$10,000 to get that kind of (inaudible).

3 SPEAKER: I guess I'm looking for when
4 is it becoming the norm that you assume that
5 everyone has their personal genome?

6 DR. KORLACH: I can't answer that.
7 (Laughter)

8 SPEAKER: Yeah.

9 DR. KORLACH: (inaudible) because that
10 goes beyond the scientific feasibility.

11 SPEAKER: Right, right.

12 DR. KORLACH: There's a whole bunch of
13 other things that needs to happen for that.

14 MR. ROSENFELD: For sure less than 10
15 years scientifically; there's no question. You
16 can almost do it now if you wanted to if you're
17 willing to pay for it. I'm not sure if it'll ever
18 be the norm in the United States. If you look at
19 market research studies of that kind of thing,
20 most people don't actually want to know those
21 things. There's actually only a small population
22 of people who are interested in getting their

1 genome sequenced, so I don't think it will be the
2 norm.

3 SPEAKER: We get the data so that we can
4 do the test later cheaper, so we're not actually
5 getting any information other than Gs and Ts, but
6 then that's something to mine and interrogate when
7 you want to know it. That would be the model.

8 MR. ROSENFELD: I understand. People
9 are nervous about that. It's a question of trust
10 of what's going to be done with those Gs and Ts.

11 DR. KORLACH: The culture could change
12 over the next 10 years possibly.

13 DR. SHAW: I think the point you just
14 made about the difference, the somatic mutation,
15 detecting it later when there's a disease state
16 means multiple times in your life potentially as
17 opposed to an early sequence that'll carry you
18 through your lifetime and the need for that.

19 DR. ZOOK: Okay, let's close by thanking
20 all the speakers and the panelists in the second
21 -- (Applause)

22 MS. VOSKANIAN-KORDI: At this point,

1 we're going to take about a 15-minute break, and
2 we'll reconvene at 3:15 for the food safety and
3 pathogen protection section.

4 (Recess)

5 MR. PETTENGILL: So I guess we'll get
6 started. Thanks everyone for coming. I guess
7 this is the final session of the two day meeting.
8 This one's on food safety and pathogen detection.
9 I'm James Pettengill and then Heike Sichtig will
10 also, will be coordinating this. I think we have
11 three talks that will focus on food safety and
12 pathogen detection, from, one from the FDA, Center
13 for Food Safety and Applied Nutrition, one from
14 NCBI, dealing with some of the database issues and
15 performing whole genome sequencing, and then a
16 final talk from somebody from FSIS.

17 I think it's sort of safe to say, at
18 least from the FDA SIPSANS perspective that next
19 generation sequencing has sort of really
20 revolutionized the food safety and pathogen
21 detection. I mean, I think we'll be hearing a lot
22 about that today from the first speaker, and then

1 I think another big part of it is to bring the
2 whole genome sequencing and food safety into
3 practice and in production, it's been a sort of
4 multi-agency collaboration. I think we'll hear a
5 little bit about that today also. So with that I
6 guess we'll get started.

7 Our first speaker today is Mark Allard
8 from the FDA Center for Food Safety and Applied
9 Nutrition. He'll be talking about FDA Genome
10 Tracker surveillance and source tracking for
11 Salmonella Listeria. So we're just zooming down
12 into a bacterial genome here. We'll get to the
13 base level here soon enough. (laughter)

14 SPEAKER: I just want to say as an
15 introduction to my talk, I hope all of the food
16 and genetics folks (inaudible), I guess what I
17 want to say is what we're doing to bacteria, will
18 be the future for (inaudible) and for food the
19 future is here and now. And so what we'll be
20 talking about microbial trace backs (inaudible)

21 MR. ALLARD: Thank you and I'll proceed.
22 Essentially, to bring next generation into the

1 regulatory lab, and I represent the research end
2 of the Center for Food and Safety and Applied
3 Nutrition, which essentially supports our
4 regulator state and federal laboratories, we did
5 timing testing to say, will next generation lab
6 response be at least as fast as traditional typing
7 procedures, which is PFGE, and of course the
8 answer is yes. The current technology is
9 relatively rapid; in three or four days to a week
10 you can have an answer. This may not be the
11 clinical level, at fast enough, but in food safety
12 and trace back of outbreaks, this is a much great
13 improvement over the traditional methods.
14 Essentially what we're leveraging with whole
15 genome sequencing of bacteria, is old ideas in
16 evolution that dated all the way back to Darwin
17 and that is that the principle of inheritance will
18 betray the original birth place of organisms. And
19 we've known this from the beaks of finches to the
20 traditional museum curation and taxonomy and the
21 genome is just additional characteristics that
22 will trace organisms right back to the source.

1 And so for FDA, it's critically
2 important to understand where an outbreak
3 occurred, where the contamination is arising in
4 the food supply and to be able to remove it from
5 the food supply. And so in 2012, we had sort of a
6 watershed event for FDA in an example of a sushi
7 outbreak. And this is the traditional view, is
8 you see a large cluster of clinicals all of
9 identical genotype. This would be done based on
10 PFGE, but you don't know where the food is from
11 and so the trouble with -- the difficulty with the
12 PFGE is for some bacteria like salmonella, they
13 are so clonal that they all share the same PFG
14 type. So it doesn't really help you very much
15 with source tracking, because your last ten years
16 of isolates could all have the same pattern.

17 But a genomics provides additional
18 resolution and in essentially actionable
19 information in helping us identify where the food
20 came from. And so, and then when you actually
21 find the food and sequence it, in this case it was
22 a single point source outbreak, the bacteria from

1 the food is directly imbedded within the clinical
2 samples. And then when you actually do an
3 inspection, all the environmental swabs will also
4 sit in the middle of that.

5 But what's particularly important with
6 this example was we wanted to say, well, are there
7 any investigational leads that can be provided by
8 just going into our freezers, diving into them and
9 sequencing all the Salmonella Bareilly that we had
10 in the freezer, or Bareilly that was of the same
11 PFG type. And we got about a dozen isolates. Six
12 of them were PFG identical, and if you look at the
13 link just above the outbreak isolates, there's a
14 sample from India from shrimp from an inspection
15 five years earlier. And genetically, it's about
16 25 snips away, so it's not the source. The
17 outbreak cluster's only two to five snips away.
18 But it is geographically a lead. And the PFGE
19 supported that it's somewhere in Southeast Asia
20 but it could be the gulf or even potentially the
21 Pacific. Okay, but genomically, it says we're on
22 the west coast of India, and in fact, if you do

1 overlay this on a geographic google earth map, so
2 you're laying the tree on top of google earth.
3 This is what an outbreak looks like along the
4 eastern seaboard and this is a picture of where
5 the contaminant is coming from. And in this case,
6 the Salmonella Bareilly from five years previous
7 had occurred from eight kilometers away from the
8 actual produce packing house, where the
9 contamination occurred. The reason we did such a
10 good job in geographically source tracking this
11 isolate was it was a repeat offender.

12 Now we don't know if the isolate is
13 reoccurring contaminants in the wild from that
14 region and the fish that come in occasionally
15 carry that, or whether this is a contamination
16 within the fish processing packing house that's in
17 the equipment that is sporadically coming out.
18 But essentially this gave us the understanding of
19 how powerful the source tracking could be with
20 doing whole genomes and comparative clustering.
21 And so the idea was, well we should sequence
22 everything in our freezer or at least the last

1 five years of outbreaks, put them in a database
2 and use them essentially to provide leads of where
3 to send inspectors. And the technology is driven
4 by the drop in cost of equipment and so now it's
5 feasible to sequence hundreds if not thousands of
6 bacterial isolates to draft sequence quality
7 levels.

8 To build at the national database, we
9 need three components. The first are labs to
10 generate the whole genome sequence, and we're
11 interested in the distributed network because
12 that's where the isolates are coming in, a local,
13 state department of health, a local clinic or a
14 federal inspection facility. We need network
15 management. And then we need sequencing storage
16 data provider and analysis and essentially, we've
17 off, put off our data storage and partnered deeply
18 with NCBI and so they're taking over all those
19 costs, data storage. And we've been following
20 them in recommendations on how to move data
21 efficiently up to the NCBI. And Bill Klinke
22 (phonetic 010:12.4) will be talking about those --

1 that partnership and the part that the NCBI does.

2 But essentially, by shifting the cost
3 over to the MIH NCBI, we have more funds to
4 provide to more state and federal laboratories, to
5 actually put the sequencers and reagents into a
6 distributed network. Though the basic model is to
7 have national and international distributed
8 networks of sequence uploaded to a common place,
9 this is the NDC which is the NCBI, EMBL and DDBJ
10 database, and we call this database Genome
11 Tracker. And then all the sequence, all the draft
12 sequences that I'm talking to you about are
13 publicly available in real time, uploaded, and
14 that means that our federal partners can get
15 access to it, our state and foreign public health
16 laboratories and of course academics and industry.
17 The data's transparent. All of the data's
18 available.

19 And so our initial network -- the
20 network as it exists today is six state labs and
21 eleven federal laboratories. But his is rapidly
22 expanding as state networks; state departments of

1 health purchase their sequencers and are ready to
2 start uploading to the database. And we're in
3 actually the second year of a pilot study. And
4 the pilot is to build a reference database to
5 source track unknown clinicals against the known
6 food and environmental that are sitting in the
7 database.

8 And you can go to this genome tracker
9 database. It's just a bio project at NCBI. I
10 think by the end of the year, right now, we have
11 10,000 preregistered strains and 6,000 genomes.
12 This slide's already out of date. I think it's
13 7,500 Salmonella and almost another 1,000
14 Listeria. So we're going to top out of 1,000
15 isolates of mostly unpublished draft genomes,
16 uploaded in real time, available to anyone in the
17 world.

18 We are also providing, and this is, not
19 we, that's the royal we -- it's Bill and his
20 colleagues at NCBI -- they have an ftp database
21 where they do initial draft analysis and upload a
22 tree that could be visualized in their software

1 called Genome Workbench, also free. And so we
2 have existing databases now of e-Coli, Listeria
3 and Salmonella that are actively growing and
4 people are collecting it, state and federal
5 partners. And essentially, I'm just showing you
6 an example where you can see the new tree, and I
7 think for some of these databases we get a new
8 tree every day that we're averaging five to six
9 hundred new isolates every month and the trees are
10 enormous, you know, when you're talking about
11 thousands. It's hard to actually see the data in
12 a coherent way. And we're working with NCBI and
13 other third party software providers to give us
14 new tools to rapidly say, is there a new cluster?
15 Because what we're really interested in at FDA is,
16 is there a growing cluster? Is it new? Is there
17 a new association within a food and environmental
18 sample? But frankly, the CDC has been using this
19 and just understanding which clinicals are
20 associated, because usually the traditional PFG
21 would say ten are in a cluster. They're all
22 identical PFG. But really there might be two or

1 three clusters there, of the independent sources
2 of contamination arising at the same time. And so
3 over time you can look for this. Like I say,
4 we're working on tools. And if you dig into the
5 tree, you can get the associated metadata, which
6 says that it's clinical or environmental, it's
7 associated with peanut butter, cantaloupes. And
8 so we're using this data. So, like all of the
9 genomics that you've been hearing about, as the
10 cost drops and the applications increase, the
11 rapid rate of increase of genomes is rapidly
12 increasing and so the Salmonella and Listeria,
13 which are what we're primarily emphasizing, are
14 also rapidly increasing.

15 I mentioned the CDC, FDA, NCBI and USDA,
16 Listeria real time sequencing project. This was a
17 project whereby the CDC decided to get involved
18 last year and they said for their first project,
19 they said for all Listeria, we'll sequence all
20 Listeria, clinical cases in the U.S. this year,
21 and FDA and USDA agreed to do all the food in the
22 environmental, as well as dig into our freezers

1 and sequence and upload them into the database.
2 And essentially, paired with this, the CDC already
3 has programs where they're actively doing the
4 traditional epidemiological questionnaires of
5 every patient, okay? So you have good epi paired,
6 at least for the patients that answer the phone
7 and will talk to a regulator, to add this. And so
8 this is real time surveillance, but we've also
9 been doing real time surveillance of some
10 Salmonella. So for example, the states wanted to
11 look very closely at another very homogenous
12 Salmonella genome, the Salmonella Enteritidis,
13 which is what's in chickens and eggs.

14 So in Listeria, it was a watershed year,
15 2014 was a watershed year, combining inspection in
16 clinical and food isolates, the FDA closed Bruce
17 Cheese based on genomic information. This wasn't
18 epidemiological information that was supported by
19 genomic information, or used at a later date or as
20 a case study a year later. This was in real time,
21 use of genomics to show the close relationship
22 with the foods that were coming out of this

1 company and the clinical cases of the people that
2 were getting sick. And this was significant in
3 that at least one person died and several babies
4 were boarded and lost. Listeria is a killer for
5 the very old and the very young. And on this is
6 so, actively we've done cases and now, once you've
7 done a successful case for the government, they
8 don't ask you to stop, they give you more and more
9 and more cases to start doing this on a regular
10 basis. And so we're now in a zone where we're
11 looking at a new case about every month and
12 continuing to essentially integrate. There's a
13 lot of integration. This is a new tool and the
14 epi's need to know who pulls the trigger. When do
15 we say we should start sequencing? There's a
16 whole sort of coordination that's directly
17 affecting regulatory pieces of FDA.

18 I was going to, as long as this is a
19 standardization meeting, I want to talk just a
20 little bit about validation, though this is
21 usually Errol's job, to do this part. We've done
22 a lot of validation. Essentially, if you're going

1 to use this as a regulatory tool, then you're
2 going to close down companies or put people in
3 jail, you have to have a validated approach that
4 it's reproducible and you get the same answer.
5 And so we've done technical performances on
6 standards. So we work very closely with Justin
7 Zuck and PICA in helping build the reference
8 material and sequencing reference material on
9 multiple platforms. We've also done several
10 intra-laboratory variation sequencing studies with
11 hundreds of runs as well as typical
12 intra-laboratory variation and then also testing
13 the actual bioinformatics pipeline because the
14 pipeline that brings you to the interpretation
15 also has to be validated, reproducible and
16 accurate, and so you also want to have good epi to
17 back it up to say you're getting the right answer,
18 or, if you're sequencing the same isolate, you
19 know what the answer is.

20 And what this shows is that essentially
21 there is variation from laboratory to laboratory
22 and people can have a bad day, but in general the

1 sequences either work or they don't work. It's
2 pretty obvious when a sequencing run fails. And
3 you can see this in measurements like the N50 or
4 the number of contigs, or actually even just the
5 number of reads in the instrument. It's
6 relatively straight forward. When you do cost
7 comparisons, and you just count up the number of
8 snips that are different in what should be
9 identical runs, you rarely see more than one or
10 two differences and this is in the same range as
11 within a single point source outbreak, where we
12 see say 1 to 5 snips in the intra-lab comparisons,
13 and ten snips even lower.

14 And then, I don't have a picture but
15 essentially when you load it into a file of
16 genetic tree, that's the other and most important
17 piece is, would you ever make a different
18 regulatory decision, to say this is part of the
19 outbreak or not part of the outbreak? And we
20 never see these cases. And so we're working to
21 publish that, those results, because this kind of
22 key validation, but I like to say, all the data is

1 up in the Genome Tracker database, or much of it
2 is. And so people can look at this data and
3 carefully study it themselves or test their own
4 validated pipelines.

5 So the submission test, NCBI also has to
6 be coordinated, integrated and standardized across
7 the network and so we're working carefully on what
8 kind of information and metadata goes up with the
9 sequence data, the level of coverage, and the
10 basic QC that's done. Here's another example
11 where a QC score, where you can see, it comes
12 right out of the top where the number of reads is
13 much smaller and you can instantly see that there
14 was some problem with this particular run. The
15 fact is is that these robots are very robust and
16 they don't make very many mistakes. The biggest
17 mistakes are always human error. It's
18 mishandling, mislabeling, cross contamination,
19 things that the bacterial community has long known
20 are issues. And so you have to do careful
21 bacteriology and microbiology to control that.

22 There has also been a large an ongoing

1 discussion about what's the minimum metadata.
2 We're releasing this to the public. And so
3 there's less data than we can -- as a regulatory
4 agency we know much more detail about an
5 inspection than the sampling. And that's behind
6 the fire wall of FDA and we would share that with
7 a legitimate regulatory body, like a state
8 department of health or a foreign body. But it's
9 not made publicly available, and so it's all
10 linked to a key. And I'm not going to give you
11 examples. Go to Genome Tracker and just randomly
12 pull out samples. There are many needs in this.
13 We need widely available commercial solutions. We
14 need customized solutions. We need more automated
15 methods for microbial, anti-microbial resistance,
16 virulence, and we need very -- we're working
17 closely on easy forms to interpret physicians,
18 epidemiologists and researchers working directly
19 with our epidemiologists at FDA and the regulatory
20 bodies in the food safety to directly support
21 this. Thank you.

22 (applause) Questions?

1 SPEAKER: Thanks, it's very interesting,
2 impactful work. Thanks for the talk. Just
3 wondering about the (inaudible) that follow
4 genetics. What kind of techniques are advocated
5 -- is it phasean, maximum likelihood or
6 (inaudible)

7 MR. ALLARD: So we've used many methods
8 and that's a slide I didn't leave in and so, but
9 the current FDA procedure, if you do a search at
10 github for SIPSAN you can see our whole published
11 procedure where it's transparent. The data
12 analysis is full transparent and described with
13 open source software. The fact is is that what
14 FDA is interested in is the very tip of the tree,
15 because you were saying, is this clinical and
16 another clinical match and does it match a food,
17 and so we don't really care if the lower parts of
18 the tree are correct or not. It's the very tip is
19 where the regulatory decision is. And because
20 these are such recent events, almost any
21 phylogenetic method will work.

22 SPEAKER: Okay.

1 MR. ALLARD: There are very low levels
2 of homoplasy.

3 SPEAKER: Okay.

4 MR. PETTENGILL: No questions from
5 online.

6 MR. ALLARD: I just want to say with my
7 -- oh go ahead.

8 SPEAKER: Sorry, quick question. Have
9 you considered including academic partners in the
10 tracker network?

11 MR. ALLARD: We already have academic
12 partners in the network. Early on, essentially,
13 most departments of health don't have
14 bioinformatic excellence; don't have any genetic
15 or genomic experience. So for example, our New
16 York lab in Wadsworth in Albany works closely with
17 Cornell food safety and in our site in Florida,
18 works closely with University of Florida in
19 Gainesville. And so, but usually we're working
20 with, directly with state partners because we're
21 trying to support and improve this for regulatory
22 action and in department of health labs. But yes,

1 not only that, anyone in the network, we've asked
2 them to be regional hubs and reach out, because
3 isolates could come from a microbial ecology, or a
4 vet path laboratory at a veterinary school, or
5 other food safety groups that are studying what's
6 in the -- what bacteria are in the tomato field or
7 the lettuce field. And we're looking for industry
8 involvement. You know, the idea is food industry
9 could learn more from this kind of database, the
10 same way that the hospital acquired infections
11 folks are saying, how is the pathogen moving
12 through the hospital? Food safety and industry
13 should want to know the same thing. How do we get
14 this bacteria out of our equipment and what's
15 imbedded reoccurring? And then the last way that
16 we think industry will use this was, it's really a
17 litigious issue. They don't want to get stuck
18 with the hot potato and be blamed for the recall.
19 They'll look upstream and say, what products came
20 in that brought contamination into their product.

21 MR. PETTENGILL: So our next speaker is
22 Bill Klimke from the NCBI and he'll be talking

1 about lessons learned from a real time Listeria
2 project.

3 MR. KLIMKE: So this is Marc's sized.
4 Can we get this more my size? It doesn't grow?
5 It doesn't --

6 SPEAKER: You can sit down if you like.

7 MR. KLIMKE: We have ones that you know,
8 push buttons and they move up and down. This is a
9 new building. They should, they should all be
10 robotocized already. So the good news is Marc
11 gave my talk already, so I'll be happy to take any
12 questions.

13 SPEAKER: Sorry about that.

14 MR. KLIMKE: He's stealing my trees. So
15 Marc gave an actually a pretty excellent
16 introduction to our involvement with this. So I'm
17 going to have to, I'm going to be able to skip
18 past a lot of this. Just to note, we are really
19 early on the technology curve. I think we're in
20 the early adopter stage. And I know a lot of
21 people say that they want standards, but I think
22 standards is a good thing to discuss but if you

1 expect standards to drop out of the sky into your
2 lap, it's going to take a lot of work. So I'll
3 come to some of the issues later on.

4 So when FDA approached us, about this
5 project, they mentioned that they have to do a lot
6 of food safety. There's an enormous amount of
7 food and the obvious solution to the food safety
8 issue that we all stop eating is somehow
9 politically untenable. So they're going to do a
10 lot of food inspections and there's increasing
11 amount of food being imported into the country.
12 And as Marc pointed out, the basic idea of
13 subtyping these molecular methods is, can you do a
14 match between what is making a person sick and
15 what is in the food source, whether it's the
16 spinach that you bought at the supermarket
17 yesterday or is it something that you bought or
18 ate in the cafeteria this morning.

19 And current methods don't have the
20 resolution that is necessary for all isolates.
21 PFG works pretty well. It's not completely
22 useless but it doesn't work in all cases. And

1 that's probably best distinguished in the sort of
2 schematic of a tree where you see a postal type
3 pattern that's sort of occurring in different
4 parts of the tree. PFG pattern you would not be
5 able to distinguish these isolates, but using
6 whole gene sequencing, you can.

7 You saw this model where NCBI is playing
8 a very central role in accepting data from all
9 across the world and then the analysis can be done
10 basically anywhere. This slide as well, and I
11 think that you saw this slide but for the first
12 part of the initial kick off, the first pilot year
13 of this project, data has been flowing from the
14 network into FDA and then into NCBI. FDA really
15 wants to get away from this model where they are
16 sort of brokering the data. They want to have all
17 this distributed network where all the partners
18 can submit to us.

19 All right, after FDA approached us, CBC
20 also discussed starting a pilot project for
21 Listeria whole time, or real time sequencing using
22 whole genes. Marc was telling you about some of

1 the success of that. Just to give you an idea of
2 why Listeria was chosen -- it's got a low
3 incidence in the population but it has high
4 mortality and morbidity. As you were mentioning,
5 the people who are old. There's a small genome
6 size that has excellent epi to do the action
7 between the sequencing surveys with people whose
8 current typing methods are problematic, though
9 that's a problem with this midget desk.

10 (laughter)

11 So you saw this template for pathogen
12 bio sample information, sort of the connection of
13 what is a sample, when the sample is collected,
14 where and who collected it? This is established
15 template in the bio sample database at NCBI. The
16 template exists. The current discussions between
17 the different agencies are that contextual level
18 of detail that can go into each one of these
19 fields and that's an ongoing discussion between
20 our different agencies, but this sort of minimum
21 metadata of what to collect is actually I think
22 well established for tens of thousands of samples

1 now.

2 The (inaudible) has this model that is
3 sort of an analog to what we're trying to do here,
4 is you have this weather. Weather is global. It
5 affects everyone and it's not responsible for
6 borders. It crosses and affects everyone, and
7 there's all these instruments out there collecting
8 information. And you can integrate that
9 information to predictions. For example, we have
10 Hurricane Sandy and its early predictions of where
11 it will cross and what state it will impact. In
12 the early hurricane season you don't really know
13 where it's going to make landfall and so you
14 wouldn't -- you can make predictions but you
15 wouldn't actually evacuate Boston for example.
16 But as you get the better resolution, you can see
17 all these models in real time, both the raw data
18 and the predictions. Then you get better
19 resolution. And doing this publicly is important
20 because scientific methods have not established
21 patterns for doing this. There's not an easy
22 button you can push and just say well which city

1 should we evacuate?

2 So how do we get involved in this or
3 what is our role? So Marc gave you some examples
4 of some outputs, but this is sort of our pipeline
5 right now. The short reads arrive through the
6 network and get deposited in SRA. And we pick up
7 those short reads and we do a little bit of
8 modification, trimming, and using camers we
9 determined the nearest reference in the reference
10 tree and this is for all bacteria, which is 26 to
11 27,000 genomes right now. So using camers
12 (phonetic, same as earlier) we identified the
13 nearest neighbor and in parallel, I have to
14 explain this in parallel. I'm going to say it
15 again -- in parallel, we do both reference
16 assisted assembly and a panel of the novo
17 assemblers. The reason to do this is the
18 reference assisted assembly is very quick, very
19 rapid and we can usually, for the close neighbors
20 or clonal outbreaks we cover almost the entire
21 genome sequence. And then the novo assemblers are
22 mostly there to pick up the mobile elements --

1 plasmas and phage. This is different what most
2 other people do. Most other people either do read
3 mapping to a reference that's been pre-selected or
4 they just do a de-novo assembly. And so we think
5 our method is better and we've done extensive
6 testing on this. We combine that final into a
7 final assembly, do read mapping to measure the
8 quality of the assembly and from there we can now
9 re-cluster that, or we can on to do snip analysis.
10 We can do annotation and deposit in a jim bank.
11 In fact, we've actually finished the engineering
12 for that step and so, the first Listeria genome
13 from the project has now been deposited jim
14 (phonetic, same as previous) bank and then very
15 quickly we'll fill in the gap for the other almost
16 10,000 samples that have been submitted. So those
17 will be in jim bank (phonetic, same as previous)
18 shortly as well, besides just being in the short
19 read archive.

20 Then Marc gave you this comparison
21 between PFG and post net (phonetic 0:32.48.0),
22 then we actually, comparison between whole genome

1 sequencing and PFG pattern typing. And we
2 actually did a real time experiment last year with
3 New York State where they were looking at whether
4 they could do the PFG pattern or the whole genome
5 sequencing typing in the same amount of time.
6 Effectively, it was the same amount of time. The
7 only thing I want to point out on this slide, a
8 pointer -- people online can't see this, but you
9 see this stack of very short parallel steps?
10 That's basically NCBI's involvement. Everything
11 upstream is all the laboratory response, the
12 collecting the sample, culturing, laboratory prep,
13 the sequencing and then submission to NCBI. So
14 the work we can turn around within approximately
15 on average 4 to 5 hours the result, which we
16 haven't even optimized yet, so we think we can
17 turn this around even more quickly in the future.

18 What have we done so far? So as Marc
19 was saying, we have -- this is based on the
20 beginning of September, so we had almost 9,000
21 samples in total. We have Listeria from CDC and
22 the FDA, and even I should point out, the Food and

1 Environmental Research Agency in the U.K., so this
2 is starting to expand into a global international
3 network. And I was just at the Global Microbial
4 Identifier Meeting in York, U.K. which is sort of
5 an ongoing attempt to build a global system like
6 this, although there's still far too much talking
7 at that meeting and not enough doing in my
8 opinion. The U.S. is a forerunner and a leader in
9 this field. And I as a Canadian feel that we
10 should expand this to a more global system. But
11 we have commitments from the public health agency
12 and England to submit all their Salmonella to
13 (inaudible) Laboratory, so this really is starting
14 to become a global system.

15 We have all the Salmonella from the FDA.
16 The first year of the project was a great deal of
17 retrospect and samples being submitted but we're
18 starting to get commitments from states to
19 sequence their clinicals to start to make those
20 connections as well. We have some e-Coli as well.

21 This is just sort of the number of
22 samples that have been submitted per day over the

1 last year and you can see the spikiness of the
2 samples. And there are a couple of downtimes that
3 are partly due to our problems, so sometimes our
4 computer network has had issues. We had a
5 thunderstorm that flooded our computer room, so
6 that went down for a while. But those are things
7 that we're going to have to be aware of in a
8 global network system. And we also get spikes
9 when I think people have gone on holidays and then
10 come back and then have to submit 100 samples on a
11 Friday.

12 But there are a number of issues that I
13 see. People have presented two types of talks at
14 this conference. One is the optimistic version
15 and the more pessimistic version. Everybody who
16 knows me knows I'm a pessimist so I'm going to
17 present some of the issues that have arisen. As
18 can be expected, when you're doing real time
19 sequencing, you can run into some problems, and
20 some of these are as Marc mentioned, human errors.
21 Some of them are cases of the wrong organism being
22 submitted and so somebody isolates what they think

1 is Salmonella but it turns out to be (inaudible).
2 We've seen mixed organism cases. We have some
3 data issues with how much data we think we can
4 trust, so less than 5x coverage, probably not
5 enough, duplicates, those are just human errors,
6 making mistakes. (inaudible) is only on here, I
7 think Jonas has already left but it's mostly
8 because we haven't built the system to handle pac
9 bio assemblies yet because we were assuming, if
10 you're sequencing pac bio you're going to do the
11 assembly yourself.

12 I'm going to go through a couple of
13 these errors in detail. So incorrect organism --
14 this can be both in what is (inaudible) right now.
15 Everyone knows that this is probably a problem.
16 We have, as I said, almost 27,000 assemblies and
17 some number are obviously going to be wrong.
18 We're working on trying to fix that, so here's an
19 example of a Campylobacter jejuni, which clusters
20 with Campylobacter e- Coli and we informed the
21 submitter of this that they probably had the wrong
22 organism name and they agreed with us and they

1 made a fix. And then here's an example of
2 Enterococcus Faecalis which clusters with
3 Enterococcus Faecium and this was supposed to be
4 an HNP reference genome but apparently it was not
5 a reference genome. We informed the Brode
6 Institute that this was a problem and they are
7 going to correct this very soon. It's amazing
8 what public outing of your mistakes does to your
9 cleaning up systems properly.

10 We have ways of labeling this data as
11 problematic so this is a couple of examples from
12 our assembly database. They are listed as
13 anomalous assembly. Two of these were
14 misassembled and this last one is actually, was
15 submitted as Bacillus Subtilis but it's actually,
16 someone has cloned a Cyanobacteria genome into
17 Bacillus Subtilis so this actually links the
18 publication describing that process. It's not
19 that the experiment was done incorrectly but the
20 fact that this is Bacillus sort of, you want to
21 give warning that there's something else going on
22 here.

1 And then a quick tangent on, we have
2 ways of labeling reference material, that you've
3 heard lots of talks about reference material at
4 this meeting. In the (inaudible) database we have
5 this field called reference material, which sort
6 of, we can give you a description of what that
7 reference material is. You probably can't read
8 that, but it's like a post snip standard reference
9 material used to do pattern typing. And we also
10 have links, for example to ETTC culture collection
11 where if it was a culture collection, if you got
12 it from DSMZ or ATCC, you could actually provide
13 that information as well.

14 So what else have we seen as errors?
15 We've seen mixed samples, so if you ignore the
16 colors, if you just look on the X axis, you see
17 the assembly size. So this is a Listeria and you
18 see the nice clustering at around three mega bases
19 for Listeria, and then on the left is, on the Y
20 axis is the N50 (inaudible) length and so in the
21 upper left you see a lot of high quality
22 assemblies, single contig or just a few contigs

1 (phonetic, same as previous). The bottom left
2 there's some bad assemblies. But the ones that
3 are problematic are everything on the right of
4 this X axis. And what is going here? We have
5 some that are almost double the size you'd expect
6 of Listeria. Well, here's an example of a
7 contaminated sample with both Listeria and
8 Staphylococcus. So we have a fairly decent
9 assembly for the Listeria but we have a lot of
10 Staphylococci for a lot of fragmented contigs.
11 So this is obviously a case where this sample
12 cannot be used for trace back because we don't
13 know what happened to this sample.

14 Then we have a problem with the
15 platform. So Illumina, most of the data that
16 we're getting from this network is Illumina. And
17 we have sometimes the case in Illumina where we
18 have carry over contamination. So what we have on
19 this slide is four samples. Histogram of the
20 coverage and the number of vases, and the sample
21 on the upper left, the upper right and the bottom
22 right are decent high quality assemblies. The one

1 on the bottom left, I don't know if you can see
2 this, but there is a little spike of low coverage
3 bases in this assembly. And so this looks like
4 this is carry over from one run to the next run in
5 your sequencing run. And this is a case where
6 NCBI is going to make a recommendation. It's not
7 a standard, but it's sort of a recommendation to
8 our partners to explore the use of combinatorial
9 barcoding to sort of separate this issue out. So
10 this is of the case where you would, for example,
11 sample one, you would use barcode A and B, and
12 sample 2, you would use barcode C and D. And so
13 next week when you're doing a sample 127, if you
14 see barcodes A and B, you would recognize it was
15 from a previous run. You could throw that data
16 away. So they're exploring this and this is where
17 the collaborative nature of our network helps.

18 And then the more insidious case of a
19 mixed sample is where it's the same species. So
20 here we have a fairly decent run. It assembles
21 very well, three point almost two mega bases, very
22 few contigs (phonetic, same as previous) but when

1 assembling map the reads back to the assembly to
2 measure the quality, we recognized a whole bunch
3 of alternate alleles. And you can recognize this
4 when you see; this is in Genome Workbench or tool
5 where you can sort of load a bound file. At the
6 bottom you see this coverage graph and then all
7 these little red lines are errors. If you zoom in
8 on that, what you see is for -- you see a set of
9 variants, a core occurring in a set of reads, and
10 these variants complement each other very well.
11 So that's unlike what you would see with a normal
12 assembly where you have -- you always see variants
13 or errors in that actual read, but they're not
14 suggesting that the -- what we think is happening
15 here is that we have a mixed population of
16 *Listeria* in the sample.

17 So you could imagine, this is easy to
18 see, right? There's a lot of nice red lines and
19 they all look like they're in the same position
20 and they have the same base. But imagine you were
21 in a laboratory and you were seeing an outbreak
22 and you had two strains that were very very very

1 similar to each other. Could we possible tease
2 apart that this one is different by 12 snips and
3 this one by 15 snips, across three mega bases?
4 They may be very problematic. So this is
5 something to be aware of when doing this type of
6 trace back investigation using the whole genome
7 sequencing.

8 All right, so the idea is that whole
9 genome sequencing would give us this resolution
10 and as Marc said, we really want to focus on the
11 tips of the tree -- what is the epidemiological
12 relevance. And we want to distinguish independent
13 events, so we want to separate things that are
14 assembly errors or repeats. And you can see this
15 in a (inaudible) comparison of two assemblies.
16 You see the variances of these sort of accumulate
17 along the graph, or sample on the X axis and
18 assemble on the Y axis. Get a uniform
19 distribution and these two samples; they're so
20 divergent that they wouldn't even be considered
21 part of the same outbreak. But if you don't
22 properly account for repeats, then you'll get a

1 spike like this in the graph, where in a single
2 part of the genome you see a huge increase in the
3 number of variants. And you want to properly
4 filter those out, so you're only getting the
5 independent events as a measurement. And we're
6 working on using an algorithm developed at the
7 Sanger Pathogen Center and we're looking at ways
8 of how we can rapidly detect these. Because you
9 really want to get down to this, as Mark said, the
10 tips of the tree, and rapidly report on these
11 cases that we think are sort of clonal or near
12 clonal (phonetic, same as previous)
13 indistinguishable. So we're working on outputs
14 for both -- right now we're doing a (inaudible)
15 tree of all the Salmonella. That's what Marc was
16 saying we produce routinely daily, depending on
17 whether we get samples or not, and we're working
18 on rapid ways to do a snip tree. We've looked at
19 both parsimony methods and rapid distance methods.

20 And so the idea is sort of to report to
21 the people, something like continuing along this
22 weather theme. You know, you can look on the

1 weather app on your iPhone or whatever phone you
2 have and you can see that there is either a
3 thunderstorm happening or there's a tornado watch,
4 which means there's something to be aware of, but
5 you know, don't be worried, or a tornado warning,
6 which means, get underground as soon as possible.
7 So this is a mockup of sort of like, you know,
8 what's happening in Salmonella today in the U.S.?
9 You sort of come and you can see how many samples
10 have been updated and then get like a rapid report
11 of a new clinical to an environmental link and you
12 can look at the tree, just within this small part
13 and see that, you know, there's a clonal link
14 between two isolates for example.

15 So I'd just to acknowledge all the
16 people in CBI -- it's a vast army of people and we
17 know there's a vast army of people at both CDC,
18 FDA, UC Davis, well okay, we'll ignore that for
19 now, USDA, state labs, and various partners and
20 collaborators that we worked on this project over
21 the years. And I'd be happy to take any
22 questions.

1 (applause)

2 SPEAKER: More a comment than a
3 question, food processing has been very successful
4 in increasing the, or actually in decreasing the
5 amount of bacterial contaminants in food, whereas
6 it's been much less successful for decreasing
7 bacterial -- sorry, viral contaminations, so if
8 you have the enteric viruses that are breaking
9 through and not decreasing in terms of food borne
10 pathogens. So do you see an increasing trend
11 towards using your NGS or setting up libraries for
12 viral based pathogens as opposed to bacterial
13 based?

14 MR. KLIMKE: Probably a question to ask
15 FDA and USDA, not us.

16 SPEAKER: No, we're planning to do
17 norovirus and those are what's included. That's
18 the plan.

19 SPEAKER: Any other questions? So
20 briefly from the online audience -- you touched
21 barcodes, but the question was, do you see barcode
22 bias specifically with Illumina and if so, how do

1 you deal with it?

2 MR. KLIMKE: Do we see barcode bias?

3 SPEAKER: That was the question, yeah.

4 MR. KLIMKE: We haven't thoroughly
5 investigated that based on the submissions that we
6 get, because we don't always -- we're not always
7 informed at submission time what barcodes have
8 been used. So we would have to actually probably
9 contact people and collaborate with them to
10 determine that.

11 SPEAKER: Any other questions?

12 MR. PETTENGILL: Okay, our third and
13 final speaker, is Kristin Holt from the Food
14 Safety and Inspection Service, which is part of
15 the USDA, and she'll probably actually not
16 overlap, maybe much with any of the previous
17 speakers. And she'll be talking about the use of
18 next generation sequencing data to enhance food
19 safety and the need for data format
20 standardization.

21 SPEAKER: All right, just a quick
22 interruption, after this talk, we will have quick

1 closing remarks, before we end the session. Thank
2 you.

3 MS. HOLT: Okay, hopefully I've got the
4 mikes correct. Great, okay. Well great. I
5 appreciate this opportunity to continue the
6 discussion of the use of next generation
7 sequencing and food safety. I know that during
8 the workshop, a few people have asked, kind of
9 polled the audience for things, and I have a
10 question for the audience, so, how many people
11 have a relationship with food or eat food? Raise
12 your hand. (laughter) Okay, well I actually had
13 a kind of a hidden motive in asking you that
14 question, because I'm the last presenter today.
15 And I just wanted to kind of get some oxygen and
16 endorphins going, so, because, and I'm going to
17 move, try to move quickly there, and I do have
18 some repeat content I think. Okay.

19 So the Food Safety and Inspection
20 Services, the public health agency in the U.S.
21 Department of Agriculture is responsible for
22 ensuring that the nation's commercial supply of

1 meat, poultry and egg products is safe, wholesome
2 and correctly labeled and packaged. So I
3 underline meat, poultry and egg products because
4 that's what we're about and FDA is about all the
5 other foods. And since we're talking about next
6 generation sequencing, I think it's obvious for me
7 to give you an overview of the FSIS office of
8 public health science laboratory services team.
9 So we have three field service laboratories that
10 are responsible for microbiological and chemical
11 analyses of foods and environmental swabs. Again,
12 this relates to meat, poultry and egg products.
13 And we have a strong focus on the control of data
14 integrity and standards. Our labs are located in
15 Alameda, St. Louis and Athens, Georgia, and Athens
16 we also have some specialty units, pathology
17 branch and an outbreak section of the eastern
18 laboratory or OSSL. Also in Athens, we have our
19 executive associate for lab services, our lab
20 quality assurance staff and I'm going to call out
21 Cathy Pence whose here, to show how we're into
22 quality. We wanted to be here at this meeting and

1 find out about data standards. And we also have
2 our food emergency response network staff.

3 The mission of the FSIS laboratories is
4 to perform analyses on meat, poultry and egg
5 products. We have a very robust food sampling
6 program, at least for food. This is very very
7 robust. We get over 100,000 samples and perform
8 over one million analyses a year, so our labs are
9 very busy. Again, a strong emphasis on quality.
10 We have a quality system and this graphic shows,
11 you know, the types of documents that we have at
12 all of our laboratories. We have quality manual,
13 SOPs et cetera. And our labs are accredited to
14 the ISO 1702.5 standard.

15 Now I'm going to drill down a little bit
16 and talk to you about bacterial isolate
17 characterizations that we do at our eastern
18 laboratory which is located within the Richard
19 Russell Research Center in Athens, and there's a
20 picture taken on a beautiful spring day. So what
21 is the current role of FSIS in bacterial isolate
22 characterization? So we do pulse field gel

1 electrophoresis, or PFGE, we've done that since
2 the beginning of PulseNet, we were an original
3 member in that network. We upload quite a bit of
4 patterns and meta data to go with it. I don't
5 know why I have just numbers back to 2006, but
6 we've got data going all the way back to the mid
7 1990's. We do MLVA; molecular serotyping of the
8 Salmonella isolates using CDC's new molecular
9 serotyping method. We do anti-microbial
10 susceptibility testing and we have storage of lots
11 of isolates.

12 Now this is really busy, so I just want
13 to orient you to a couple things and make really a
14 point about the work load that the folks have
15 there in Athens, in characterizing the bacterial
16 isolates. So the top left graphic is e-Coli 0157
17 and the top right is Listeria. And those programs
18 have been pretty stable -- not a lot of changes
19 there. But if you look at the Salmonella in the
20 middle, which is on a different order of
21 magnitude, you see kind of a gap in Salmonella in
22 the early years and you say well, they do poultry,

1 surely they've got a lot of poultry isolates.
2 Well actually, our sister agency, the USDA or as
3 was doing the characterizations and we took those
4 over recently, so then we had the big jump in the
5 numbers. And then the bottom row is for non 0157
6 (inaudible) and enterococcus and those big jumps
7 are new programs that we actually have recently
8 implemented. So to give you kind of a sense of
9 the jump, in numbers we had 492 isolates
10 characterized by this crew of people in 2010. And
11 just this first six months of this year, they did
12 over 3,000. So this is a very busy group.

13 So why the increase in isolates? We
14 just have new regulatory testing programs, we are
15 getting isolates from baseline studies that we're
16 going to characterize all those isolates and
17 sometimes that can be a lot of isolates and we
18 picked up the responsibilities that ARS had
19 before. We're doing NARM sequel sample testing
20 and we do special projects and outbreak
21 investigations in the recent Salmonella Heidelberg
22 investigation associated with chicken. We had a

1 lot of isolates that we did characterizations of.

2 So why is FSIS becoming involved in
3 whole genome sequencing? Well it fits well with
4 our FSIS strategic goals, our strategic plan and
5 this is just a sampling of some that seem to
6 relate, but the one at the bottom is, you know, we
7 have goals to use innovative methodologies,
8 processes and tools and this of course is a
9 perfect fit and we have high level support within
10 our agency to move forward.

11 The benefits of whole genome sequencing,
12 now this is a really busy slide, and I actually
13 curtailed some information. I think I could have
14 given you four, five slides with tiny font. But
15 let me walk you through a couple things here. So
16 the first thing is really, I think, visionary.
17 FSIS anticipates that whole genome sequencing will
18 eventually replace PFGE, serotyping and other
19 testing and we will quickly be able to quickly
20 tell us about a strain's virulence and resistance
21 potential, it will streamline our lab methods.
22 And in the clinical world, we're looking forward

1 to the fact that we might get resistance and
2 virulence data and information really really
3 quickly. And of course, it's a powerful and
4 discriminatory tool, compared to PFGE. It's fast
5 and cheap. We're going to be more efficient with
6 our resources. It's going to improve
7 surveillance. You know, for epi investigations
8 that we think we're going to continue to see that
9 we have to do trace back and epi in
10 investigations, not just rely on matches from
11 whole genome sequencing between a clinical isolate
12 or food. We got to bring that extra investigation
13 data to the consideration, because we wouldn't use
14 lab evidence alone to take regulatory action we
15 have to, in a food borne illness investigation.
16 It's going to improve the classification of ill
17 patients, meaning that it's more likely that a
18 case is a case and we're doing the trace back on
19 the right person and the right food.

20 So how have we been involved? We've
21 been involved in the global microbial identifier
22 that was mentioned earlier. We've been

1 contributing isolates for sequencing to the
2 Hundred K Genome Project. We are part of that
3 Listeria real time surveillance project that you
4 heard earlier. We're not actually doing the
5 sequencing yet ourselves, but we're close. So to
6 participate in the project we have contributed a
7 168 isolates, sent them to FDA and they've kindly
8 sequenced them and uploading the data to NCBI.
9 And we've also sent Salmonella isolates over.

10 We've received two sequencers in April
11 2014. We have two aluminum (inaudible), so that
12 was very exciting for us. We may -- taken our
13 scientists and done (inaudible) to CDC, NIH and
14 FDA and they graciously opened their doors and
15 showed us everything, answered all our questions
16 and continue to collaborate and help us move our
17 whole genome sequencing activities along. We have
18 had successful sequencing runs with Listeria,
19 Salmonella and STEC and there was perfect timing
20 when we got our sequencing, we had in our
21 experience, and FDA had a proficiency test that we
22 participated in. So that was really great.

1 We're developing our knowledge and
2 application within the agency. We're all very
3 used to PFGE terms. So we're trying to bring
4 people along to hold genome sequencing terms. And
5 what's missing here is that we're not yet
6 uploading directly to NCBI yet, but we hope that
7 will come very soon.

8 And then I won't spend much time on this
9 because this was already mentioned. We're
10 participating in this project, and calls and
11 conferences, et cetera where this activity is
12 discussed.

13 And where do we get our *Listeria*
14 isolates? We have a regulatory testing program of
15 ready to eat meat and poultry products. We do
16 environmental swabs with these manufacturing
17 facilities, and if we had a food borne illness
18 investigation, we collected samples. We may have
19 some isolates from there. We are planning to do
20 the real time sequencing at our lab in fiscal year
21 2015 and we think we'll have about a capacity to
22 do about 10 to 20 isolates every two weeks. And

1 we do have banked isolates, so we will be pulling
2 some of those out and doing those.

3 Just some extra considerations for shiga
4 toxin, producing e-Coli 0157 and now we're doing
5 testing for the 6 non 0157 serum groups. We're
6 going to do the real time sequencing there. And
7 if you're -- the graphic there is I guess if
8 you're curious, what sero groups we're seeing are.
9 We're seeing a lot of 0103. And we have banked
10 isolates. Another consideration is for Salmonella
11 and Campylobacter we have a lot of isolates and I
12 think this is the same challenge on the clinical
13 side. So we're going to definitely do that -- you
14 know, whole genome sequencing for food borne
15 illness investigation samples. But we're waiting
16 to kind of see how CDC, FDA, NIH and FSIS come
17 together on a strategy to take on Salmonella and
18 Campylobacter because it's a little different than
19 Listeria. There's a lot more isolates.

20 So what are the data considerations for
21 us? They're basically the same as for everybody.
22 Data storage and transmission, of course, for us,

1 we have the quality standards that we have to live
2 up to. And then also, the epidemiological
3 understanding of how we're using the new data --
4 that's a learning curve for us, but I think it's a
5 learning curve for everybody. But we think it's
6 really a great tool. Successes -- you know,
7 commercial software is there, our IT
8 infrastructure's going, getting better and the
9 agency is learning more about sequencing. And I
10 guess everybody kind of had a graphic like this.
11 This is just basically; we've got our flow mapped
12 out.

13 Meta data -- this has been pretty
14 straight forward for us. We're going to provide
15 and we are providing with the isolates we're
16 sending the FDA, the product type with you know,
17 what kind of food item it is, the year the sample
18 was collected, the state where it was collected,
19 any subtyping information we have and the meta
20 data is available as soon as we pass the isolate
21 or we're going to upload. We're not going to have
22 any lag on there. It will go with the sample or

1 the raw reads.

2 And the assessment of sequence quality
3 -- this is what we have. And I guess one of the
4 draws for me to come to the meeting was also to
5 learn, you know, what is the scientific and public
6 health community using for quality standards?
7 Because we're regulatory agency and for us to take
8 an action, we really need to be doing what
9 everybody else is doing, that there's consensus.
10 And I guess we're learning that I guess there's
11 not some hard numbers and hard agreement and we're
12 really just kind of getting there and I appreciate
13 FDA hosting this meeting.

14 So for the scientific and public health
15 approach, you know, we know there are a lot of
16 benefit from whole genome sequencing. For us, we
17 do recognize that the Information Quality Act
18 pertains to us and other agencies. And USDA
19 created guidelines on the information quality. We
20 of course, have to have that strong data quality
21 data standardization, so that we can use that
22 information and defend and take an action.

1 So I'm going to end with just kind of a
2 little bit of a story with the focus of how the
3 sequence data from just one person's isolate could
4 be very very important to crack open a food borne
5 illness investigation. SO CDC estimates that each
6 year roughly one in six Americans get sick,
7 128,000 are hospitalized and 3,000 die of food
8 borne disease. And I'm going to talk now about
9 Salmonella causing about 1.2 million illnesses in
10 a year. Most of the time you're thinking about
11 poultry products, but when you have an
12 investigation that starts, you don't really know
13 it's really going to be the poultry. There have
14 been Salmonella outbreaks connected to vine
15 vegetables like tomatoes, fruits, nuts and nut
16 butters and of course, sometimes it's not in food
17 at all. It's backyard poultry, or somebody's pet
18 reptile or amphibian.

19 And there are many benefits we see with
20 whole genome sequencing too -- the epidemiologic
21 and the trace back investigations. Food borne
22 illness investigations are multifaceted, they're

1 multi-disciplinary and they're collaborative
2 undertakings. We believe in the three legged
3 stool of investigations, with the legs being
4 laboratory, epidemiology and environmental health,
5 and I highlighted in blue, you know I guess the
6 clear connects with whole genome sequencing. For
7 a laboratory, I think it's pretty obvious, but for
8 epidemiology, we think we're going to have better
9 cluster protection, better case classification, so
10 we're going to rule out some of those people who
11 probably just kind of sporadic illnesses that are
12 not part of that outbreak. That's going to be
13 really powerful for us. And for trace back,
14 again, we interview people and we're going to be
15 tracing back the right food.

16 So to walk through, I'll just kind of
17 walk through quickly here. So you know, if you
18 eat a contaminated food item, for Salmonella it
19 takes a couple of days for you to exhibit the
20 symptoms and be ill, and then you decide to go to
21 the doctor. The doctor asks for a stool specimen.
22 That stool specimen then gets tested at the

1 clinical lab, Salmonella is identified and then a
2 clinical lab will send the isolate to the public
3 health lab for further characterization such as
4 PFG. And then the case is confirmed as part of
5 the outbreak. So it can actually take maybe two
6 to four weeks after you ate that contaminated food
7 item to being confirmed as part of the outbreak.
8 So the significance of that is, when the health
9 department calls you up, they're going to ask you
10 what you ate, maybe four weeks ago. And they're
11 going to ask you the seven days, or three to seven
12 days before you got sick. So what were you doing
13 five weeks ago and what did you eat and where did
14 you go? It's pretty tough, so that's one of the
15 things I wanted to just keep that in your mind, as
16 we move through the rest of my slides. And the
17 other thing is again, I think whole genome
18 sequencing can kind of cut some of this timeline
19 down and give us some information really quickly.

20 So now we're using PFGE and I'm just
21 going to use that as an example because you know,
22 this was a new technology 20 years ago. It was

1 innovative. It got awards. But you know, it's
2 been used for 20 years. And it's strong, it's
3 powerful and maybe it's not as discriminatory but
4 as much as whole genome sequencing, but when that
5 isolate gets to the health department, we've got
6 it, the labs do standard methods, you can roll the
7 data all up together. It's not a perfect tool,
8 but it is a great tool and it's worked really
9 well. And we're probably going to continue to use
10 it for a few more years now. So basically, so
11 that person who got sick, and had that isolate, go
12 to the health department, this will be compared
13 either at the state, you know, looking at the
14 central database or at the CDC and roll up the
15 national data.

16 So this is kind of a hypothetical
17 investigation that I put together. It's kind of
18 got hallmarks of some real investigations. So you
19 imagine there's a Salmonella Heidelberg with a PFG
20 pattern A, it's a sub-cluster because it's maybe a
21 common clone of Heidelberg. But we get down to
22 those people who four or five weeks later, you

1 interview them and you have a good food history.
2 For some reason, they really remember what they
3 did, where they ate, what they shopped, where they
4 shopped etcetera. And so we might have a big
5 outbreak, but we're only working with a handful of
6 people who actually have what we call good food
7 histories or good exposure history. So every
8 person here is really important and we don't want
9 to lose that data because somebody said it doesn't
10 meet a quality mark. So for this scenario, we've
11 got to interview these four patients -- oh, sorry
12 -- so by interviewing them, there's some kind of a
13 poultry item we're worried about -- chicken or
14 something else. And we're finding from the trace
15 back going to the grocery stores, the shopper card
16 loyalty data, that it looks like these are big
17 chicken eaters. And believe me, this is how this
18 really works. People eat a lot of chicken in the
19 United States. And they eat it from all kinds of
20 different places.

21 So we get okay, establishment A is
22 showing up on three out of four, but yet

1 establishment B, that their chicken's showing up
2 on three out of four and establishment C. So I'm
3 really excited to see this fifth case come on from
4 Minnesota to kind of triangulate and move it out
5 of the area. What did they eat? And they had a
6 good history, establishment A. So establishment A
7 essentially kind of lights up. And we are then
8 like okay, we have a lot of confidence here. And
9 then, for case patient two, say the shopper
10 loyalty card information was slow to come in. It
11 comes in and voila, there establishment A lights
12 up, and we really have it firm.

13 But I guess this is just to illustrate
14 that if we lose the data quality or we don't have
15 standards, we're not all together, we can lose one
16 person and it can really break up our ability to
17 then take some kind of action like a product
18 recall, do a public health alert, go to that
19 facility and check them out, which we probably
20 would be doing anyway, regardless of that we'd be
21 checking out several of these.

22 So anyway, just want to kind of walk you

1 through the, this is kind of where the rubber
2 meets the road, and we need good data quality
3 standards with things all coming together. But
4 thank you.

5 (Applause)

6 MS. HOLT: We're not officially part of
7 the Genome Tracker network. But FDA has been
8 really just incredibly gracious. Send us your
9 isolates, send us your isolates, send us your
10 isolates. So we have sent lots of isolates and
11 they have done the sequencing and I mentioned
12 before, we visited them and they walked us through
13 all their protocols et cetera. So we appreciated
14 FDA's work.

15 SPEAKER: We appreciate your
16 (inaudible).

17 MS. BOLT: All right great.

18 SPEAKER: I have a quick question. So
19 with the increase in sort of Next Generation
20 sequencing being used and given that you pointed
21 it out, sort of the importance of you know, good
22 exposure, history and sort of the epi and sort of

1 that linkage information. Do you know if there's
2 a sort of concurrently novel method being
3 developed, to sort of better you know track people
4 that are sick and get information from them, that
5 sort of would complement the NGS input?

6 MS. HOLT: So really that's probably a
7 question for CDC and the state health departments.
8 So I'll just say though, yes, enhancements to
9 getting the questionnaires standardized, getting
10 questionnaires collected, you know, people
11 interviewed quicker. There's a lot of, I mean,
12 there would be a whole 'nother presentation on
13 that, so yeah. On the epi side, there's a lot of
14 activities and there's a lot of close, close work
15 between the epi's and the lab folks with this
16 Listeria whole genome sequencing, kind of learning
17 together as we go.

18 SPEAKER: You know, in that example that
19 you gave, there's also the possibility that all of
20 those A, B, C and D, received chicken from the
21 same points or it's going backwards one level.
22 And so what -- do you immediately go out and

1 investigate them anyway, even if you found that
2 commonality of A? Because it could have been one
3 layer backwards?

4 MS. HOLT: Yeah, and that does happen.
5 So when we do a trace back on a person, it starts
6 largely say, like with the grocery store. Then it
7 goes to the distributor, then it goes to that next
8 point. And sometimes that's a processing plant,
9 and then it goes back another step to a slaughter
10 plant. So yeah, it carries all the way through.
11 And then triggers to actually go to a facility.
12 We can look at their data just remotely, their
13 micro profile from testing over the years. But we
14 can, you know, really, we don't have to go do a
15 full formal investigation. We can start, you
16 know, going back and tracing and looking at them
17 and talking to them even and saying, what can you
18 tell us to help? So we have a lot of latitude.
19 And then we have some of the more formal things
20 that are, you know, are kind of listed at the
21 bottom there, where we have described protocols.
22 So we have a lot of latitude.

1 MR. PETTENGILL: That concludes the
2 final session, but there's going to be a few
3 closing remarks.

4 MR. SIMONYAN: I'll take like five
5 minutes of the time that is left. We today got
6 all of the stakeholders. Oh, sorry. Today and
7 yesterday, we got representatives from all of the
8 stakeholders. We got device manufacturers, we
9 heard from them. We heard from consumers of the
10 (inaudible) information, those being hospitals,
11 research institutes, regulating organizations or
12 consultants, or like national health
13 organizations. We didn't hear actually from
14 hardware manufacturers and we didn't hear from
15 software manufacturers, from industry. And that
16 is also important. We also want their input into
17 the standardization effort and I'll make sure, I
18 mean, I know they are in our lists with just
19 perhaps, short sightedly didn't give them a chance
20 to talk or maybe the time was not enough, but I'll
21 make sure that they also have a talk when we are
22 working on standards, and when we are working on

1 bioinformatics validation as well.

2 So what our vision is, what the
3 procedure would be after this, we have developed
4 the draft documents which will be introducing what
5 we are trying to do. But after this meeting, we
6 got a lot of input from all of you, by talking to
7 you, we want to get a few days on editing this
8 draft, trying to incorporate the advice, the
9 ideas, the meaning which we have put into our
10 brains, and then, release those drafts after we
11 rotate (inaudible) maybe a day or two, we will
12 release this draft by clearly stating that this is
13 a draft. There is nothing written in stone. It
14 is a suggestion where we can start our
15 conversation, and how we can continue the
16 conversation. At that point we will send an email
17 to all of the participants in this conference and
18 all of those people who registered for continued
19 workshop participation and then we will outline
20 what is the procedure, when do we meet, when is
21 the first meeting, what is the agenda of the first
22 meeting and ask you for agenda items on it.

1 Hopefully this will work out and important thing
2 is that all of us are here for one and only goal,
3 with different perhaps means but there is only one
4 ultimate goal behind all of this, that's the
5 health of the humans. And the impact we can have
6 to health of the humans. People who use and reuse
7 this work, but it is important to remember those
8 humans are our children, our mothers, our fathers
9 and us. So by having that beautiful idea, I want
10 to close this conference by having that beautiful
11 idea, however we do and whatever we do, we are
12 doing it for humans. Thank you. And let's thank
13 all of the organizers.

14 (applause) And Carolyn Wilson,
15 Carolyn Wilson wanted to say a
16 Couple of words from FDA perspective.

17 DR. WILSON: Okay. I will be quick. I
18 know it's been a long two days. But I do want to
19 also just echo the thanks to all the speakers, all
20 the people who were here for a portion or all of
21 the two days, and offered up their suggestions,
22 their questions. I think what I came away from

1 these two days understanding is that there's a lot
2 of concerns, challenges, and there's a lot of
3 shared understanding of what those problems are in
4 the field, in order to make the potential promise
5 of this technology really give the full benefit,
6 as Vahan said to you, the public health. And I
7 think that at least now that we have a shared
8 understanding of what we need to do, and I assume
9 the next steps are really to prioritize what those
10 needs are -- what I got from it is some will be
11 easier to tackle than others. And that you know,
12 some may take longer and maybe we don't even know
13 enough yet at this point to know how to tackle
14 them. So I think those are additional things that
15 people need to think about going forward, but
16 again, thank you everyone, and we look forward to
17 continuing this dialog with the community.

18 (Whereupon, the PROCEEDINGS were
19 adjourned.)

20 * * * * *

21

22

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

CERTIFICATE OF NOTARY PUBLIC

STATE OF MARYLAND

I, Mark Mahoney, notary public in and for the State of Maryland, do hereby certify that the forgoing PROCEEDING was duly recorded and thereafter reduced to print under my direction; that the witnesses were sworn to tell the truth under penalty of perjury; that said transcript is a true record of the testimony given by witnesses; that I am neither counsel for, related to, nor employed by any of the parties to the action in which this proceeding was called; and, furthermore, that I am not a relative or employee of any attorney or counsel employed by the parties hereto, nor financially or otherwise interested in the outcome of this action.

(Signature and Seal on File)

Notary Public, in and for the State of Maryland

My Commission Expires: November 1, 2014

