

PRE- AND POST-MARKET EVALUATION OF AUTONOMOUS AI/ML: LESSONS LEARNED FROM PRIOR CAD DEVICES

Nicholas Petrick, Ph.D.

Division of Imaging, Diagnostics and Software Reliability (DIDSR)
Office of Science and Engineering Labs (OSEL)
Center for Devices and Radiological Health (CDRH)
U.S. Food and Drug Administration (FDA)
nicholas.petrick@fda.hhs.gov

OUTLINE

- CAD devices types
- CAD assessment framework
- Some lessons learned
 - and how they might apply to autonomous AI/ML

COMPUTER-AIDED DIAGNOSIS (CAD)

<p>CADe</p>	<ul style="list-style-type: none"> • <u>Computer aided detection</u>: Aids in localizing/markings regions of that may reveal specific abnormalities
<p>CADx</p>	<ul style="list-style-type: none"> • <u>Computer aided diagnosis</u>: Aids in characterizing/assessing disease, disease type, severity, stage, progression
<p>CADe/x</p>	<ul style="list-style-type: none"> • <u>Computer aided detection/diagnosis</u> : Aid in localizing and characterizing conditions
<p>CADt</p>	<ul style="list-style-type: none"> • <u>Computer aided triage</u>: Aids in prioritizing/triaging time sensitive patient detection and diagnosis
<p>CADa/o</p>	<ul style="list-style-type: none"> • <u>Computer-aided acquisition/optimization</u>: Aid in the acquisition/optimization of images/diagnostic signals

These CAD device types are general not organ/disease specific

INITIAL SCREENER

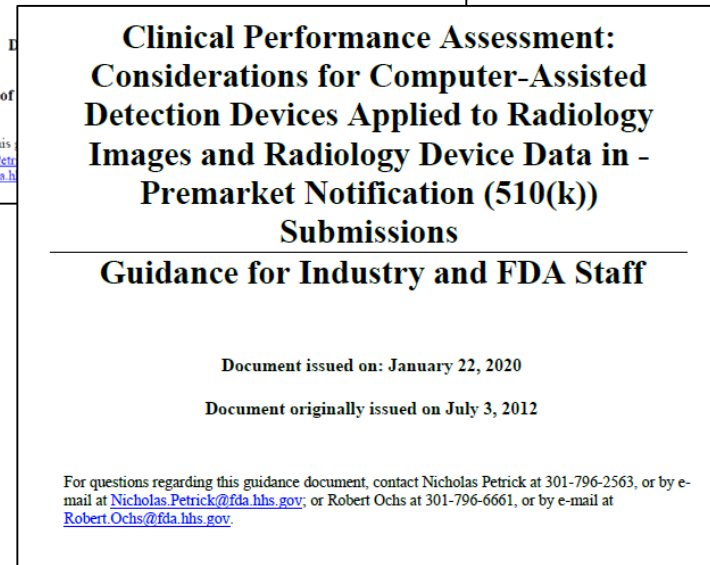
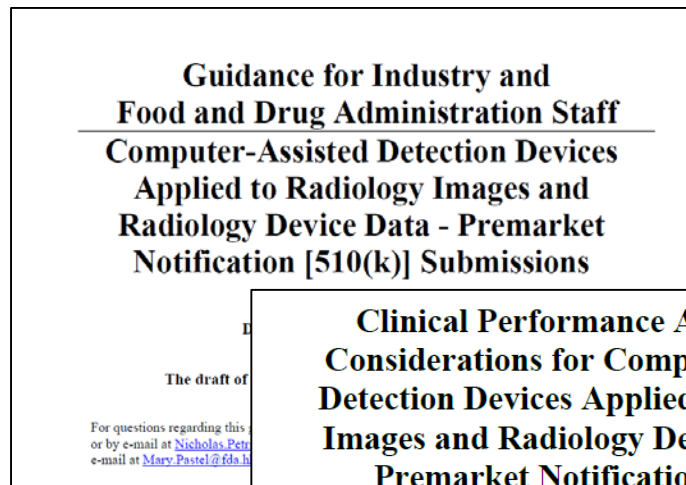
- Initial Screener
 - Automated cervical cytology slide reader
 - BD FocalPoint GS (P950002-S002)
 - Selects up to 25% of Pap smear slides that need no further review



*<https://www.bd.com/en-us/offerings/capabilities/cervical-cancer-screening/cytology-instruments/focalpoint-gs-imaging-system>

CAD DEVICE ASSESSMENT

- Two CADe guidance documents
 - Clinical assessment guidance
 - Updated 1/22/2020 to reflect CADe reclassifications





BASIC CADE ASSESSMENT FRAMEWORK

- Detailed Descriptions
- Verification and Validations (V&V)
 - Standalone performance assessment
 - Clinical performance assessment
 - ...
- Labeling
- ...

DETAILED DESCRIPTIONS

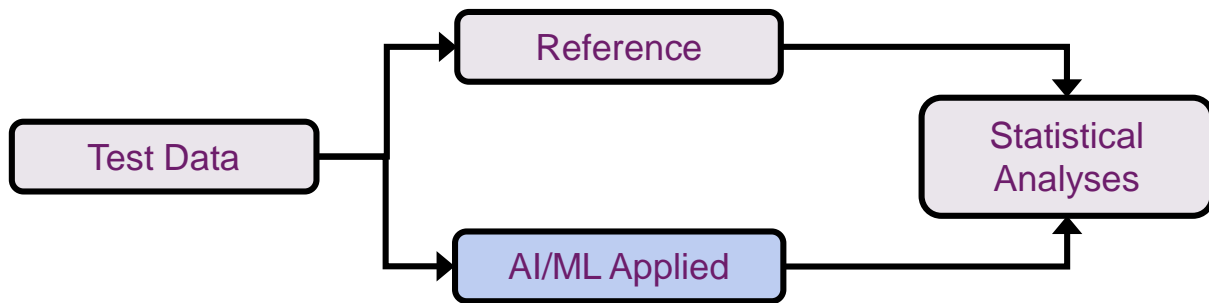
- CAD algorithm and component
 - Device usage and patient population
 - Algorithm design and function
 - Training process
 - Training/tuning datasets
 - ...

- Performance testing protocols and dataset(s)
 - Statistical analysis plan
 - Validation dataset(s)
 - Independent of training/tuning datasets & only accessed after algorithm finalized
 - Reference standard
 - Scoring method
 - ...

STANDALONE PERFORMANCE ASSESSMENT



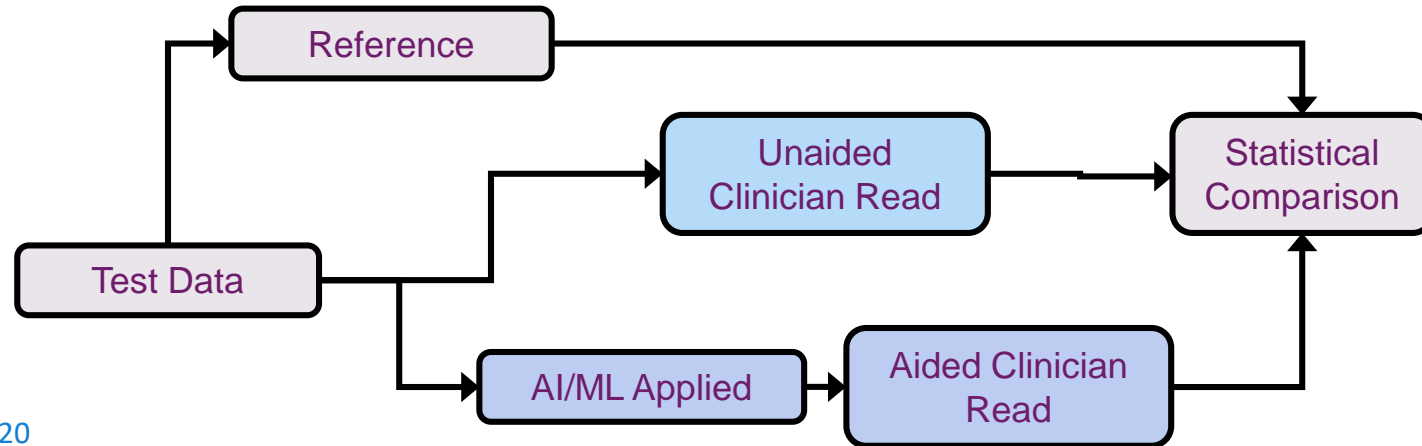
- Benchmark performance (includes CIs)
- Facilitate generalizability analysis



CLINICAL PERFORMANCE ASSESSMENT



- Clinicians' performance utilizing the device
 - Multi-reader multi-case designs
 - Prospective/retrospective





OTHER CAD ASSESSMENT FRAMEWORKS

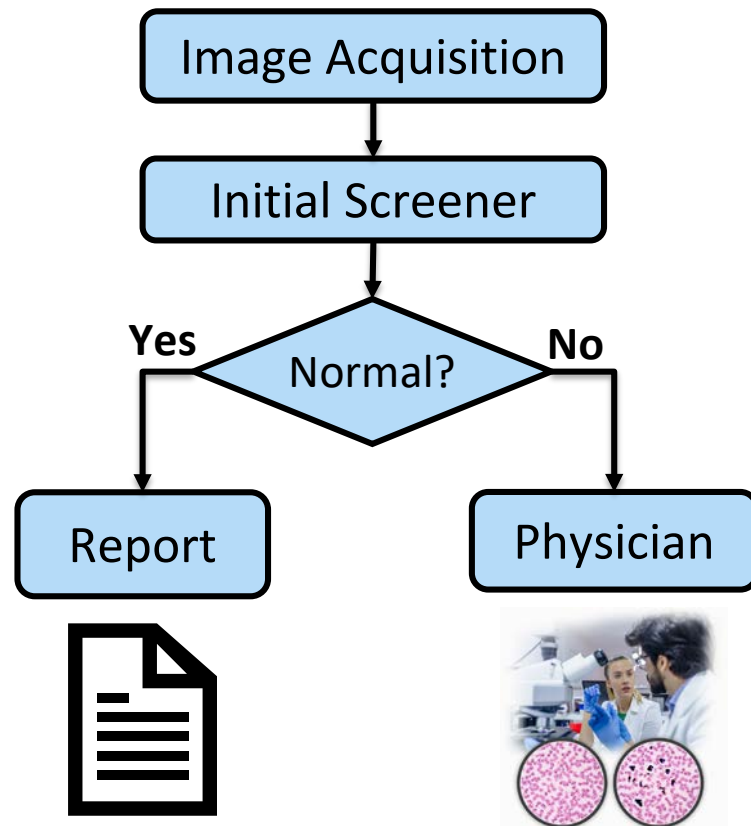
- Special Controls
 - CADx: DEN170022 Quantitative Insights QuantX
 - CADe+CADx: DEN180005 Imagen OsteoDetect
 - Generally consistent with CADe Guidance approaches
 - CADt: DEN170073 Viz.AI ContraCT
 - Detailed Descriptions
 - ...
 - V&V
 - Standalone performance testing
 - Performance testing that demonstrate device will provide effective triage
 - ...
 - Labeling
 - CADa/o: DEN190040 Bay Labs Caption Guidance

LESSON 1

- CAD assessment framework largely applicable to autonomous AI/ML
 - Important differences
 - CADe/CADx
 - Includes standalone and clinical reader performance
 - Generally utilizes enriched dataset
 - Autonomous AI/ML
 - Include standalone/benchmark performance
 - Include clinical comparison with standard of care (?)
 - Representative dataset likely required

INITIAL SCREENER

- Initial Screener (AI/ML+Reader)
 - Normal: AI/ML or reader determine “Neg”
 - Positive: Both AI/ML & reader determine “Pos”
- Properties
 - $Se_{AI/ML+Reader} \leq Se_{SOC}$
 - AI/ML may miss cancers
 - $Sp_{AI/ML+Reader} \geq Sp_{SOC}$
 - SOC reader may misdiagnose normal cases determined correctly by AI/ML
- Assuming reader performance same before/after introducing initial screener



LESSON 2

- Makeup of validation dataset impacts safety and effectiveness assessment

COMMON TEST DATASET TYPES

- Enriched datasets
 - Challenging cases to probe differences in modalities
 - Enriched subgroups to assess generalizability
 - Technological subgroups
 - Patient subgroups

- Representative datasets
 - Representative of clinical population
 - Representative of imaging acquisition landscape

CAD TEST DATASETS

- Evaluating CADs in FDA submissions
 - Generally not measuring aid under clinical conditions
 - Generally not comparing devices across different datasets

 - Standalone
 - Representative (somewhat) to benchmark performance
 - Limited number of sites, likely some selection bias
 - Enriched subgroups/imaging technologies
 - Informs labeling and facilitates generalizability analyses

 - Clinical reader study
 - Typically enriched with challenging cases
 - Clinical readers generally not truly representative of reader population

AUTONOMOUS AI/ML TEST DATASETS



- Evaluating autonomous AI/ML in FDA submissions
 - Likely want ability to direct compare algorithms across different datasets
 - Standalone
 - Representative of patient population and clinical image acquisitions
 - More sites, minimize potential selection bias
 - Enriched patient or imaging technologies subgroups
 - Inform labeling, determine limitations, generalizability analysis
 - Standard-of-Care comparison study
 - Representative of clinical population, clinical image acquisitions, clinical readers
 - Automated Pap smear reader device used a two-arm prospective intended use study at five cytology laboratories

DATASET SIZE

- How much data is enough?

CAD has generally uses modest sized enriched validation datasets

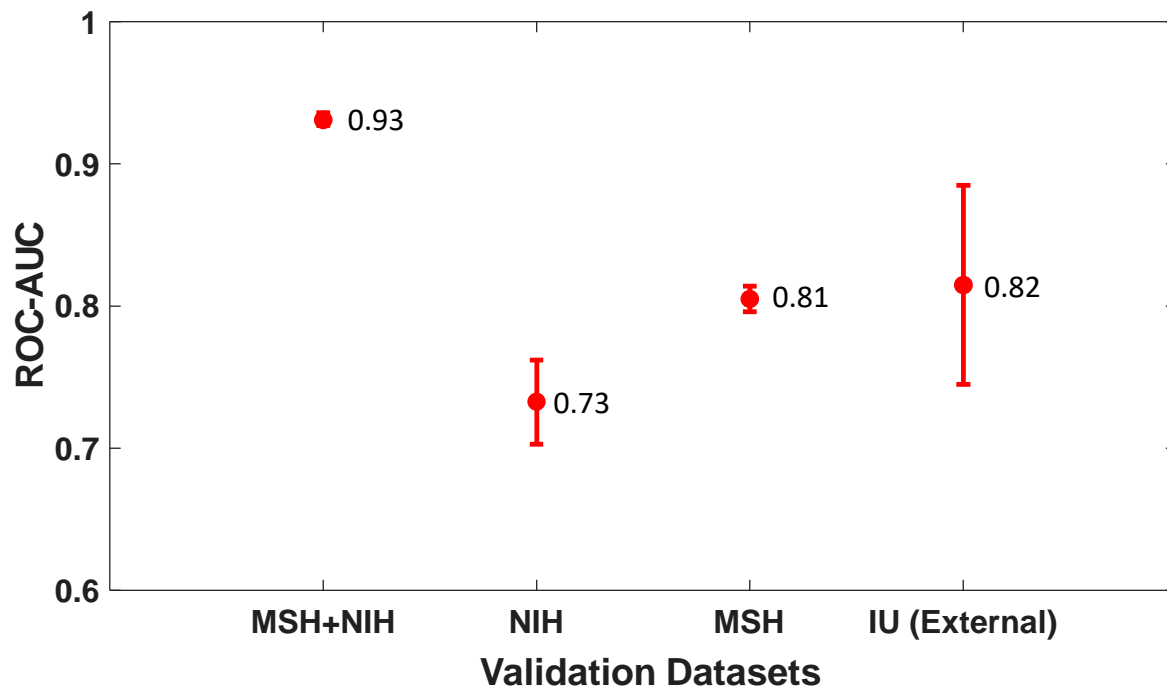


Automatous AI/ML likely to require larger representative validation datasets



LIMITS OF AUTONOMOUS AI/ML TESTING

- Zech et al*, CNN pneumonia screener for chest X-rays
 - Training MSH & NIH
 - Train: 109,922 cases
 - Tune: 15,144
 - Validation
 - NIH: 22,062
 - MSH: 8,388
 - IU: 3,807
 - Pneumonia prevalence
 - NIH: 1.2%
 - MSH: 34.2%
 - IU: 1.0%

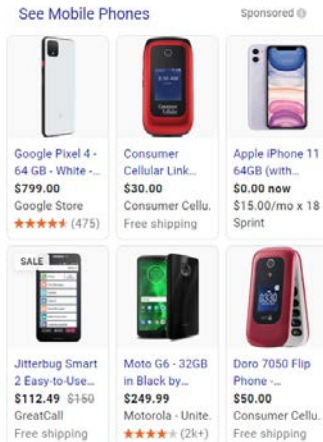


LIMITS OF AUTONOMOUS AI/ML TESTING

- Limits of benchmarking AI/ML performance

Extend to new imaging devices without new studies?

Extend to other patient groups without new studies?



Adults



Children

LESSON 3

- Quality control (QC) is critical for automated AI/ML
 - QC as safety control

AUTOMATED PAP SMEAR INITIAL SCREENER



- Automated QC checks on inputs & AI/ML processing
 - To ensure real-time adequacy of algorithm processing
- Clinical QC
 - To confirm ongoing AI/ML functionality
 - System identifies up to 15% of successfully processed slides for full manual rescreening

SUMMARY OF CAD LESSONS

- For autonomous AI/ML
 - CAD assessment framework largely applicable
 - Large representative validation dataset likely needed
 - Ongoing quality control may be a useful safety control



U.S. FOOD & DRUG
ADMINISTRATION