- Electronic Health Records
- patient
- Clinical decision support systems
- research
- ...

- Clinical information
- Results from other laboratories

interpretation

ordering

reporting

collection

analysis

transport

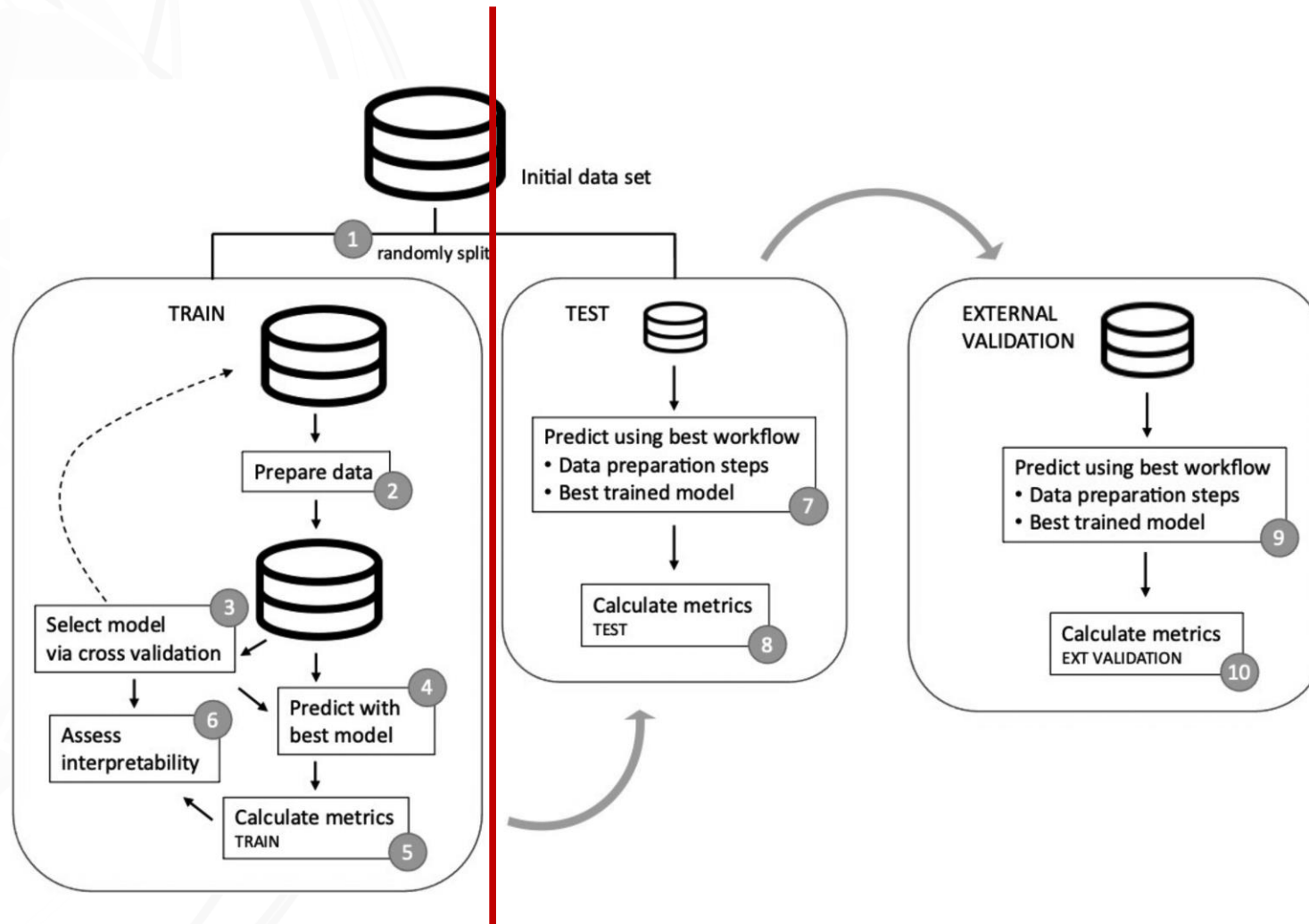Bietenbeck, A., & Streichert, T. (2021). Preparing laboratories for interconnected health care. Diagnostics, 11(8), 1487.

Formulate problem → Collect and prepare data → Validate and select model → Explain and interpret model → Implement model

In the manuscript:

- 30.000 features from novel Omics method
- 60 samples (30 healthy, 30 ill)
- Principal component analysis
- Support vector machine for classification
- Leave-one-out cross validation
- No external validation
- Code available
- AUC: .69

🚩 Insufficient cases for number of features

🚩 Weak validation

What might have happened…

- Nobody understands the data so let's do machine learning
- First approach (e.g. removal of correlated features, random forest): AUC .56
- Next approaches: try out other pre-processing pipelines, algorithms… (> 100 permutations …)
- Report only the best result

1. New cases not similar enough to any of the training examples – failure to generalize
2. Similar inputs associated with different outputs
3. Defined outcomes are controversial because of an ill-defined gold standard
4. Insufficient infrastructure or resources (data scientists) for machine learning
5. Unreliable outcome labelling, lack of in-house expertise to provide training diagnoses.
6. No clear strategy or understanding of the operational context
7. Traditional rule-based software methods are equivalent/better (simple or well-characterized problem)
8. Insufficient data (quantity or quality)
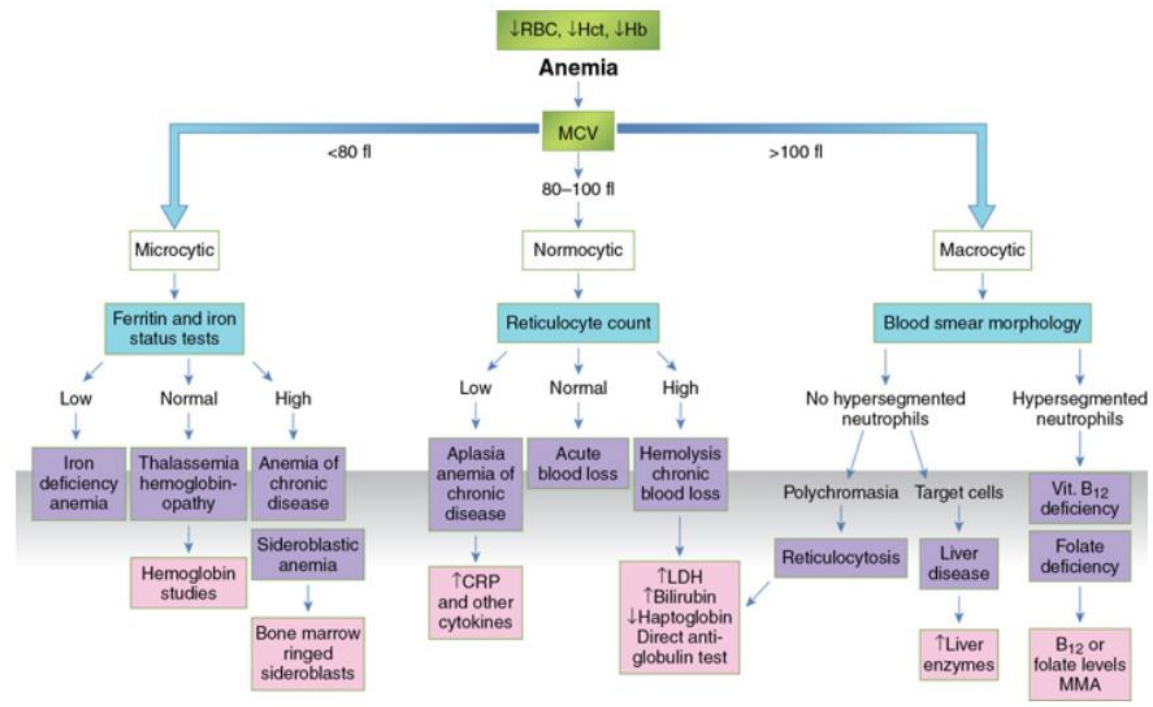
Features:

MCV

HB

Ferritin

Reticulocytes

Haptoglobin

Outcome:
- Iron deficiency anemia
- Renal anemia
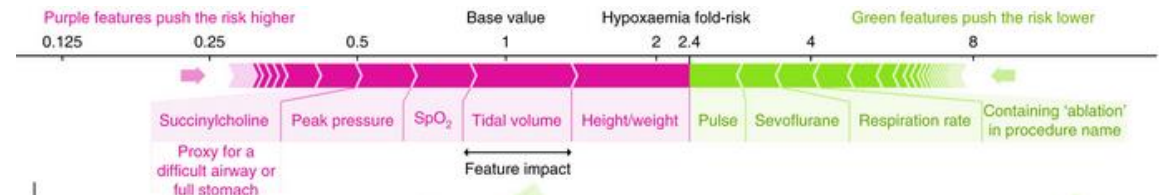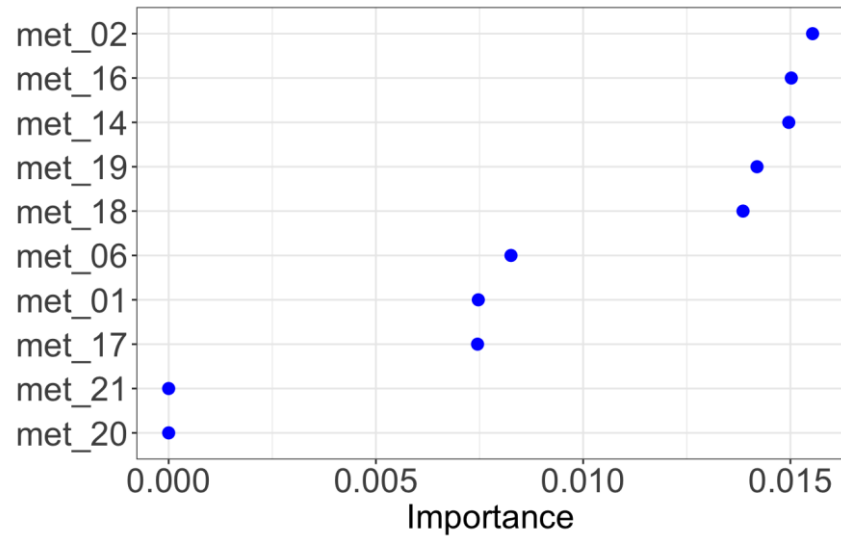- Hemolytic anemia
- Other forms of anemia

"deep neural network"

Rifai, Nader. Tietz textbook of clinical chemistry and molecular diagnostics-e-book. Elsevier Health Sciences, 2017.

1. New cases not similar enough to any of the training examples – failure to generalize
2. Similar inputs associated with different outputs
3. Defined outcomes are controversial because of an ill-defined gold standard
4. Insufficient infrastructure or resources (data scientists) for machine learning
5. Unreliable outcome labelling, lack of in-house expertise to provide training diagnoses.
6. No clear strategy or understanding of the operational context
7. Traditional rule-based software methods are equivalent/better (simple or well-characterized problem)
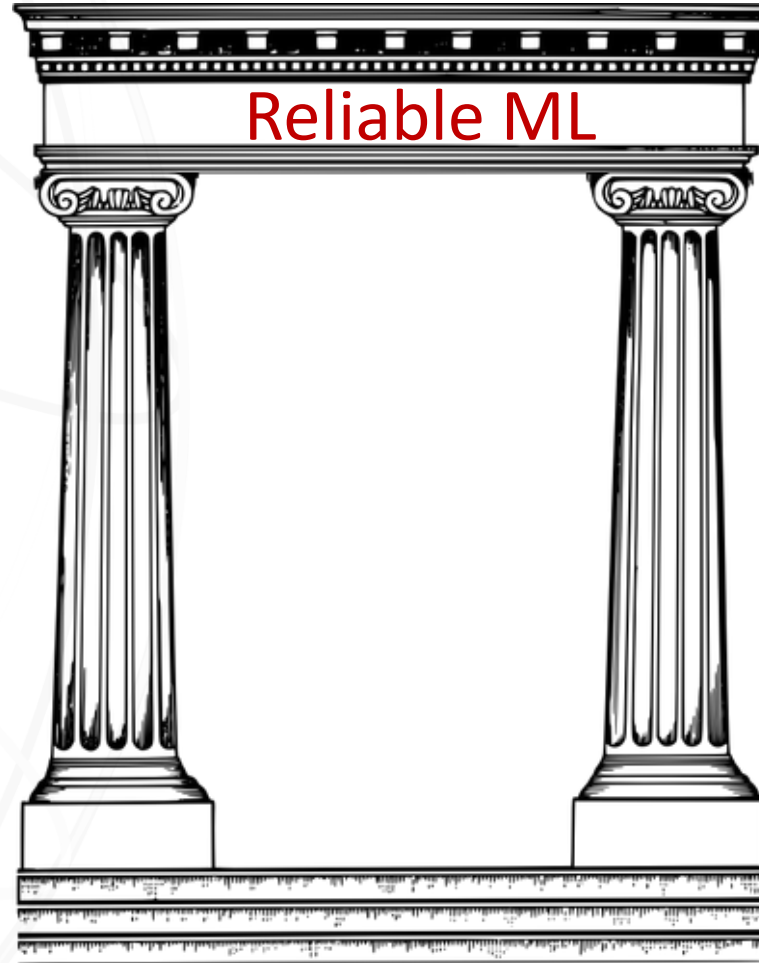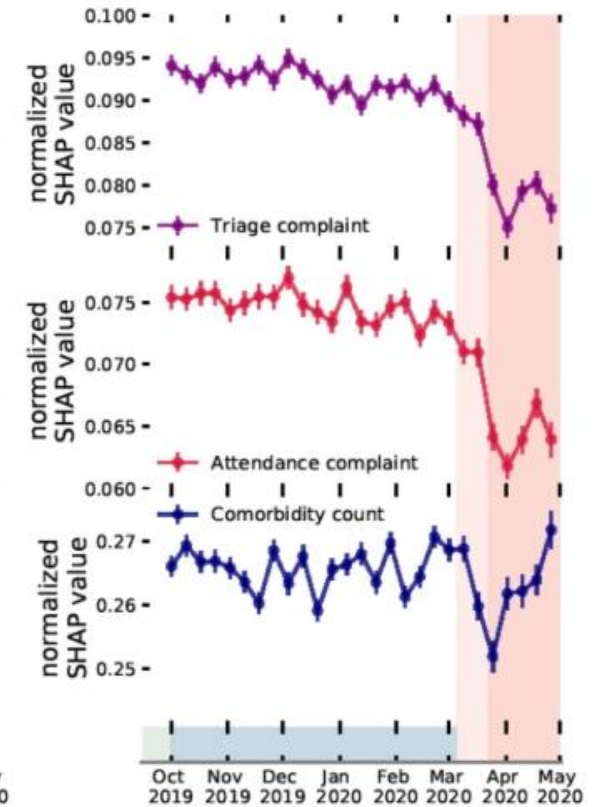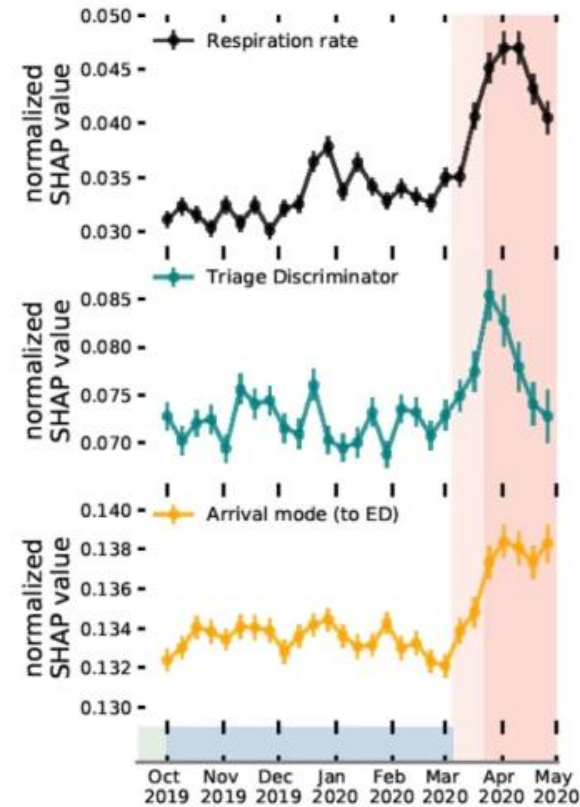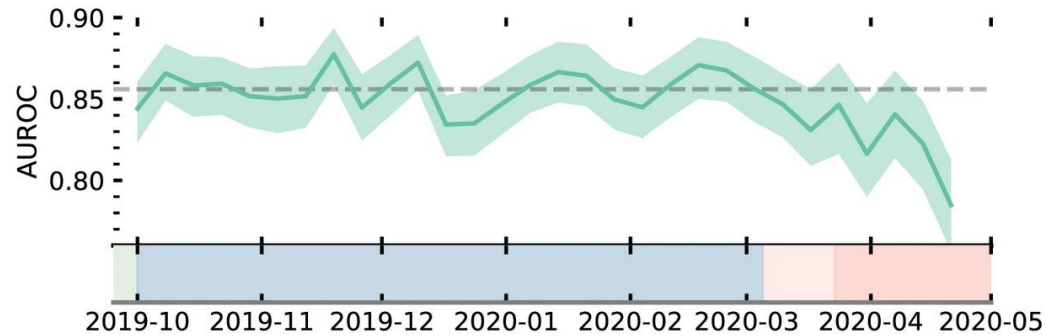8. Insufficient data (quantity or quality)

# Feature importance analysis





Lundberg, Scott M., et al. "Explainable machine-learning predictions for the prevention of hypoxaemia during surgery." Nature biomedical engineering 2.10 (2018): 749-760.

*Recommendation 13: <span style="color:red">Interpret the results and performance</span> of the selected model using suitable global and/or local interpretability methods.  <span style="color:red">Address performance and potential harms</span> in relevant subgroups and clinical scenarios.*
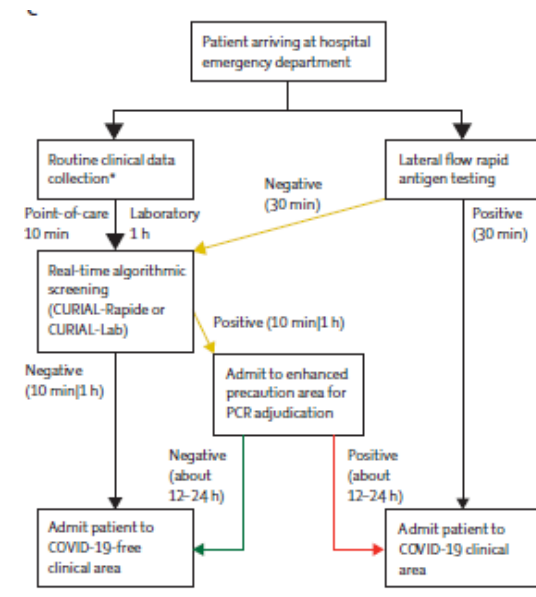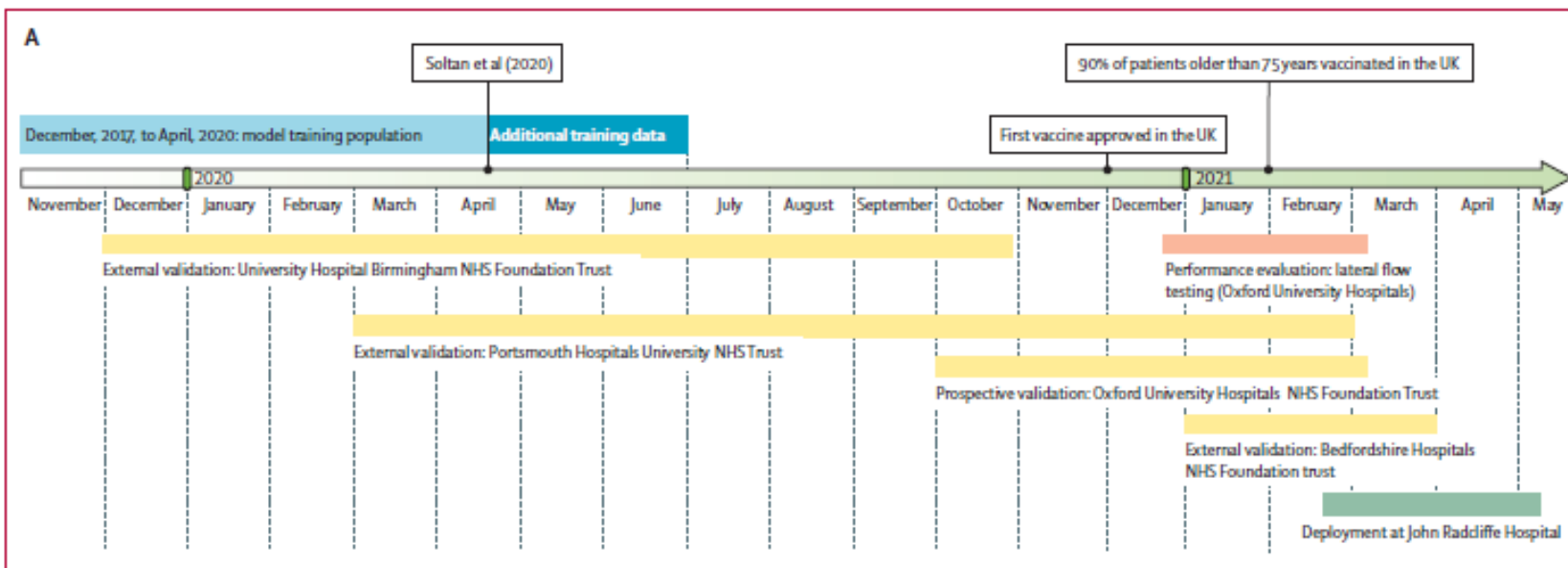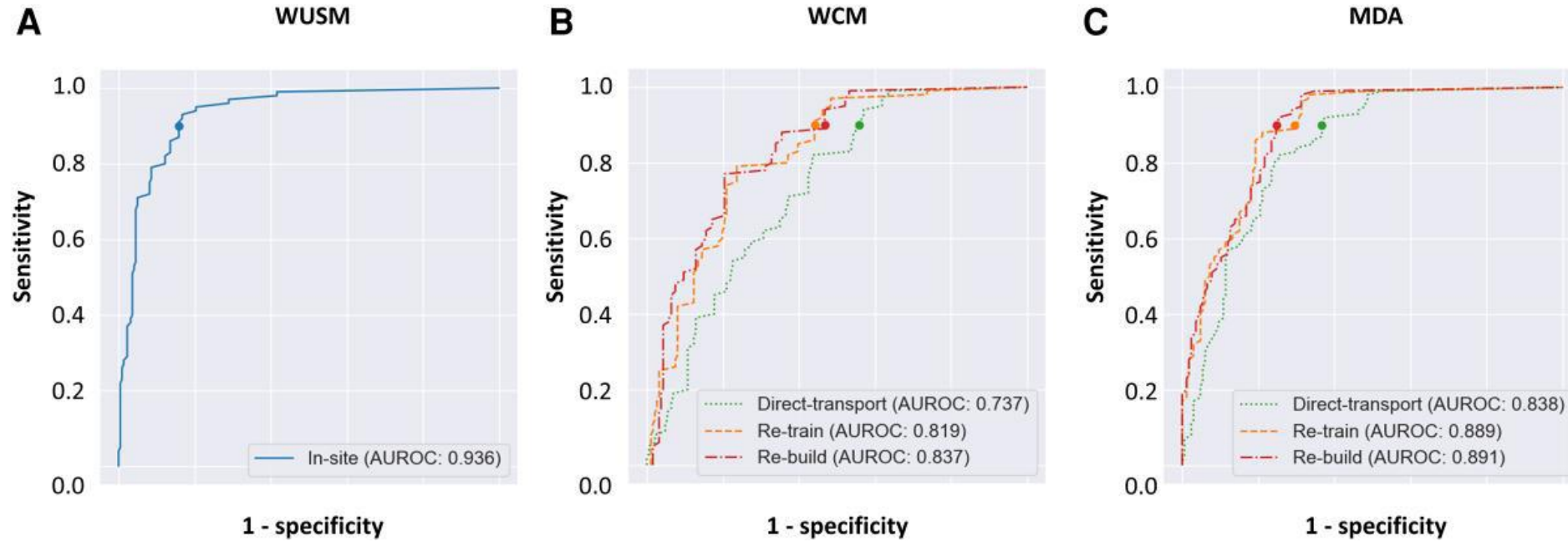
Reliable ML

Data

Interpretability, (knowledge)

Duckworth, Christopher, et al. "Using explainable machine learning to characterise data drift and detect emergent health risks for emergency department admissions during COVID-19." Scientific reports 11.1 (2021): 23017.

Soltan, Andrew AS, et al. The Lancet
Digital Health 2022

"This difference could be partially attributed to the fact that both WUSM and MDA laboratories use the same analyzers to conduct routine chemistry tests [...]"

Karger, Amy B., et al. "Long-term longitudinal stability of kidney filtration marker measurements: Implications for epidemiological studies and clinical care." Clinical chemistry 67.2 (2021): 425-433.

- Keep it simple: Only use Machine Learning when you have to.

- Beware of data leakage!

- Machine Learning is no "magic bullet":  Insufficient data (quantity and quality) cannot lead to convincing results.

- Play to your strengths: Use interpretability methods to evaluate machine learning models.

- Only stable, traceable measurements can guarantee stable, transferable ML models.

**Clinical Chemistry 69:7**
**690–698 (2023)**

**Special Report**

**Machine Learning in Laboratory Medicine:**
**Recommendations of the IFCC Working Group**

Stephen R. Master [iD],[a,b,*] Tony C. Badrick,[c] Andreas Bietenbeck [iD],[d] and Shannon Haymond[e,f,*]

- https://area9lyceum.com/laboratorymedicine/

- IFCC Webinar part 2

- lab@bietenbeck.net

**NEJM Knowledge+** | **AACC Learning Lab**

**MACHINE LEARNING (ADVANCED)**
AACC Learning Lab Advanced

Author information
Shannon Haymond, PhD, MSPA
Stephen R. Master, MD, PhD
Li Zha, PhD

Learning objectives:
✓ Explain principles of machine learning
✓ Describe machine learning process

START LEARNING NOW

**IFCC**
International Federation
of Clinical Chemistry
and Laboratory Medicine

For further information, visit
**www.ifcc.org | eacademy.ifcc.org**

*eAcademy*

**ifcc**
International Federation
of Clinical Chemistry
and Laboratory Medicine