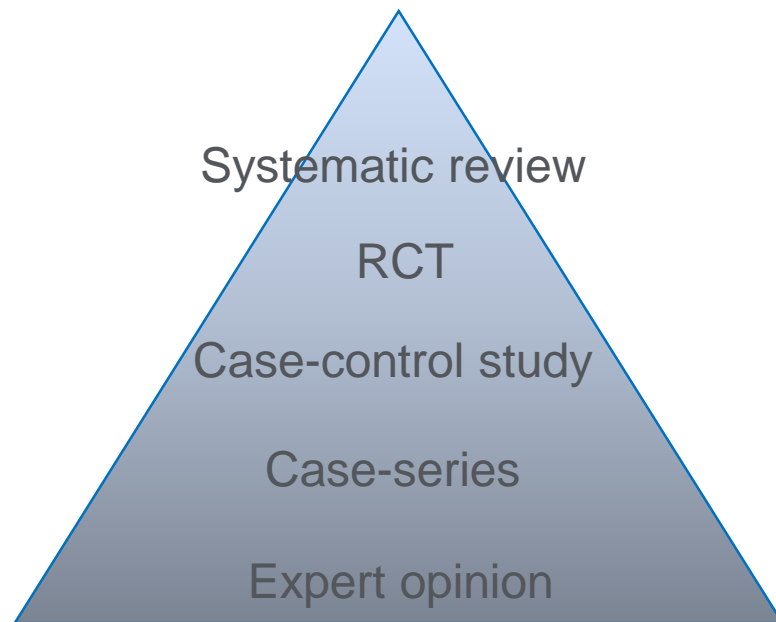# Evaluation of medical algorithms

## AI for Health, NY

November 2018
Naomi Lee. The Lancet.

# Evidence based medicine

- Eminence based medicine

- 1980/1990s

- Medical statistics: the RCT and meta-analysis

- Critical analysis

Systematic review

RCT

Case-control study

Case-series

Expert opinion

# Medical journals

Top journals are trusted sources of information:

- Quality assurance
- Peer review
- Standards

- Select practice changing research

# What are the standards applied?
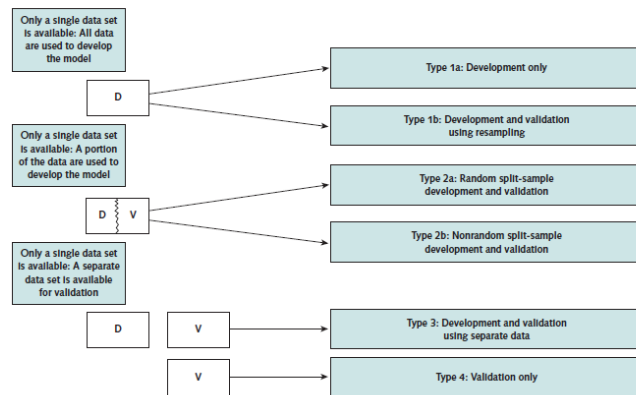
ICMJE

Helsinki declaration

EQUATOR NETWORK

Author guidelines from journals

# EQUATOR NETWORK

- Reporting guidelines for health research

- Transparent reporting of a multivariate prediction model for an Individual Prognosis or Diagnosis (TRIPOD)

- "Gives keys details of how prediction models were developed and validated in order to assess generalizability and risk of bias"

- External validation in a separate dataset



Figure 1. Types of prediction model studies covered by the TRIPOD statement.

# What is 'practice changing'?

- Accuracy of diagnosis/prediction

- Evidence of efficacy
    - Clinically meaningful endpoint
    - Compared again current standard

- Cost effectiveness

THE LANCET

# Why is that a problem?

- Adoption of unassessed technology causes patient harm

# What are the pitfalls for AI?

- A lot of health AI research isn't externally validated

- It doesn't demonstrate clinical efficacy or cost effectiveness

- # How do we transition to the mainstream?


- Standards

  - Quality standards

  - Reporting Guidelines


- Framework for assessing efficacy and cost effectiveness

Thank you….