



Automated quality assessment of medical evidence



—
Simon Šuster

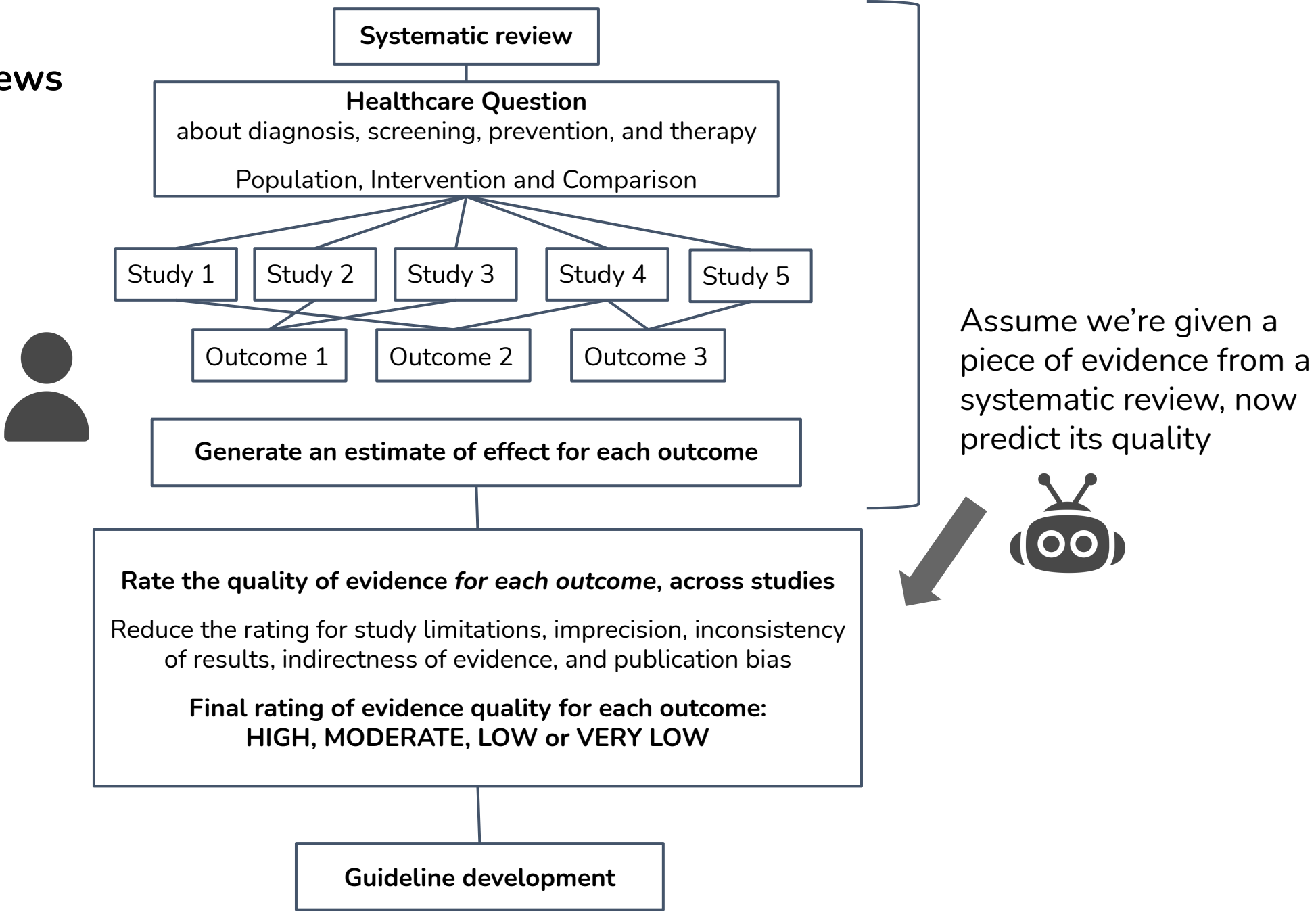
with Tim Baldwin, Antonio Jimeno Yepes, Jey Han Lau,
David Martinez Iraola, Yulia Otmakhova, Karin Verspoor

Stream 4

29/1/2021



Constructing systematic reviews and quality assessment



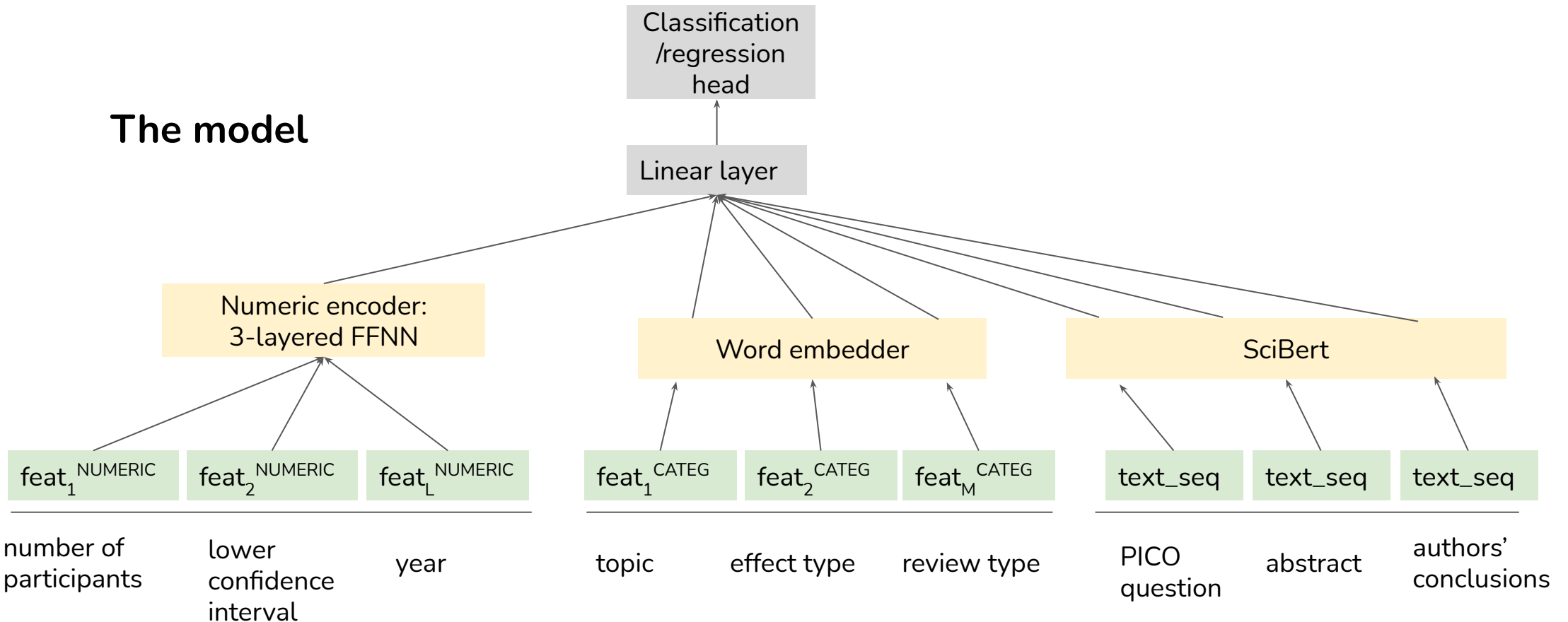
Modified from Guyatt et al. 2011

Obtained from Cochrane Systematic Reviews by extracting relevant data from Summaries of Findings and other parts of reviews.

The data

- ~7,000 reviews
- ~13,500 data points (pieces of evidence), split into train/dev/test sets
- Author-assigned quality scores represent our gold standard (labels)

The model



10-fold cross-validated results

<i>problem:</i>	4-level grading	binarised grading	multi-labeling of reasons																						
<i>labels:</i>	<table><tr><td>high</td><td>F1: 0.47</td></tr><tr><td>moderate</td><td>0.53</td></tr><tr><td>low</td><td>0.51</td></tr><tr><td>very low</td><td>0.46</td></tr></table>	high	F1: 0.47	moderate	0.53	low	0.51	very low	0.46	<table><tr><td>high/moderate</td><td>F1: 0.72</td></tr><tr><td>low/very low</td><td>0.76</td></tr></table>	high/moderate	F1: 0.72	low/very low	0.76	<table><tr><td>risk of bias</td><td>F1: 0.7</td></tr><tr><td>imprecision</td><td>0.6</td></tr><tr><td>inconsistency</td><td>0.2</td></tr><tr><td>indirectness</td><td>0.3</td></tr><tr><td>publication bias</td><td>0.1</td></tr></table>	risk of bias	F1: 0.7	imprecision	0.6	inconsistency	0.2	indirectness	0.3	publication bias	0.1
high	F1: 0.47																								
moderate	0.53																								
low	0.51																								
very low	0.46																								
high/moderate	F1: 0.72																								
low/very low	0.76																								
risk of bias	F1: 0.7																								
imprecision	0.6																								
inconsistency	0.2																								
indirectness	0.3																								
publication bias	0.1																								
<i>overall result:</i>	<p>F1: 0.5 MAE: ~0.6 random: 0.25 F1, 1.2 MAE majority: 0.13 F1, 0.8 MAE</p>	<p>F1: 0.74 random: 0.52 F1 majority: 0.35 F1</p>	<p>F1: 0.4 random: 0.3 F1 majority: 0.3 F1</p>																						

Ongoing work and open questions

Incorporating information from primary studies when assessing overall quality

- problem of retrieval and efficient encoding
- *possible solution*: attentive module that selects the studies based on similarity between questions in the studies and the piece of evidence

Using different label types in a single model

- instances annotated with quality scores *and* downgrading reasons
- *possible solution*: a stacked/multi-task model

If grading reliability of human reviewers is poor, how does it affect learning?

- *possible solution*: obtain multiply graded pieces of evidence and report interrater agreement; are there categories with higher reliability?

In adults without cardiovascular disease, does Mediterranean diet (compared to no dietary intervention) help reduce the risk of cardiovascular disease?

*CVD mortality
stroke
myocardial infarction
total cholesterol change
...*

*Myocardial infarction as outcome:
Risk: 12 per 1000 (Intervention)
16 per 1000 (Control)
...*

*GRADE: ⊕⊕⊕⊖ (low)
Downgraded by one level for imprecision. Confidence interval is wide enough to include both an important increase or decrease in the outcome.
Downgraded by one level for risk of bias. The only included study was the PREDIMED trial retracted due to methodological issues with randomisation [...]*

