



# Automated quality assessment of medical evidence



—  
Simon Šuster

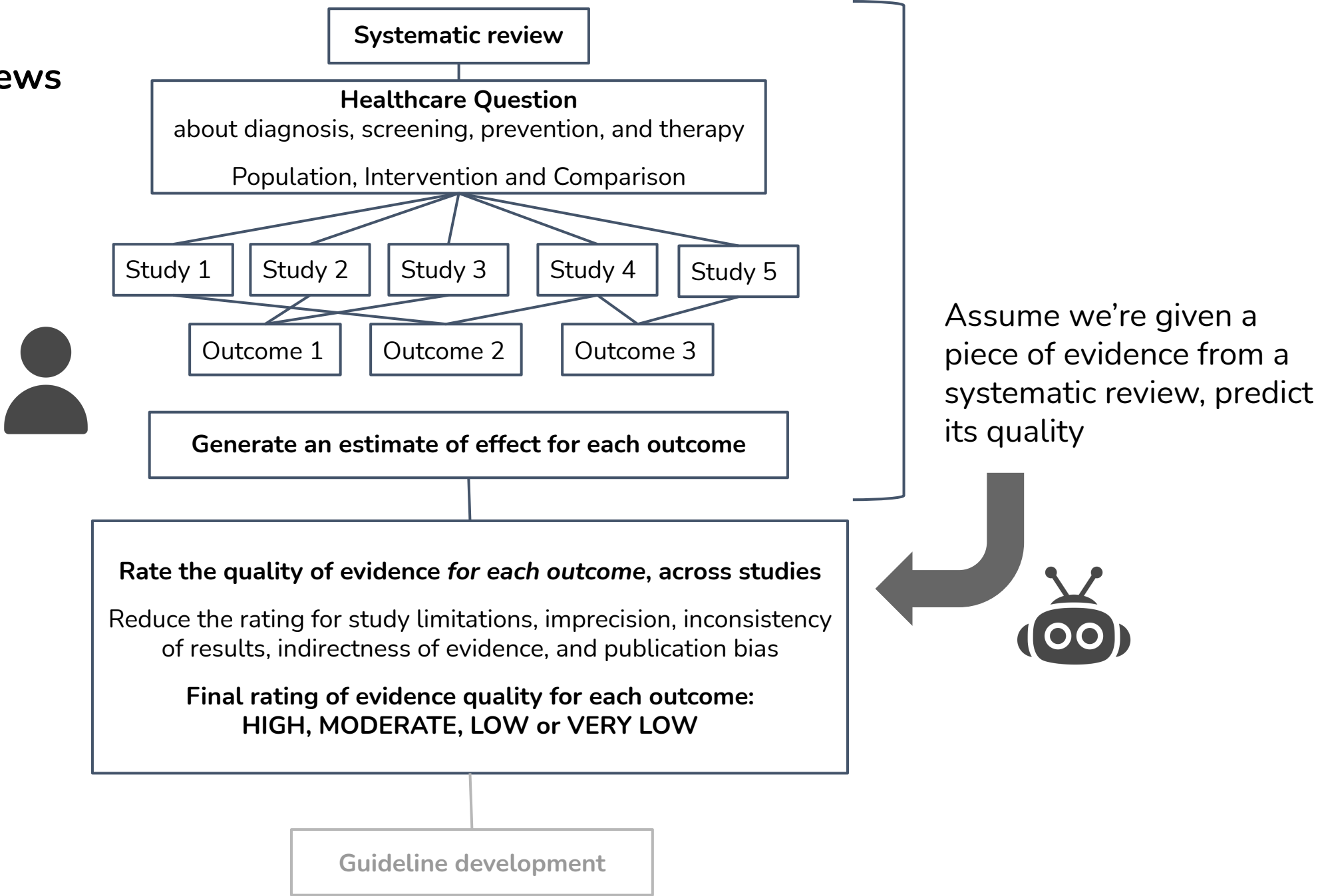
with Tim Baldwin, Antonio Jimeno Yepes, Jey Han Lau,  
David Martinez Iraola, Yulia Otmakhova, Karin Verspoor

Stream 4

22/3/2021



# Constructing systematic reviews and quality assessment



# Landscape

## Predicting strength of recommendation of a body of evidence (Sarker et al., 2015)

- 1,100 abstracts, 3 levels according to SORT (Ebell et al., 2004)
- Publication metadata features and word n-grams: 64% accuracy

### **Limitations:**

- Unclear what the score measures (strongly reflects the publication types)
- Loosely defined SORT criteria and inclusion criteria (doesn't follow PICO)
- Cohen's kappa of around 0.5 for human annotators

## Grading individual studies with isolated criteria

- Risk-of-bias assessment in RobotReviewer and TrialStreamer (Marshall et al., 2017 & 2020)

### **Limitations:**

- Does not grade the body of evidence but individual studies

## Semi-automated quality assessment (SAQAT; Stewart et al., 2015)

- Human reviewers answer checklist questions
- Final score assigned by a Bayesian network

### **Limitations:**

- Still largely manual

# Data creation



Pharmacological interventions for cognitive decline in people with Down syndrome
 Interventions for restoring patency of catheter lumens
 Mediterranean-style diet for the primary and secondary prevention of cardiovascular disease
 Wheat flour fortification for improving iron status

~7,000 systematic reviews (majority from 2010-)

Extract data related to quality appraisal from summaries of findings

Outcome	Relative effect* (95% CI)	Number of participants (studies)	Quality of the evidence (GRADE)	Comments
Cognitive abilities (Severe impairment Battery, SIB)	The mean change in cognitive abilities in the intervention groups was 5.32 points higher (0.27 lower to 1.13 higher)	165 (3 RCTs)	⊕⊕⊕⊖ Low <sup>†</sup>	
Behavioural problems (business scales)	The mean change in behavioural problems the intervention groups was 0.42 points higher (0.68 lower to 0.83 higher)	157 (3 RCTs)	⊕⊕⊕⊖ Low <sup>†</sup>	
Adverse events	Risk with placebo	OR 0.32 192 per 1000 (0.16 to 0.62)	⊕⊕⊕⊖ Low <sup>†</sup>	
	Risk with dapspeil	Study population 313 per 1000 (167 to 768)		
Carer stress	No data available			
	No data available			
Institutional/home care	No data available			
	No data available			



13,500 outcomes rated for quality using GRADE framework (with justifications)

- Validation of extraction procedures against human-verified data (Conway et al., 2017)
- Prepared data for 10-fold CV, with train/dev/test splits
  - instances built from same SR kept in the same split
- Author-assigned quality scores represent our gold standard (labels)

## Levosimendan compared with placebo for cardiogenic shock or low cardiac output syndrome

**Patient or population:** adults with cardiogenic shock or low cardiac output syndrome

**Settings:** hospital

**Intervention:** levosimendan

**Comparison:** placebo

Outcomes	Anticipated absolute effects (95% CI)		Relative effect (95% CI)	No of participants (studies)	Quality	Comments
	Risk with placebo	Risk with levosimendan				
<b>All-cause short-term mortality:</b> range 4 to 6 months	<b>Moderate<sup>1</sup></b>		<b>RR 0.48</b> (0.12 to 1.94)	55 (2)	⊕⊕⊕⊖ <b>very low<sup>3,4</sup></b>	Studies included participants with LCOS or CS due to HF or AMI
	187 per 1000	<b>90 per 1000</b> (22 to 363)				
	<b>High<sup>2</sup></b>					
	500 per 1000	<b>240 per 1000</b> (60 to 970)				

<sup>3</sup>Downgraded one step due to study limitation because of lack of blinding of participants and physicians, and missing information on randomisation in the larger study.

<sup>4</sup>Downgraded two steps for imprecision due to few events and the confidence interval crosses the line of no difference and includes possible benefit from both approaches.

# Tasks



Quality score  
(GRADE)

very low	0
low	1
moderate	2
high	3



Binary quality  
score

very low-low
moderate-high



Downgrading  
reasons




imprecision
inconsistency
indirectness
RoB
pub. bias



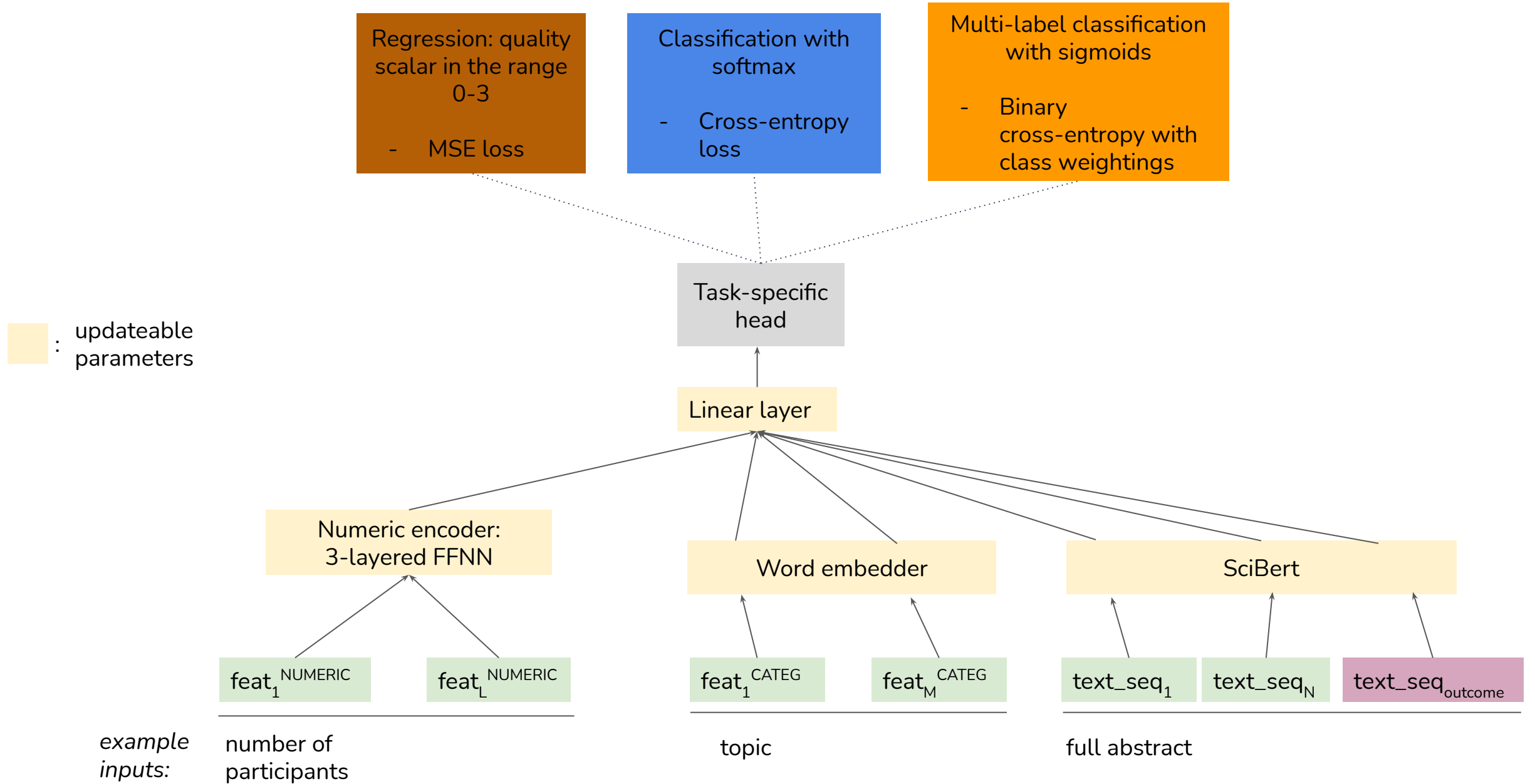
Number of reason  
types

0
1
2
3



-  classification
-  regression
-  multi-label classification

# Base model



# Feature space

## Textual

- parts of SRs that are likely to discuss quality
- impose little assumptions, open-ended solution
- (4 in total) *Authors' conclusions, plain language summary, abstract conclusion, full abstract*

## Categorical

- meta-data about the review and non-numerical statistical information
- (3) *Review type, topics, type of effect*

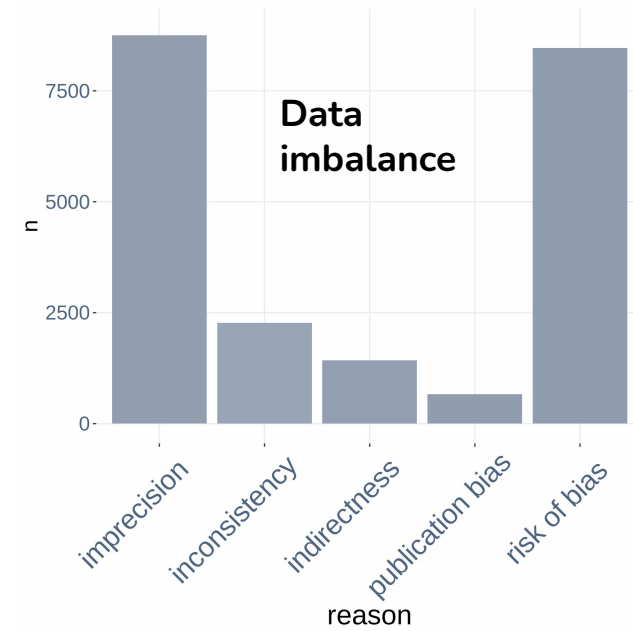
## Numerical

- meta-data about the review and statistics
- (13) *Num. of included studies, year, num. of outcomes, relative effect, lower CI, upper CI, ...*

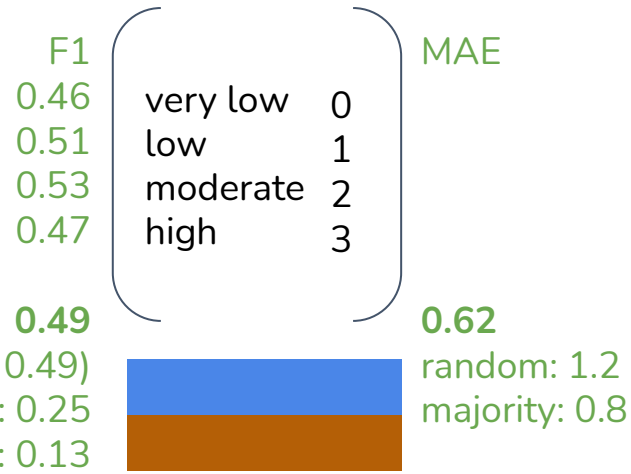


# Results (averaged over 10 folds)

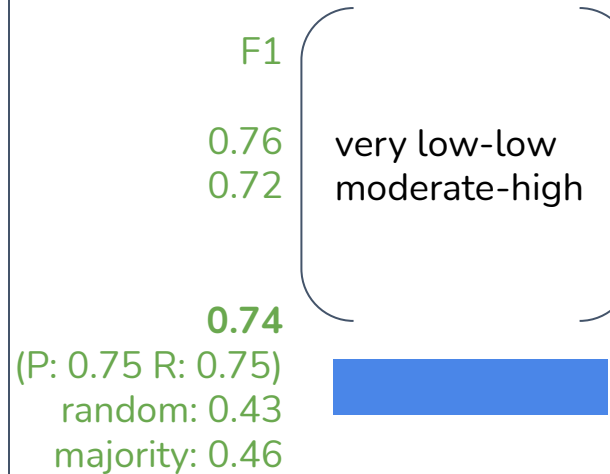
- classification
- regression
- multi-label classification



## Quality score (GRADE)



## Binary quality score



## Downgrading reasons



# Reliability of ratings

Existing user studies (Meader et al., 2014, Berkman et al., 2013, Hartling et al., 2013, Mustafa et al., 2013, Atkins et al., 2005)

- Limited by small sample sizes and datedness
- Poor to almost perfect agreement
- RoB and imprecision (risk of random errors)

What we know so far:

- 4-level quality annotation perhaps too granular/fine distinctions somewhat arbitrary
- Binarisation makes the task easier
- Fairly good performance on some reason classes
  - data imbalance likely a problem (less common reasons predicted less well)
- Our small-scale reliability study for risk-of-bias on recurring primary studies

Possible future work:

- Expert independently grades evidence and provides justifications; then measure agreement with the Cochrane authors →hard and time-consuming.
- Expert only judges the validity of assigned quality grades and justifications; identifies support within the review for authors' decision
- Multiple reviews for the same PICO question -> trouble defining the equivalent questions; unlikely to have the same set of primary studies

# Cochrane's RoB2 framework for RCTs

## Judgement per type

judgement<sub>0</sub>  
judgement<sub>1</sub>  
...

- Bias arising from the randomisation process
- Bias due to deviations from intended interventions
- Bias due to missing outcome data
- Bias in measurement of the outcome
- Bias in selection of the reported result

## Risk-of-bias types

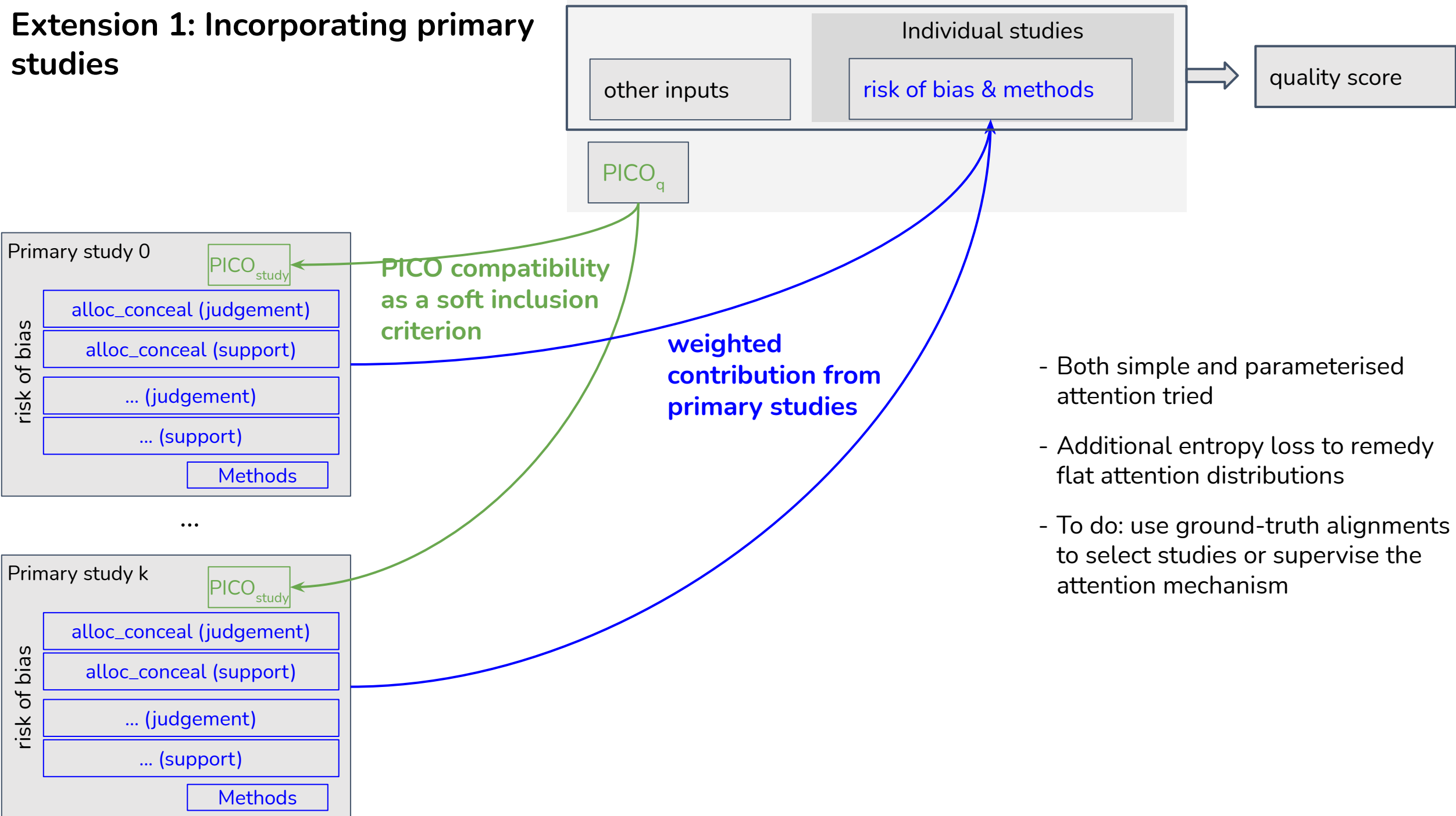
## Signalling questions

- is the allocation sequence random?
- is the allocation sequence adequately concealed?
- do baseline differences between intervention groups suggest a problem with the randomization process?

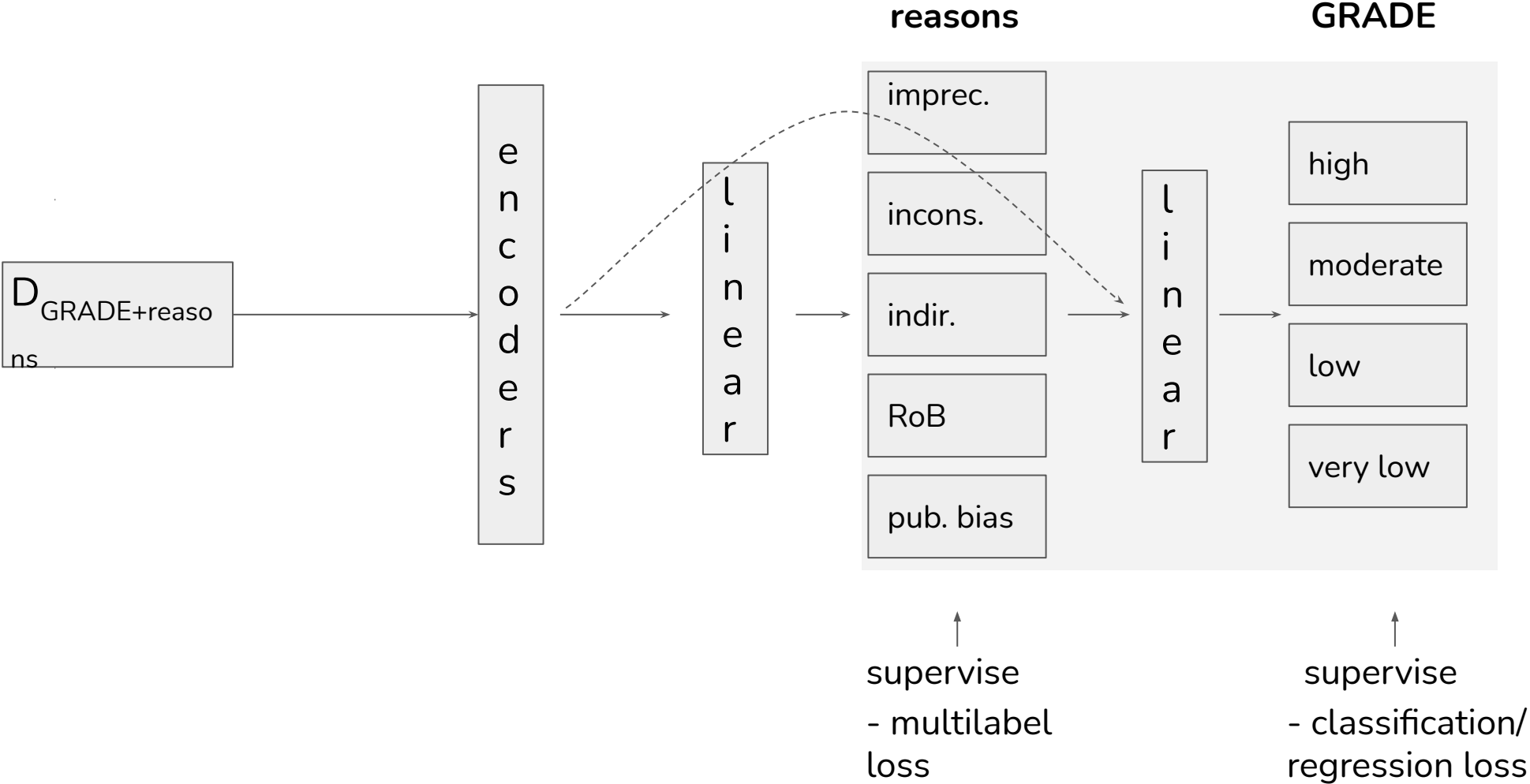
- participants were aware of their assigned intervention during the trial;
- carers and people delivering the interventions were aware of participants' assigned intervention during the trial.
- ...

Overall judgement: low, some concern, high

# Extension 1: Incorporating primary studies



# Extension 2: Joint modeling of downgrading reasons and quality score



# Ongoing work and open questions

Obtain data about included studies per PICO

- extract from review editing software files
- observe effect of augmented input on predictions

Augment data for poorly represented reason types

Reliability study

More insights into dataset:

- e.g. are some interventions more likely to yield high-quality evidence?  
(e.g. pharmacological vs. surgical)
- cluster/label the PICO criteria and relate to the quality of evidence

# Systematic review

## Healthcare Question

about diagnosis, screening, prevention, and therapy  
Population, Intervention and Comparison

Study 1

Study 2

Study 3

Study 4

Study 5

Outcome 1

Outcome 2

Outcome 3

Generate an estimate of effect for each outcome

Rate the quality of evidence for each outcome, across studies

Reduce the rating as needed (study limitations, imprecision, inconsistency of results, indirectness of evidence, publication bias)

Increase the rating (e.g. large effect size)

Final rating of evidence quality for each outcome:  
**high, moderate, low or very low**

Guideline development

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 [...]