



Evidence Selection in Long Documents

Subendhu Rongali, Rajarshi Das
University of Massachusetts Amherst

UMassAmherst
College of Information & Computer Sciences

Problem Definition

- Given a medical trial report and a structured question prompt about the findings in the report, find the evidence and answer the question.
- There is a need to automate this, given the volume of biomedical evidence being released everyday.
- The size of these reports (up to thousands of words) causes significant computational challenges.
- Dataset – Annotated Pub-med articles and prompts (10k).
- Prior approaches don't use transfer learning, evidence attention is scattered.



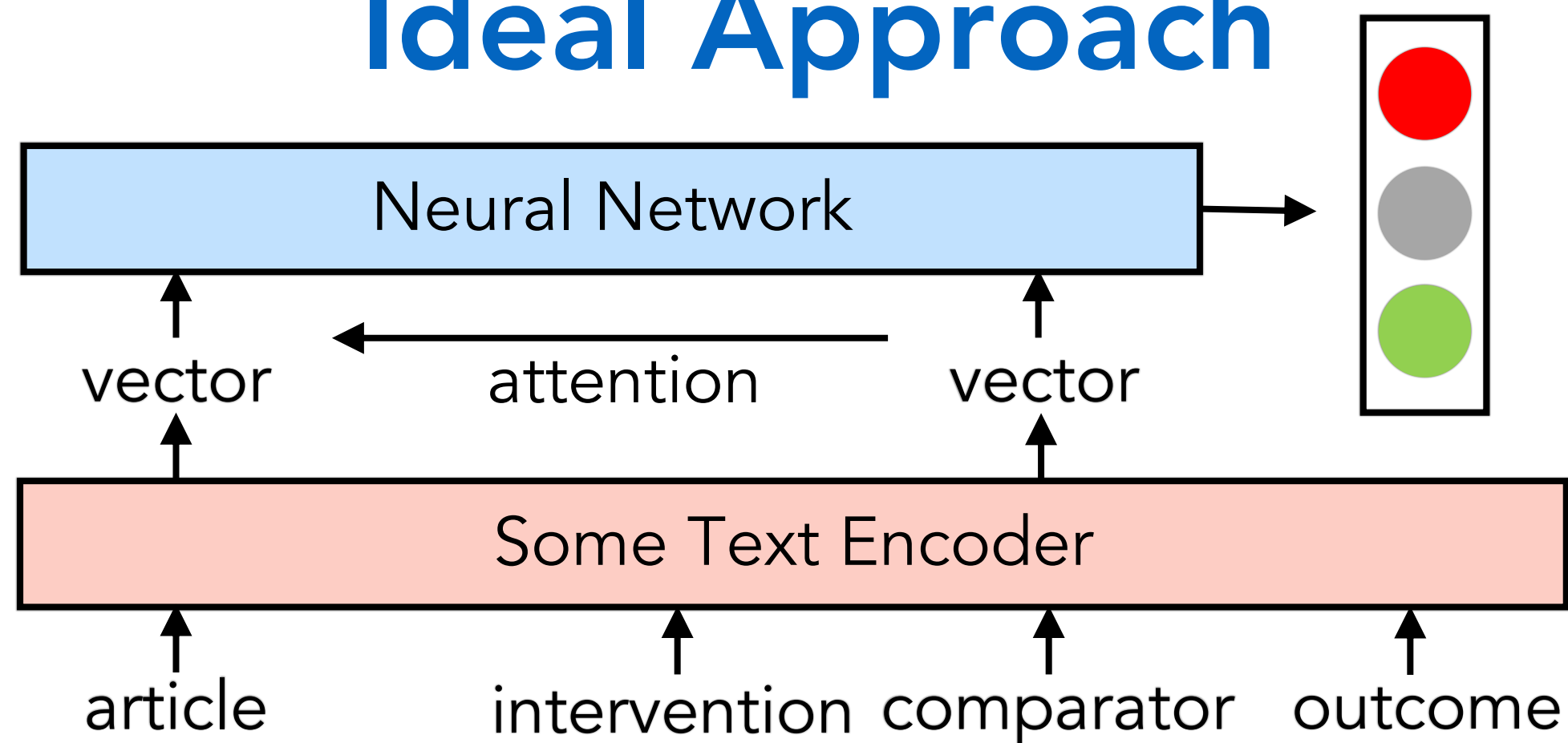
Intervention: *metronidazole*
Comparator: *placebo*
Outcome: *pre-term birth*

prompt

Patients receiving metronidazole experienced significantly fewer pre-term births than those in the comparison group.

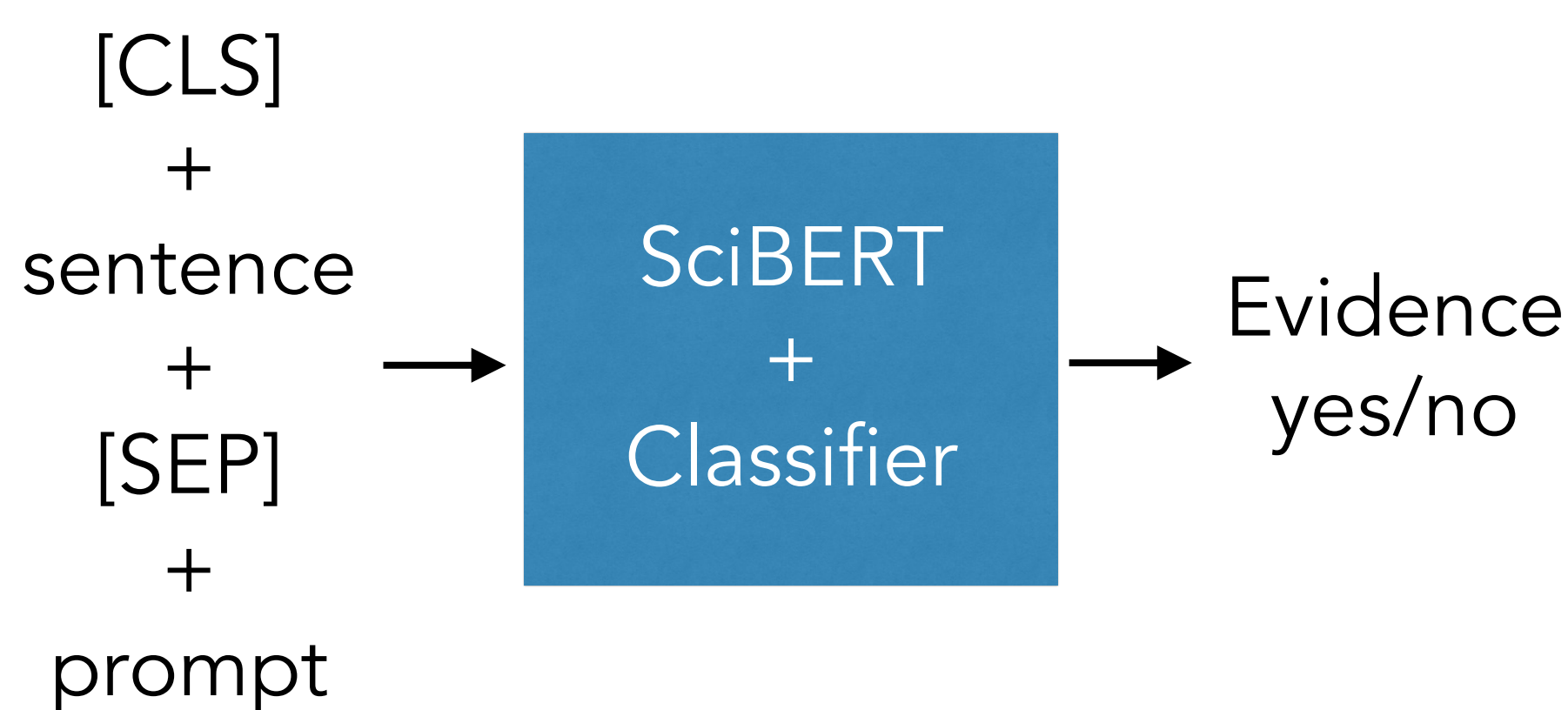


Ideal Approach

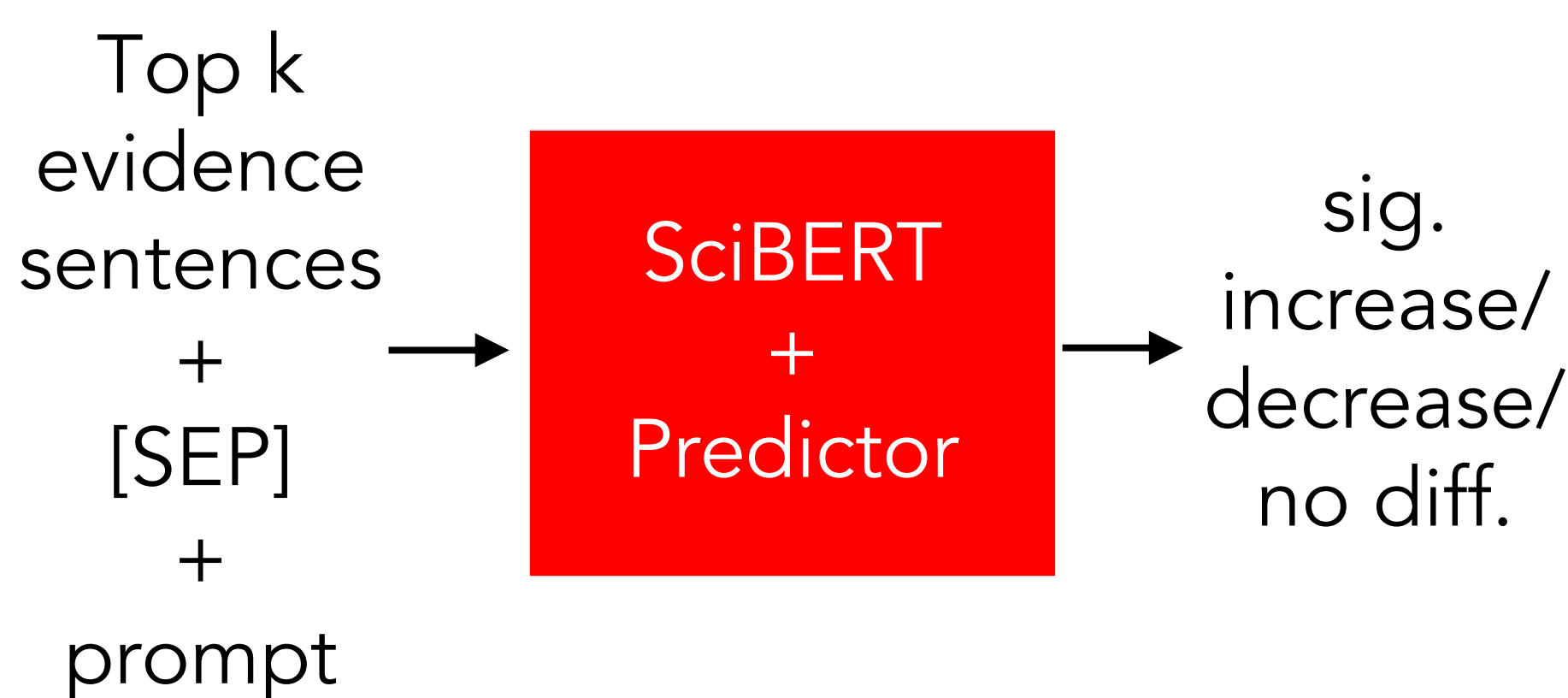


Our Approach

- We use a pipelined approach to handle processing the entire article through a BERT encoder.
- We first train an evidence classifier to pick the possible evidence sentences.
- We then train a model to predict the answer based on the top evidence sentences.



- We train by sampling positive and negative evidence sentences and maximizing their margin.
- We then run this model on all sentence to get their evidence scores and pick the top k.



Results and Findings

- The evidence classifier is trained using pairwise positive and negative samples. The negative samples include random non-evidence sentences, as well as evidence sentences for other prompts.
- We experimented with hinge loss with margins in the range (0.5, 0.8), and binary cross entropy loss.
- We later switched to using paragraphs instead of sentences to include more context.

Evidence Classifier	
	Accuracy
SciBERT + Hinge	84.5
SciBERT + BCE	81.2
Predictor (oracle)	
	F-score
Ours - SciBERT	81.6
Previous best	73.9
Overall	
Heuristics	43.1
Previous best	53.1
Ours - SciBERT Pipeline	42.2

Conclusions & Next Steps

- The evidence classifier and predictor seem to be working well on their own. But the overall pipeline is broken. This needs to be fixed.
- A lot of sentences lack context. The marked evidence sentences don't have the prompt words in them. We tried to add some context but they are still lacking.
- We will try to bridge the two models in the pipeline using a weak signal - something on the lines of negatively reinforcing evidence sentences that don't help the downstream predictor.