

## INTRODUCTION

An Outcome is a measurement or an observation used to capture and assess the effect of a treatment [1]. Automating Outcome Detection (OD) could **speed up** access to evidence necessary in health care decision making. Given a sentence, "There was no significance between group difference in the incidence of **wheezing** or **shortness of breath**", OD extracts outcomes such as those underlined and in bold font.

OD has however previously been hindered by an absence of a consensus on how outcomes should be reported and classified. Moreover, datasets like EBM-NLP [2] supporting OD have been found erroneous [3] with flaws like,

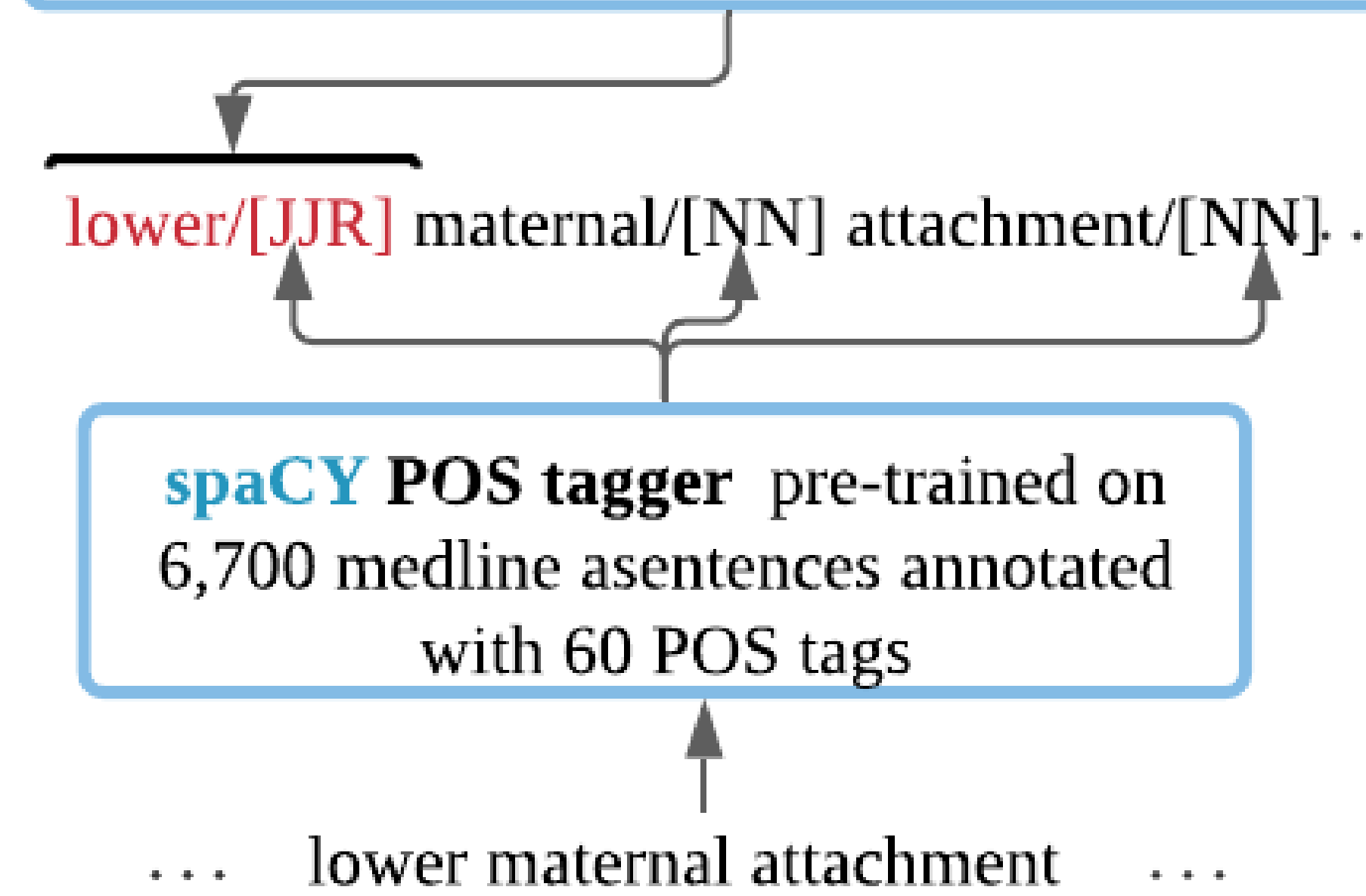
- Flaw 1: Inclusion of unnecessary text.
  - statistical metrics e.g. "mean arterial pressure".
  - Clinical measurement tools e.g. "Quality of life Questionnaire".
- Flaw 2: Failure to identify independent and granular outcomes.
  - e.g. "cardiac arrest and heart failure"
- Flaw 3: Imprecise outcome annotations.
  - e.g. "Suicidal Ideations" annotated as a **Mortality** outcome rather than Mental outcome.

## EBM-COMET (ANNOTATION & EVALUATION)

Annotation category	Annotated text
Simple	<p>...Tai Chi may alleviate &lt;P 0, 28&gt;depression&lt;/&gt;of the elderly through modulating autonomous nervous system or &lt;P 0&gt;heart rate variability&lt;/&gt;...</p> <p><b>depression</b> - [0:Physiological, 28:Emotional Functioning] <b>heart rate variability</b> - [Physiological]</p>
Complex	<p>...The objective of this study was to evaluate &lt;P 0&gt;(S2)right heart size &lt;P 0&gt;and function&lt;/&gt;assessed by echocardiography during long term treatment</p> <p><b>right heart size</b> - [0:Physiological] <b>right heart function</b> - [0:Physiological]</p>

## CORRECTING FLAWED OUTCOME ANNOTATIONS

**Rule based syntactic chunking** will eliminate un-necessary text including metrics, contextually comperative POS, punctuations



**Figure 1:** Part-Of-Speech (POS) Tagging and Rule-based Chunking to build EBM-NLP<sub>rev</sub>

Model	EBM-NLP	EBM-NLP <sub>rev</sub>
biLSTM	72.2	<b>80.3</b>
biLSTM - Flaw 2	-	<b>74.3</b>

**Table 1:** F1 (%) for OD on original and revised version of EBM-NLP, including when only Flaw 2 is corrected.

**Data set statistics:** 300 RCT PubMed Abstracts, 5193 sentences, an average of 0-4 outcome phrases/sentence.

### Full Outcome phrase detection

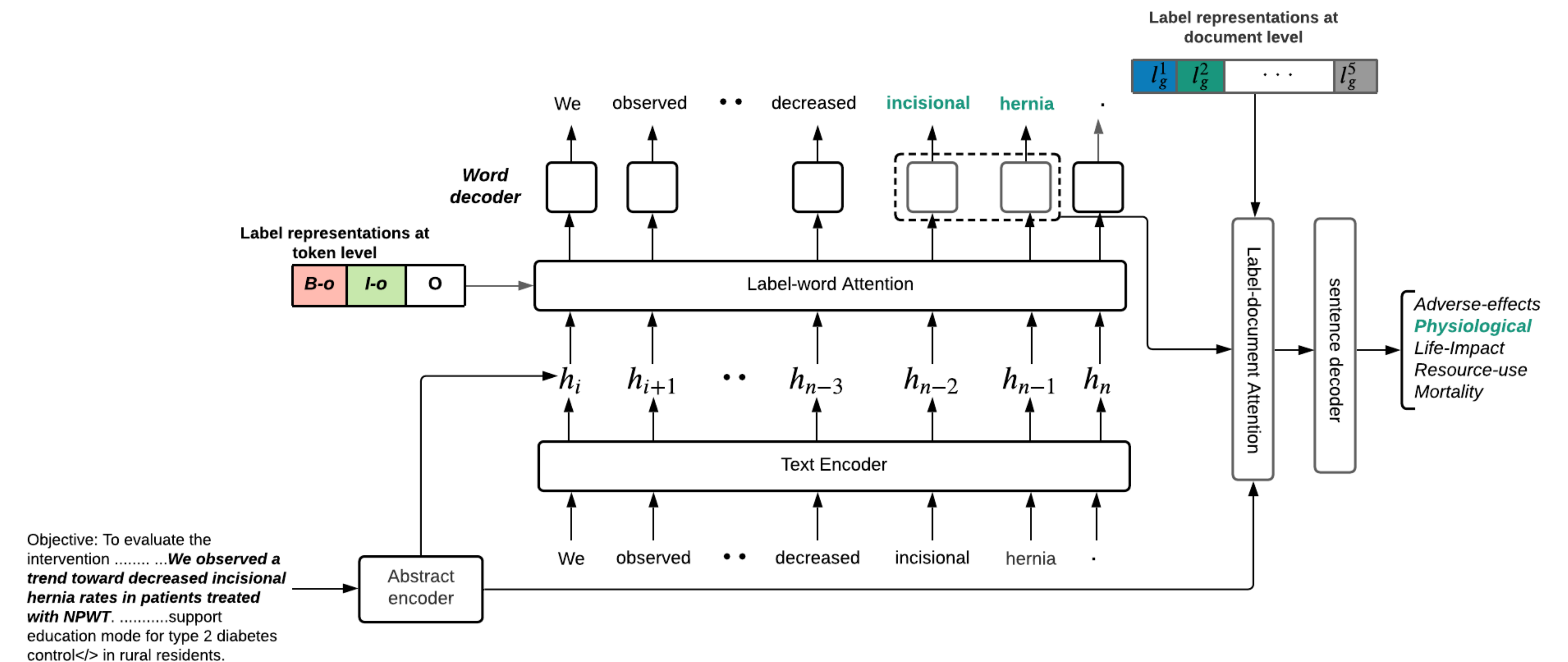
Ground truth:- Systolic blood pressure  
Predicted:- **Systolic** blood pressure

	P	R	F
Traditional NER evaluation:-	100	66.7	80.2
Full outcome phrase evaluation:-	0	0	0

Model	EBM-NLP <sub>rev</sub>	EBM-COMET
BioBERT [4]	<b>53.1</b>	<b>81.3</b>
BioELMo [5]	52.0	75.0
SciBERT [6]	52.8	77.6
ClinicalBERT [7]	51.0	68.5
BioFLAIR [8]	51.4	68.5

**Table 2:** F1 (%) for OD using in-domain CLMs

## JOINT OUTCOME SPAN DETECTION (OSD) & OUTCOME CLASSIFICATION (OC)



**Figure 2:** Label-word context aware attention framework (LCAM) for joint OSD and OC [9]

## OSD & OC

Given a sentence  $s = \{w_i\}_{i=1}^M$ , OSD identifies an outcome span  $o_d = \{b_i\}_{i=1}^N$ , and OC predicts an outcome type  $t(o_d) \in \mathcal{Y}$  for  $o_d$ , where,  $N \leq M$

$$h_n^c = \text{BioBERT}(w_n) + \frac{1}{|a|} \sum_{n=1}^{|a|} (\text{BioBERT}(w_n)) \quad (1)$$

Label-word attention representation (OSD)

$$A_n^{(1)} = \text{softmax}(\mathbf{W} \tanh(\mathbf{V}h_n)) \quad \& \quad A_n^{(2)} = \mathbf{U}h_n \quad (2)$$

$$E^{t_l} = A^{(1)}h_n^\top + A^{(2)}h_n^\top \quad (3)$$

$$L_{osd} = - \sum_{n=1}^N \sum_{i=1}^{|l_w|} y_{n,i} \log(\hat{y}_{n,i}). \quad (4)$$

Label-word attention representation (OC)

$$L_{oc} = - \sum_{i=1}^{|L_s|} y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \quad (5)$$

$$\text{Combined loss } L = L_{osd} + L_{oc} \quad (6)$$

Model	OSD	OC
LCAM	68.0	83.0
LCAM - Abstract	65.0	78.0
LCAM - Attention	58.0	71.0

**Table 3:** OSD and OC performance F1 (%) on EBM-COMET

## LABEL-ALIGNMENT FOR DATA AUGMENTATION

### Algorithm 1 Label Alignment

- 1: **Input:** comparable datasets  $\mathcal{S}$  &  $\mathcal{T}$
- 2: **for** each label  $l$  in  $\mathcal{S}$ :
- 3: Create an **embedding**  $l_s$  by  $l_s = \frac{1}{|l_s|} \sum_i^{l_s} O_{l_s}$
- 4: where  $O_{l_s} = \frac{1}{d} \sum_i^{i+(d-1)} \text{BioBERT}(w_i)$
- 5: and  $i$  &  $i + (d - 1)$  are the first and last words
- 6: of an outcome span labelled  $l_s$  i.e.  $O_{l_s}$
- 7: **for** each label  $l$  in  $\mathcal{T}$
- 8: Compute cosine\_similarity (cos) of  $l_s$  &  $l_t$
- 9: Reannotate  $l_t$  outcomes with most similar  $l_s$ .

## REFERENCES

- [1] Williamson et al. Comet handbook: version 1.0. 2017.
- [2] Nye et al. EBM-NLP corpus. *ACL*, 2018.
- [3] Abaho et al. Correcting crowdsourced outcome annotations. In *CEUR Workshop Proceedings*, 2019.
- [4] Lee et al. Biobert. *Bioinformatics*, 36(4):1234–1240, 2020.
- [5] Jin et al. Probing biomedical embeddings. *NAACL*, 2019.
- [6] Beltagy et al. Scibert. *EMNLP*, 2019.
- [7] Alsentzer et al. Clinicalbert. *NAACL*, June 2019.
- [8] Sharma and Daniel. Bioflair. *arXiv:1908.05760*, 2019.
- [9] Abaho et al. Joint span detection and classification for health outcomes. *arXiv preprint arXiv:2104.07789*, 2021.
- [10] Abaho et al. Assessment of contextualised representations in outcome detection. *Manuscript under review*, 2021.