

# What do you learn from context?

## Probing for sentence structure in contextualized word representations

Ian Tenney<sup>\*1</sup>, Patrick Xia<sup>2</sup>, Berlin Chen<sup>3</sup>, Alex Wang<sup>4</sup>, Adam Poliak<sup>2</sup>, R. Thomas McCoy<sup>2</sup>, Najoung Kim<sup>2</sup>, Benjamin Van Durme<sup>2</sup>, Samuel R. Bowman<sup>4</sup>, Dipanjan Das<sup>1</sup>, and Ellie Pavlick<sup>1,5</sup>

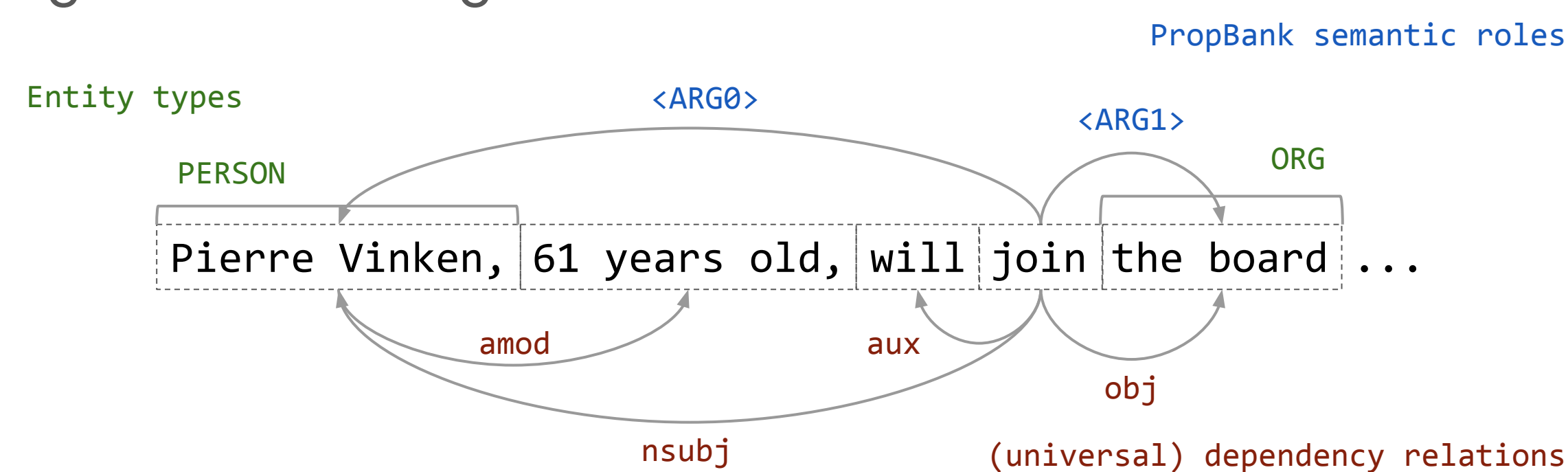
<sup>1</sup>Google AI Language, <sup>2</sup>Johns Hopkins University, <sup>3</sup>Swarthmore College, <sup>4</sup>New York University, <sup>5</sup>Brown University

### Overview

**Question:** Contextual encoders like ELMo and BERT outperform traditional word embeddings - why?

- What kind of linguistic information is modeled?
- Mostly syntax, mostly semantics, or both?
- Long range, or just local?
- Are there qualitative differences between encoders?

**Our approach:** use traditional intermediate tasks - tagging, parsing, etc. - to *probe* contextual representations for linguistic knowledge.



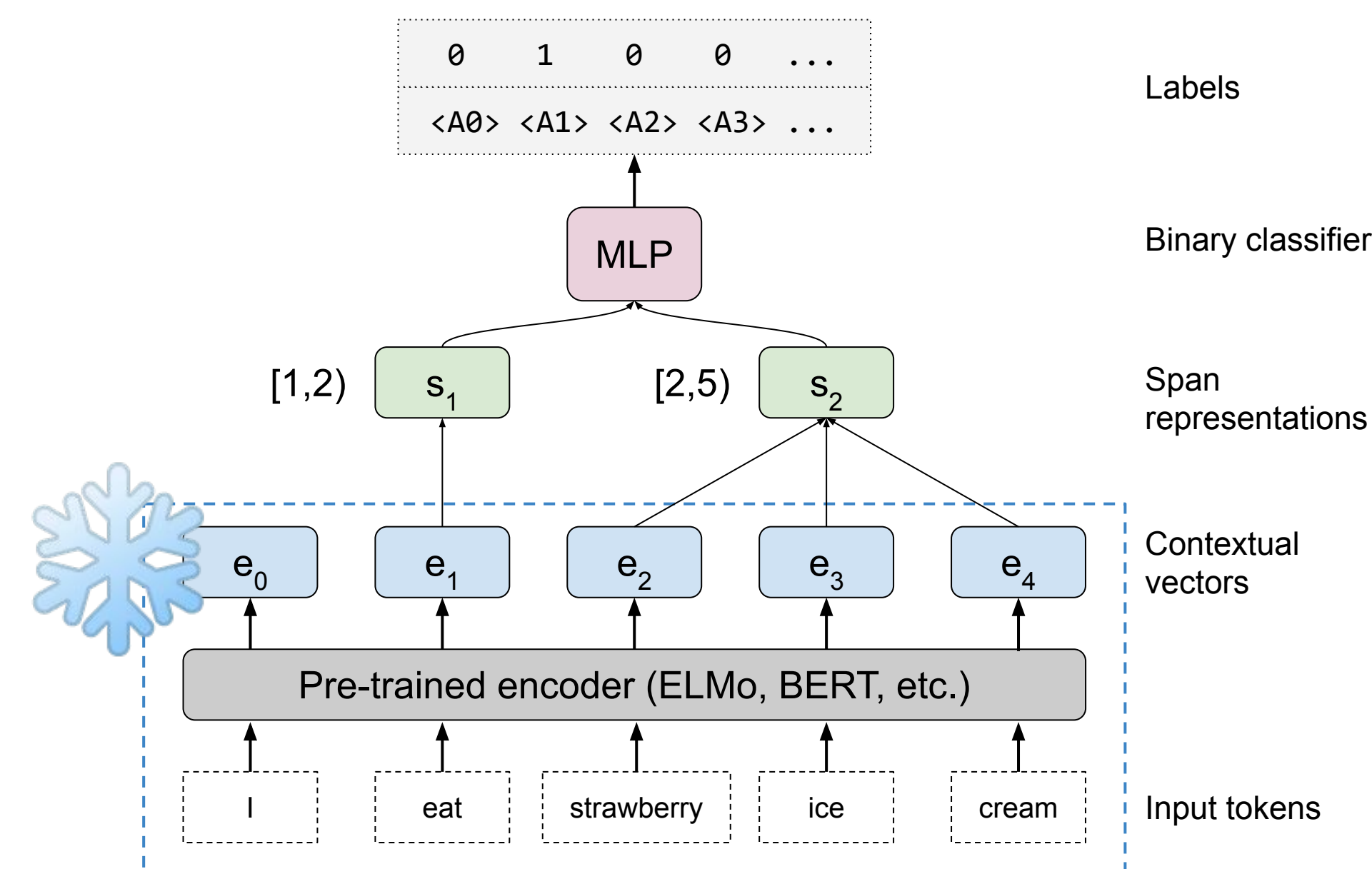
### Takeaways:

- Pretrained LM encoders (ELMo, BERT, GPT) give strong representations for syntax, but weaker for semantic phenomena.
- Encoders are able to capture long-range dependencies, and trained models outperform random baselines.

### Edge Probing

Multi-label classifier over pairs of spans. Given contextual vectors  $E = [e_0, e_1, \dots, e_n]$ , predict:

- **Unary:** label(s) for span1 =  $[i_1, j_1)$
- **Binary:** label(s) for (span1 =  $[i_1, j_1)$ , span2 =  $[i_2, j_2)$ )



Classifier sees *inside* of spans only - no other context. Must use encoder representations to infer role of phrase(s) in a sentence.

- Encoder weights are frozen - don't learn new features.
- ELMo-style scalar mixing across encoder layers.
- Common metric (Micro F1) for comparison across tasks.

### Results

Labeling versions of 10 tasks, ranging from syntax to semantics. Lexical baseline gives performance from word priors, independent of the rest of the sentence.

	ELMo			BERT (base)		BERT (large)	
	Lexical (char CNN)	Full context (mix all)	Abs. Δ vs. Lexical	Full context (mix all)	Abs. Δ vs. ELMo full	Full context (mix all)	Abs. Δ vs. BERT base
Part-of-Speech	90	97	6.3	97	0.0	97	0.2
Constituents	69	85	15.4	87	2.1	87	0.4
Dependencies	80	94	13.6	95	1.1	95	0.3
Entities	92	96	3.5	96	0.6	96	0.3
SRL (all)	74	90	16.0	91	1.2	92	1.0
Core	74	93	19.0	94	1.0	95	1.0
Non-core	75	84	8.8	86	1.8	87	1.0
OntoNotes coref.	75	84	8.7	90	6.3	91	1.2
SPR1	80	85	4.7	86	1.3	86	-0.3
SPR2	82	83	1.0	84	0.7	84	0.3
Winograd coref.	54	54	-0.8	55	1.4	61 ± 6	6.5
Relations	58	78	22.1	82	4.2	82	0.5

- ELMo is great at syntactic tasks, less so at semantics...
- BERT improves on ELMo by around 20% (relative) on most tasks, over 40% on coreference
- GPT is approximately on-par with ELMo; CoVe somewhat lower.

### Probing Data

Task	Example target	
Part-of-Speech	The important thing about Disney is that it is a global [brand] <sub>1</sub> .	NN (noun)
Constituents	The important thing about Disney is that it [is a global brand] <sub>1</sub> .	VP (verb phrase)
Dependencies (UD EWT)	The important [thing] <sub>2</sub> about Disney [is] <sub>1</sub> that it is a global brand.	nsubj (nominal subject)
Entities	The important thing about [Disney] <sub>1</sub> is that it is a global brand.	Organization
Semantic roles (SRL)	[The important thing about Disney] <sub>2</sub> [is] <sub>1</sub> that it is a global brand.	ARG1 (agent)
Semantic proto-roles (SPR)	[It] <sub>1</sub> [endorsed] <sub>2</sub> the White House strategy ...	{awareness, existed_after, ...}
OntoNotes coreference	The important thing about [Disney] <sub>1</sub> is that [it] <sub>2</sub> is a global brand.	True
Winograd coreference (DPR)	[Characters] <sub>2</sub> entertain audiences because [they] <sub>1</sub> want people to be happy. Characters entertain [audiences] <sub>2</sub> because [they] <sub>1</sub> want people to be happy.	True False
SemEval relations	The [burst] <sub>1</sub> has been caused by water hammer [pressure] <sub>2</sub> .	Cause-Effect (e <sub>2</sub> , e <sub>1</sub> )

### Ablations on ELMo

How much is just due to architecture? (versus training)

- ELMo outperforms random-ELMo (orthonormal LSTM weights)

What about local features?

- ELMo outperforms lexical baseline + local CNN ( $\pm 1$  or  $\pm 2$  tokens)
- Real ELMo performance drops only slowly on long-range edges

