

**BERT**



# BERT

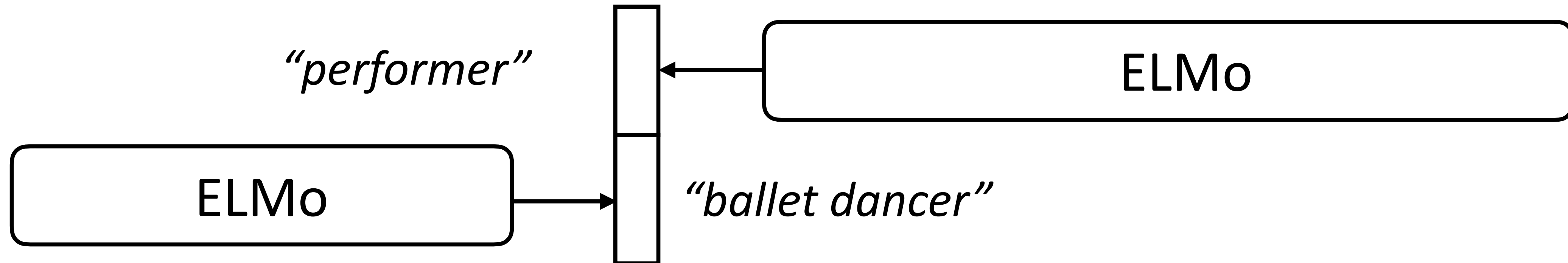
---

- ▶ AI2 made ELMo in spring 2018, GPT (transformer-based ELMo) was released in summer 2018, BERT came out October 2018
- ▶ Four major changes compared to ELMo:
  - ▶ Transformers instead of LSTMs
  - ▶ Bidirectional model with “Masked LM” objective instead of standard LM
  - ▶ Fine-tune instead of freeze at test time
  - ▶ Operates over word pieces (byte pair encoding)

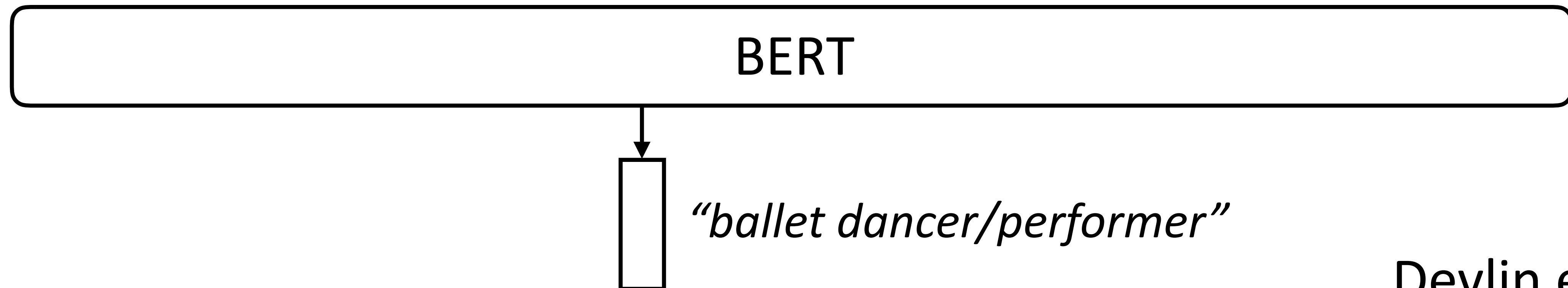


# BERT

- ▶ ELMo is a unidirectional model (as is GPT): we can concatenate two unidirectional models, but is this the right thing to do?
- ▶ ELMo reprs look at each direction in isolation; BERT looks at them jointly



*A stunning ballet dancer, Copeland is one of the best performers to see live.*

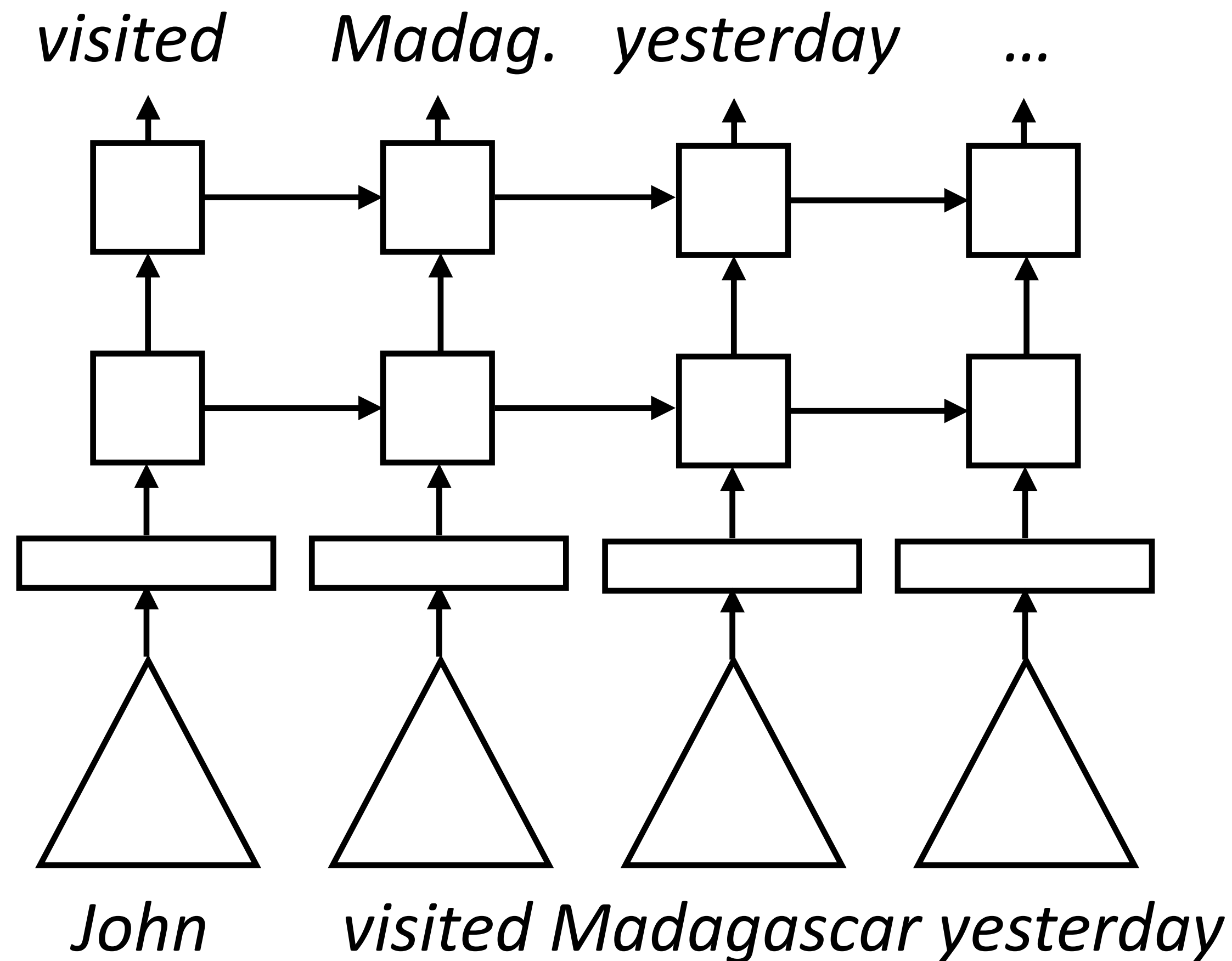




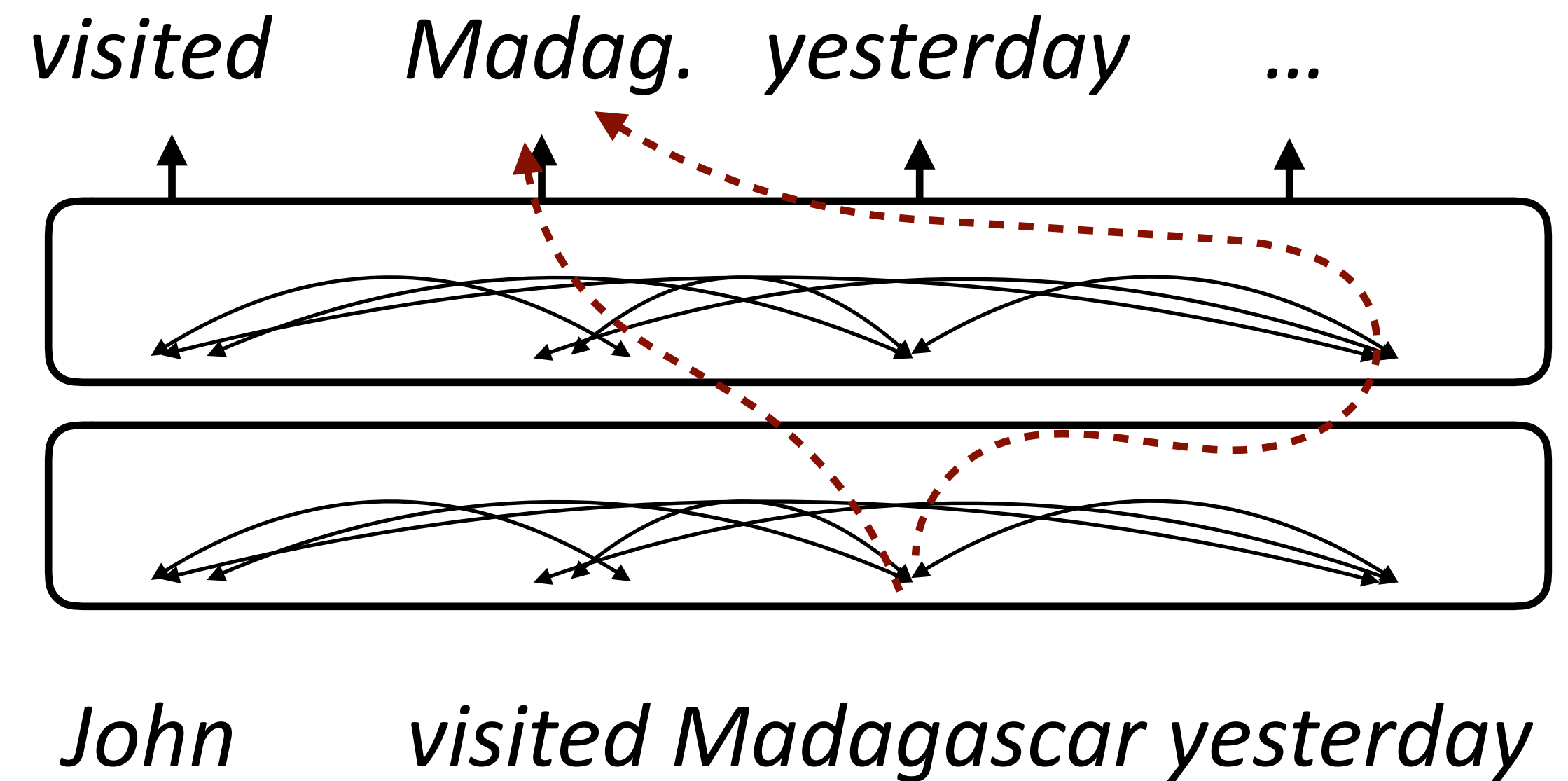
# BERT

- ▶ How to learn a “deeply bidirectional” model? What happens if we just replace an LSTM with a transformer?

## ELMo (Language Modeling)



## BERT

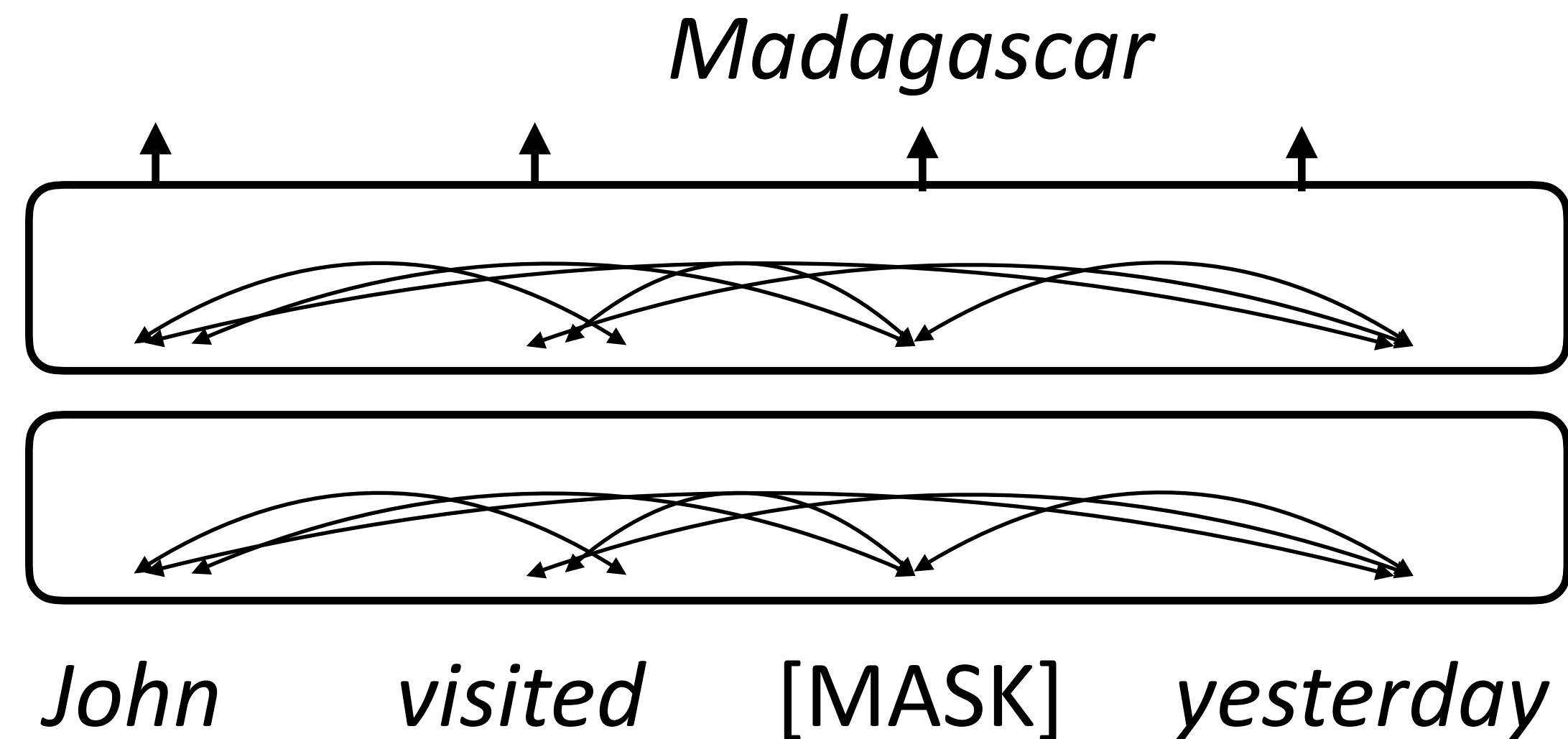


- ▶ You could do this with a “one-sided” transformer, but this “two-sided” model can cheat



# Masked Language Modeling

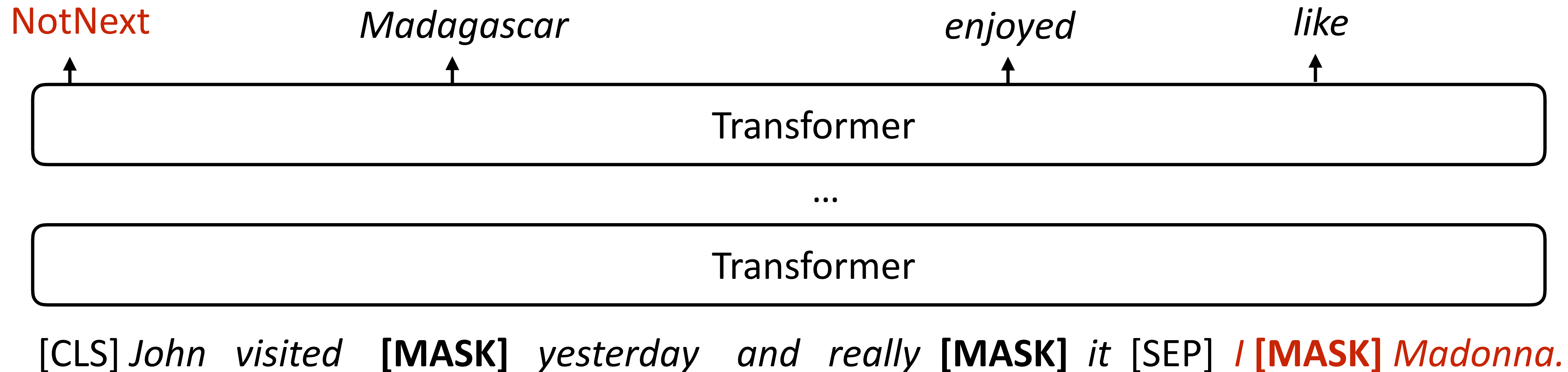
- ▶ How to prevent cheating? Next word prediction fundamentally doesn't work for bidirectional models, instead do *masked language modeling*
- ▶ BERT formula: take a chunk of text, mask out 15% of the tokens, and try to predict them





# Next “Sentence” Prediction

- ▶ Input: [CLS] Text chunk 1 [SEP] Text chunk 2
- ▶ 50% of the time, take the true next chunk of text, 50% of the time take a random other chunk. Predict whether the next chunk is the “true” next
- ▶ BERT objective: masked LM + next sentence prediction

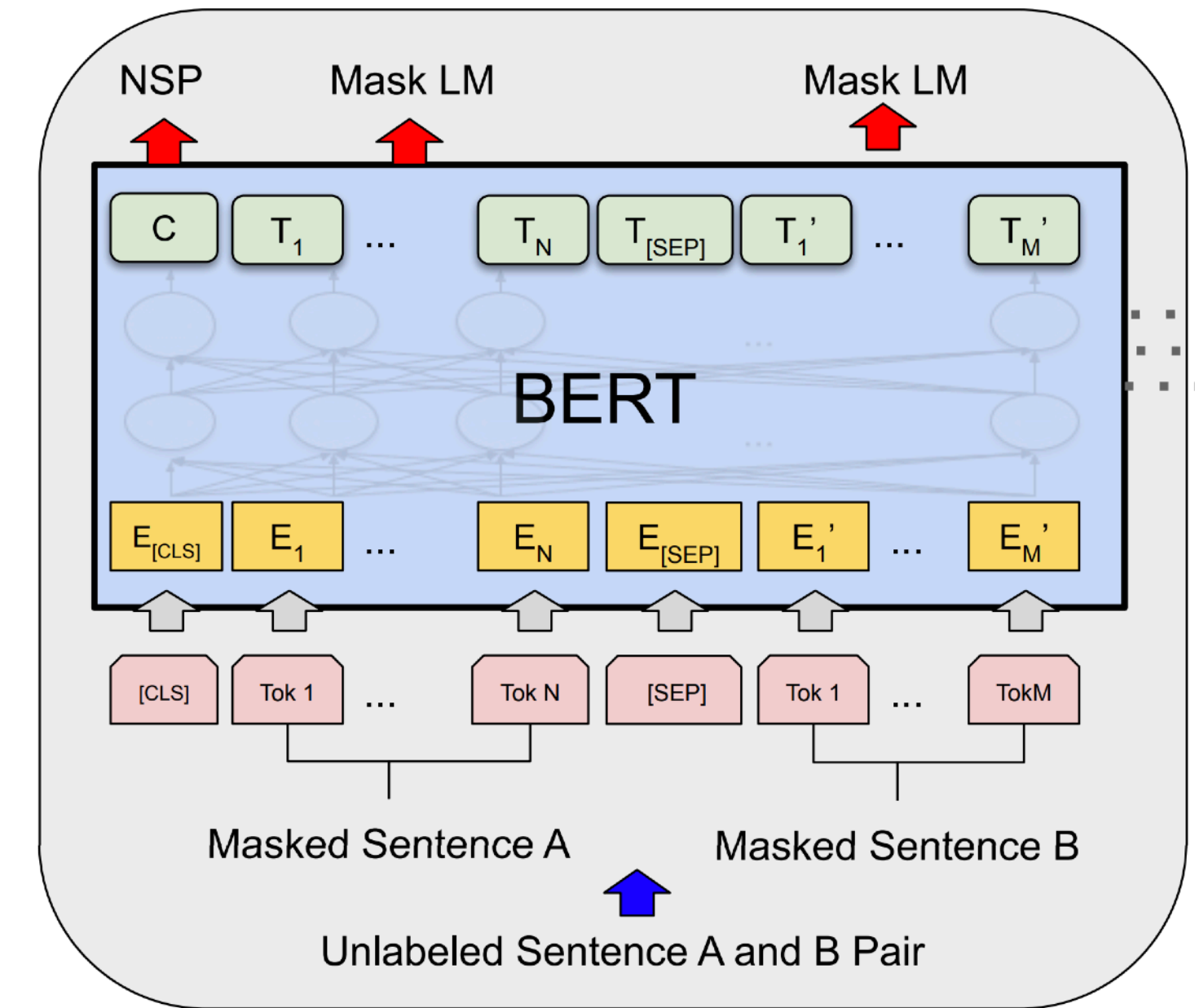






# BERT Architecture

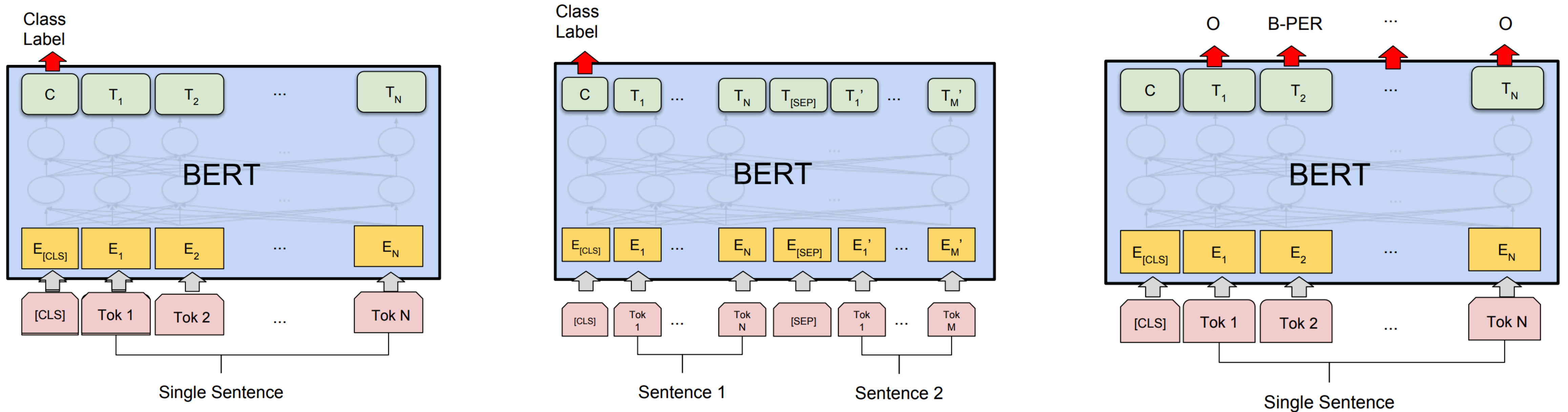
- ▶ BERT Base: 12 layers, 768-dim per wordpiece token, 12 heads. Total params = 110M
- ▶ BERT Large: 24 layers, 1024-dim per wordpiece token, 16 heads. Total params = 340M
- ▶ Positional embeddings and segment embeddings, 30k word pieces
- ▶ This is the model that gets **pre-trained** on a large corpus



Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	E <sub>[CLS]</sub>	E <sub>my</sub>	E <sub>dog</sub>	E <sub>is</sub>	E <sub>cute</sub>	E <sub>[SEP]</sub>	E <sub>he</sub>	E <sub>likes</sub>	E <sub>play</sub>	E <sub>##ing</sub>	E <sub>[SEP]</sub>
Segment Embeddings	E <sub>A</sub>	E <sub>A</sub>	E <sub>A</sub>	E <sub>A</sub>	E <sub>A</sub>	E <sub>A</sub>	E <sub>B</sub>	E <sub>B</sub>	E <sub>B</sub>	E <sub>B</sub>	E <sub>B</sub>
Position Embeddings	E <sub>0</sub>	E <sub>1</sub>	E <sub>2</sub>	E <sub>3</sub>	E <sub>4</sub>	E <sub>5</sub>	E <sub>6</sub>	E <sub>7</sub>	E <sub>8</sub>	E <sub>9</sub>	E <sub>10</sub>



# What can BERT do?



(b) Single Sentence Classification Tasks:  
SST-2, CoLA

(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG

(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

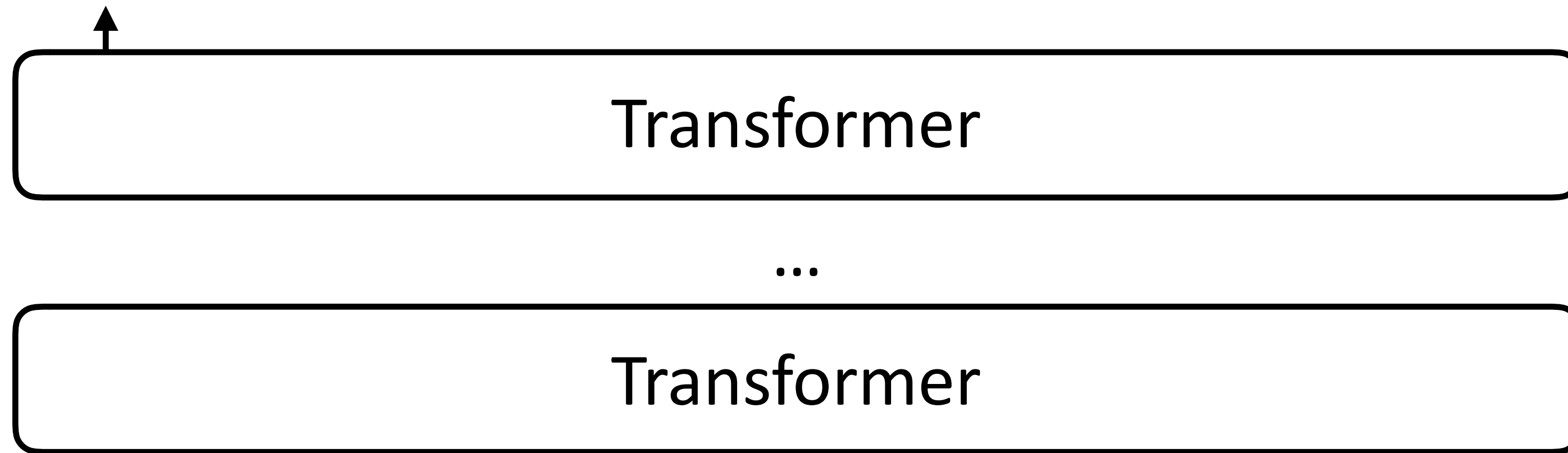
- ▶ Artificial [CLS] token is used as the vector to do classification from
  - ▶ Sentence pair tasks (entailment): feed both sentences into BERT
  - ▶ BERT can also do tagging by predicting tags at each word piece
- Devlin et al. (2019)



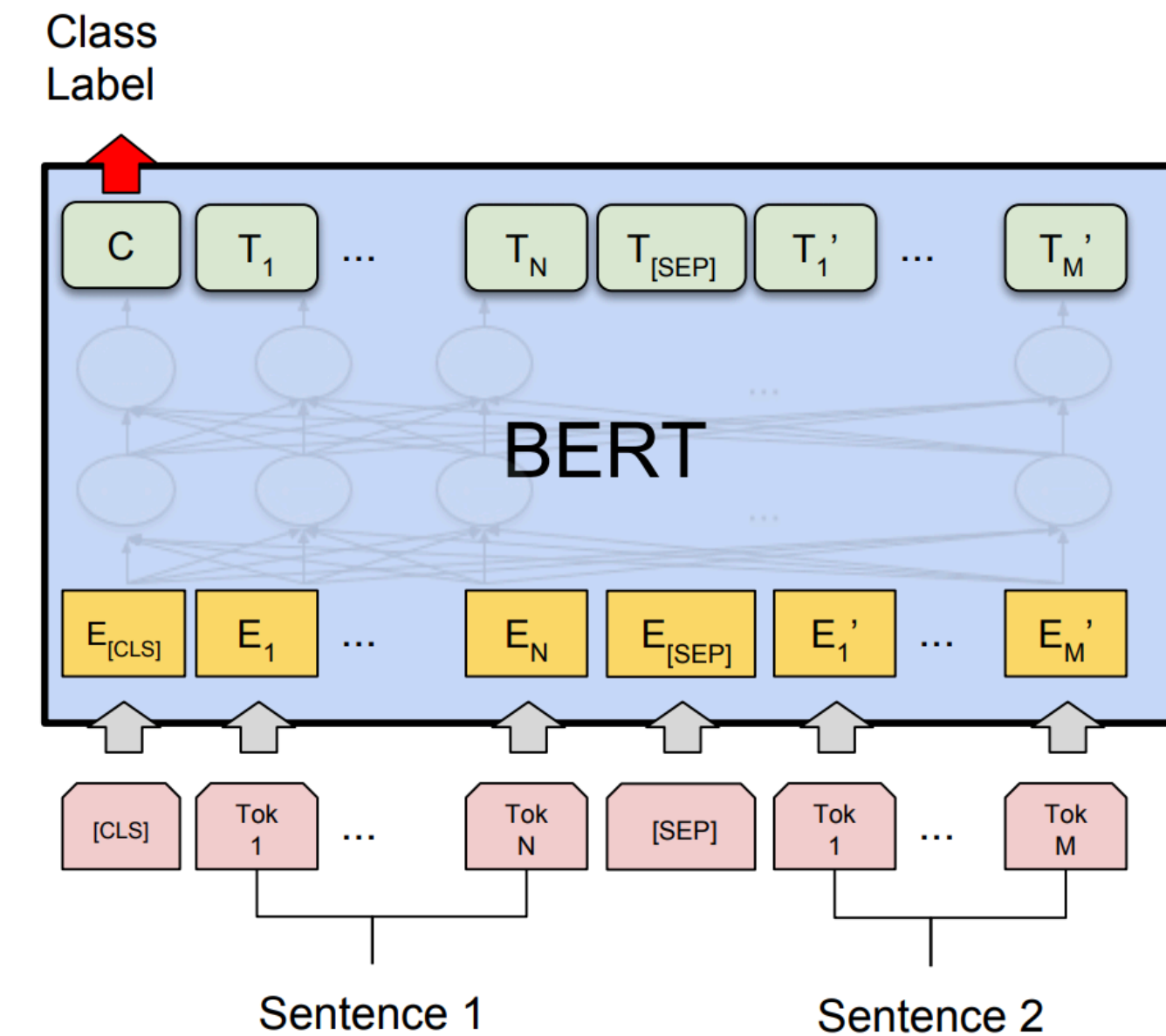


# What can BERT do?

Entails (first sentence implies second is true)



[CLS] A boy plays in the snow [SEP] A boy is outside



(a) Sentence Pair Classification Tasks: MNL, QQP, QNLI, STS-B, MRPC, RTE, SWAG

- ▶ How does BERT model this sentence pair stuff?
- ▶ Transformers can capture interactions between the two sentences, even though the NSP objective doesn't really cause this to happen



# What can BERT NOT do?

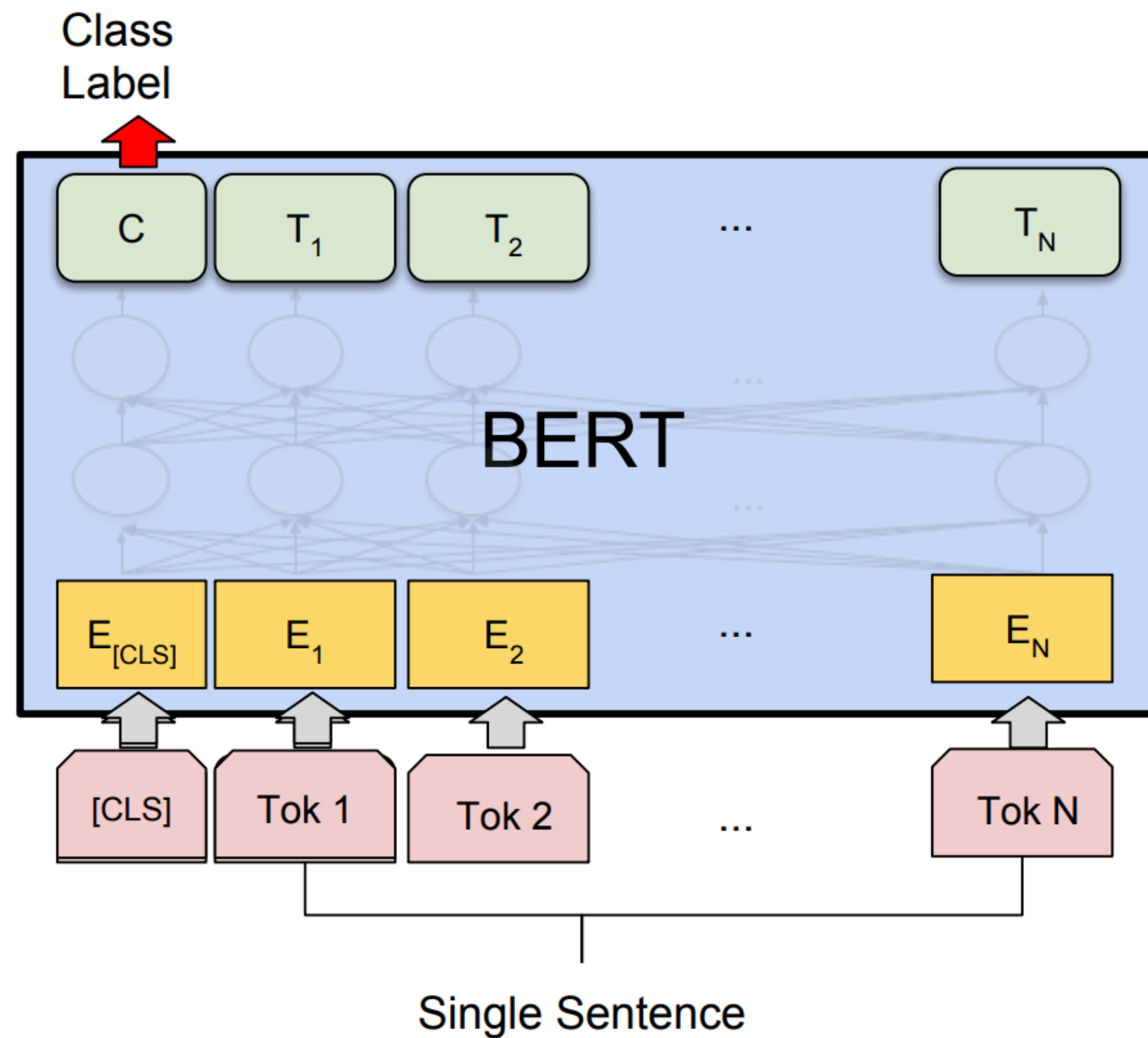
---

- ▶ BERT **cannot** generate text (at least not in an obvious way)
  - ▶ Can fill in MASK tokens, but can't generate left-to-right (well, you could put MASK at the end repeatedly, but this is slow)
- ▶ Masked language models are intended to be used primarily for “analysis” tasks



# Fine-tuning BERT

- ▶ Fine-tune for 1-3 epochs, batch size 2-32, learning rate  $2e-5$  -  $5e-5$



(b) Single Sentence Classification Tasks:  
SST-2, CoLA

- ▶ Large changes to weights up here (particularly in last layer to route the right information to [CLS])
- ▶ Smaller changes to weights lower down in the transformer
- ▶ Small LR and short fine-tuning schedule mean weights don't change much
- ▶ More complex “triangular learning rate” schemes exist



# Fine-tuning BERT

Pretraining	Adaptation	NER	SA	Nat. lang. inference		Semantic textual similarity		
		CoNLL 2003	SST-2	MNLI	SICK-E	SICK-R	MRPC	STS-B
Skip-thoughts	❄️	-	81.8	62.9	-	86.6	75.8	71.8
ELMo	❄️	91.7	<b>91.8</b>	<b>79.6</b>	<b>86.3</b>	<b>86.1</b>	<b>76.0</b>	<b>75.9</b>
	🔥	<b>91.9</b>	91.2	76.4	83.3	83.3	74.7	75.5
	$\Delta = \text{🔥} - \text{❄️}$	0.2	-0.6	<b>-3.2</b>	<b>-3.3</b>	<b>-2.8</b>	<b>-1.3</b>	-0.4
BERT-base	❄️	92.2	93.0	<b>84.6</b>	84.8	86.4	78.1	82.9
	🔥	<b>92.4</b>	<b>93.5</b>	<b>84.6</b>	<b>85.8</b>	<b>88.7</b>	<b>84.8</b>	<b>87.1</b>
	$\Delta = \text{🔥} - \text{❄️}$	0.2	0.5	0.0	1.0	<b>2.3</b>	<b>6.7</b>	<b>4.2</b>

- ▶ BERT is typically better if the whole network is fine-tuned, unlike ELMo





# Evaluation: GLUE

Corpus	Train	Test	Task	Metrics	Domain
Single-Sentence Tasks					
CoLA	8.5k	<b>1k</b>	acceptability	Matthews corr.	misc.
SST-2	67k	1.8k	sentiment	acc.	movie reviews
Similarity and Paraphrase Tasks					
MRPC	3.7k	1.7k	paraphrase	acc./F1	news
STS-B	7k	1.4k	sentence similarity	Pearson/Spearman corr.	misc.
QQP	364k	<b>391k</b>	paraphrase	acc./F1	social QA questions
Inference Tasks					
MNLI	393k	<b>20k</b>	NLI	matched acc./mismatched acc.	misc.
QNLI	105k	5.4k	QA/NLI	acc.	Wikipedia
RTE	2.5k	3k	NLI	acc.	news, Wikipedia
WNLI	634	<b>146</b>	coreference/NLI	acc.	fiction books





# Results

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>91.1</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>81.9</b>

- ▶ Huge improvements over prior work (even compared to ELMo)
- ▶ Effective at “sentence pair” tasks: textual entailment (does sentence A imply sentence B), paraphrase detection

Devlin et al. (2018)



# RoBERTa

- ▶ “Robustly optimized BERT”
- ▶ 160GB of data instead of 16 GB
- ▶ Dynamic masking: standard BERT uses the same MASK scheme for every epoch, RoBERTa recomputes them
- ▶ New training + more data = better performance

Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
RoBERTa						
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1
+ pretrain even longer	160GB	8K	500K	<b>94.6/89.4</b>	<b>90.2</b>	<b>96.4</b>
BERT <sub>LARGE</sub>						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7





# Using BERT

▶ Huggingface Transformers: big open-source library with most pre-trained architectures implemented, weights available

▶ Lots of standard models...

Model architectures

🤗 Transformers currently provides the following NLU/NLG architectures:

1. **BERT** (from Google) released with the paper [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#) by Jacob Devlin, Ming-Wei Chang, Kenton Lee and Kristina Toutanova
2. **GPT** (from OpenAI) released with the paper [Improving Language Understanding by Generative Pre-Training](#) by Radford, Karthik Narasimhan, Tim Salimans and Ilya Sutskever.
3. **GPT-2** (from OpenAI) released with the paper [Language Models are Unsupervised Multitask Learners](#) by Jeffrey Wu\*, Rewon Child, David Luan, Dario Amodei\*\* and Ilya Sutskever
4. **Transformer-XL** (from Google/CMU) released with the paper [Transformer-XL: Fixed-Length Context](#) by Zihang Dai\*, Zhilin Yang\*, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Noam Shazeer
5. **XLNet** (from Google/CMU) released with the paper [XLNet: Generalized Autoregressive and Causal Language Modeling](#) by Zhilin Yang\*, Zihang Dai\*, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Noam Shazeer
6. **XLM** (from Facebook) released together with the paper [Cross-lingual Language Model Pre-training](#) by Lample, Guillaume, Alexis Conneau, and Alexis Conneau.
7. **RoBERTa** (from Facebook), released together with the paper [RoBERTa: A Robustly Optimized BERT Base](#)

...

and “community models”

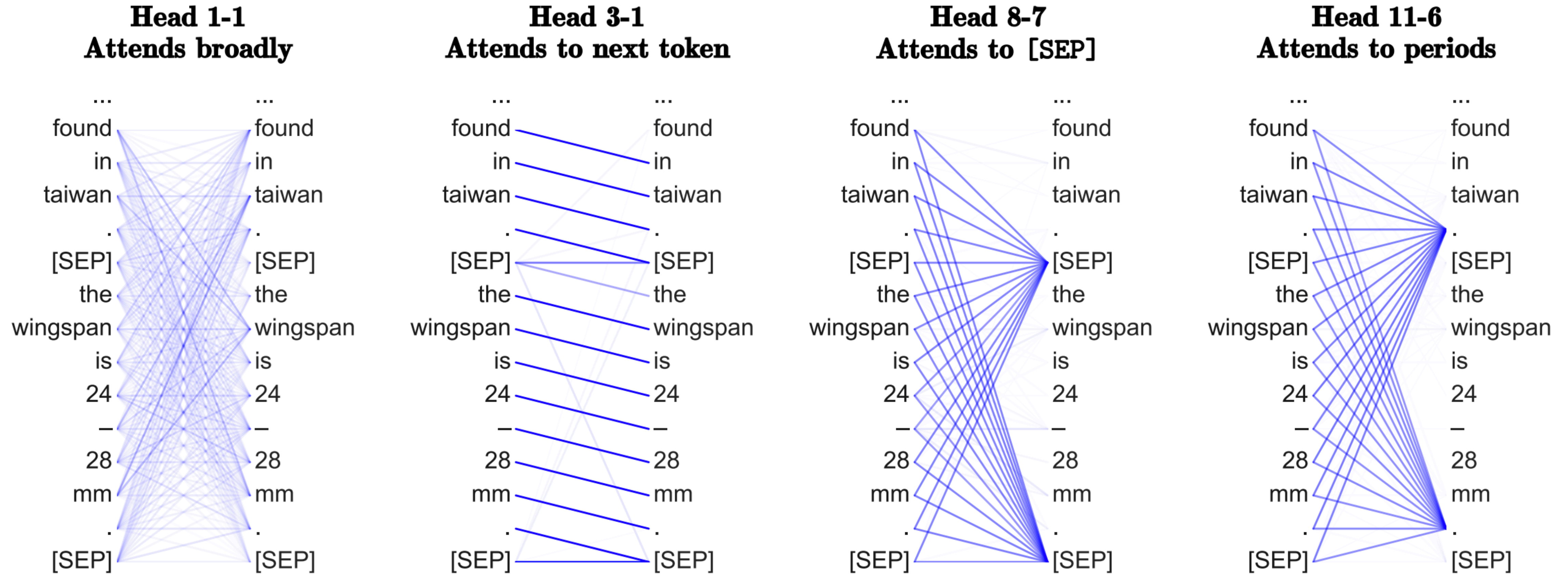
- [mrm8488/spanbert-large-finetuned-tacred](#) ★
- [mrm8488/xlm-multi-finetuned-xquadv1](#) ★
- [nlpaueb/bert-base-greek-uncased-v1](#) ★
- [nlptown/bert-base-multilingual-uncased-sentiment](#) ★
- [patrickvonplaten/reformer-crime-and-punish](#) ★
- [redewiedergabe/bert-base-historical-german-rw-cased](#) ★
- [roberta-base](#) ★
- [severinsimmler/literary-german-bert](#) ★
- [seyonec/ChemBERTa-zinc-base-v1](#) ★

...





# What does BERT learn?



- ▶ Heads on transformers learn interesting and diverse things: content heads (attend based on content), positional heads (based on position), etc.

Clark et al. (2019)

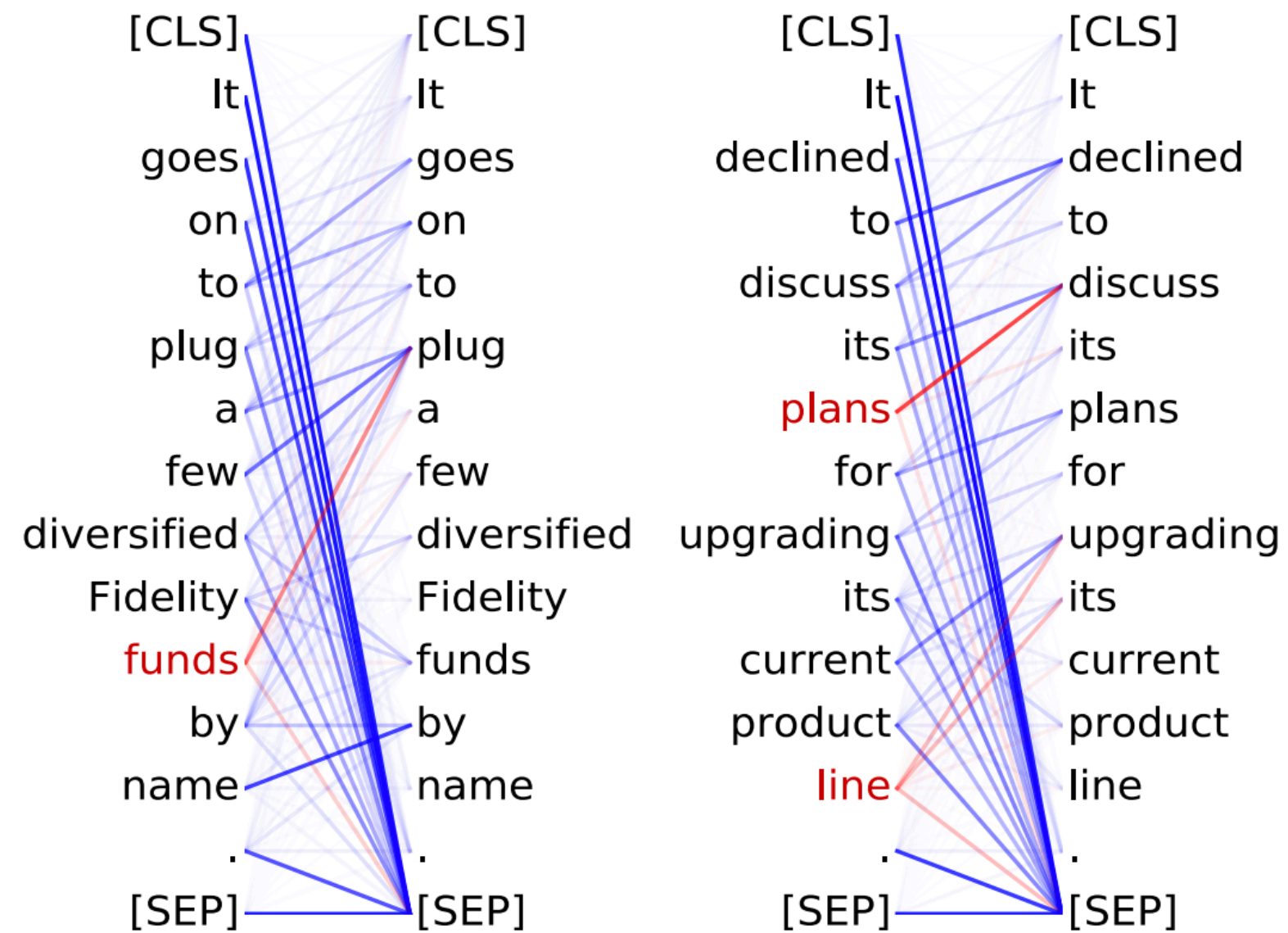




# What does BERT learn?

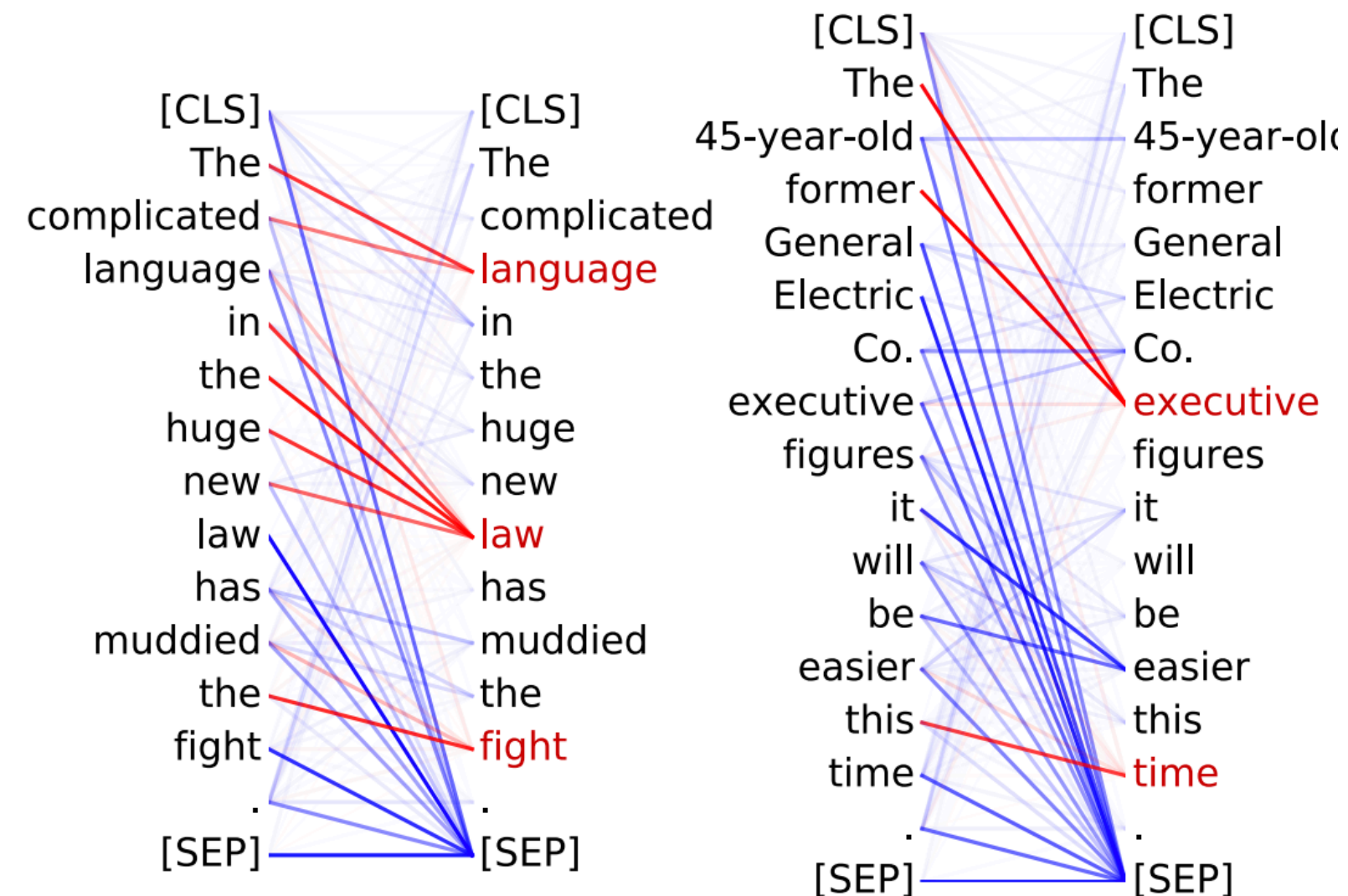
## Head 8-10

- **Direct objects** attend to their verbs
- 86.8% accuracy at the dobj relation



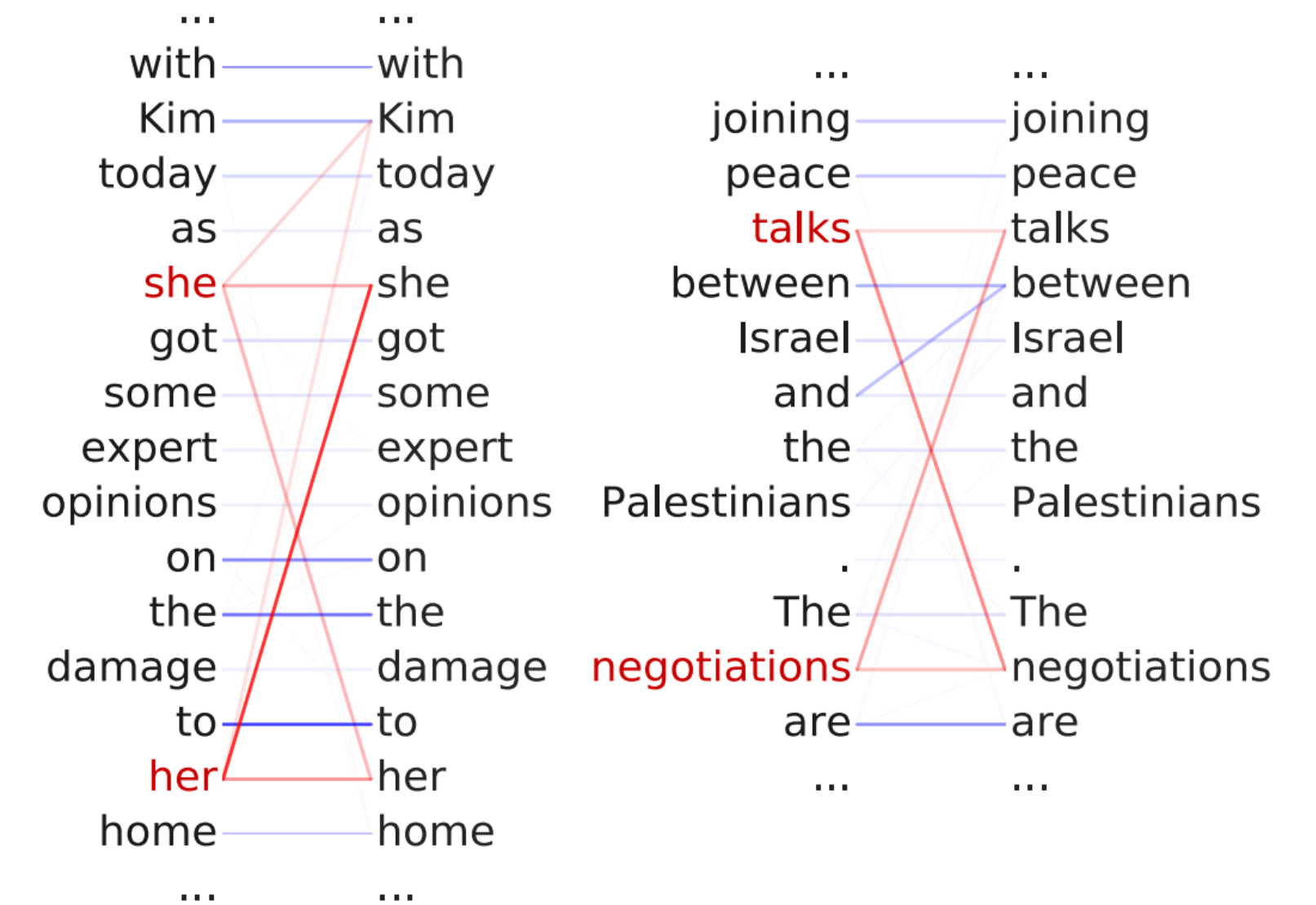
## Head 8-11

- **Noun modifiers** (e.g., determiners) attend to their noun
- 94.3% accuracy at the det relation



## Head 5-4

- **Coreferent** mentions attend to their antecedents
- 65.1% accuracy at linking the head of a coreferent mention to the head of an antecedent



- ▶ Still way worse than what supervised systems can do, but interesting that this is learned organically