

# CS378: Natural Language Processing

## Lecture 20: Pre-training, BERT

Greg Durrett



TEXAS

The University of Texas at Austin





# Announcements

---

- ▶ A5 due Tuesday
- ▶ Final project out Tuesday

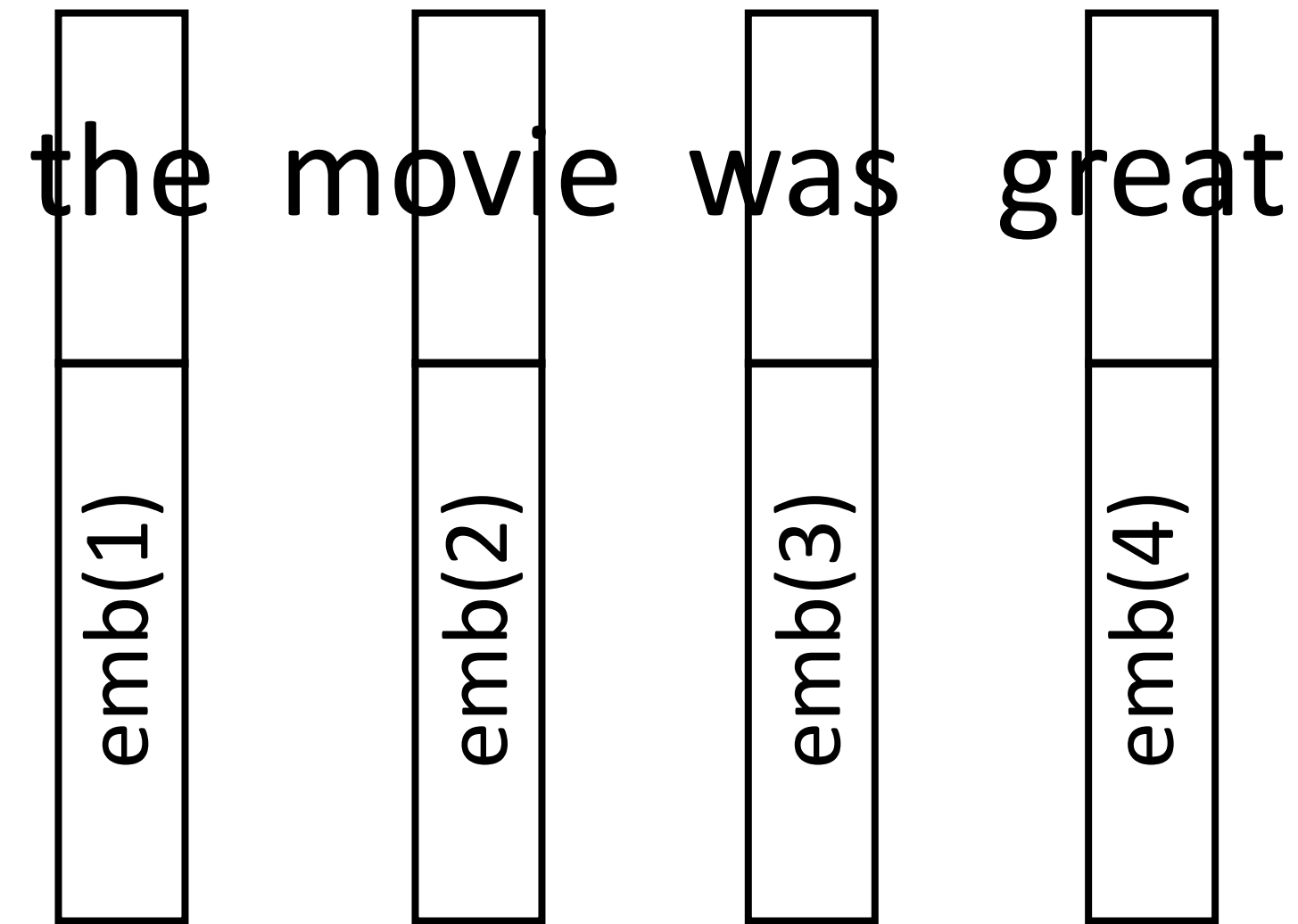
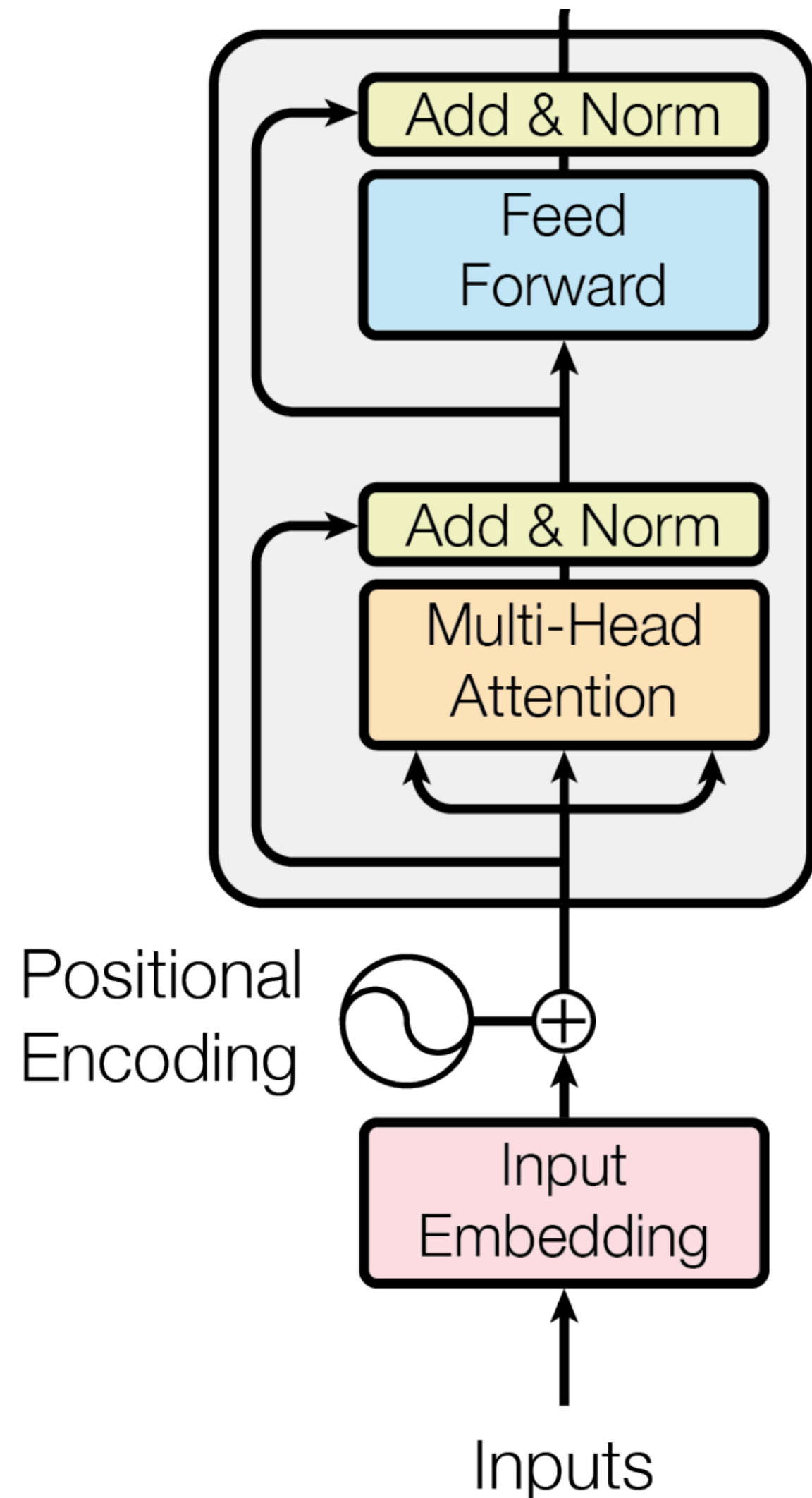


# Recap: Self Attention

---



# Recap: Transformers



- ▶ Augment word embedding with position embeddings of some kind
- ▶ Position embeddings, then repeated layers of multi-head attention and feedforward



# Today

---

- ▶ ELMo
- ▶ BERT
- ▶ BERT results
- ▶ Applying BERT

ELMo



# What is pre-training?

---

- ▶ “Pre-train” a model on a large dataset for task X, then “fine-tune” it on a dataset for task Y
- ▶ Key idea: X is somewhat related to Y, so a model that can do X will have some good neural representations for Y as well
- ▶ ImageNet pre-training is huge in computer vision: learn generic visual features for recognizing objects
- ▶ GloVe can be seen as pre-training: learn vectors with the skip-gram objective on large data (task X), then fine-tune them as part of a neural network for sentiment/any other task (task Y)



# GloVe is insufficient

---

- ▶ GloVe uses a lot of data but in a weak way
- ▶ Having a single embedding for each word is wrong

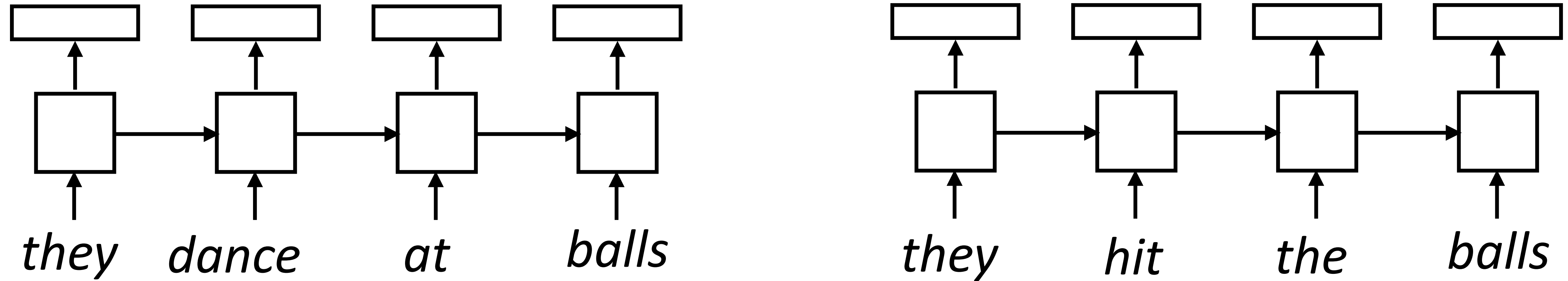
*they dance at balls*      *they hit the balls*

- ▶ Identifying discrete word senses is hard, doesn't scale. Hard to identify how many senses each word has
- ▶ How can we make our word embeddings more *context-dependent*?





# Context-dependent Embeddings

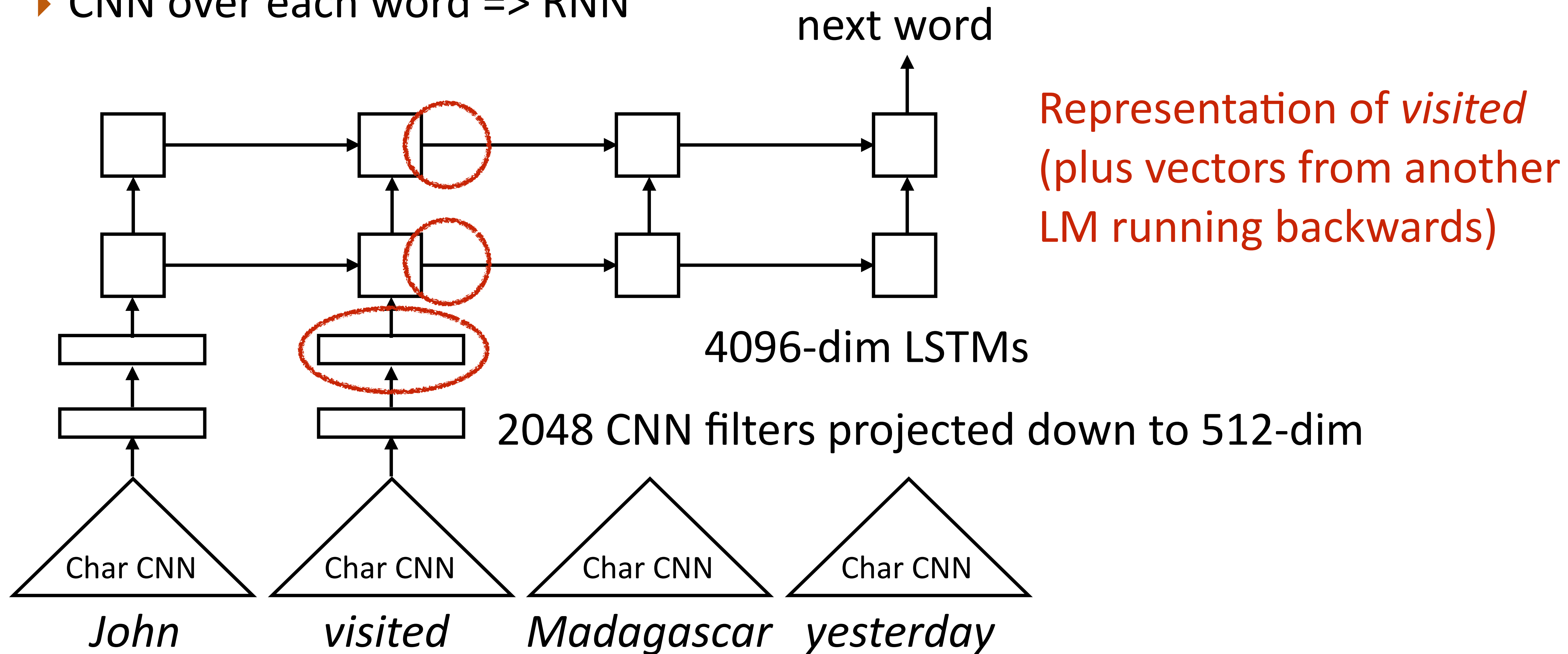


- ▶ Train a neural language model to predict the next word given previous words in the sentence, use the hidden states (output) at each step *as word embeddings*
- ▶ This is the key idea behind ELMo: language models can allow us to form useful word representations in the same way word2vec did



# ELMo

- CNN over each word => RNN





# ELMo

---

- ▶ Use the embeddings as a drop-in replacement for GloVe
- ▶ Huge gains across many high-profile tasks: NER, question answering, semantic role labeling (similar to parsing), etc.
- ▶ But what if the pre-training **isn't only the embeddings?**

BERT



# BERT

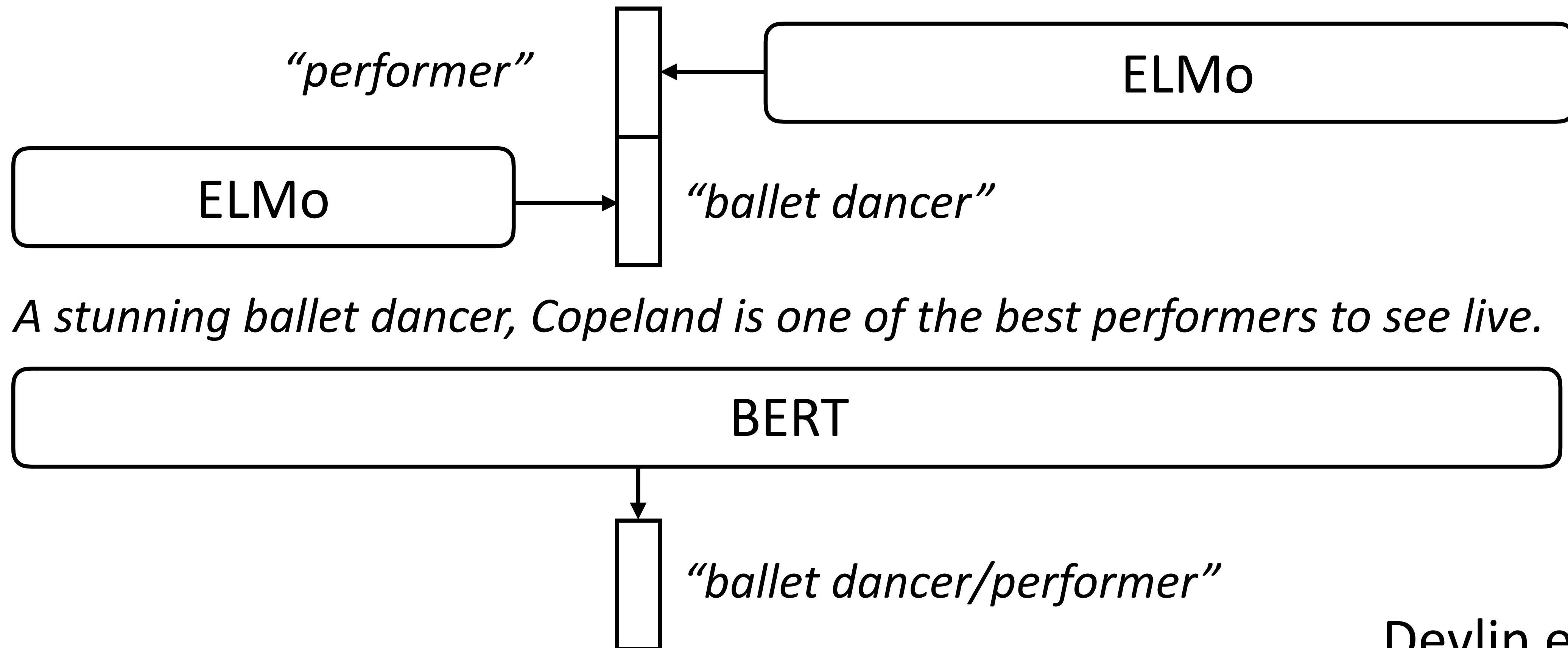
---

- ▶ 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

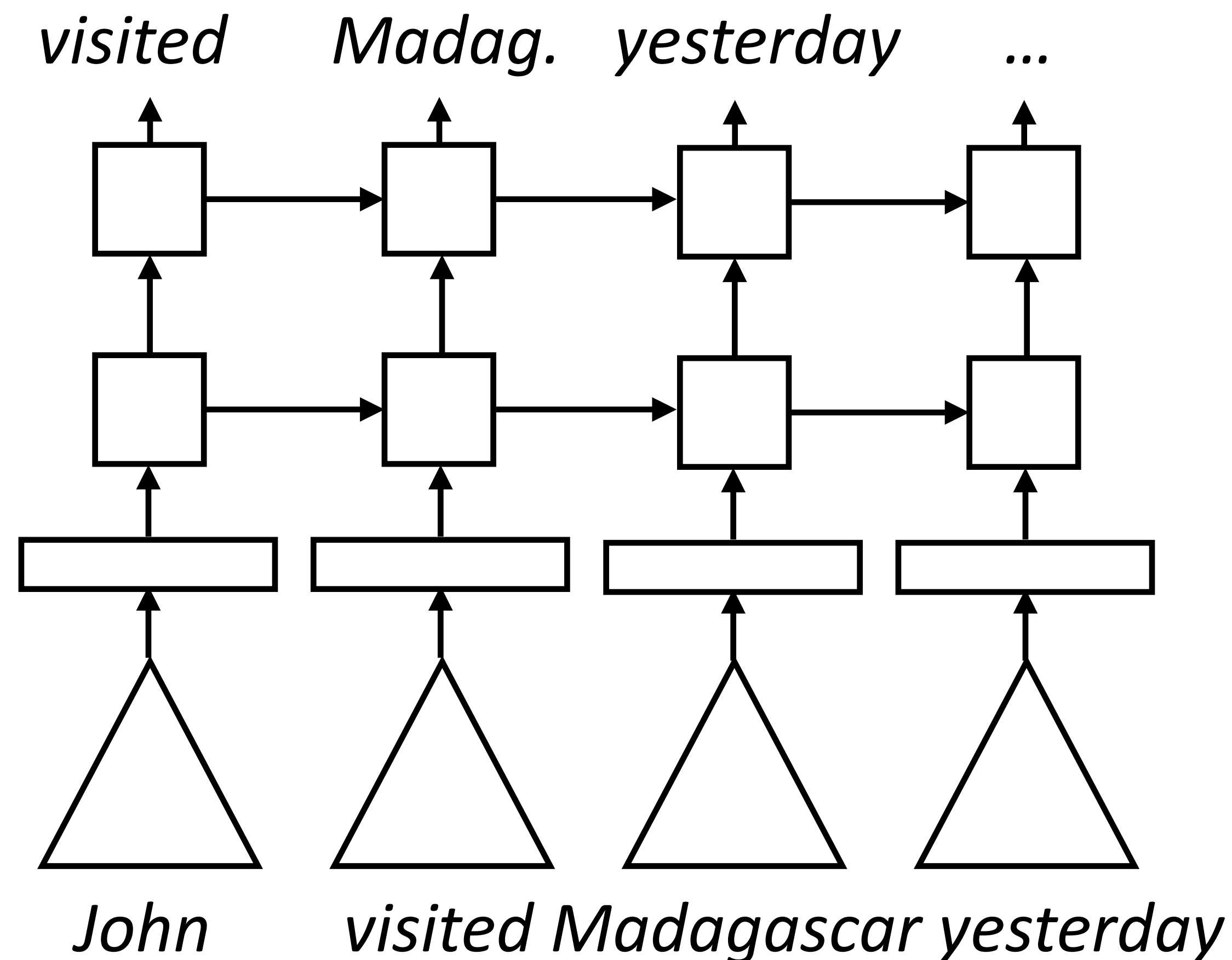




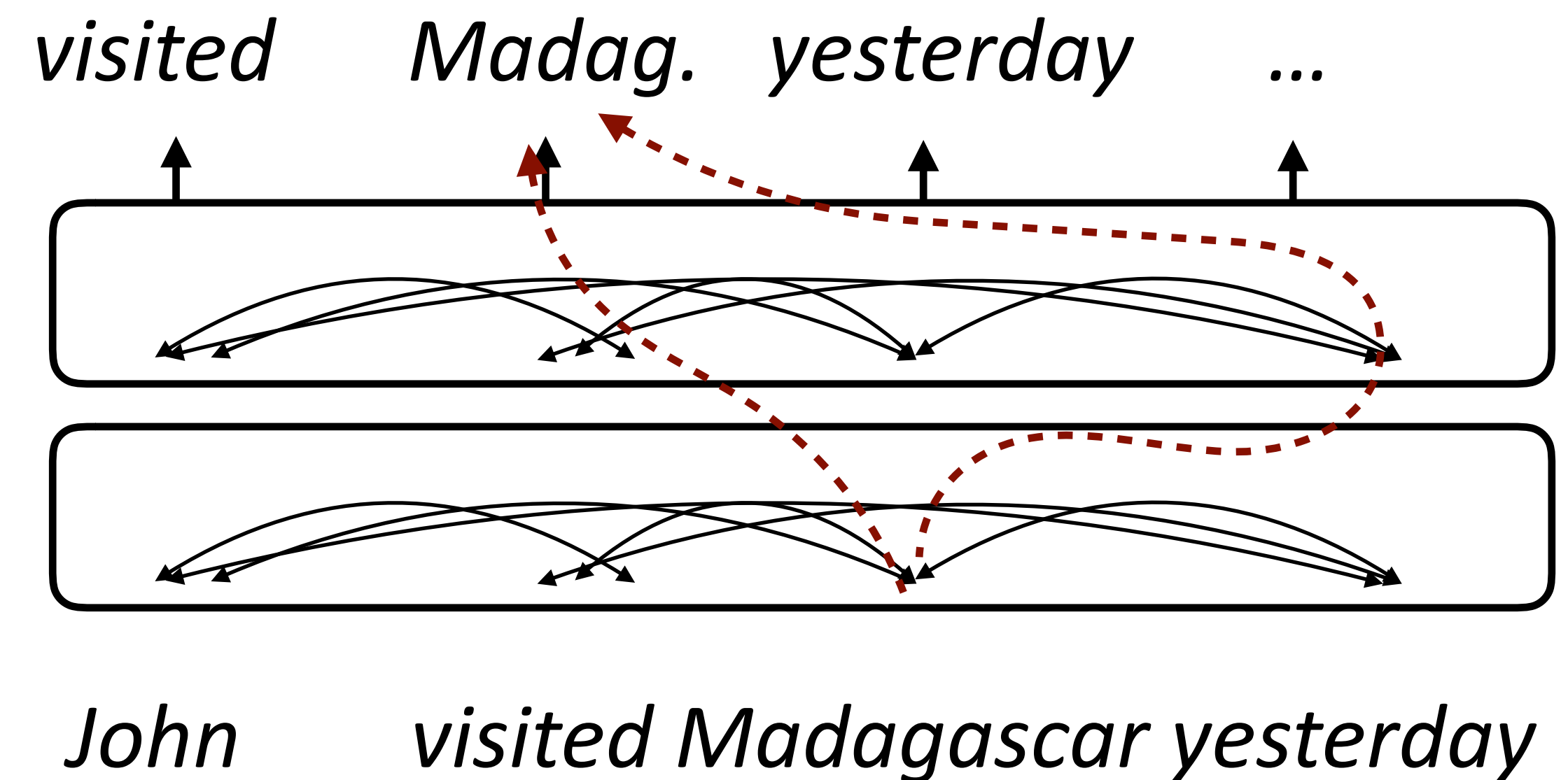
# BERT

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

## ELMo (Language Modeling)



## BERT



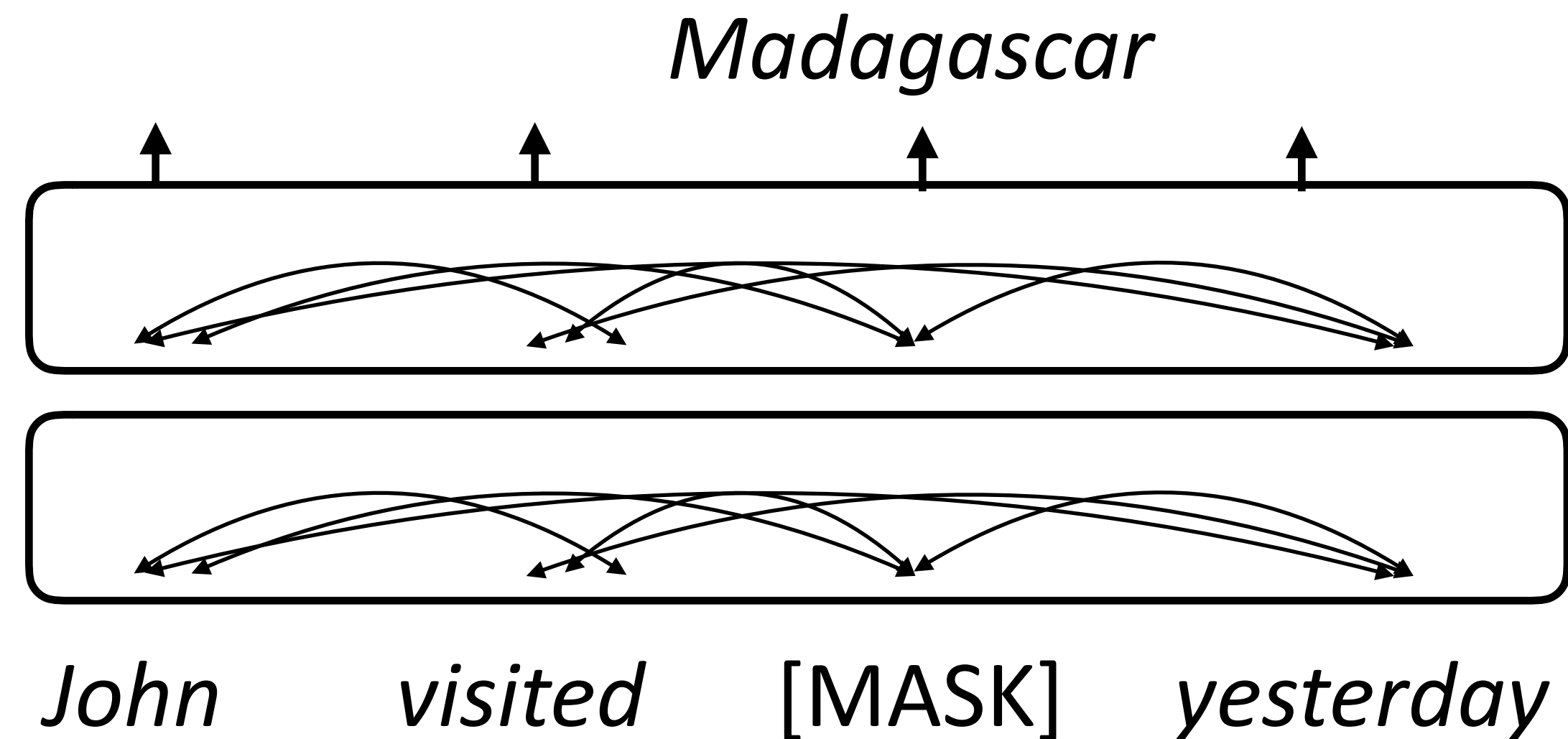
- ▶ You could do this with a “one-sided” transformer, but this “two-sided” model can cheat





# 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

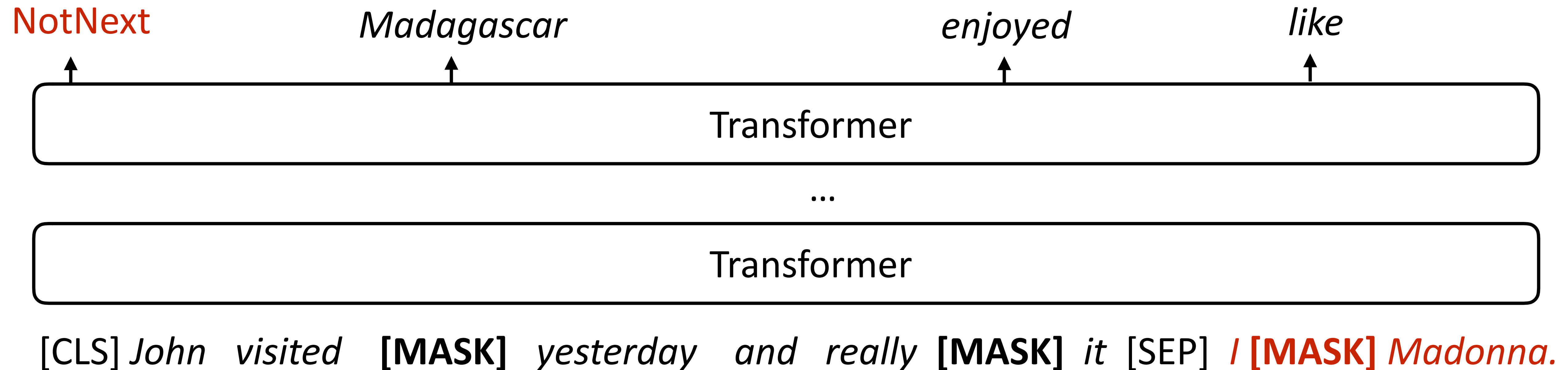






# Next “Sentence” Prediction

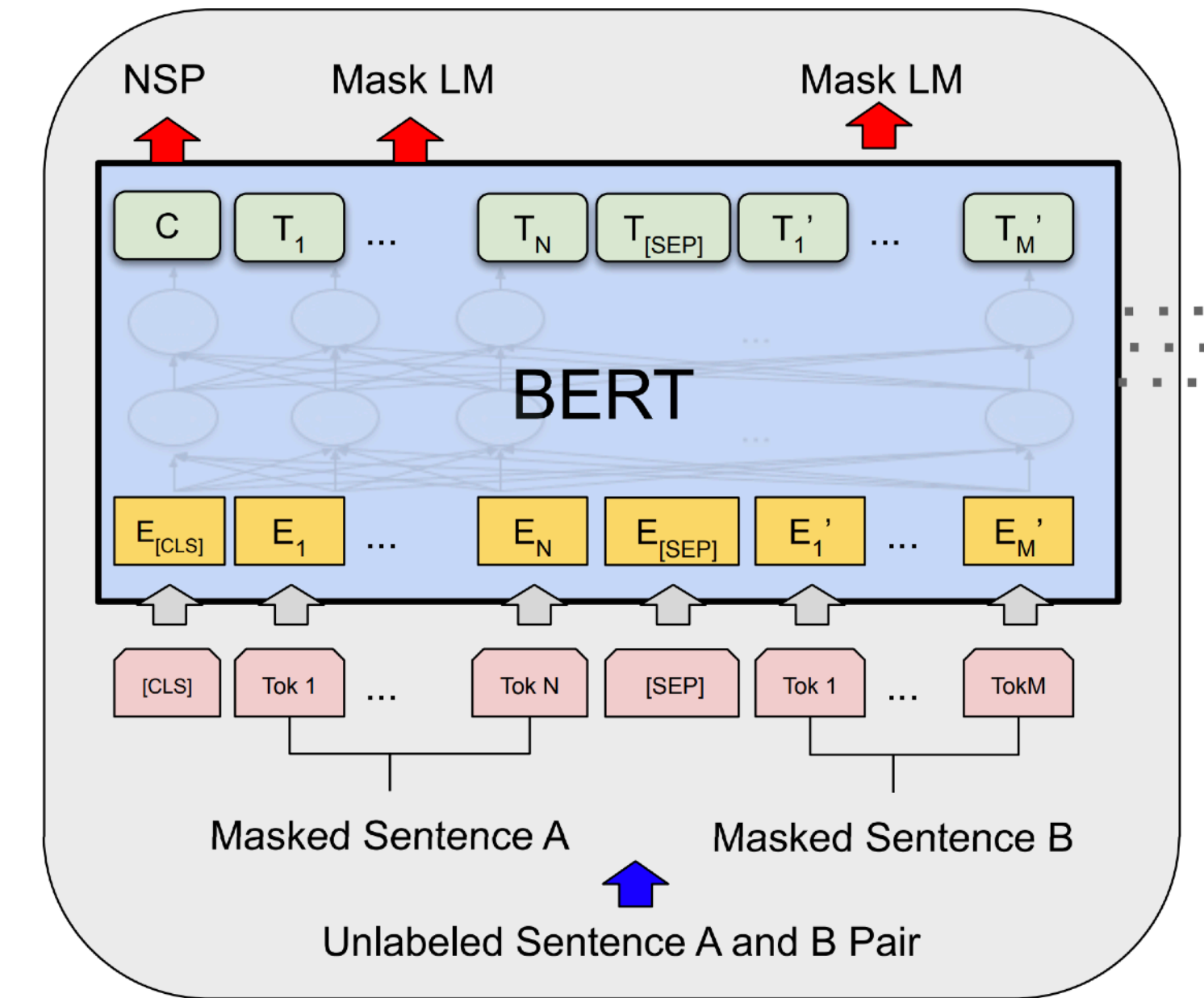
- ▶ Input: [CLS] Text chunk 1 [SEP] Text chunk 2
- ▶ 50% of the time, take the true next chunk of text, 50% of the time take a random other chunk. Predict whether the next chunk is the “true” next
- ▶ BERT objective: masked LM + next sentence prediction





# 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

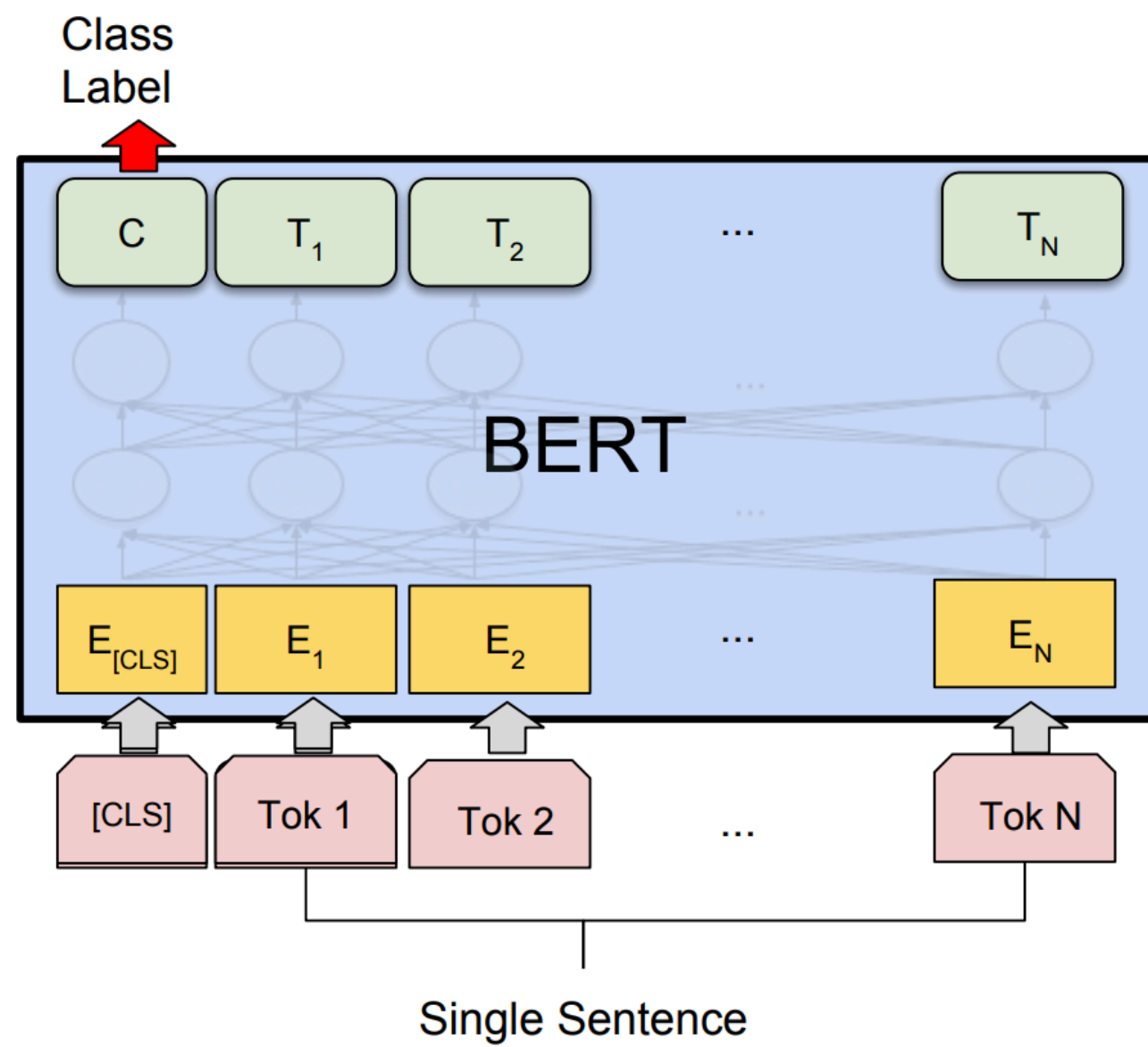


Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{my}$	$E_{dog}$	$E_{is}$	$E_{cute}$	$E_{[SEP]}$	$E_{he}$	$E_{likes}$	$E_{play}$	$E_{##ing}$	$E_{[SEP]}$
	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

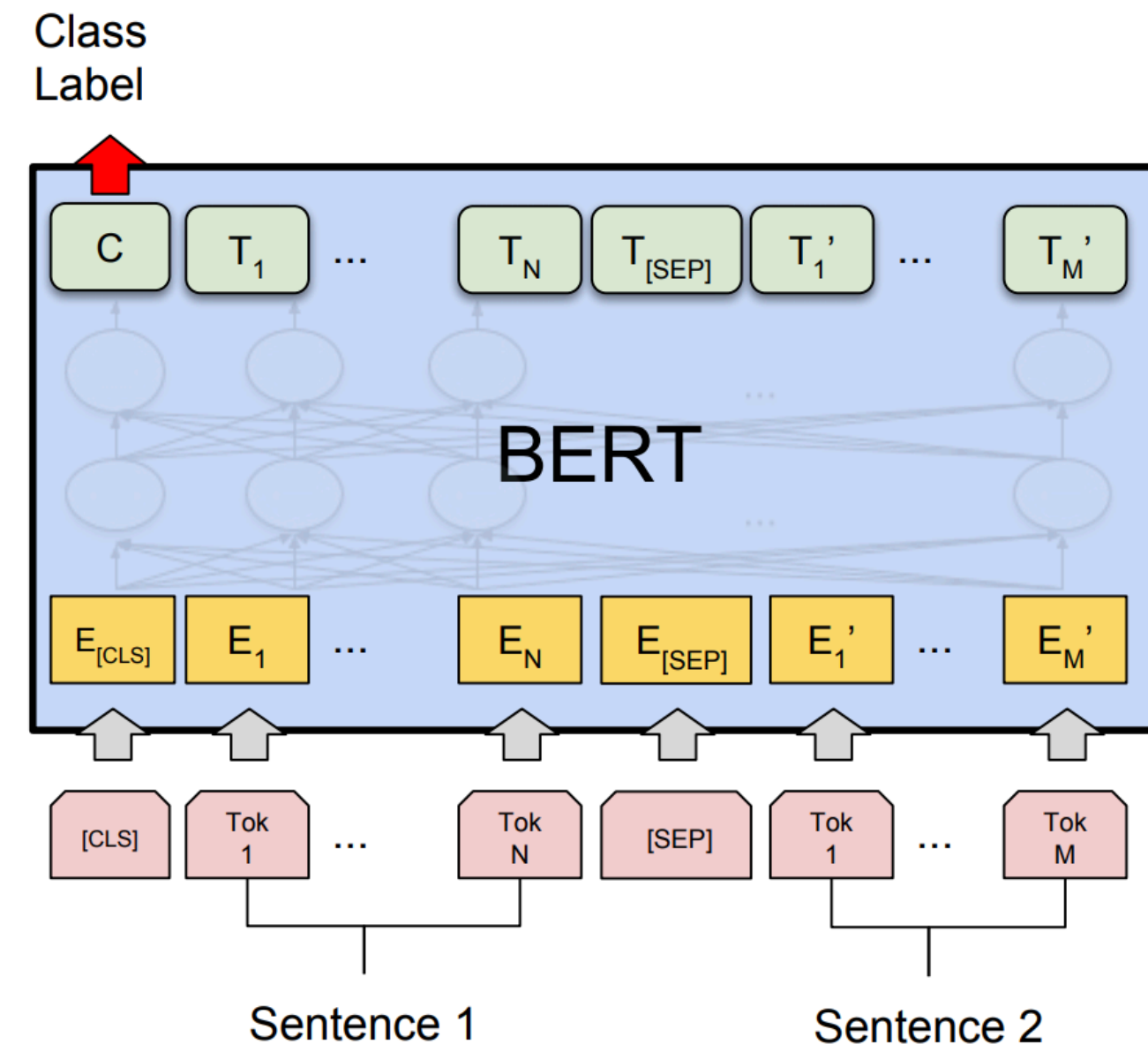
Devlin et al. (2019)



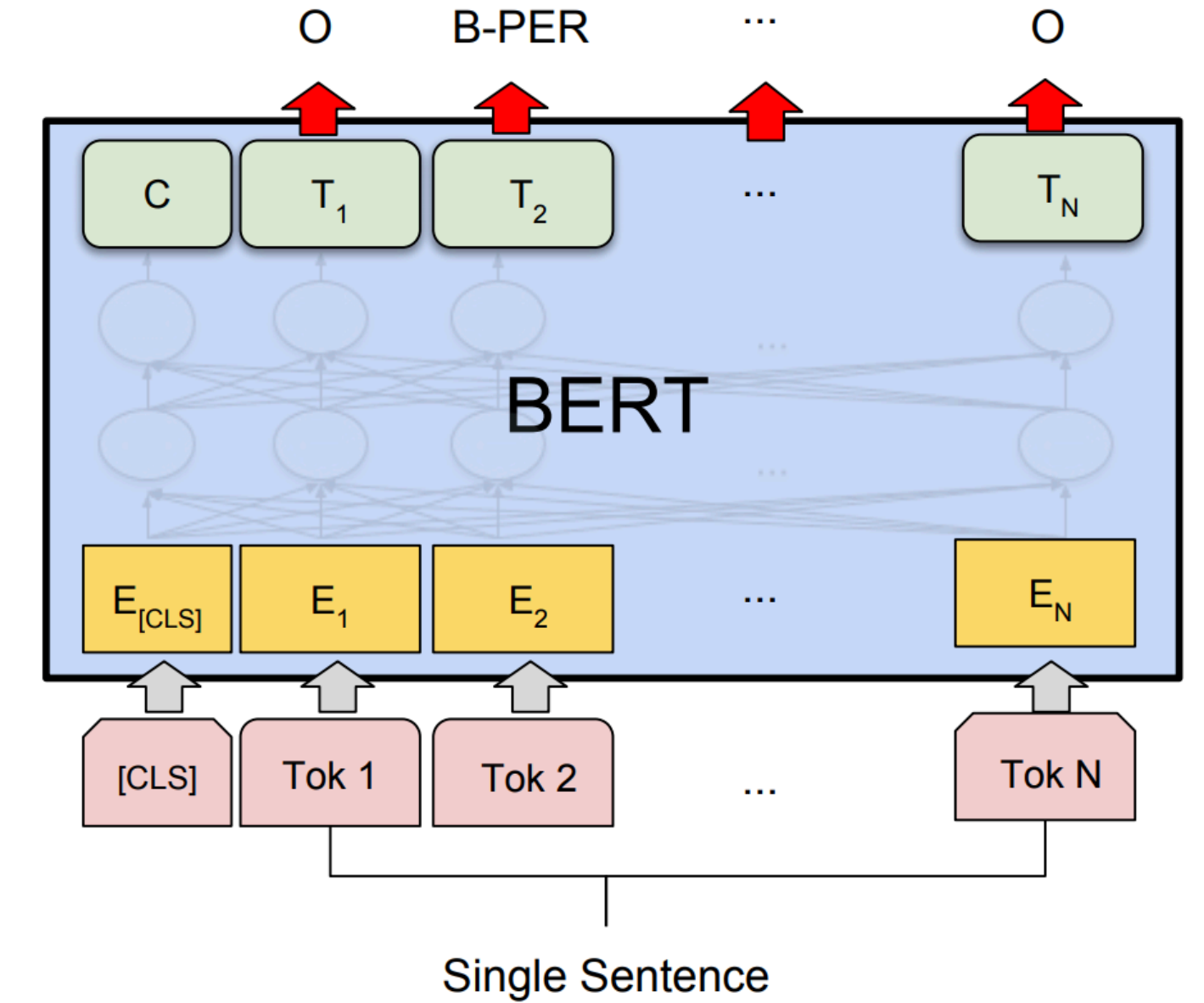
# What can BERT do?



(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

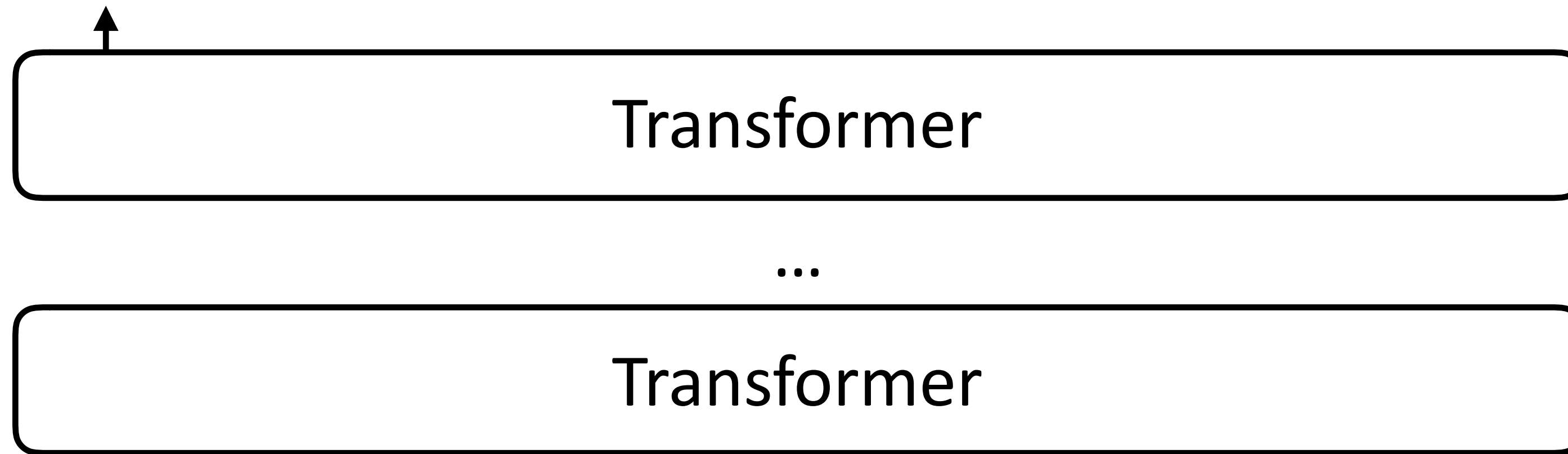
- ▶ Artificial [CLS] token is used as the vector to do classification from
  - ▶ Sentence pair tasks (entailment): feed both sentences into BERT
  - ▶ BERT can also do tagging by predicting tags at each word piece
- Devlin et al. (2019)



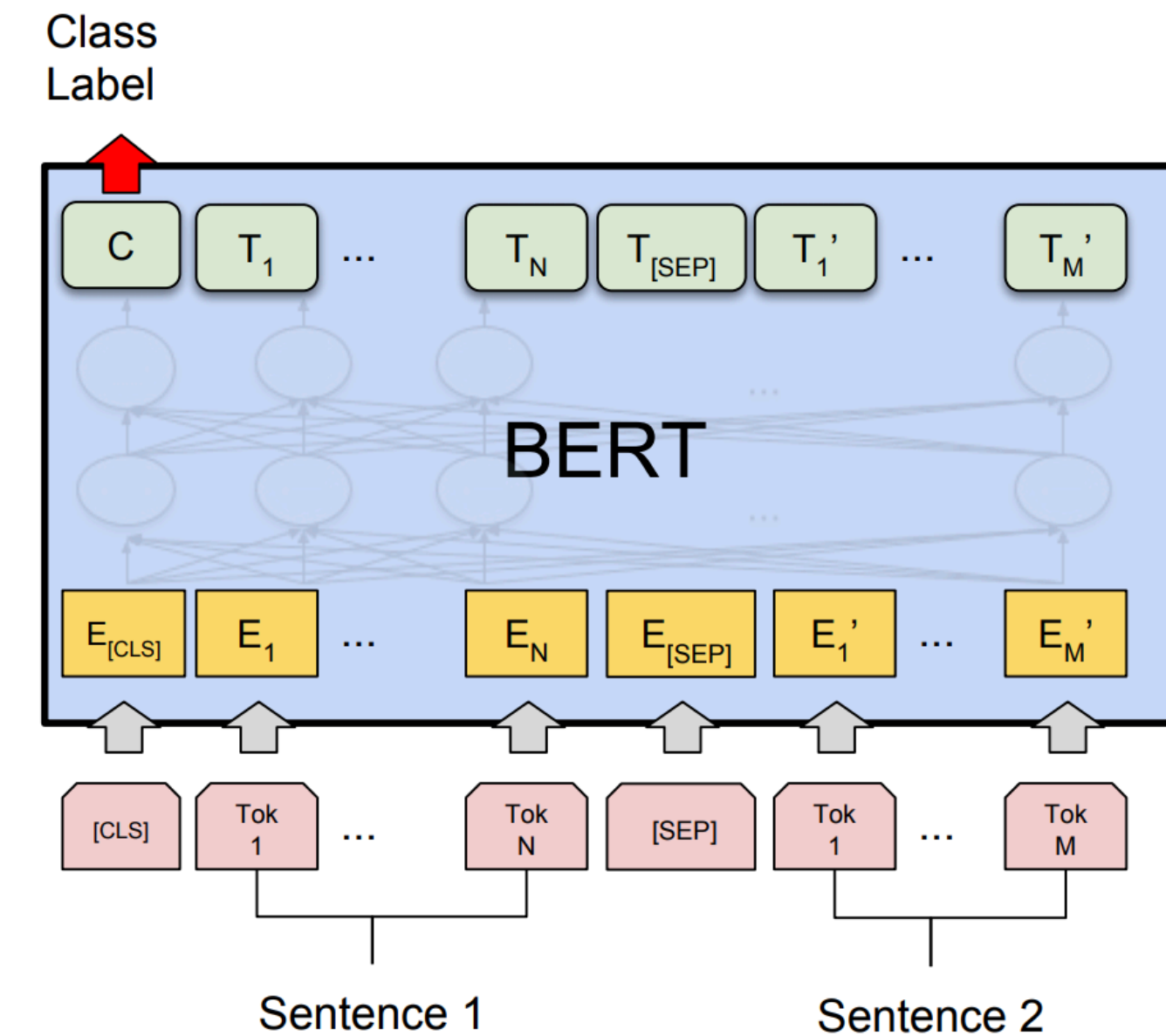


# What can BERT do?

Entails (first sentence implies second is true)



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



(a) Sentence Pair Classification Tasks:  
MNLI, 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?

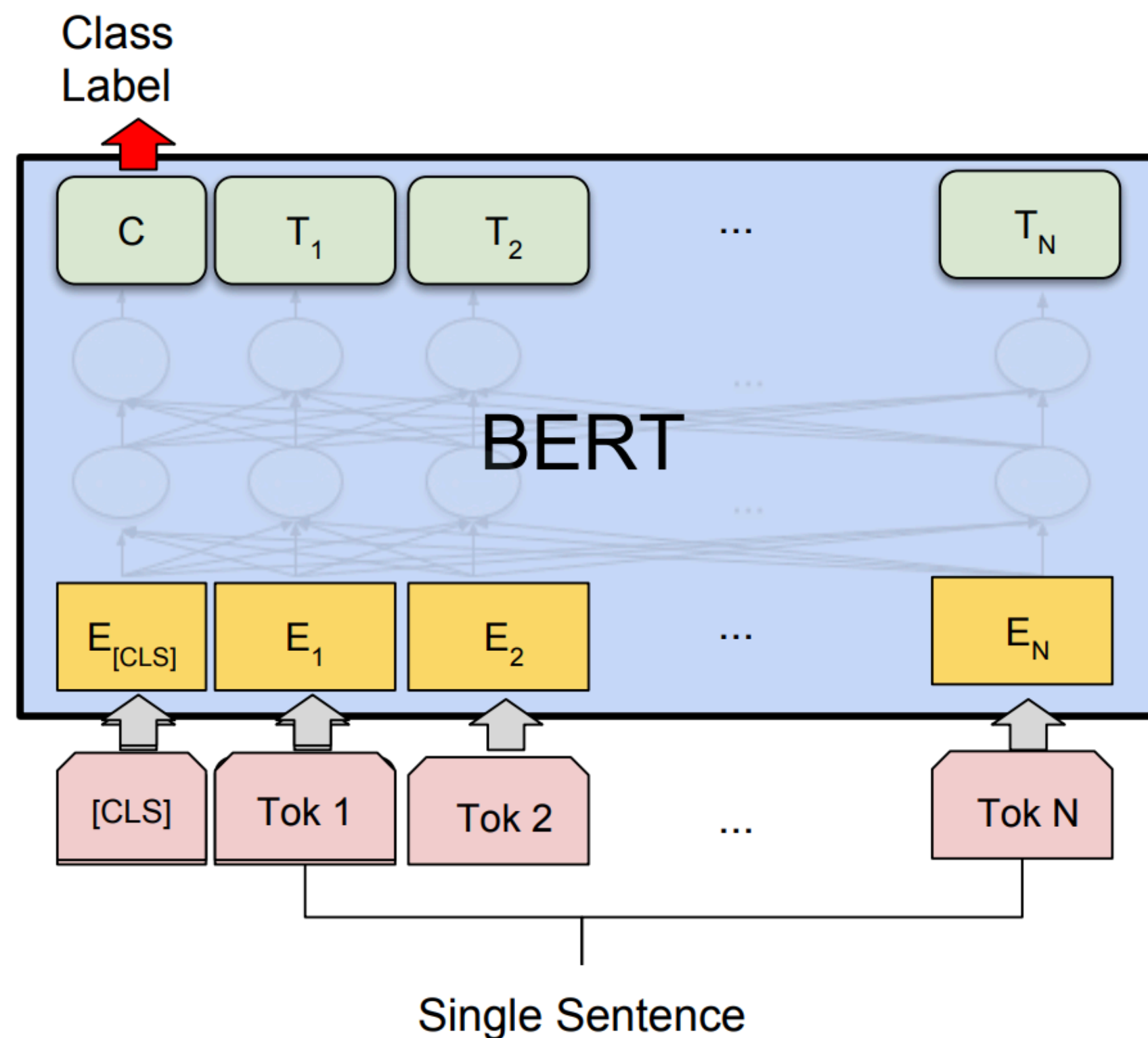
---

- ▶ 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

# BERT Results



# 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

Wang et al. (2019)





# Results

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average -
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)



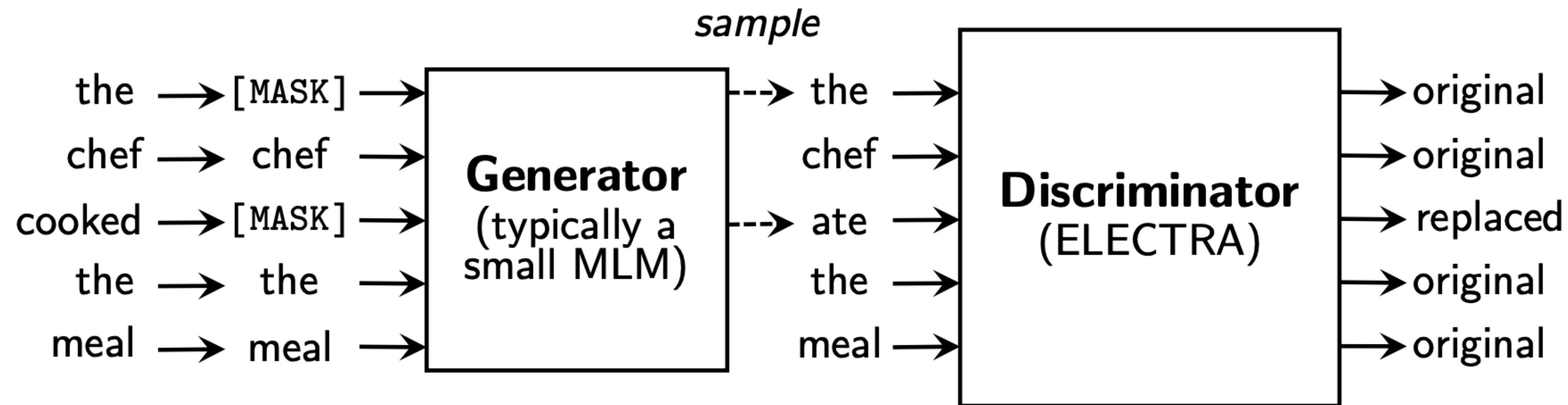
# RoBERTa

- ▶ “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

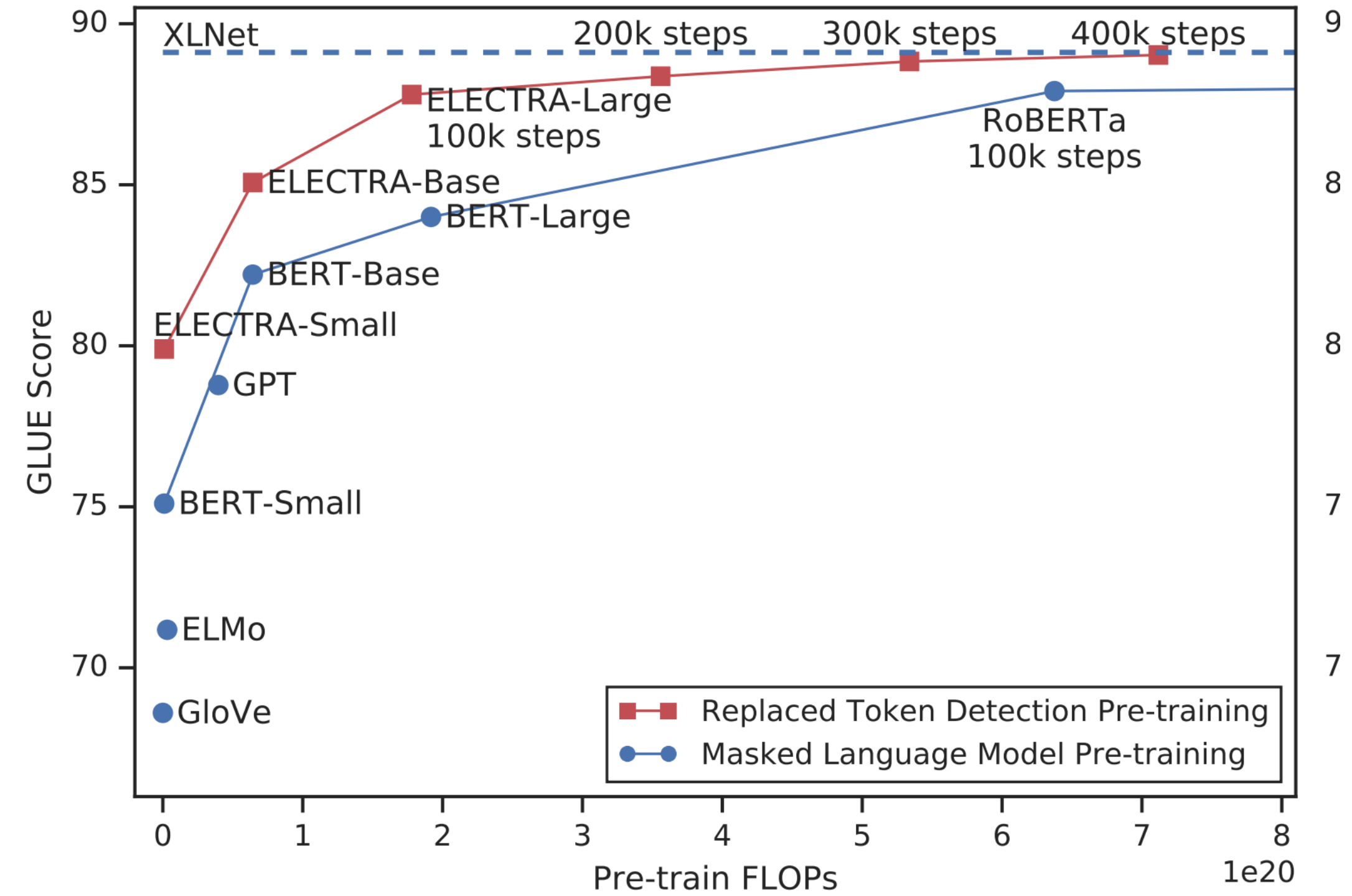


# ELECTRA



Clark et al. (2020)

- ▶ Discriminator to *detect* replaced tokens rather than a generator to actually *predict* what those tokens are
- ▶ More efficient, strong performance







# 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 Ed Chi
5. **XLNet** (from Google/CMU) released with the paper [XLNet: Generalized Autoregressive and Causal Modeling for Language Understanding](#) by Zhilin Yang\*, Zihang Dai\*, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Ed Chi
6. **XLNet** (from Facebook) released together with the paper [Cross-lingual Language Understanding](#) by Lample, Guillaume, Alexis Conneau, and Alexis Conneau.
7. **RoBERTa** (from Facebook), released together with the paper [Robustly Optimized BERT Pre-training for Natural Language Understanding](#)

...

## and “community models”

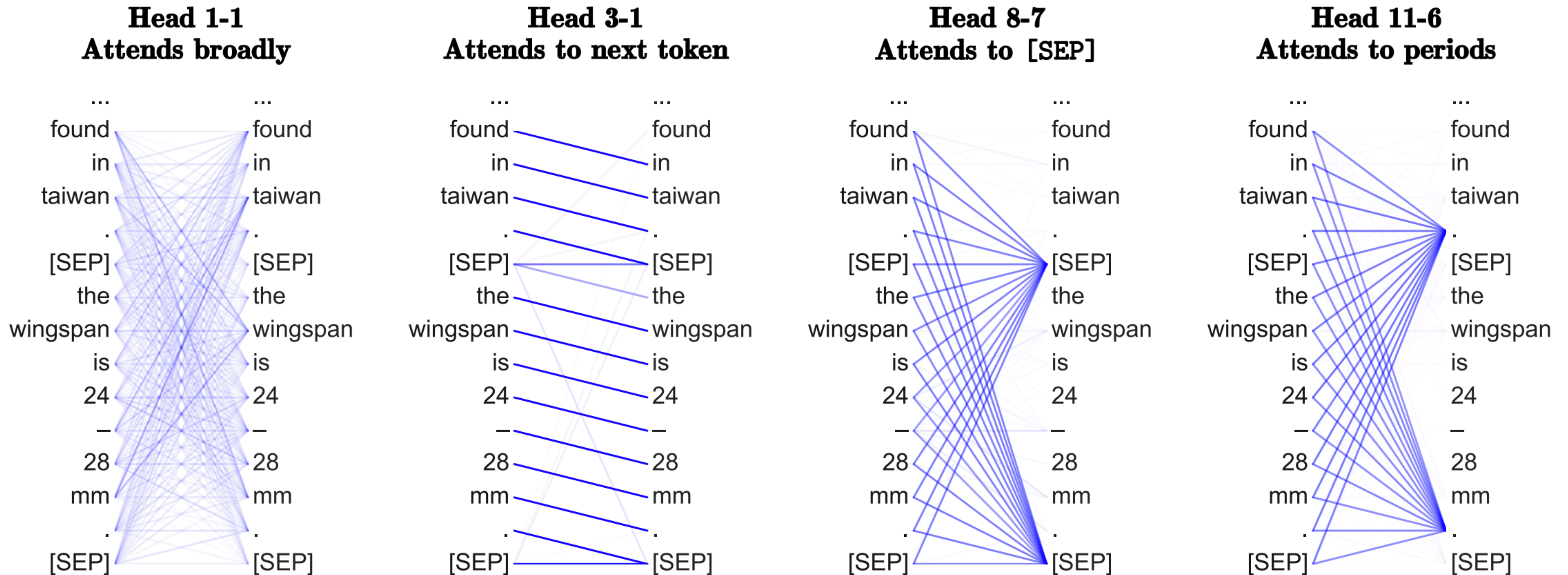
- [mrm8488/spanbert-large-finetuned-tacred](#) ★
- [mrm8488/xlm-multi-finetuned-xquadv1](#) ★
- [nlpaueb/bert-base-greek-uncased-v1](#) ★
- [nlptown/bert-base-multilingual-uncased-sentiment](#) ★
- [patrickvonplaten/reformer-crime-and-punish](#) ★
- [redewiedergabe/bert-base-historical-german-rw-cased](#)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

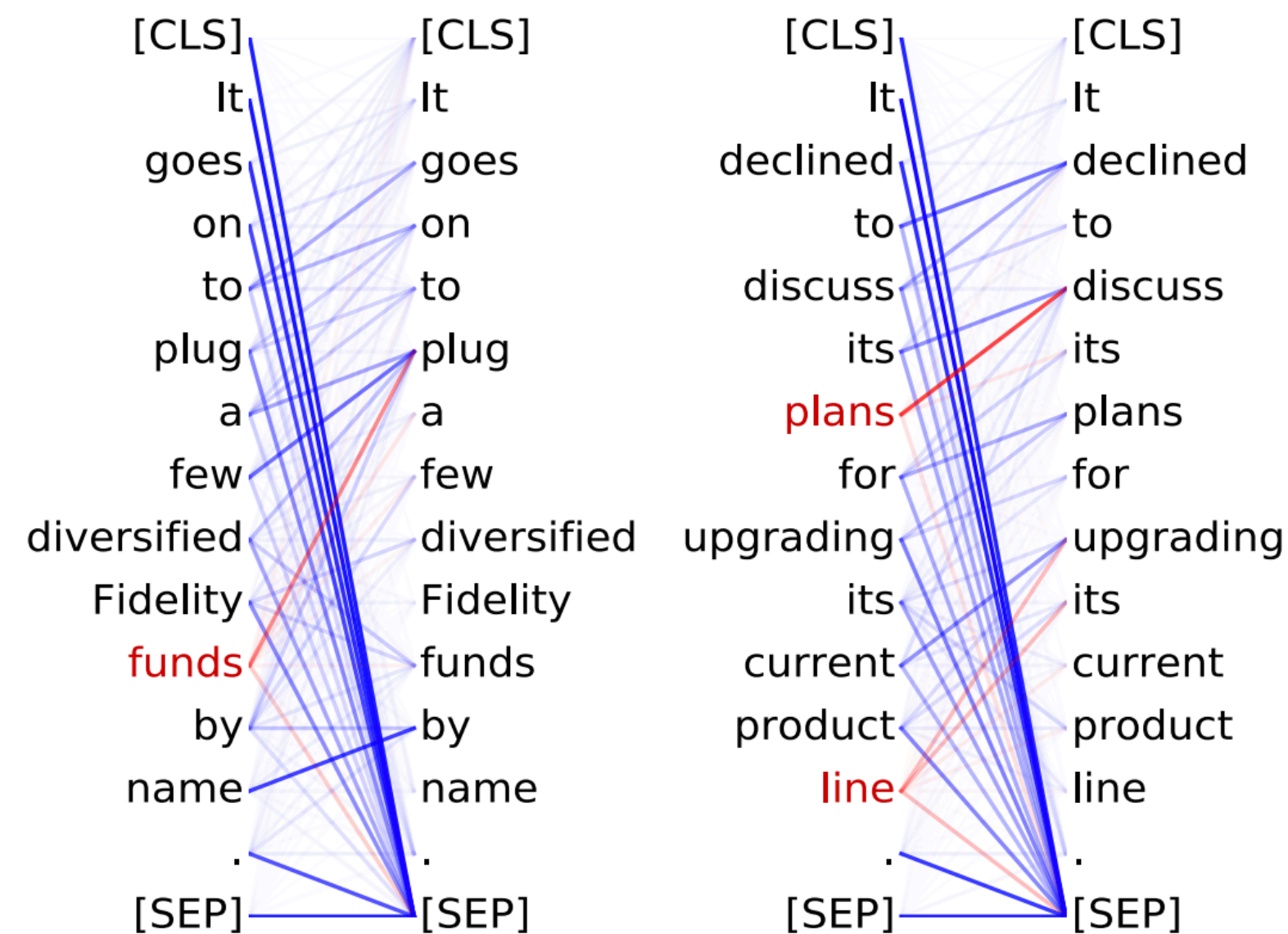




# What does BERT learn?

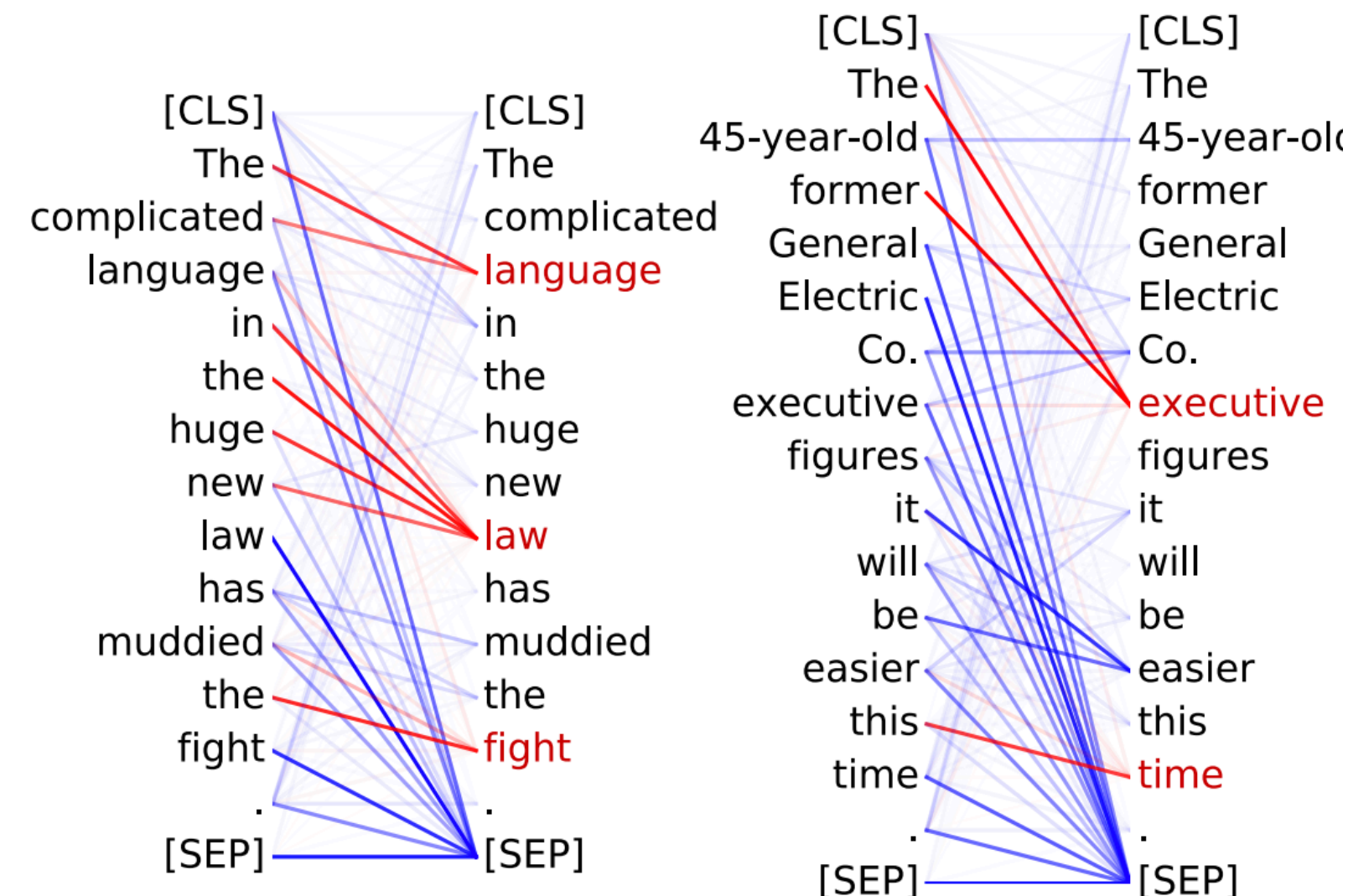
## Head 8-10

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



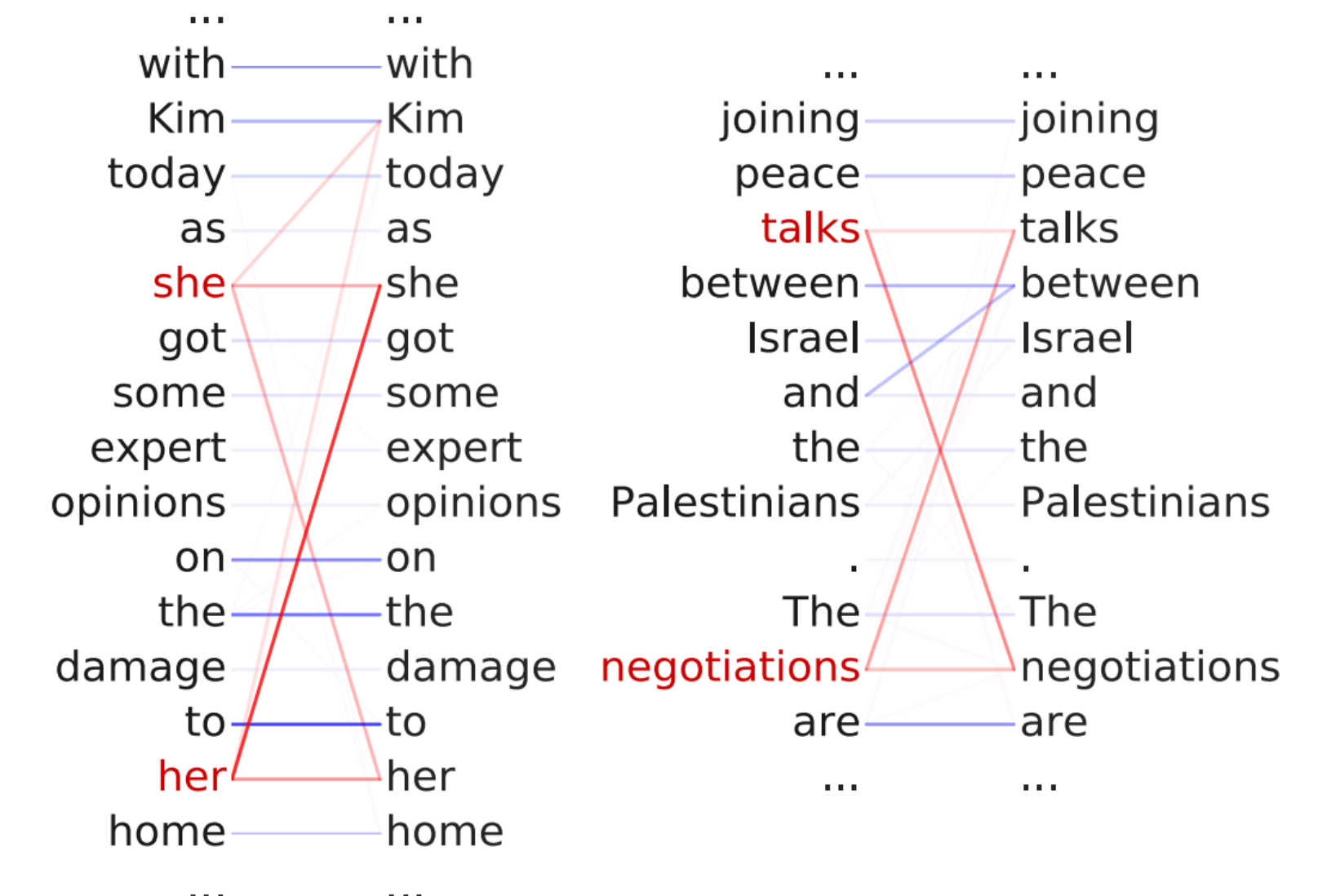
## Head 8-11

- **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

# Applying BERT



# Two Tasks

---

- ▶ Compared to ELMo, BERT is very good at **sentence-pair** tasks
  - ▶ Paraphrase detection
  - ▶ Semantic textual similarity
  - ▶ **Textual entailment**
  - ▶ **Question answering** (not really a sentence pair, but it's a pair of inputs)
- ▶ The final project will focus on when these models fail to learn the right things on these tasks. For now: crash course on these tasks + datasets





# Natural Language Inference

---

Premise

Hypothesis

A boy plays in the snow

*entails*

A boy is outside

A man inspects the uniform of a figure

*contradicts*

The man is sleeping

An older and younger man smiling

*neutral*

Two men are smiling and  
laughing at cats playing

- ▶ Long history of this task: “Recognizing Textual Entailment” challenge in 2006 (Dagan, Glickman, Magnini)
- ▶ Early datasets: small (hundreds of pairs), very ambitious (lots of world knowledge, temporal reasoning, etc.)



# SNLI Dataset

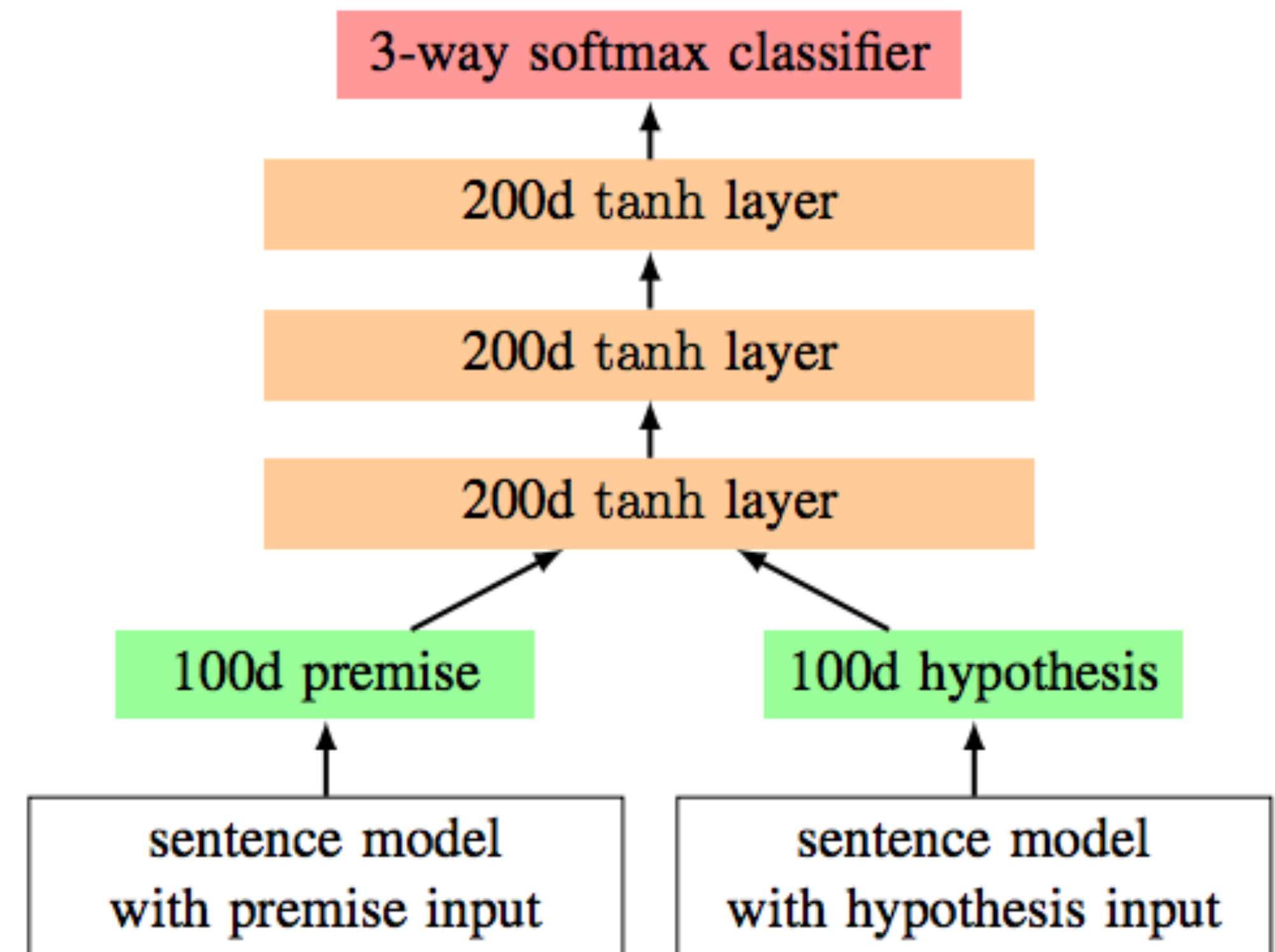
- ▶ Show people captions for (unseen) images and solicit entailed / neural / contradictory statements

- ▶ >500,000 sentence pairs

- ▶ One possible architecture:

300D BiLSTM: 83% accuracy  
(Liu et al., 2016)

- ▶ One of the first big successes of LSTM-based classifiers (sentiment results were more marginal)



Bowman et al. (2015)



# MNLI Dataset

- Drawn from multiple genres of text

Premise	Label	Hypothesis
<b><i>Fiction</i></b>		
The Old One always comforted Ca'daan, except today.	<i>neutral</i>	Ca'daan knew the Old One very well.
<b><i>Letters</i></b>		
Your gift is appreciated by each and every student who will benefit from your generosity.	<i>neutral</i>	Hundreds of students will benefit from your generosity.
<b><i>Telephone Speech</i></b>		
yes now you know if if everybody like in August when everybody's on vacation or something we can dress a little more casual or	<i>contradiction</i>	August is a black out month for vacations in the company.
<b><i>9/11 Report</i></b>		
At the other end of Pennsylvania Avenue, people began to line up for a White House tour.	<i>entailment</i>	People formed a line at the end of Pennsylvania Avenue.

Williams et al. (2018)



# How do models do it?



A **man** is eating a sandwich [SEP] A **person** is eating a sandwich



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

- ▶ Transformers can easily learn to spot words or short phrases that are transformed
- ▶ **But**, models are often overly sensitive to lexical overlap





# Question Answering

---

- ▶ Many types of QA:
- ▶ We'll focus on **factoid questions** being answered **from text**
  - ▶ E.g., “What was Marie Curie the first female recipient of?” — unlikely you would have this answer in a database
  - ▶ Not appropriate: “When was Marie Curie born?” — probably answered in a DB
  - ▶ Not appropriate: “Why did World War II start?” — no simple answer



# SQuAD

---

Q: What was Marie Curie the first female recipient of?

Passage: One of the most famous people born in Warsaw was Marie Skłodowska-Curie, who achieved international recognition for her research on radioactivity and was the first female recipient of the **Nobel Prize**. Famous musicians include Władysław Szpilman and Frédéric Chopin. Though Chopin was born in the village of Żelazowa Wola, about 60 km (37 mi) from Warsaw, he moved to the city with his family when he was seven months old. Casimir Pulaski, a Polish general and hero of the American Revolutionary War, was born here in 1745.

Answer = Nobel Prize

- ▶ Assume we know a passage that contains the answer. More recent work has shown how to retrieve these effectively (will discuss when we get to QA)



# SQuAD

Q: What was Marie Curie the first female recipient of?

Passage: One of the most famous people born in Warsaw was Marie Skłodowska-Curie, who achieved international recognition for her research on radioactivity and was the first female recipient of the **Nobel Prize**. ...

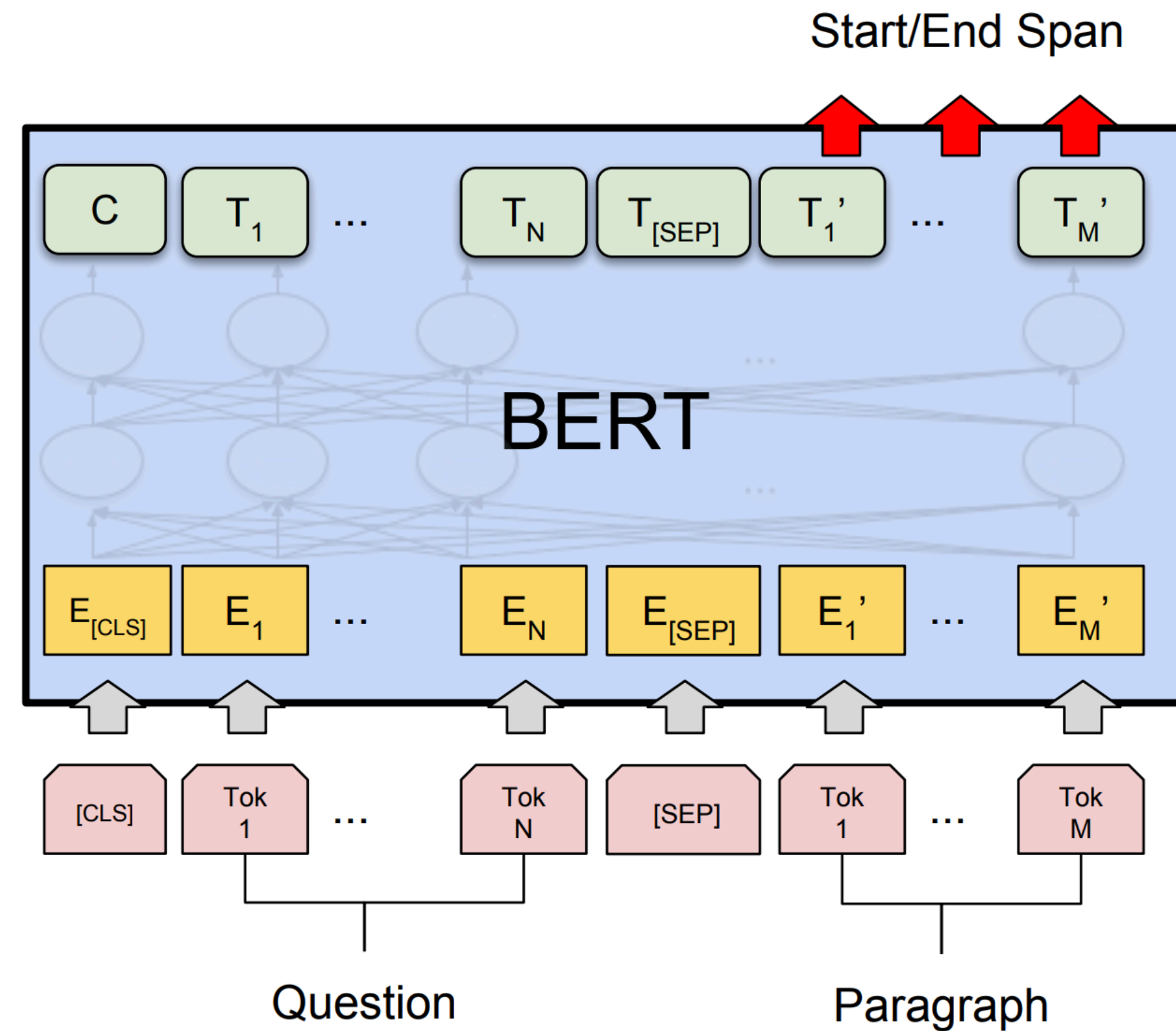
- Predict answer as a pair of (start, end) indices given question q and passage p; compute a score for each word and softmax those

$$P(\text{start} \mid q, p) = \begin{array}{ccccc} 0.01 & 0.01 & 0.01 & 0.85 & 0.01 \\ \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ & \text{recipient of the } \mathbf{Nobel Prize} . \end{array}$$

$P(\text{end} \mid q, p)$  = same computation but different params



# QA with BERT



What was Marie Curie the first female recipient of ? [SEP] One of the most famous people born in Warsaw was Marie ...





# Takeaways

---

- ▶ Pre-trained models and BERT are very powerful for a range of NLP tasks
- ▶ These models have enabled big advances in NLI and QA specifically
- ▶ Next time: final project introduction. Idea of dataset artifacts (“bad” patterns memorized by the model that hurt its ability to generalize) and what we can do about them