

Pretrained Representations: **GPT, ELMo, BERT, BART & T5**

Antoine Bosselut



Today's Outline

- **Lecture**
 - **Pretraining Transformers:** GPT
 - **Bidirectional Pretraining:** ELMo + BERT
 - **Sequence-to-sequence Pretraining:** BART + T5

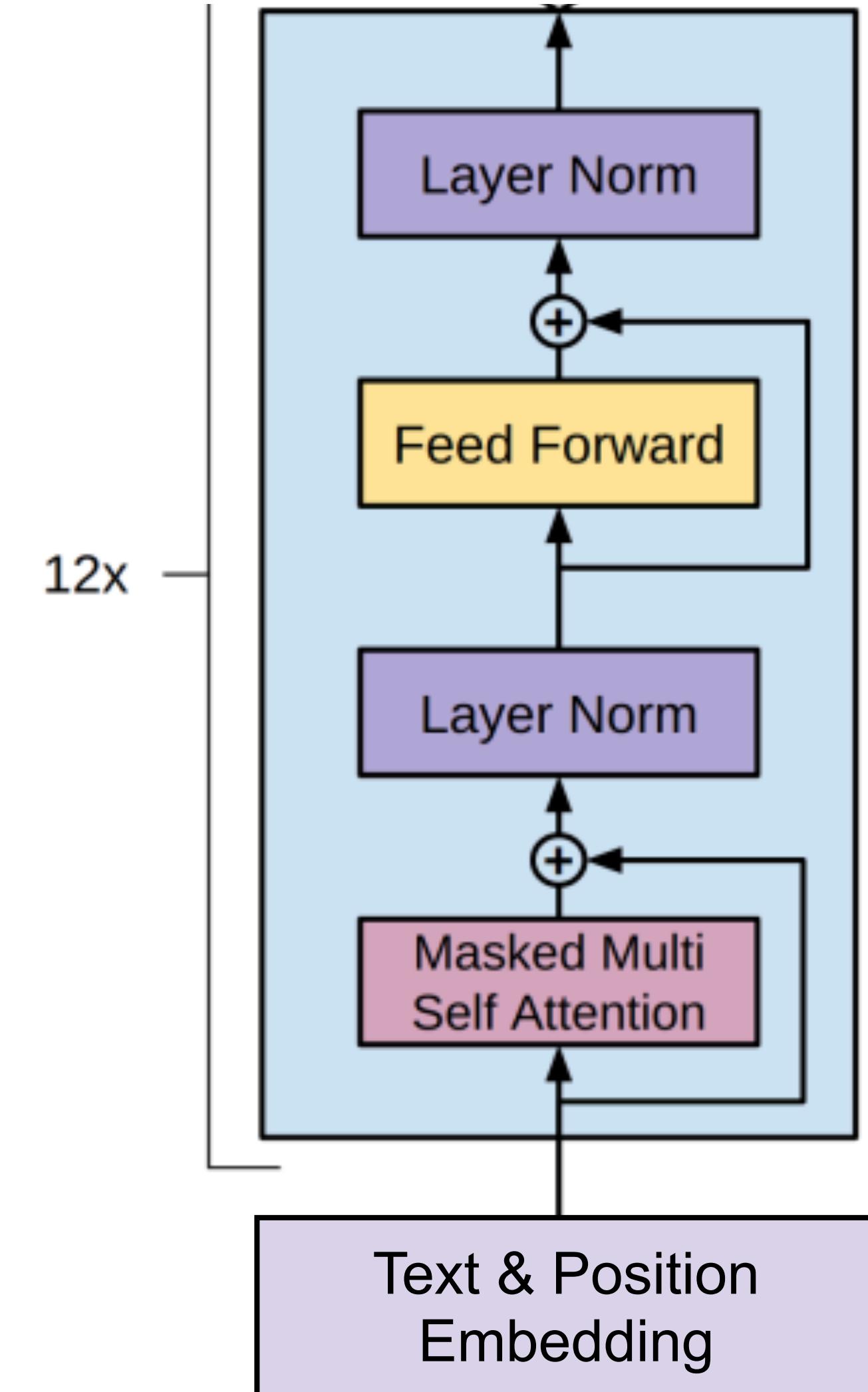
Pretrained Representations: **GPT, GPT-2**

Antoine Bosselut



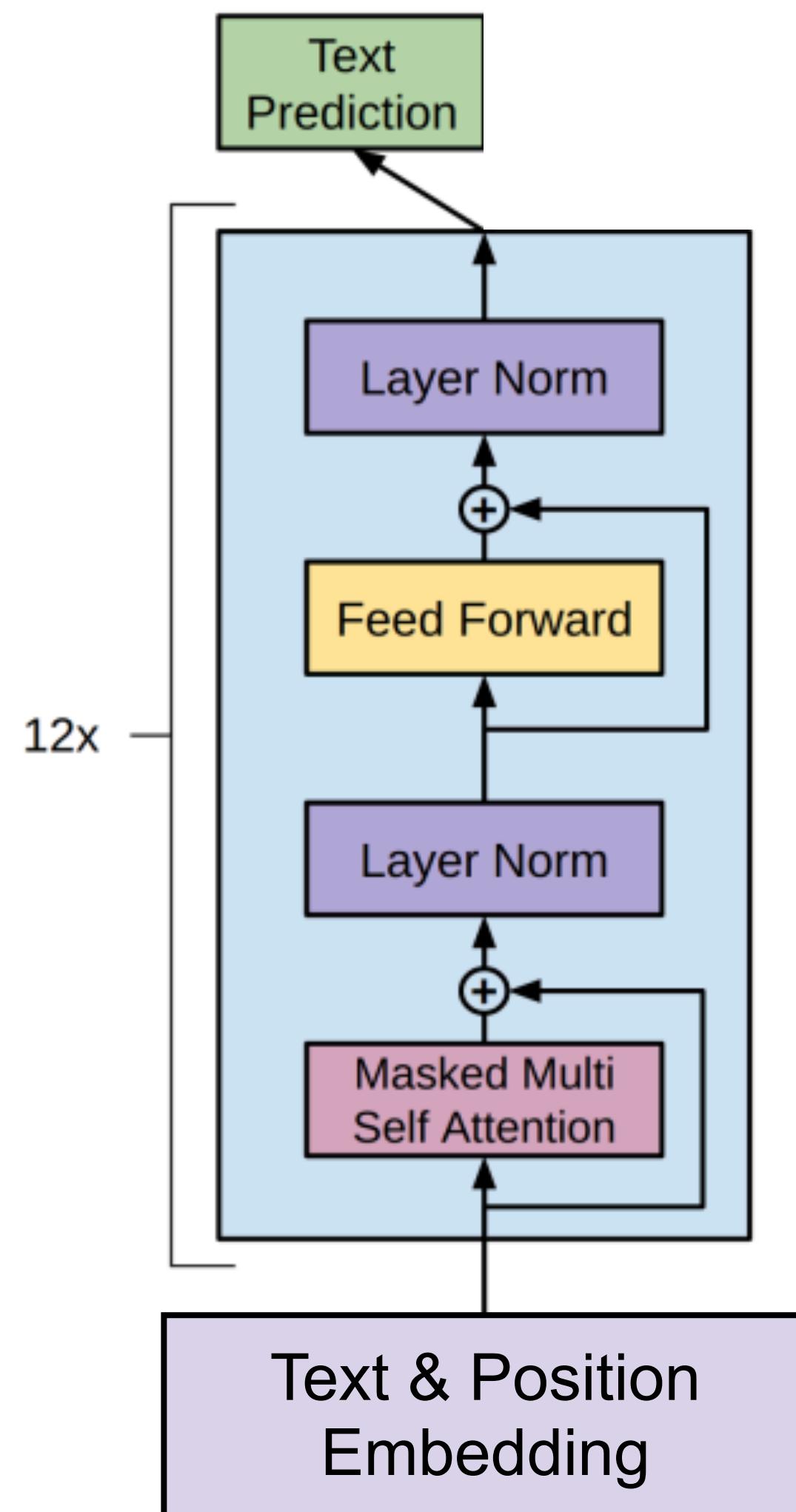
GPT: Generative Pretrained Transformer

- Called a *decoder* transformer
- **But**, actual GPT block mixes design of encoder and decoder from original transformer
- Uses masked multi-headed self-attention (**decoder**)
 - Token at each step can't attend to future tokens
- No cross-attention; only computes a self-attention over its history (**encoder**)



Training GPT

- Pretrained on TorontoBooks corpus:
7000 unpublished books (~13 GB)
- Corpus segmented broken up into
windows of 512 tokens
 - can model long-range context during
pretraining up to 512 tokens
- **Pretraining task:** next word prediction
(i.e., causal language modelling)

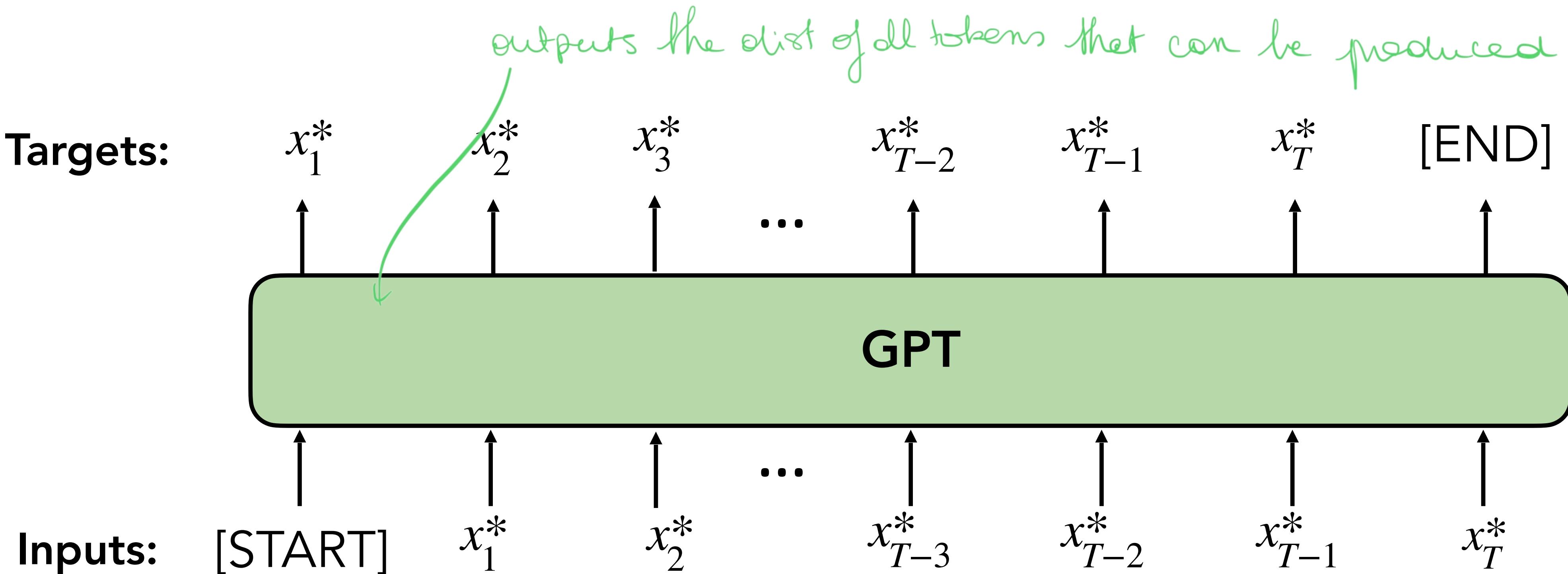


can be parallelized

Pretraining

- Minimize the negative log probability of the **gold*** sequences in your dataset

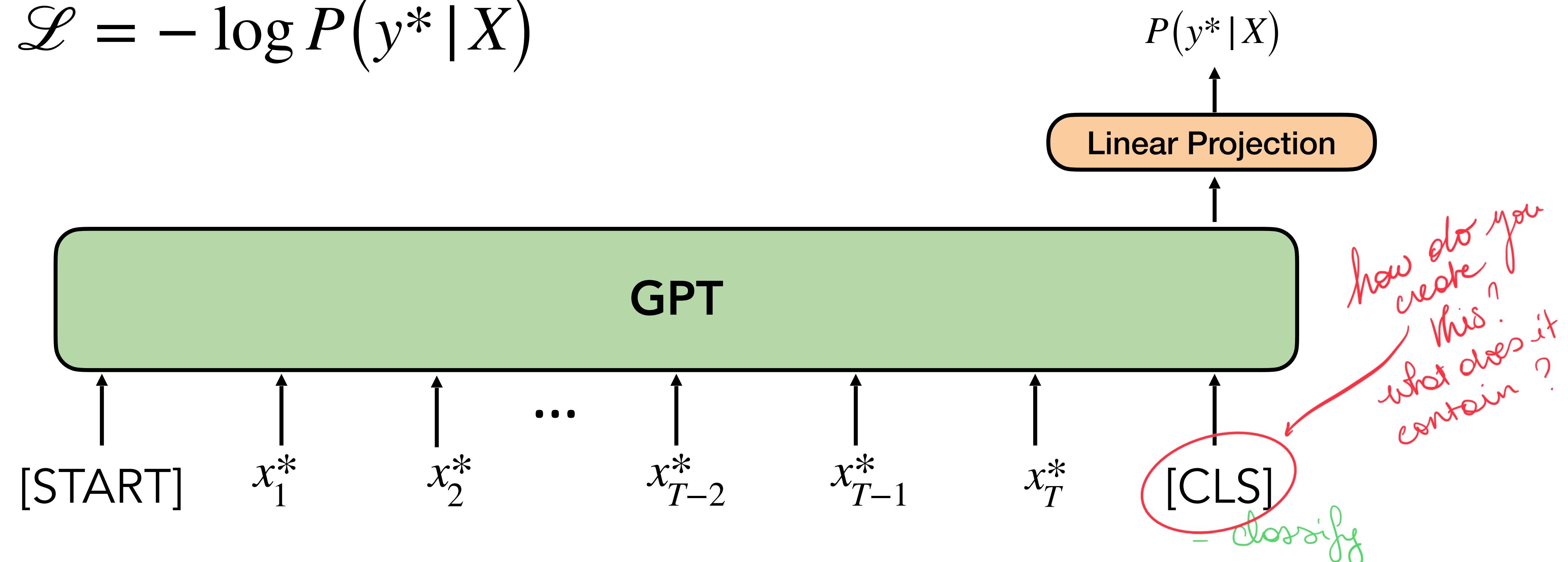
$$\mathcal{L} = - \sum_{t=1}^T \log P(x_t^* | \{x_s^*\}_{s < t})$$



Finetuning

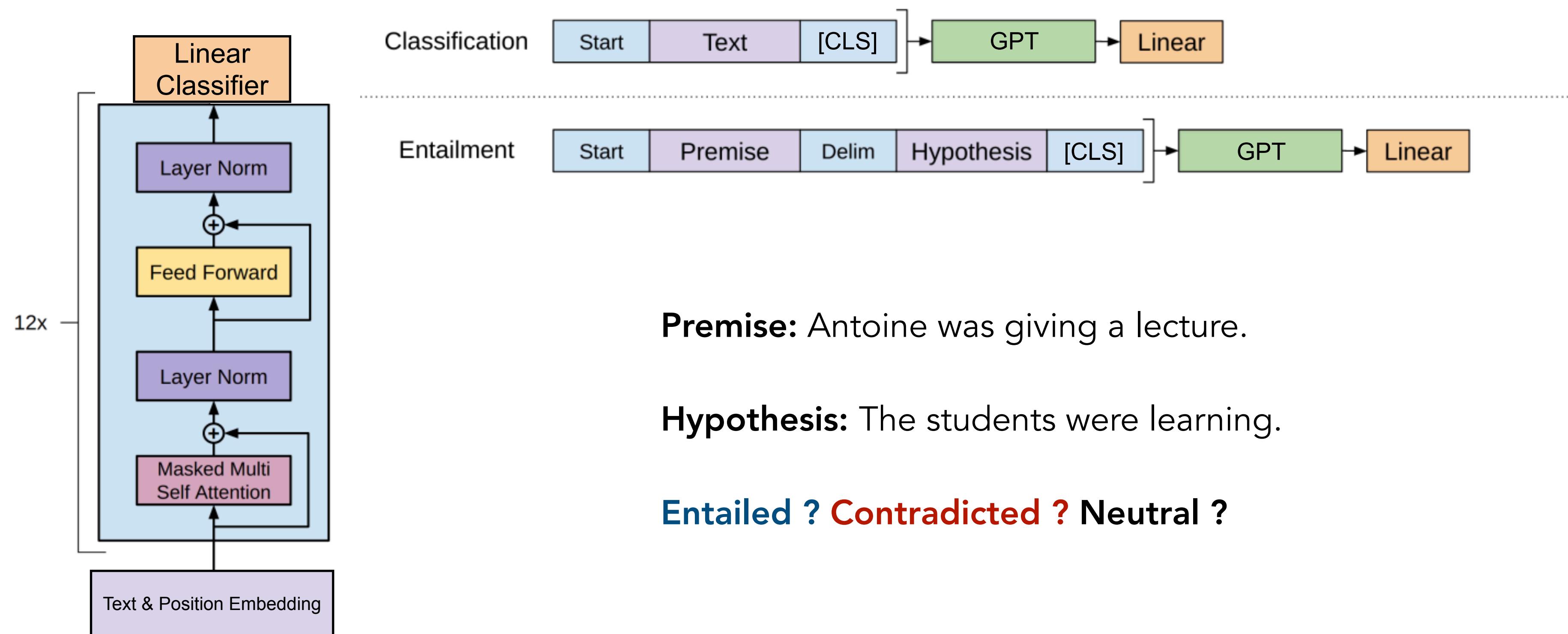
- Minimize the cross-entropy of the **gold** label of examples in your dataset

$$\mathcal{L} = -\log P(y^* | X)$$



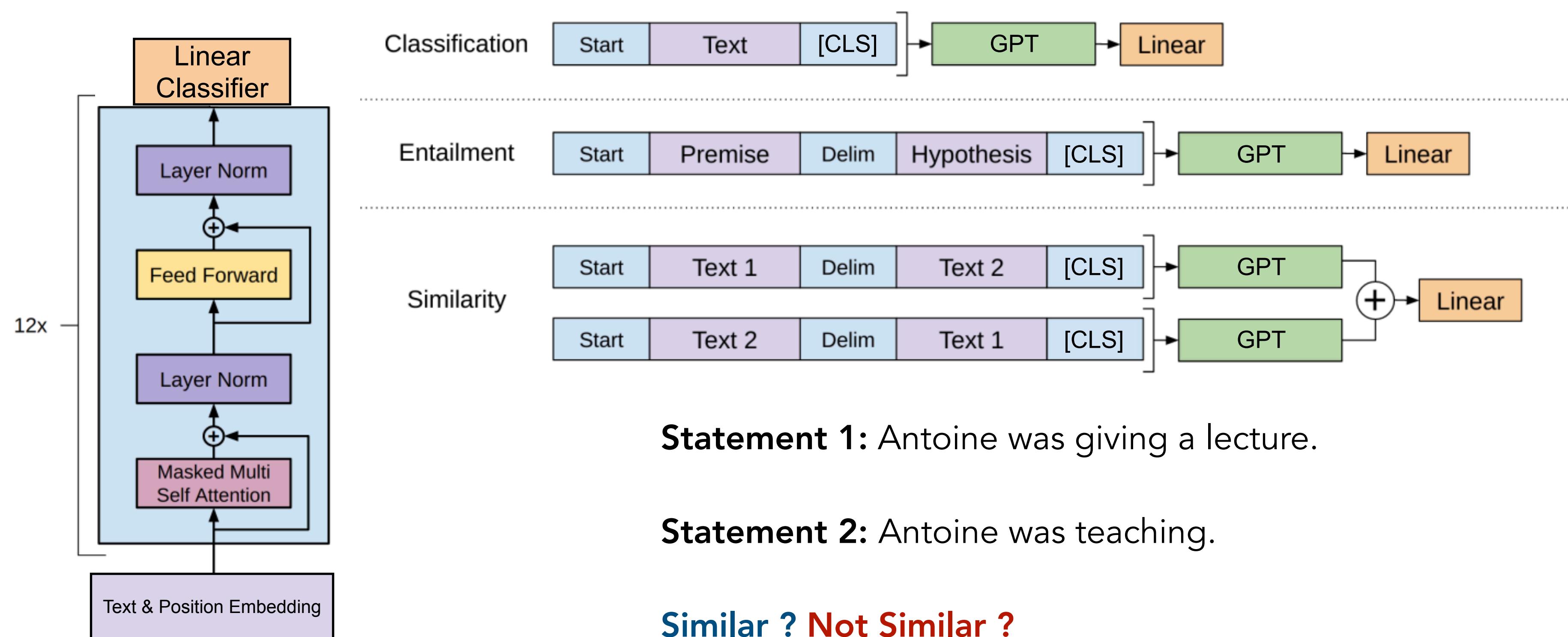
Fine-tuning

- After pre-training, model can be fine-tuned by training on individual datasets
- Pretrained model used as initialisation for training on individual tasks



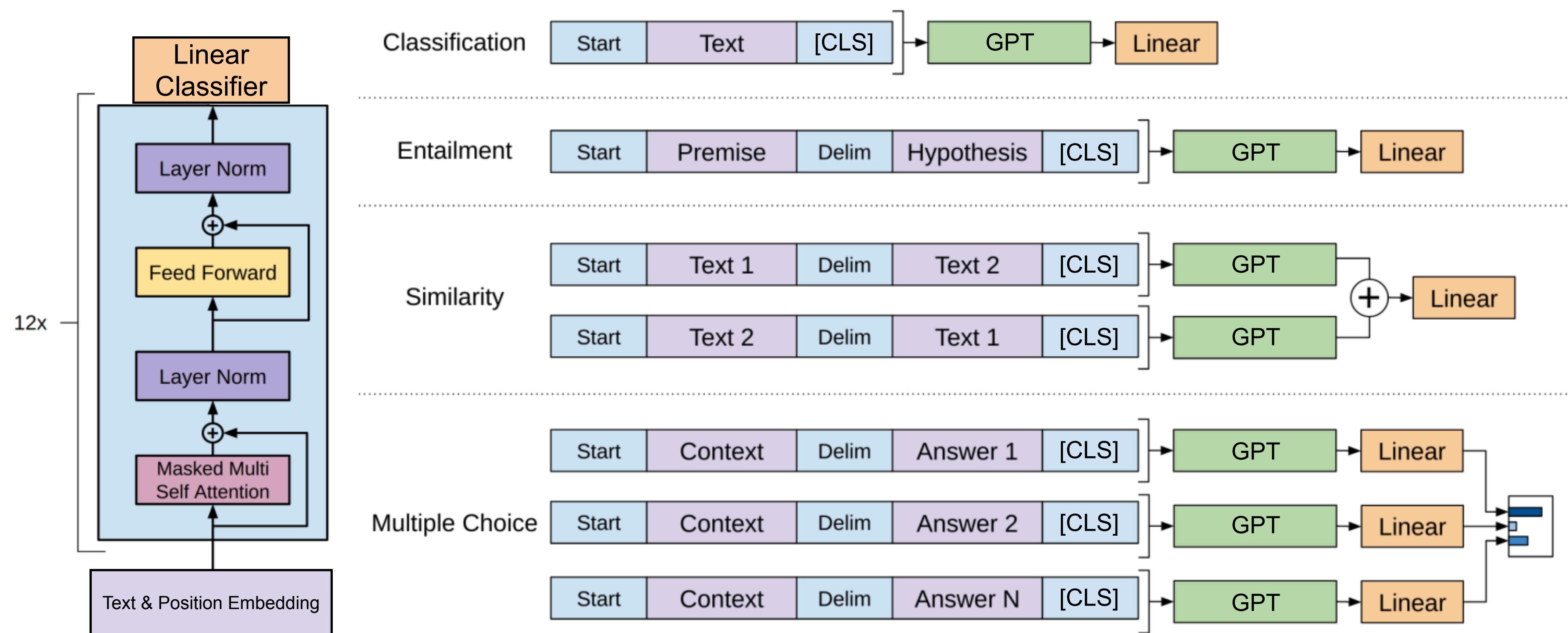
Fine-tuning

- After pre-training, model can be fine-tuned by training on individual datasets
- Pretrained model used as initialisation for training on individual tasks



Fine-tuning

- After pre-training, model can be fine-tuned by training on individual datasets
- Pretrained model used as initialisation for training on individual tasks



Massive Improvements (back then)

Dataset	Task	SOTA	Ours
SNLI	Textual entailment	89.3	89.9
MNLI matched	Textual entailment	80.6	82.1
MNLI mismatched	Textual entailment	80.1	81.4
SciTail	Textual entailment	83.3	88.3
QNLI	Textual entailment	82.3	88.1
RTE	Textual entailment	61.7	56.0
STS-B	Semantic similarity	81.0	82.0
QQP	Semantic similarity	66.1	70.3
MRPC	Semantic similarity	86.0	82.3
RACE	Reading comprehension	53.3	59.0
ROCStories	Commonsense reasoning	77.6	86.5
COPA	Commonsense reasoning	71.2	78.6
SST-2	Sentiment analysis	93.2	91.3
CoLA	Linguistic acceptability	35.0	45.4
GLUE	Multi task benchmark	68.9	72.8

1 Pretraining Phase (General Training)

Goal: Learn general language patterns by predicting the next token in large amounts of text.

How it Works:

1. Input Processing:

- Large amounts of text (e.g., books, articles, web pages) are tokenized into subwords.
- Each input sequence gets a **start token** but no **CLS token** (CLS is used in BERT, not GPT).

2. Causal Masked Self-Attention:

- GPT uses **causal masking**, meaning at each step, it can **only attend to previous tokens** (not future tokens).
- This forces the model to learn **sequential text generation** rather than bidirectional understanding.

3. Next-Token Prediction Objective:

- The model is trained to **predict the next token** given previous tokens.
- Example:

- Input: "The cat sat on the"
- Target Output: "mat"
- The model learns to output "mat" based on "The cat sat on the"

4. Loss Function:

- Cross-entropy loss is used to compare the predicted next token with the actual next token.
- The model updates its weights via backpropagation.

Outcome:

- After pretraining, GPT has learned **general language structures**, common word associations, and grammatical patterns.
- However, it **hasn't specialized** in specific tasks (e.g., translation, summarization, coding).

4. (Optional) Reinforcement Learning with Human Feedback (RLHF):

- Used in models like ChatGPT to make outputs more aligned with human preferences.
- Involves:
 - Humans ranking model responses.
 - A reward model learning these preferences.
 - Reinforcement learning adjusting GPT's behavior.

Outcome:

- GPT now **specializes** in the fine-tuning task.
- It still retains **general pretraining knowledge** but adapts to new patterns.

2 Fine-Tuning Phase (Task-Specific Training)

Goal: Adapt the pretrained model to perform a **specific task** (e.g., summarization, Q&A, chatbot dialogue).

How it Works:

1. New Dataset:

- The model is trained on a smaller, **task-specific dataset**.
- Example: If fine-tuning GPT for sentiment analysis, we train it on labeled reviews ("positive" or "negative").

2. Fine-Tuning with Supervised Learning:

- Unlike pretraining, where the model **just predicts the next word**, fine-tuning can involve:
 - **Custom formatting** (e.g., adding instruction prompts).
 - **Supervised learning** (model output is compared against labeled data).
 - **Loss function modifications** (depending on the task).

3. Parameter Updates:

- The pretrained GPT weights are **further updated** based on the new dataset.
- This allows GPT to specialize in the new task while **retaining general language understanding**.

- **Fundamental change:** No more training architectures (e.g., LSTM, GRU) initialized with pretrained embeddings (Word2vec, GloVe, fastText, etc.).
- **Finetuning pretrained full models**

Example: Pretraining vs. Fine-Tuning in Action

Pretrained GPT-3 (General Model)

💬 Input: "The capital of France is"
🤖 GPT-3: "Paris" (General knowledge from pretraining)

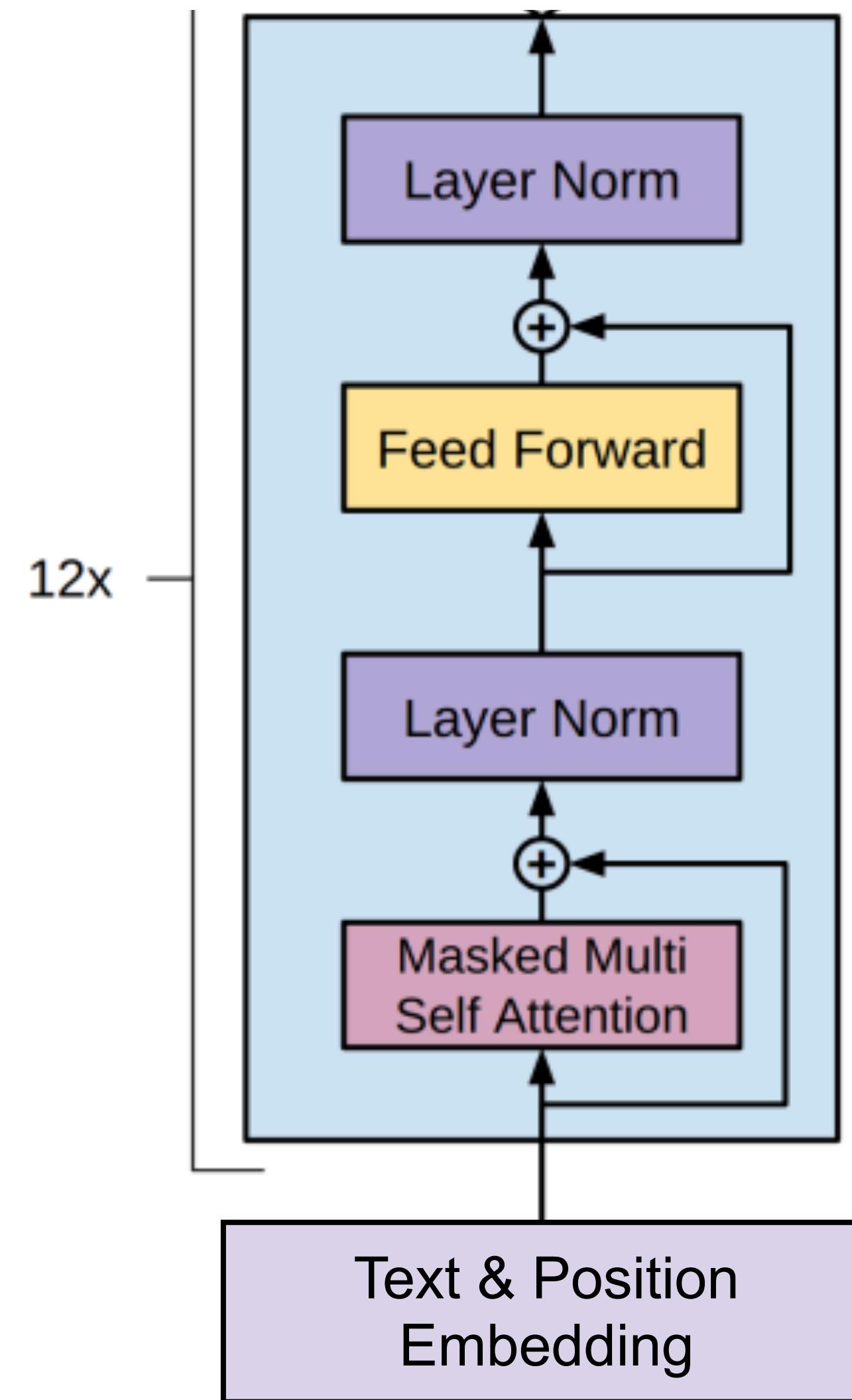
Fine-Tuned GPT-3 (QA System)

💬 Input: "Question: What is the capital of France? Answer:"
🤖 Fine-tuned GPT: "The capital of France is Paris." (More structured response because of fine-tuning)

GPT2

- Roughly the same architecture as GPT
- Trained on 40 GB of data: OpenWebText
 - Articles linked to from reddit threads
- Four model sizes:

Parameters	Layers	d_{model}
117M	12	768
345M	24	1024
762M	36	1280
1542M	48	1600



GPT2

Dataset	Metric	Our result	Previous record	Human
Winograd Schema Challenge	accuracy (+)	70.70%	63.7%	92%+
LAMBADA	accuracy (+)	63.24%	59.23%	95%+
LAMBADA	perplexity (-)	8.6	99	~1-2
Children's Book Test Common Nouns (validation accuracy)	accuracy (+)	93.30%	85.7%	96%
Children's Book Test Named Entities (validation accuracy)	accuracy (+)	89.05%	82.3%	92%
Penn Tree Bank	perplexity (-)	35.76	46.54	unknown
WikiText-2	perplexity (-)	18.34	39.14	unknown
enwik8	bits per character (-)	0.93	0.99	unknown
text8	bits per character (-)	0.98	1.08	unknown
WikiText-103	perplexity (-)	17.48	18.3	unknown

Less focus on fine-tuning; more focus on zero-shot prediction

GPT2

**SYSTEM PROMPT
(HUMAN-WRITTEN)**

Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today.

**MODEL COMPLETION
(MACHINE-WRITTEN,
SECOND TRY)**

The 19-year-old singer was caught on camera being escorted out of the store by security guards.

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back.

Scroll down for video

Shoplifting: Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today (pictured)

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back

The singer was also wearing a pair of black-rimmed glasses, a black jacket, black jeans and black sandals.

She was carrying a pair of black and white striped gloves and a small black bag.

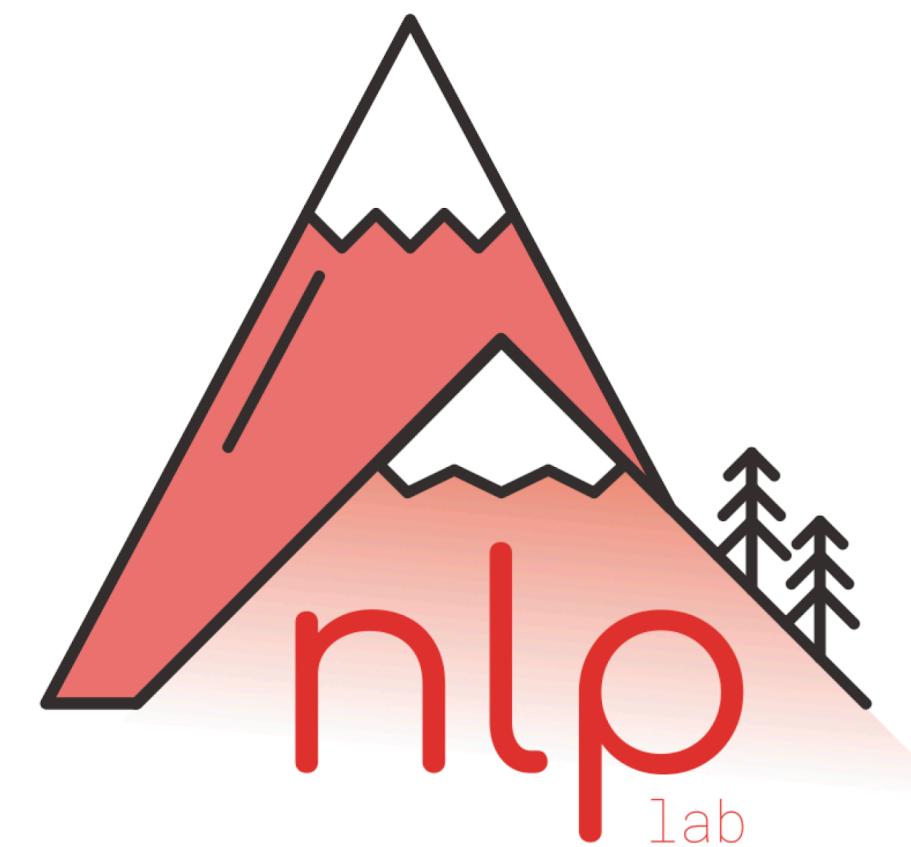
Increase in Scale

- GPT (Jun 2018): 117 million parameters, trained on 13GB of data (~1 billion tokens)
- GPT2 (Feb 2019): 1.5 billion parameters, trained on 40GB of data
- GPT3 (July 2020): 175 billion parameters, ~500GB data (~300 billion tokens)
- GPT4 (March 2023): 1.8T parameters (rumoured), ????? data

More to come in Weeks 7, 11!

Pretrained Representations: **ELMo, BERT**

Antoine Bosselut



Question

**Was GPT the first large-scale pretrained
neural representation?**

Word Embeddings: word2vec!

Question

What's an issue with word embeddings?

- 1) Chico Ruiz made a spectacular **play** on Alusik's grounder {. . . }
- 2) Olivia De Havilland signed to do a Broadway **play** for Garson {. . . }
- 3) Kieffer was commended for his ability to hit in the clutch , as well as his all-round excellent **play** {. . . }
- 4) {. . . } they were actors who had been handed fat roles in a successful **play** {. . . }
- 5) Concepts **play** an important role in all aspects of cognition {. . . }

Question

Words have different meanings in different contexts!

- 1) Chico Ruiz made a spectacular **play** on Alusik's grounder {. . . }
- 2) Olivia De Havilland signed to do a Broadway **play** for Garson {. . . }
- 3) Kieffer was commended for his ability to hit in the clutch , as well as his all-round excellent **play** {. . . }
- 4) {. . . } they were actors who had been handed fat roles in a successful **play** {. . . }
- 5) Concepts **play** an important role in all aspects of cognition {. . . }

Question

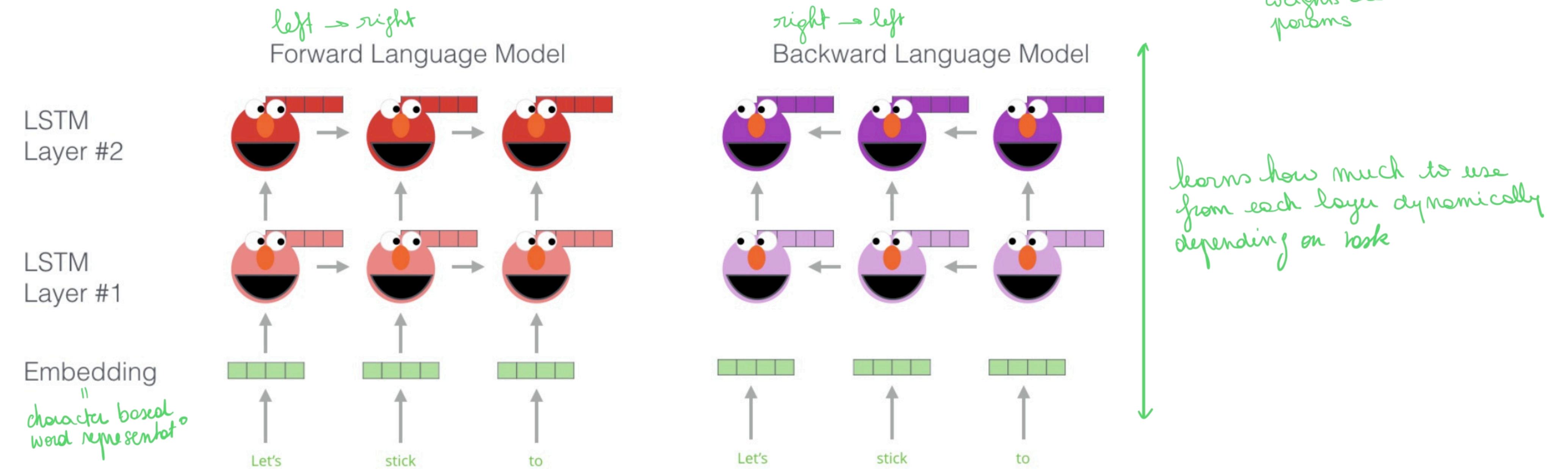
How might we integrate contextual information into word embeddings?

2018: Use bidirectional LSTMs!

ELMo

some vector to word regardless of its meaning
 \neq word2Vec, GloVe = static word representation
 $=$ dynamic, context-dependent embeddings
 provides \neq embeddings for same word in \neq contexts

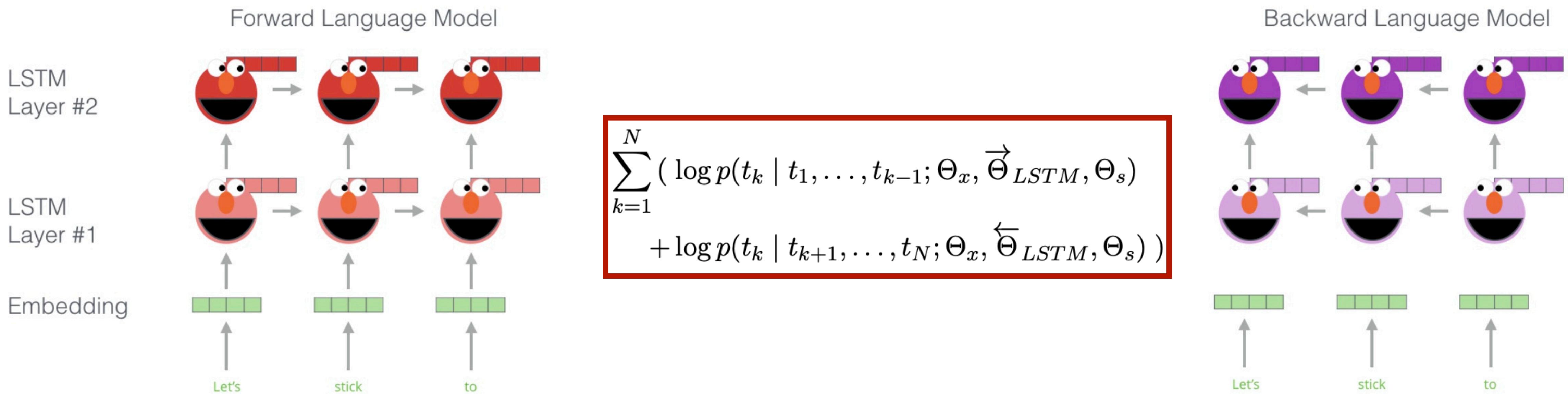
- Train two-layer LSTM-based language model on a **large** corpus
- Use **hidden states** of the LSTMs for each token to compute an embedding of each word
- LSTM should be bidirectional



1) PRETRAINED on large text corpus
2) FINETUNED for specific task by plugging them into
models for sentiment analysis, answering ...

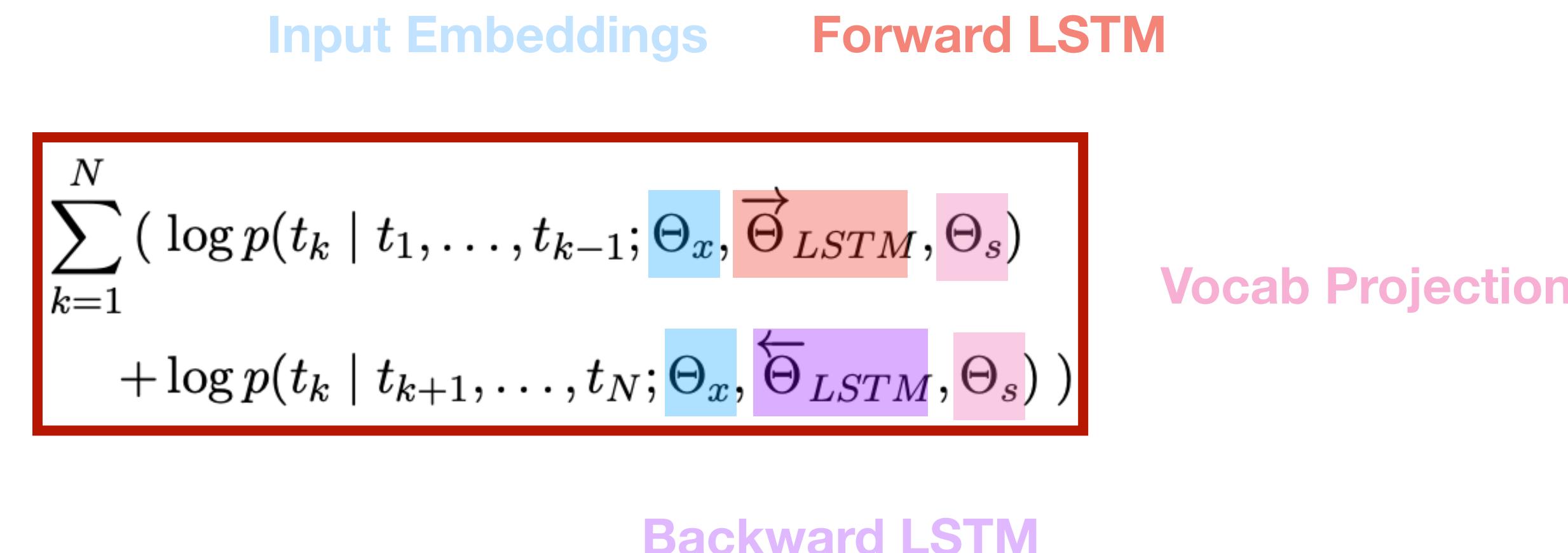
ELMo

- Train two-layer LSTM-based language model on a **large** corpus
- Use **hidden states** of the LSTMs for each token to compute an embedding of each word
- LSTM should be bidirectional

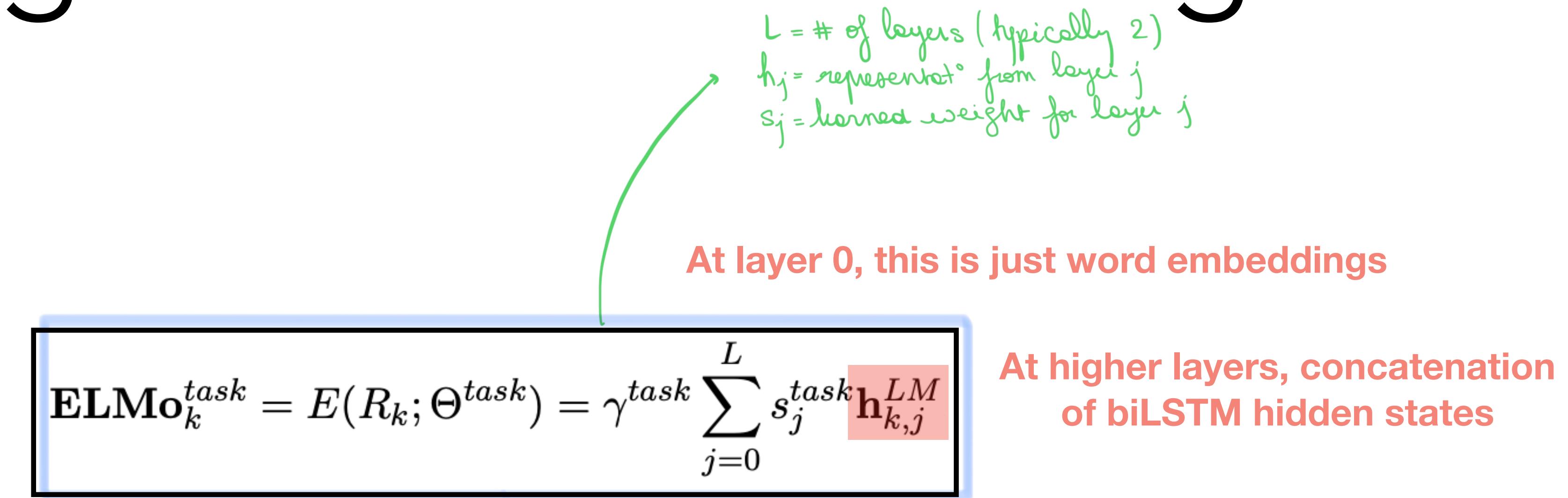


ELMo

- Train two-layer LSTM-based language model on a **large** corpus
- Use **hidden states** of the LSTMs for each token to compute an embedding of each word
- LSTM should be bidirectional
- Use 1B word benchmark (**single sentences** — why might this be a problem?)



Using ELMo Embeddings



- γ^{task} : allows the task model to scale the entire ELMo vector
- s_j^{task} : softmax-normalized weights across layers

Learn both of these parameters

Why average the representation at each layer as opposed to the final one?
For different tasks, useful representations may be at different layers

ELMo Improvements

TASK	PREVIOUS SOTA	OUR BASELINE	ELMo + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8
SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17
SRL	He et al. (2017)	81.7	81.4	84.6
Coref	Lee et al. (2017)	67.2	67.2	70.4
NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10
SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5

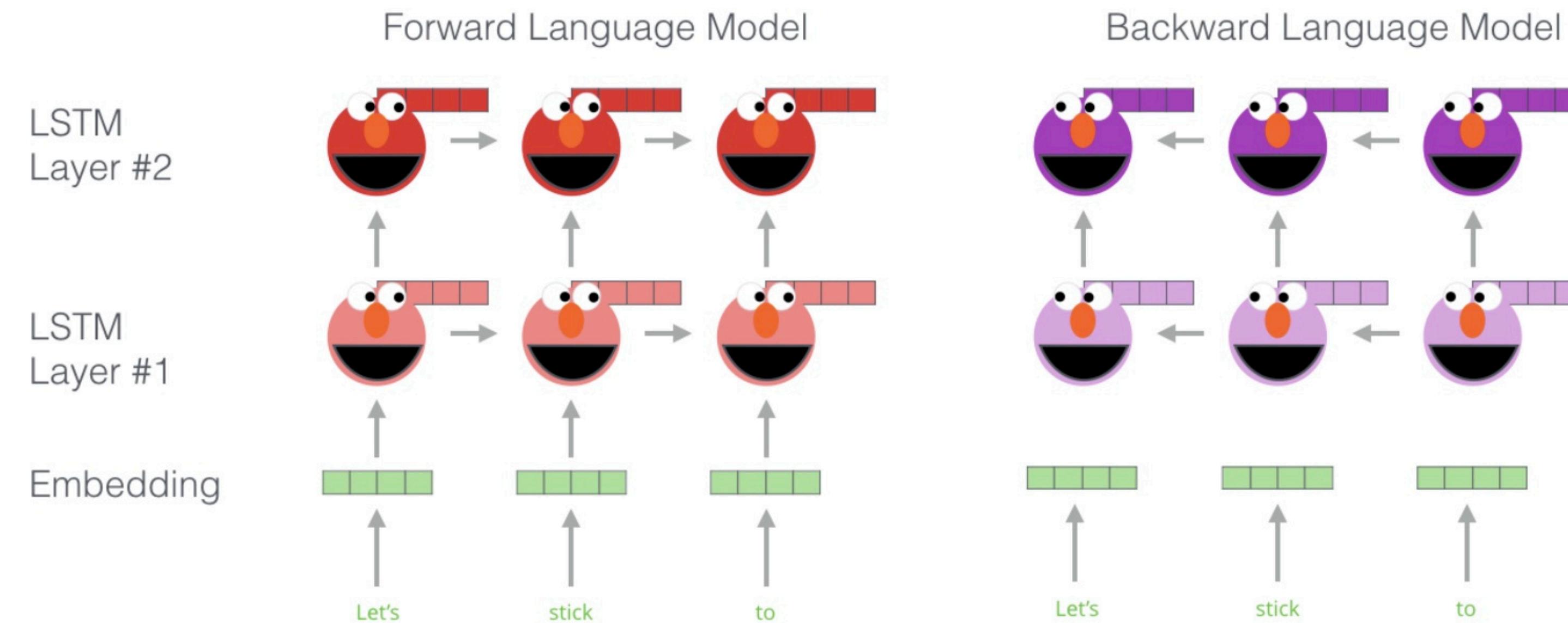
Across the board improvements over SOTA when introduced

Question

**What's a problem with ELMo's notion of
bidirectionality?**

ELMo Issue

- ELMo (and bidirectional LSTMs / RNNs, in general) are **unidirectional** models masquerading as **bidirectional**
- Separate language model encodes the forward and backward sequence (and the representations are concatenated)

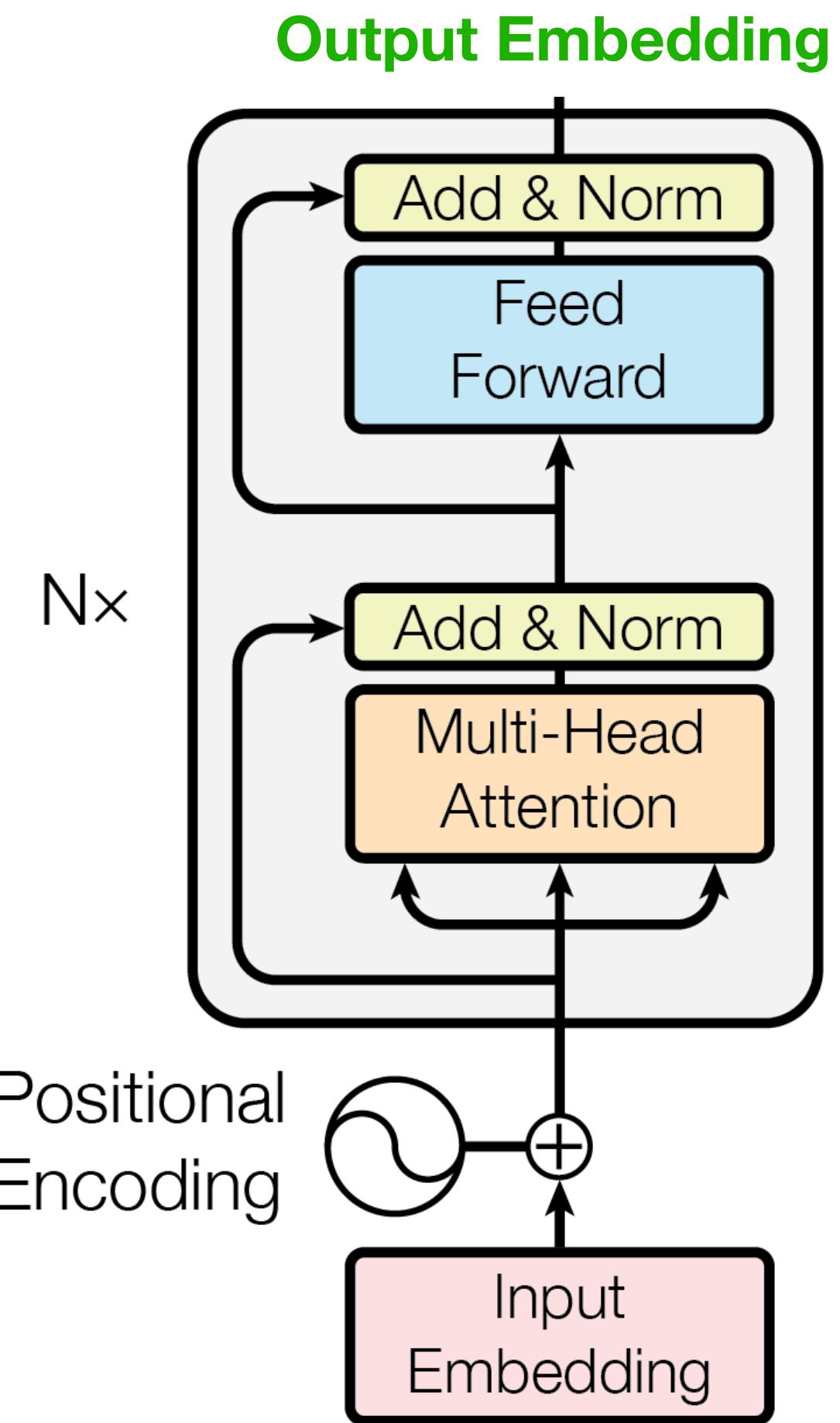
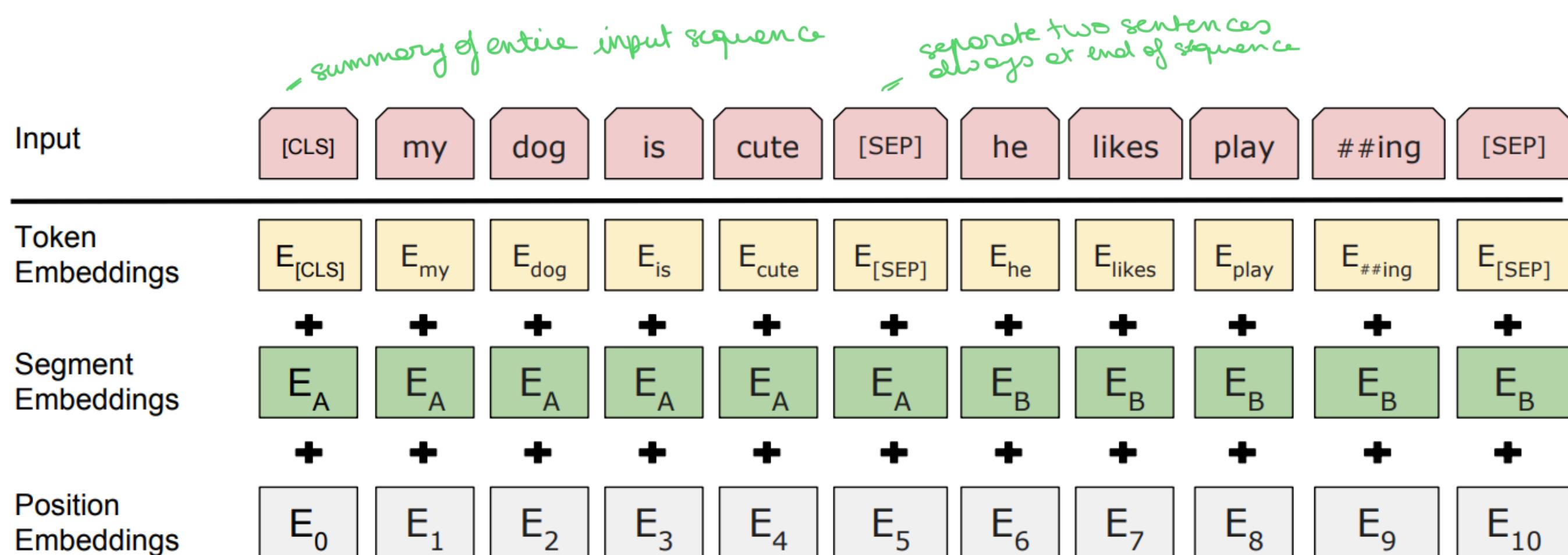


What is BERT ?



BERT Architecture

- Transformer Encoder as we saw previously!
- BERT-Base: 12 layers, $d=768$, 12 heads.
- Total params = **110M**
- BERT-Large: 24 layers, $d=1024$, 16 heads.
Total params = **340M**



- Input embeddings for $V = 30k$ word pieces
- Positional & segment embeddings

How is BERT trained?



Pretraining: Before

(Causal, Left-to-right)
Language Modeling

*I really enjoyed the movie we
watched on _____*



OpenAI

Pretraining: Two Approaches

(Causal, Left-to-right)
Language Modeling

*I really enjoyed the movie we
watched on _____*

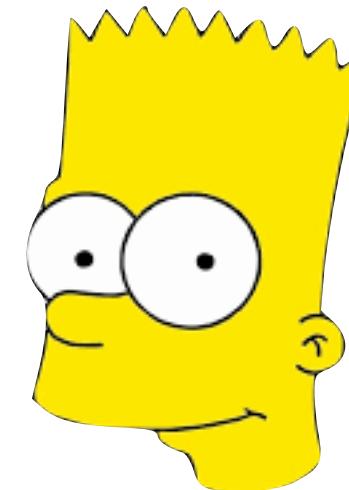


OpenAI

(Radford et al., 2018, 2019, many others)

Masked
Language Modeling

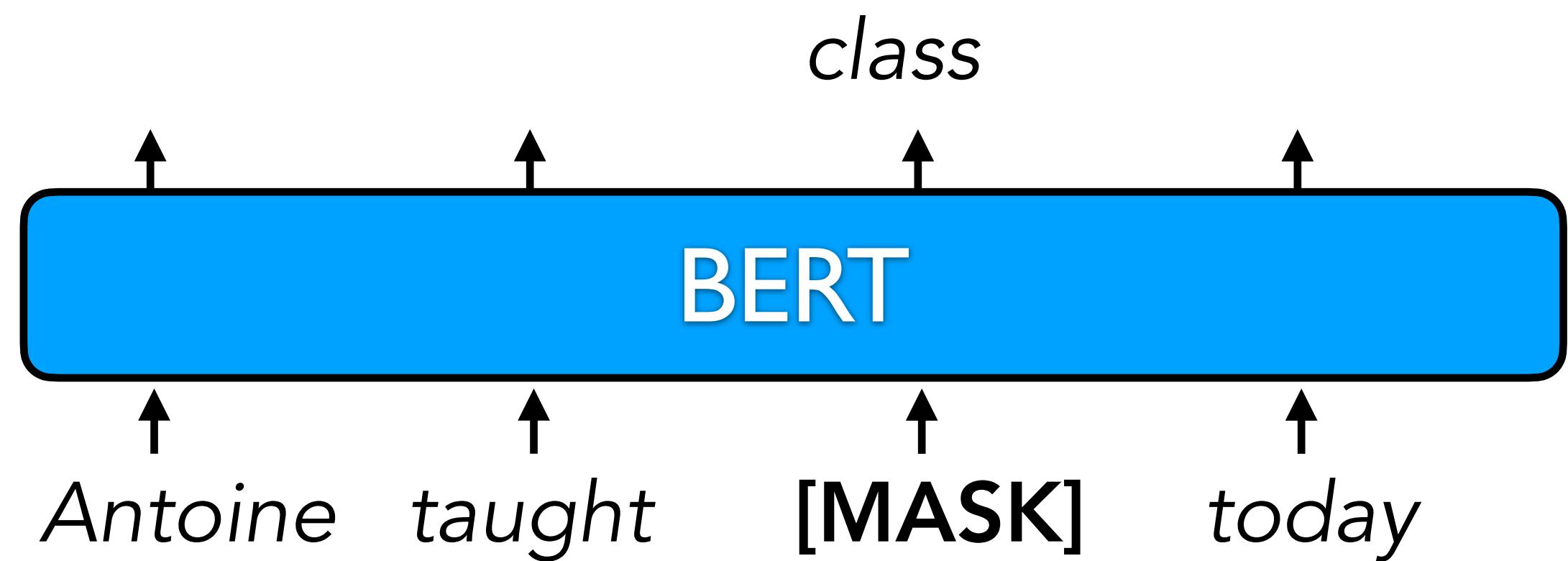
*I really enjoyed the _____ we
watched on Saturday!*



(Devlin et al., 2018; Liu et al., 2020)

Masked Language Modeling (BERT)

- **Pretraining (self-supervised learning):**
 - Done at scale on natural occurring sequences of text (any large corpus of raw text)

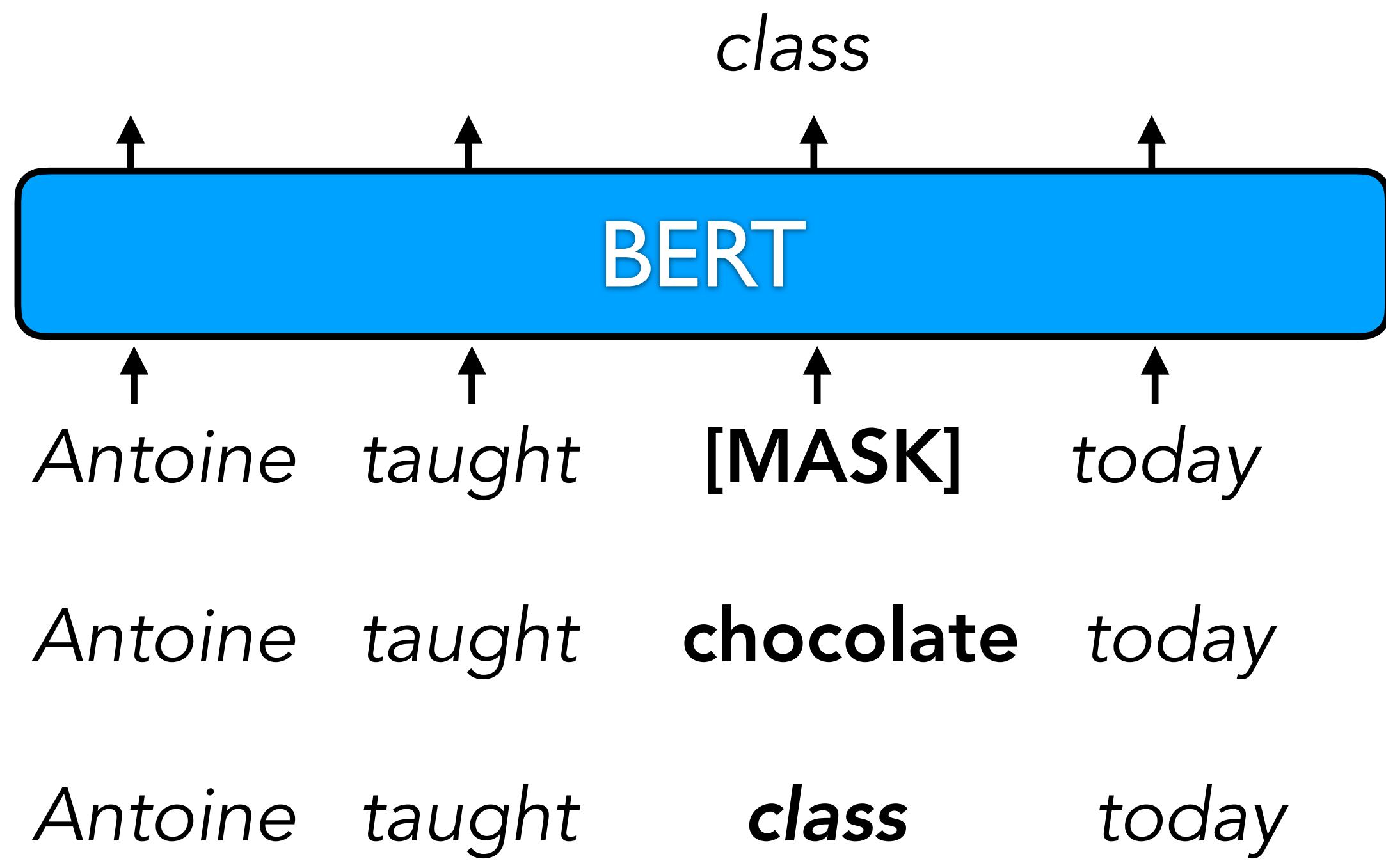


Masked Language Modeling (BERT)

- **Pretraining:** take a sequence of text, and predict 15% of the tokens

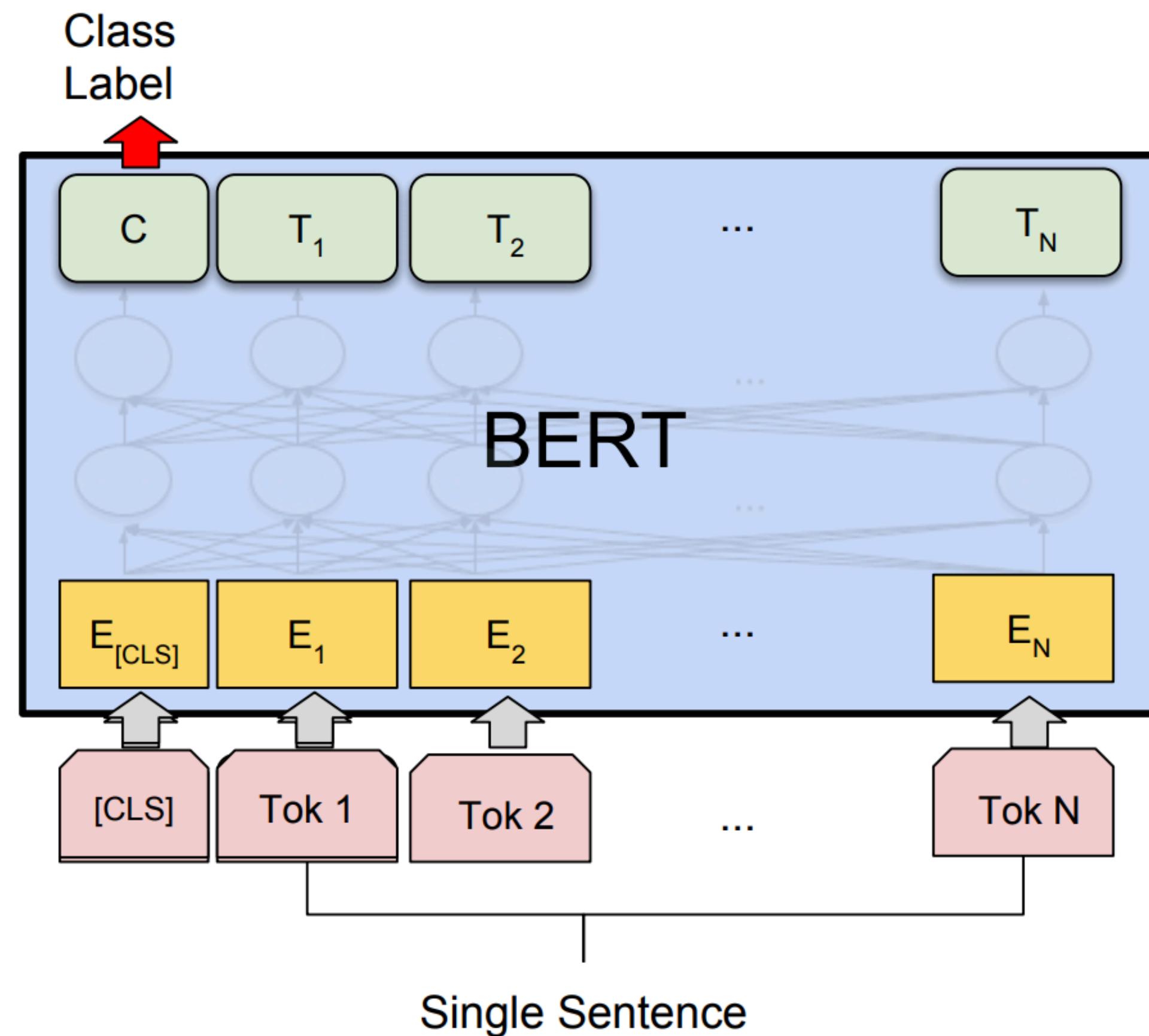
- **For 15% of tokens:**

- Replace input token with [MASK] (80% of predictions)
- Replace input token with a random token (10% of predictions)
- Keep the same input token (10% of predictions)



Fine-tuning BERT

Fine-tuning BERT for classification

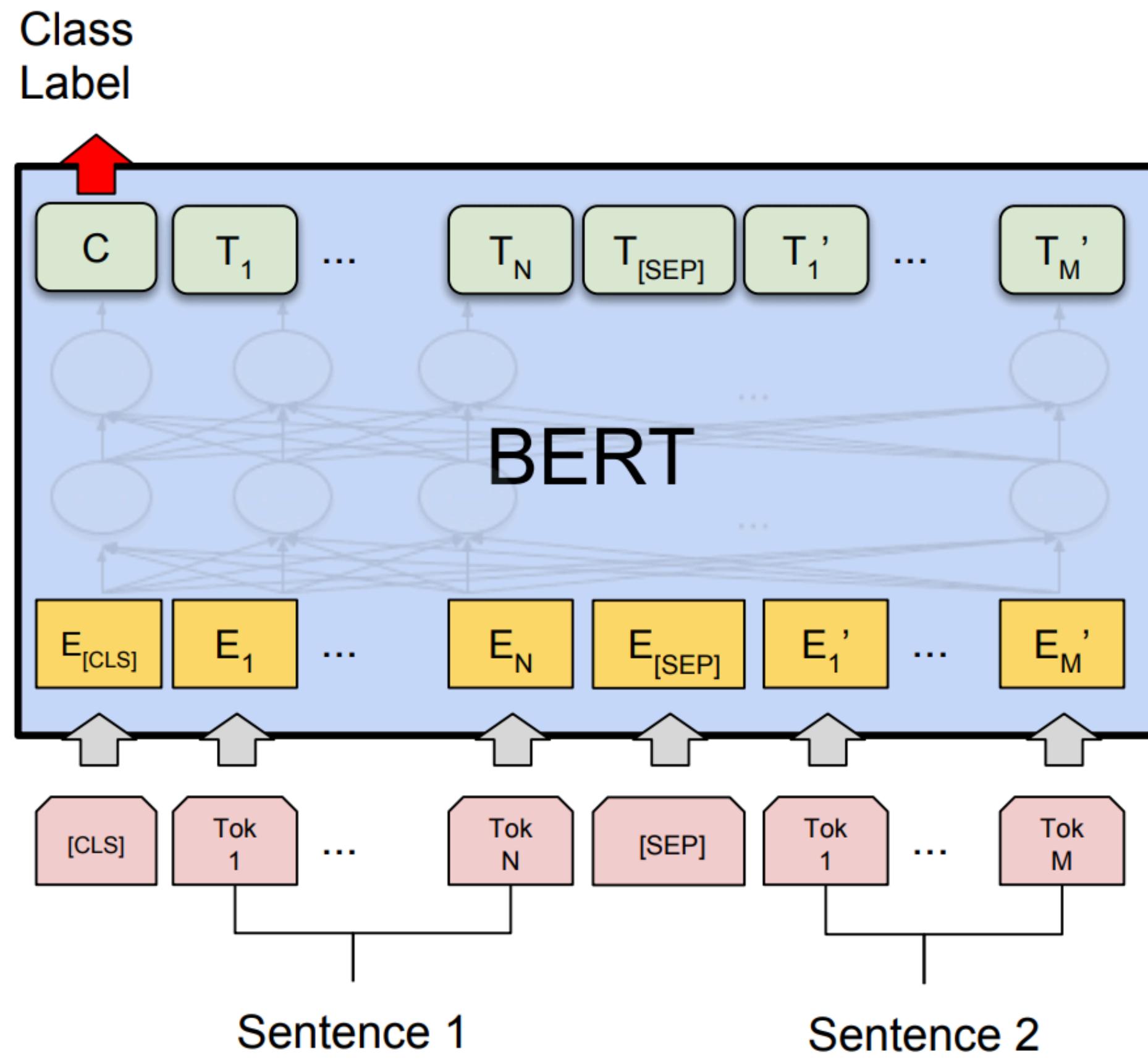


- **Fine-tuning:** Done after BERT has been pretrained (no more pretraining objectives)
- **During fine-tuning,** we update the parameters of the BERT model to learn the task
- A **contextual embedding** is outputted for each **token** in the input
- Prepend a special **token** **[CLS]** to the front of the sequence
- **Learn to classify** the **output embedding** for this **token**

How do we classify the output embedding for this token? Logistic Regression!

Fine-tuning BERT for classification

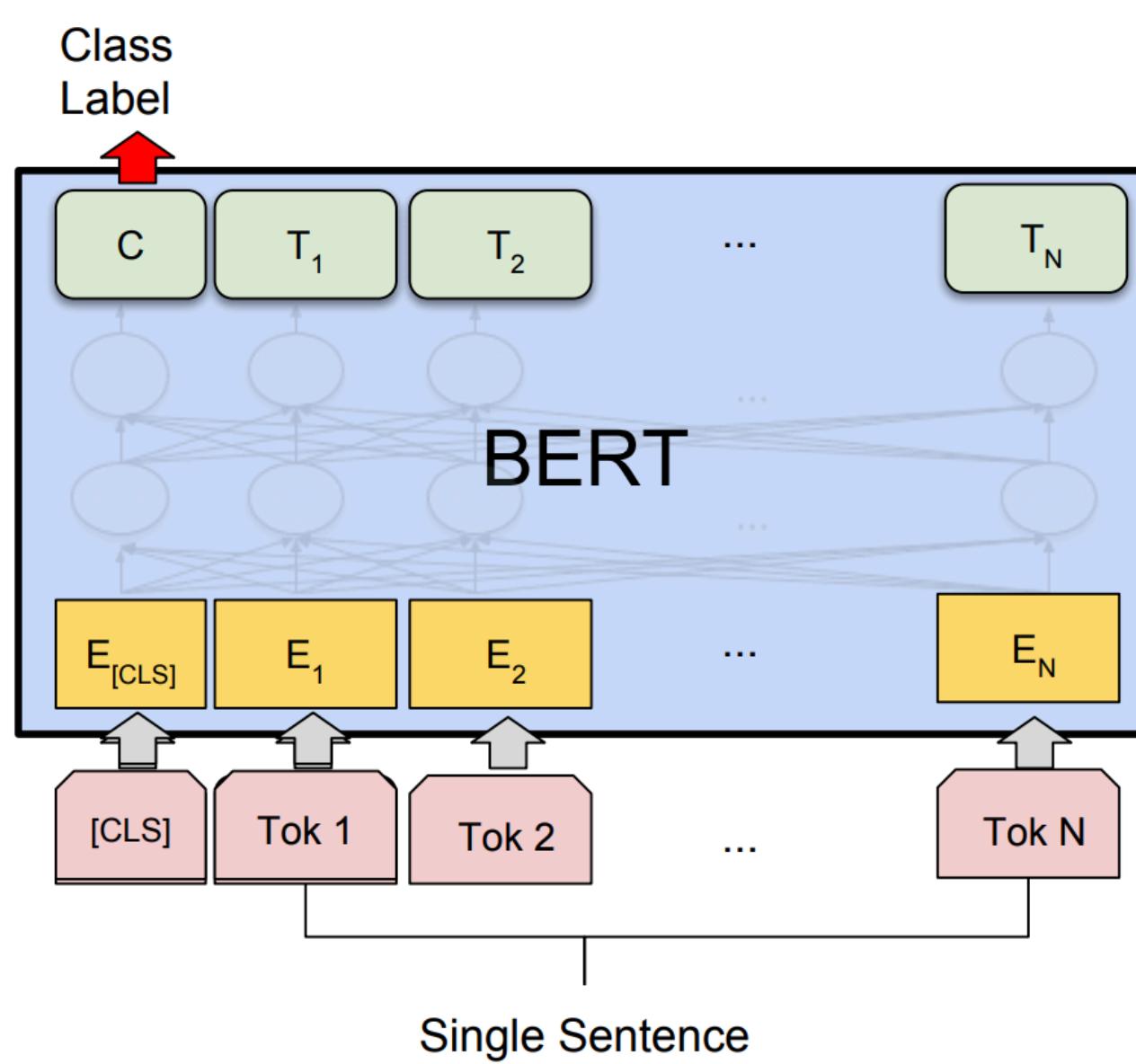
After pretraining



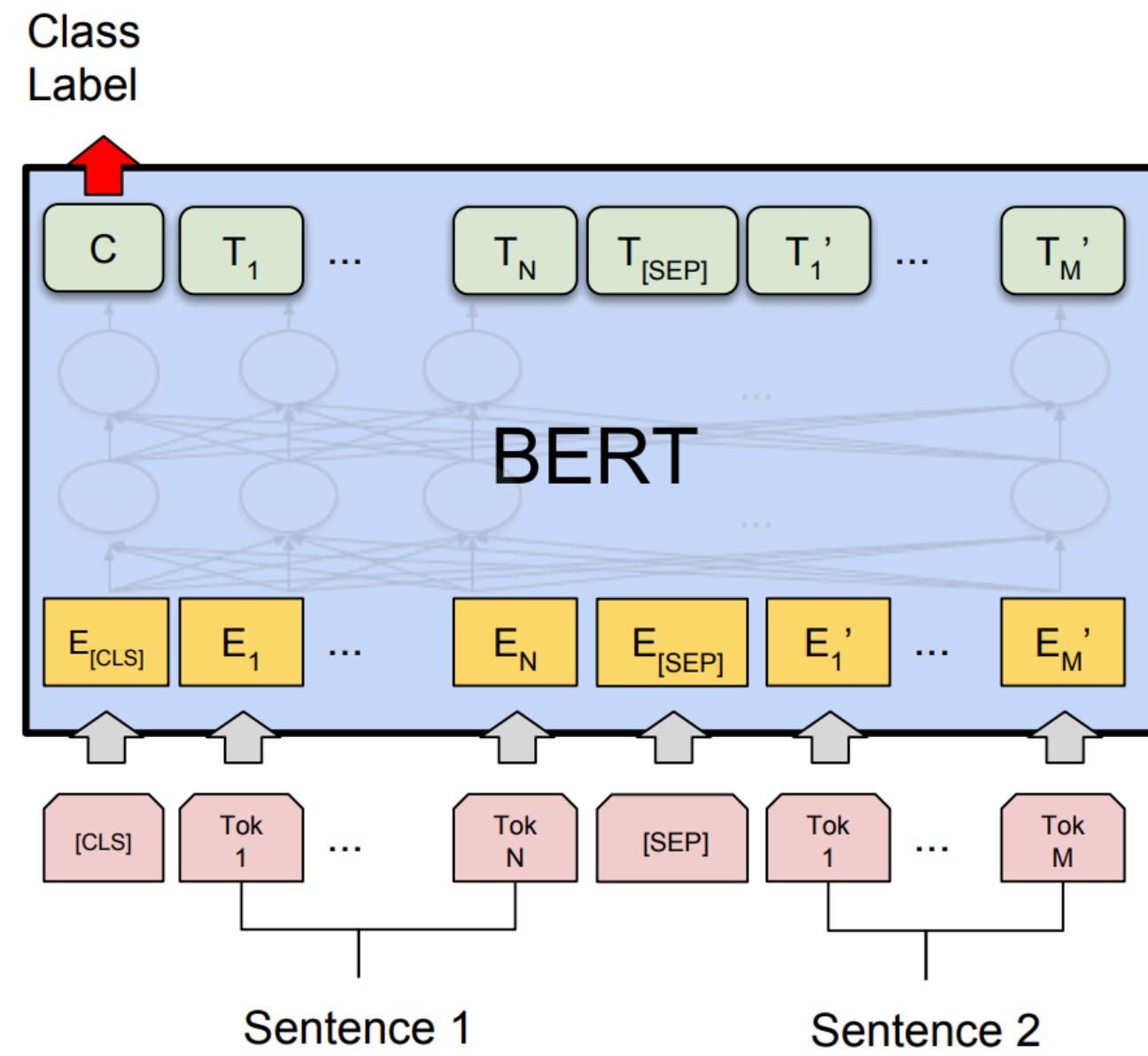
- A **contextual embedding** is outputted for each **token** in the input
- Prepend a special **token** **[CLS]** to the front of the sequence
- **Learn to classify the output embedding** for this **token**
- Separate sequences with special **[SEP]** token

Can add special meta-tokens your vocabulary when they're needed!

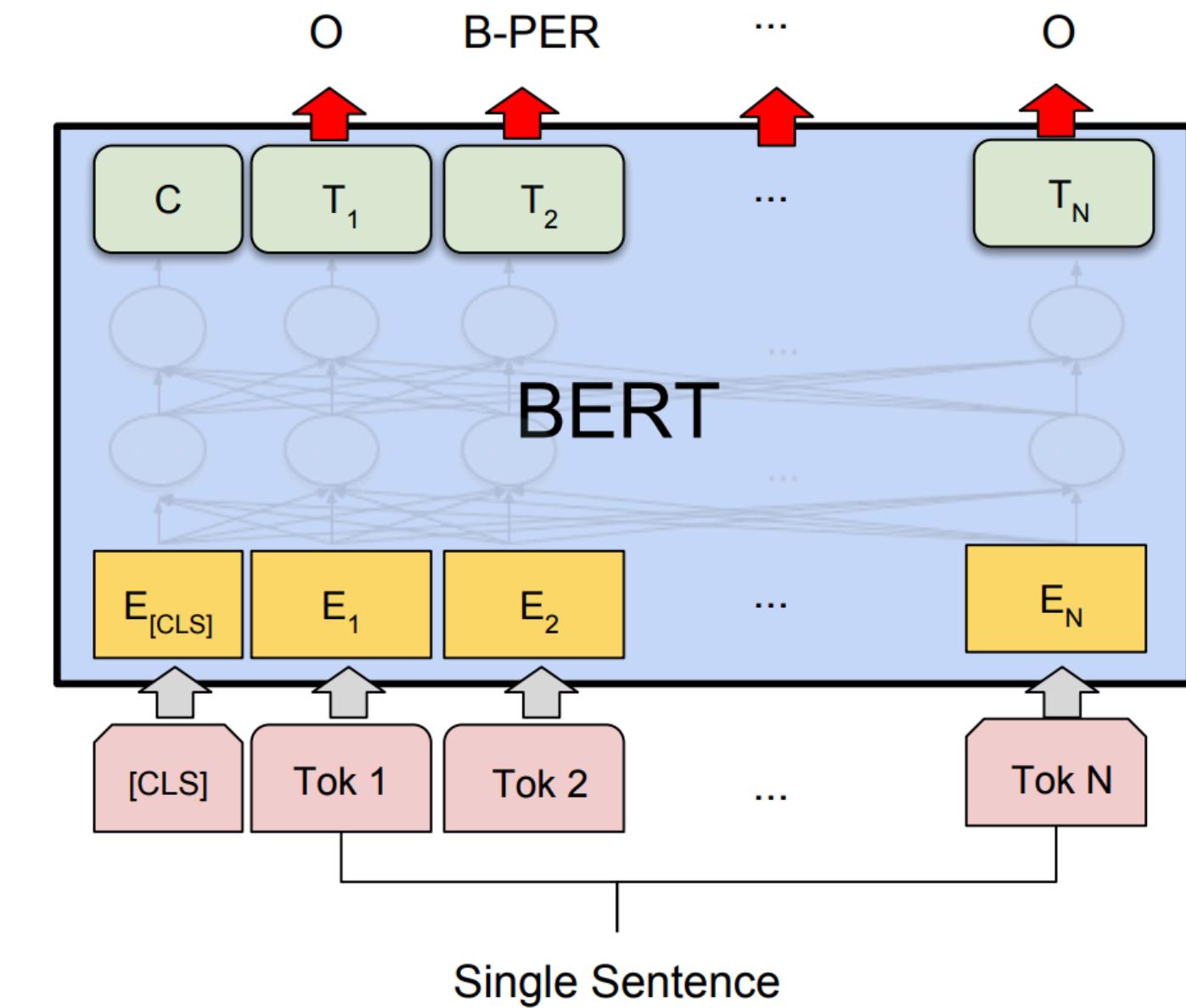
Single model starting point for many tasks



(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

- Re-using the same pretrained BERT model for fine-tuning on many tasks:
 - **Classification:** Take [CLS] output embedding as input features to classification model
 - **Sequence labeling:** Take output embedding for each token and classify individually

Question

Why do we put the [CLS] token at the front of the sequence?

Easiest place to put it. Bidirectionality ensures it attends to representations of all other tokens

GLUE: Prototypical NLU Benchmark

Corpus	Train	Test	Task	Metrics	Domain
Single-Sentence Tasks					
CoLA	8.5k	1k	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	391k	paraphrase	acc./F1	social QA questions
Inference Tasks					
MNLI	393k	20k	NLI	matched acc./mismatched acc.	misc.
QNLI	105k	5.4k	QA/NLI	acc.	Wikipedia
RTE	2.5k	3k	NLI	acc.	news, Wikipedia
WNLI	634	146	coreference/NLI	acc.	fiction books

BERT on GLUE

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 _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

For each of these tasks, a different BERT model is fine-tuned on the task data

Not the same fine-tuned BERT model that gets the same performance

BERT on GLUE

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 _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

Performance increases across the board

BERT on GLUE

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 _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

Performance increases across the board

Big increase over OpenAI GPT highlights importance of bidirectionality

Ablation Studies

Tasks	Dev Set				
	MNLI-m (Acc)	QNLI (Acc)	MRPC (Acc)	SST-2 (Acc)	SQuAD (F1)
BERT _{BASE}	84.4	88.4	86.7	92.7	88.5
No NSP	83.9	84.9	86.5	92.6	87.9
LTR & No NSP	82.1	84.3	77.5	92.1	77.8
+ BiLSTM	82.1	84.1	75.7	91.6	84.9

Table 5: Ablation over the pre-training tasks using the BERT_{BASE} architecture. “No NSP” is trained without the next sentence prediction task. “LTR & No NSP” is trained as a left-to-right LM without the next sentence prediction, like OpenAI GPT. “+ BiLSTM” adds a randomly initialized BiLSTM on top of the “LTR + No NSP” model during fine-tuning.

#L	#H	#A	Dev Set Accuracy			
			LM (ppl)	MNLI-m	MRPC	SST-2
3	768	12	5.84	77.9	79.8	88.4
6	768	3	5.24	80.6	82.2	90.7
6	768	12	4.68	81.9	84.8	91.3
12	768	12	3.99	84.4	86.7	92.9
12	1024	16	3.54	85.7	86.9	93.3
24	1024	16	3.23	86.6	87.8	93.7

Table 6: Ablation over BERT model size. #L = the number of layers; #H = hidden size; #A = number of attention heads. “LM (ppl)” is the masked LM perplexity of held-out training data.

Importance of bidirectionality

More parameters = better!

Question

Should we use BERT to generate embeddings or fine-tune the full model?

Fine-tuning vs. Embeddings

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

- BERT outputs embeddings that can be frozen and provided to a different model, but BERT performs better when fully fine-tuned

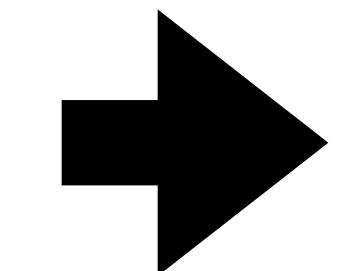
Transfer Learning

Pretraining

Learn embeddings that can be used to seed a downstream model (ELMo)

-or-

Learn a model that can be fine-tuned for many downstream tasks (GPT, BERT)



Fine-tuning

Design a new model architecture whose embeddings are initialised with pretrained embeddings. Train this model on a task of interest

- or -

Take a pretrained model and train it further on data from a task of interest

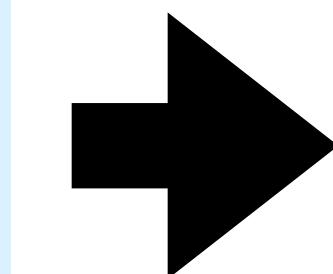
Transfer Learning

Pretraining

Learn embeddings that can be used to seed a downstream model (ELMo)

-or-

Learn a model that can be fine-tuned for many downstream tasks (GPT, BERT)



Fine-tuning

Design a new model architecture whose embeddings are initialised with pretrained embeddings. Train this model on a task of interest

- or -

Take a pretrained model and train it further on data from a task of interest

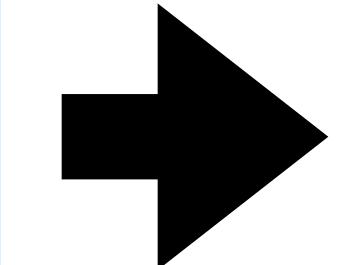
Transfer Learning

Pretraining

Learn embeddings that can be used to seed a downstream model (ELMo)

-or-

Learn a model that can be fine-tuned for many downstream tasks (GPT, BERT)



Fine-tuning

Design a new model architecture whose embeddings are initialised with pretrained embeddings. Train this model on a task of interest

- or -

Take a pretrained model and train it further on data from a task of interest

Transfer Learning

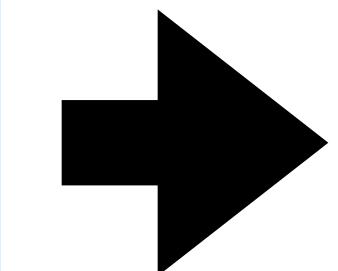
Pretraining

Uses simple training objectives

Requires tons of data

Resultant model often not useful yet

Slow & expensive; can often only do once



Fine-tuning

Done on smaller datasets

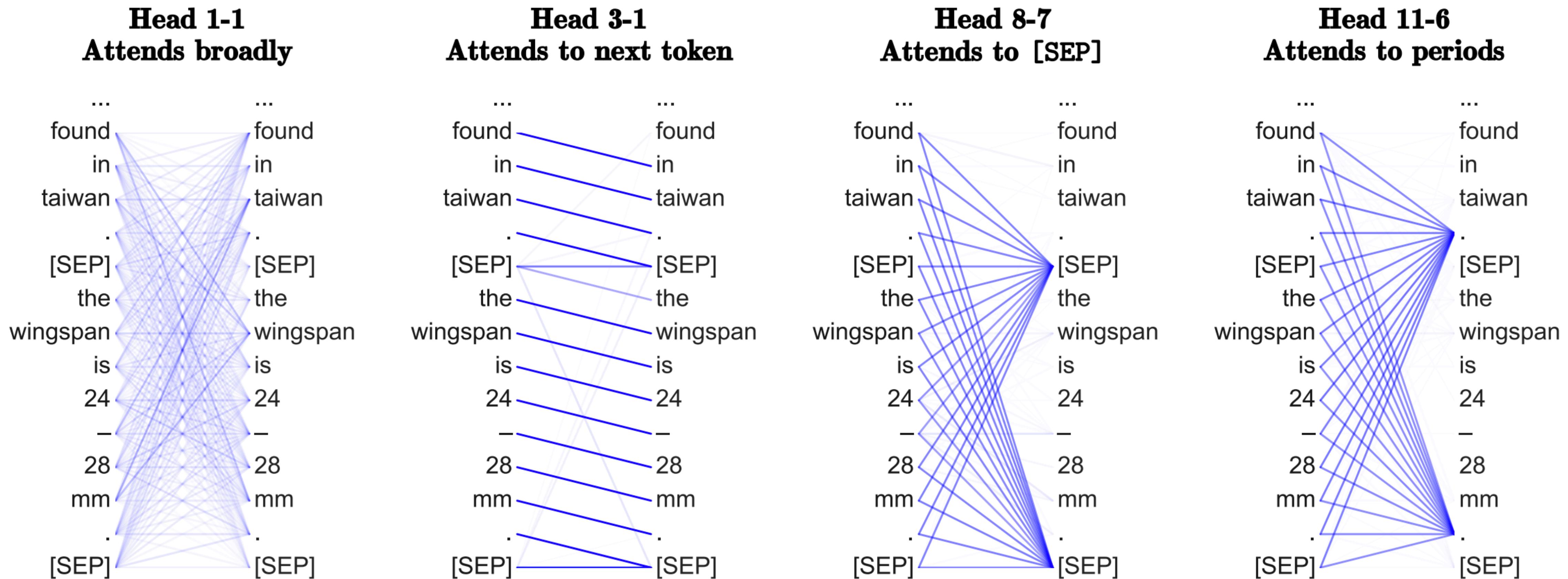
Trained on data with a more complex structure

Resultant model applied to task of interest

Typically cheaper; can afford multiple runs, hyper parameter tuning, etc.

What does BERT Learn?

What does BERT learn?

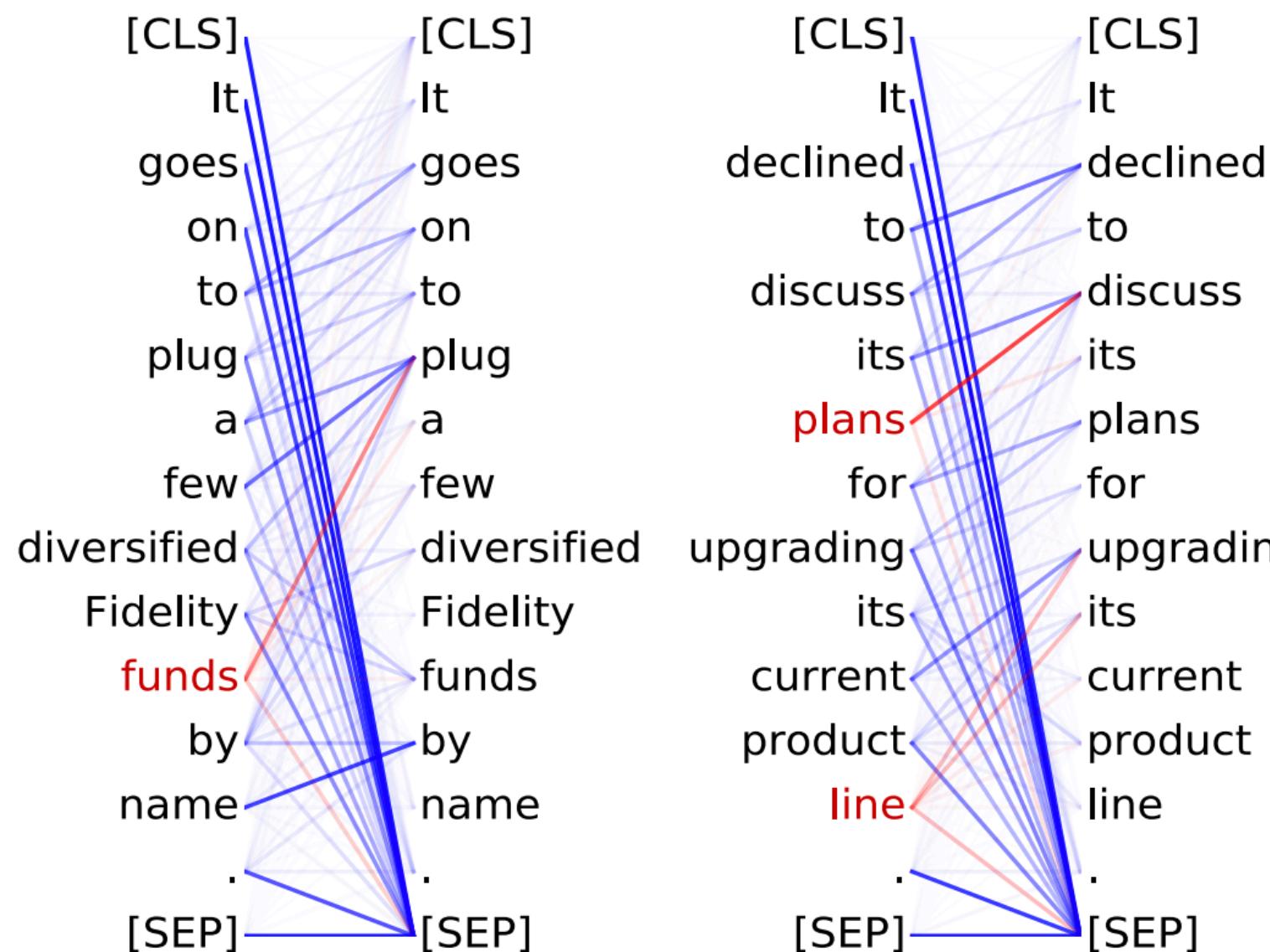


Transformer heads learn diverse concepts that map to positional, semantic, and syntactic relationships

What does BERT learn?

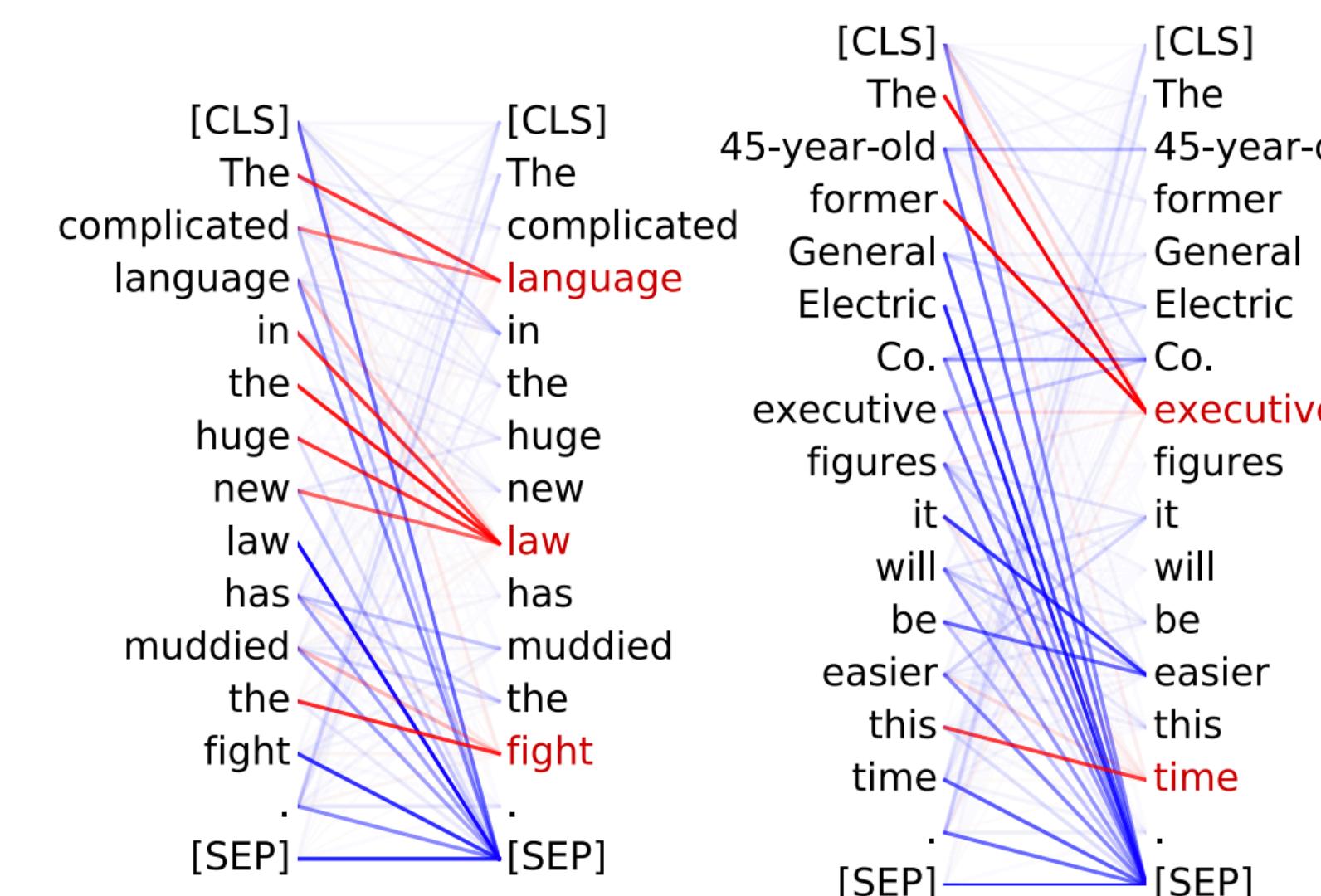
Head 8-10

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



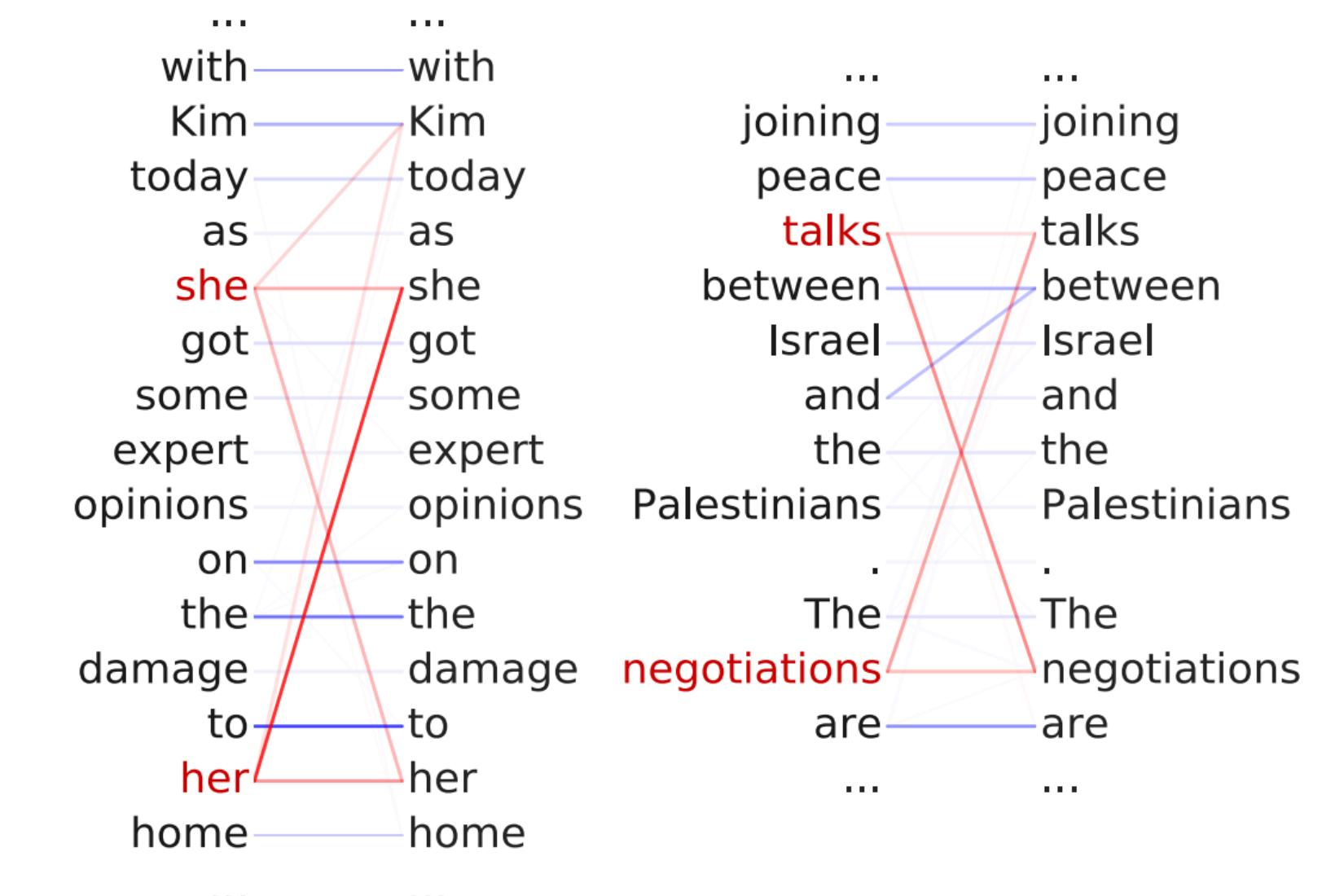
Head 8-11

- 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



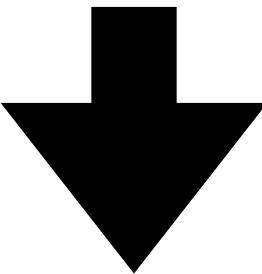
Transformer heads learn diverse concepts that map to positional, semantic, and syntactic relationships

Improvements to BERT

Feature	ELMo	BERT	GPT
Architecture	biLSTM	Transformer (encoder)	Transformer (decoder)
Contextualized Embeddings?	<input checked="" type="checkbox"/> Yes	<input checked="" type="checkbox"/> Yes	<input checked="" type="checkbox"/> Yes
Directionality	Bidirectional (LSTM)	Bidirectional (Self-Attention)	Unidirectional (Left-to-Right)
Pretraining Task	Bi-directional LM	Masked LM (MLM) + NSP	Causal LM (Next-word prediction)
Output Type	Word embeddings for each token, useful for downstream tasks	Sentence representation for classification & per-token embeddings	Continuously generates text token by token

Whole word masking

Obama was the president of the United States in 2010



[MASK] _bama _was _the _president _of _the _United _States _in _2010

vs.

[MASK] [MASK] _was _the _president _of _the _United _States _in _2010

Why might whole word masking be important?

**Too easy to predict masked subwords if rest of word is in context
— model doesn't learn as well**

RoBERTa

- “Robustly Optimised BERT” — a collection of improvements to BERT
- Same architecture as BERT
- 160 GB of training data, rather than only 13 GB in BERT

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	94.6/89.4	90.2	96.4
BERT_{LARGE}						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7

DistilBERT

- Do we need all parameters of BERT, which require lots
- What if BERT was a much smaller?
- Train with distillation over soft target probabilities of BERT (and MLM)

$$\mathcal{L}_{distil} = - P_{BERT}(y_t^* | [M], \{y_s^*\}_{s \neq t}) \log P_{dbert}(y_t^* | [M], \{y_s^*\}_{s \neq t})$$

$$\mathcal{L}_{mlm} = - \log P_{dbert}(y_t^* | [M], \{y_s^*\}_{s \neq t})$$

$$\mathcal{L}_{tot} = \gamma_1 \mathcal{L}_{distil} + \gamma_2 \mathcal{L}_{mlm}$$

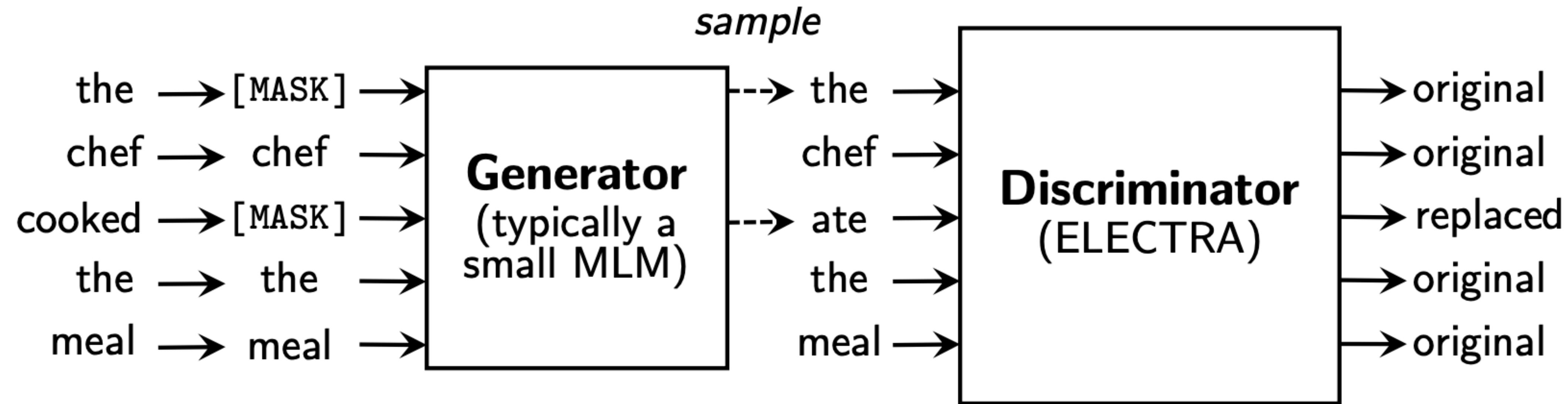
- Allows you to train much smaller DistilBERT with ~97% performance of BERT

How Does DistilBERT Work?

Instead of training a new model from scratch, DistilBERT "compresses" BERT using a technique called knowledge distillation:

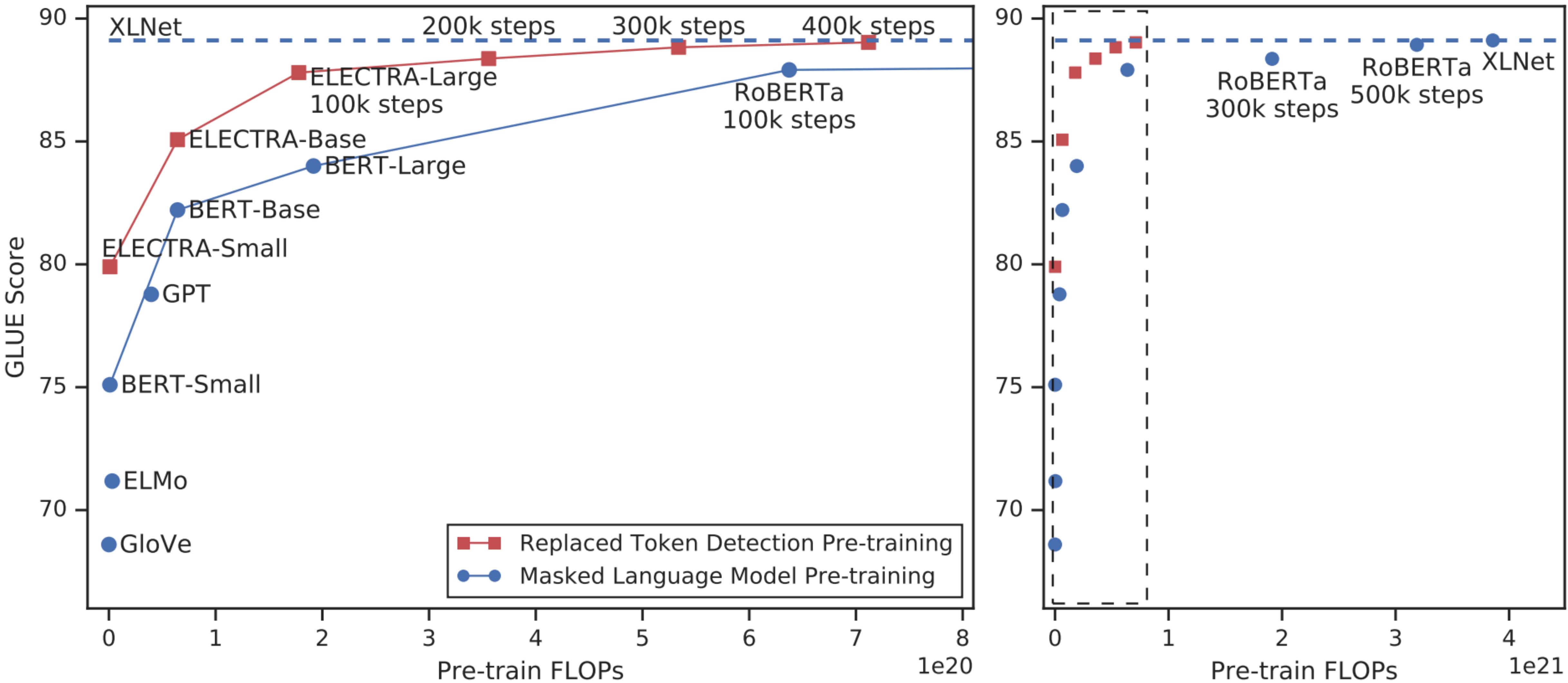
1. Teacher-Student Learning:
 - The full BERT (teacher model) trains a smaller DistilBERT (student model).
 - DistilBERT learns to mimic BERT's outputs, making it more efficient.
2. Removes Next Sentence Prediction (NSP):
 - BERT uses Masked Language Modeling (MLM) + NSP, but DistilBERT only uses MLM (since NSP isn't that useful).
3. Reduces Model Size:
 - Uses half the number of layers (12 → 6).
 - Keeps attention weights similar to BERT's.

ELECTRA



- Recall: BERT only learns from 15% of masked tokens (**quite inefficient!**)
- Instead, predict whether a token is corrupted on the discriminator
- Learning from all tokens drastically speeds up training

ELECTRA



Question

**What was the main improvement of
BERT-style models over GPT?**

Bidirectionality allows for more expressive representations to be learned

Question

**What can BERT NOT due as a result of its
masked LM training objective?**

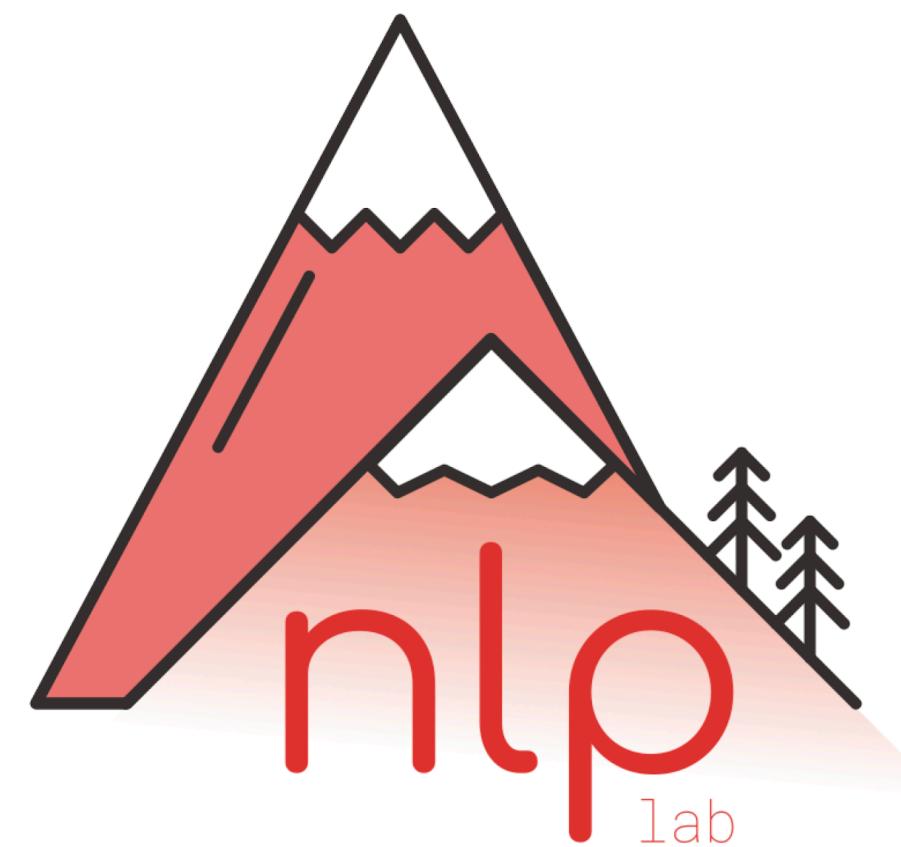
Generate text!

Recap

- **Contextualised embeddings:** Let us model words and sequences conditioned on the context around them
- **ELMo:** One of the first models for contextualized embeddings based on bidirectional LSTMs *good for pre-trained embeddings*
- **GPT:** First model for contextualised embeddings using transformer models
good for generating text as a language model
- **BERT:** Improving transformer-based contextual representations using masked language models and bidirectional encoding *good for classification + sequence labelling*
- Many variants of BERT in recent years!
- Much work on analysing information learned in contextual representations (Week 13)

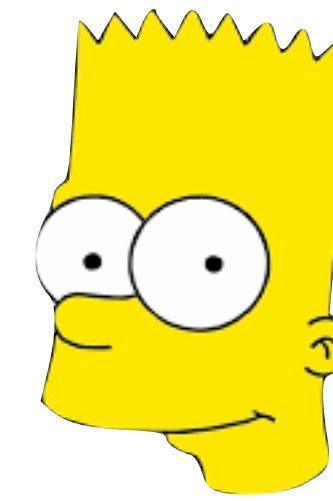
Pretrained Representations: **BART & T5**

Antoine Bosselut

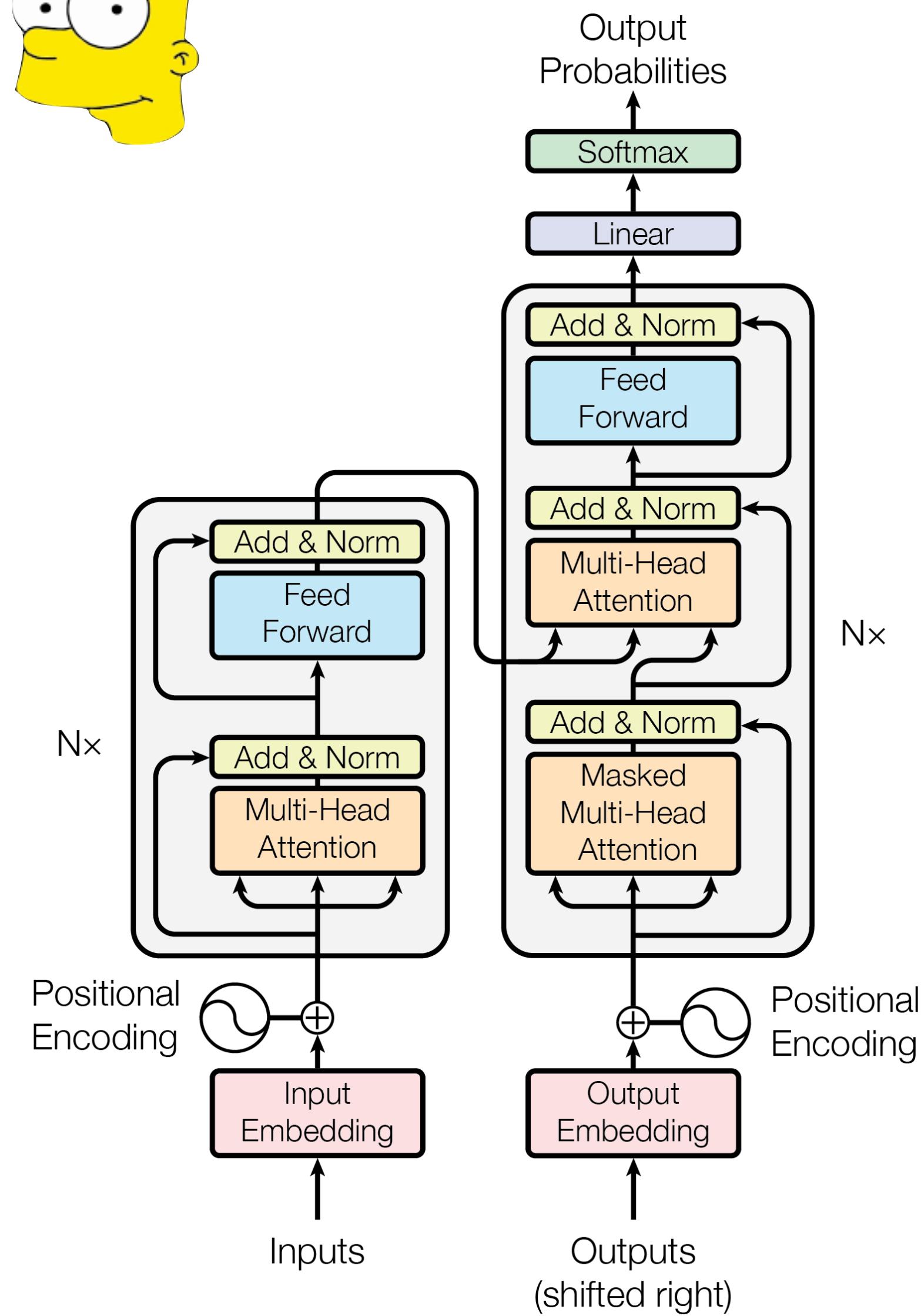


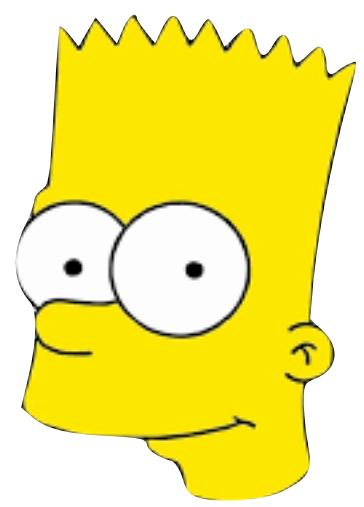
How should we pretrain
sequence-to-sequence models?

BART



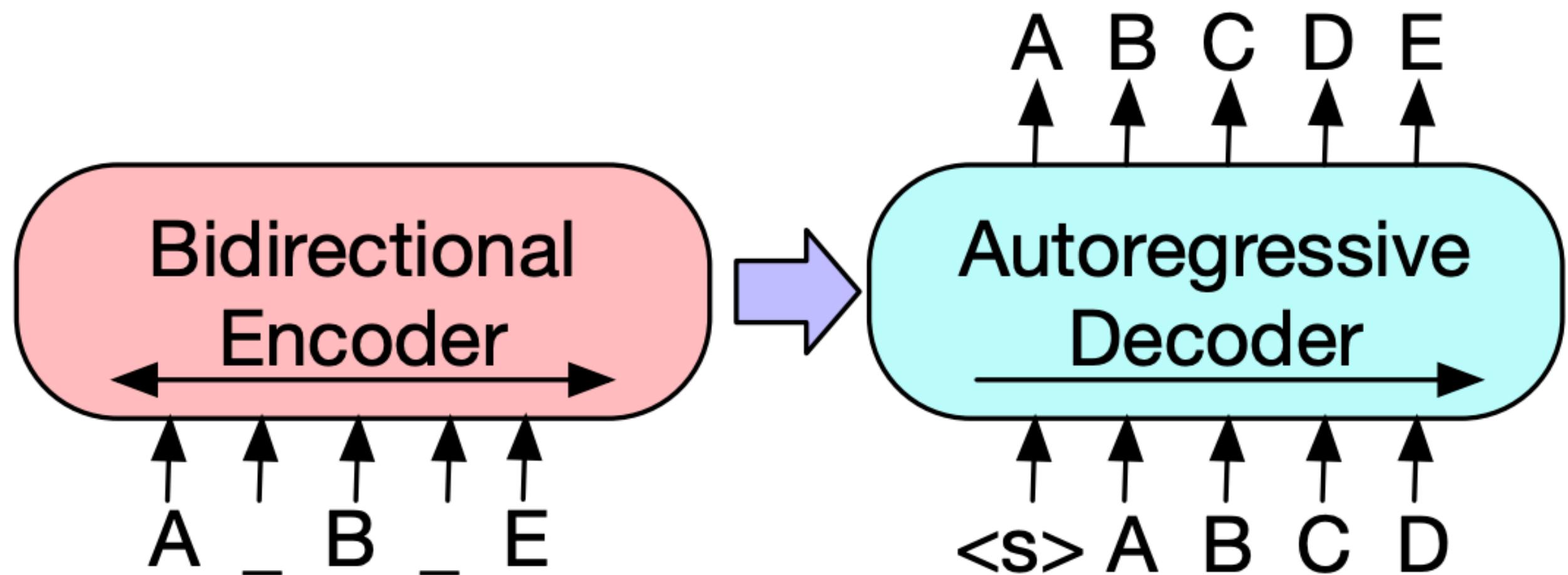
- Classic transformer architecture
- Bidirectional encoder feeds into autoregressive decoder
- Cross-attention layers in decoder are back!
- **BART-base**: 6-layers each in encoder and decoder; 140M parameters
- **BART-large**: 12 layers each in encoder and decoder; 400M parameters

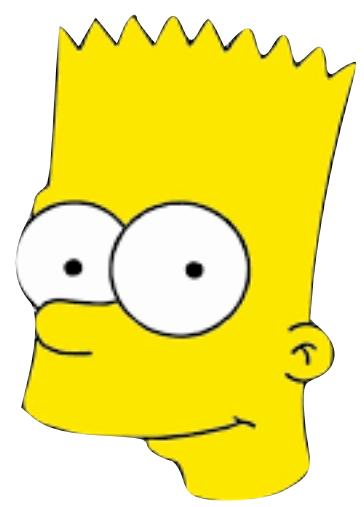




BART Pretraining

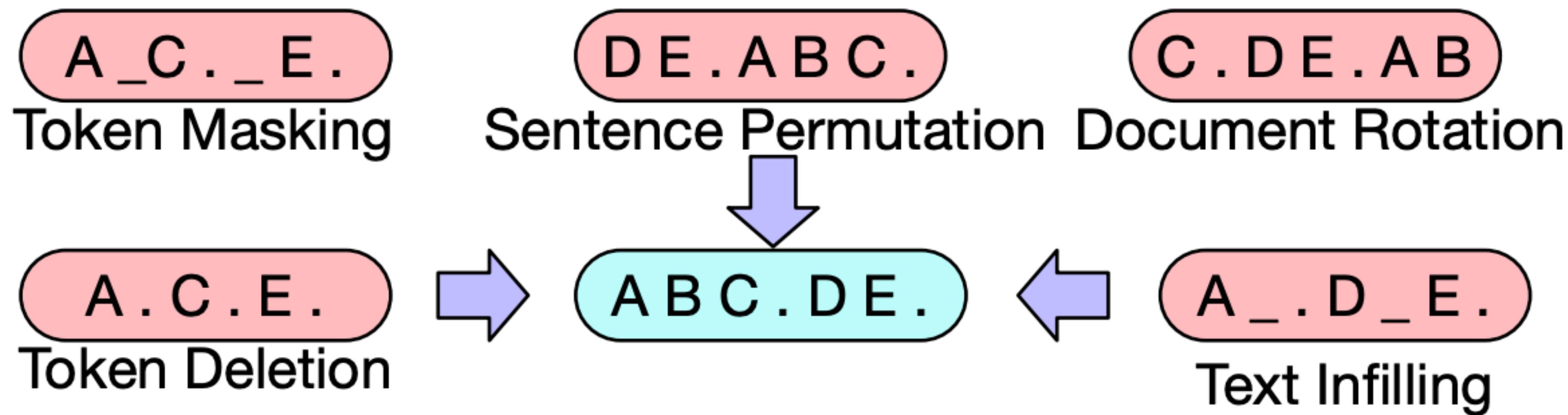
- Pretraining BART combines elements of BERT and GPT!
- **BERT-style:** input texts corrupted before they are passed to bidirectional encoder
- **GPT-style:** model is trained with a language modelling objective in the decoder: predict the next word!





BART Pretraining

- We're not reconstructing the input the same way as BERT, so can we corrupt the input in different ways?
- Many corruption strategies can be used on the encoder side



Can do all the same tasks

- BART can also do all the tasks that BERT does!
*Fine tuning = works well
for text generat°, summarizat°,
Machine translat°, text complet°*
- **Classification:**
 - Give input to both encoder AND decoder (input the sequence twice)
 - Append [CLS] token to decoder sequence and classify its output
- **Sequence Labeling:**
 - Give input to both encoder AND decoder (input the sequence twice)
 - Classify decoder output representations for each token

Can do all the same tasks

	SQuAD 1.1	SQuAD 2.0	MNLI	SST	QQP	QNLI	STS-B	RTE	MRPC	CoLA
	EM/F1	EM/F1	m/mm	Acc	Acc	Acc	Acc	Acc	Acc	Mcc
BERT	84.1/90.9	79.0/81.8	86.6/-	93.2	91.3	92.3	90.0	70.4	88.0	60.6
UniLM	-/-	80.5/83.4	87.0/85.9	94.5	-	92.7	-	70.9	-	61.1
XLNet	89.0/94.5	86.1/88.8	89.8/-	95.6	91.8	93.9	91.8	83.8	89.2	63.6
RoBERTa	88.9/94.6	86.5/89.4	90.2/90.2	96.4	92.2	94.7	92.4	86.6	90.9	68.0
BART	88.8/94.6	86.1/89.2	89.9/90.1	96.6	92.5	94.9	91.2	87.0	90.4	62.8

Almost as good as RoBERTa

Way better than BERT! Why ?

Trained on way more data!

Results: Summarization

This is the first time anyone has been recorded to run a full marathon of 42.195 kilometers (approximately 26 miles) under this pursued landmark time. It was not, however, an officially sanctioned world record, as it was not an "open race" of the IAAF. His time was 1 hour 59 minutes 40.2 seconds. Kipchoge ran in Vienna, Austria. It was an event specifically designed to help Kipchoge break the two hour barrier.

PG&E stated it scheduled the blackouts in response to forecasts for high winds amid dry conditions. The aim is to reduce the risk of wildfires. Nearly 800 thousand customers were scheduled to be affected by the shutoffs which were expected to last through at least midday tomorrow.

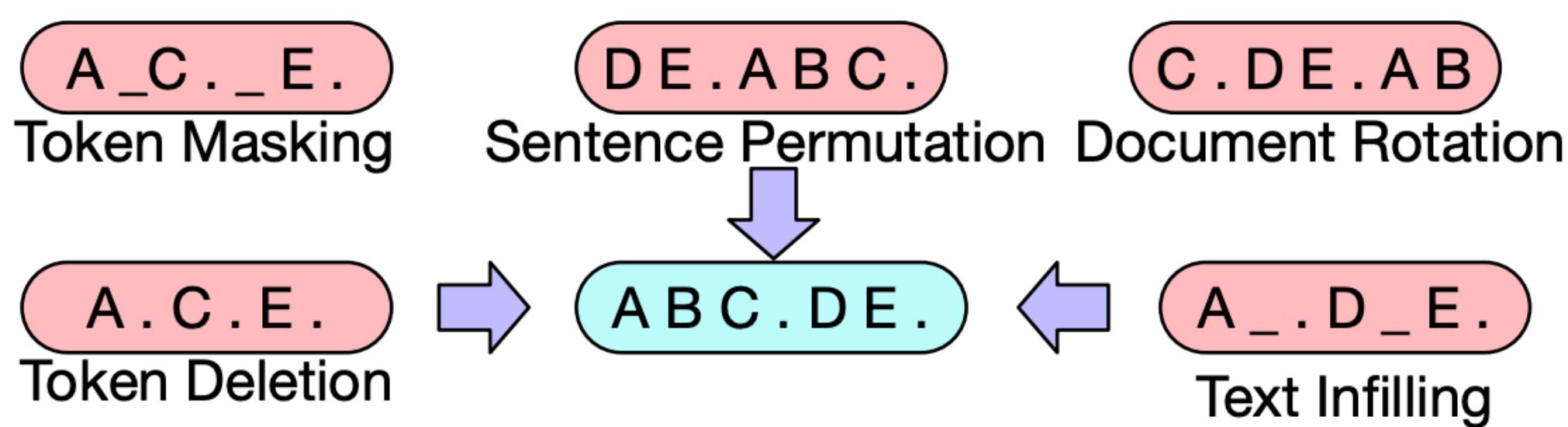
Kenyan runner Eliud Kipchoge has run a marathon in less than two hours.

Power has been turned off to millions of customers in California as part of a power shutoff plan.

**However, BART can do generation tasks too
Decoder is autoregressive!**

Which encoder-side corruption?

Model	SQuAD 1.1	MNLI	ELI5	XSum	ConvAI2	CNN/DM
	F1	Acc	PPL	PPL	PPL	PPL
BART Base						
w/ Token Masking	90.4	84.1	25.05	7.08	11.73	6.10
w/ Token Deletion	90.4	84.1	24.61	6.90	11.46	5.87
w/ Text Infilling	90.8	84.0	24.26	6.61	11.05	5.83
w/ Document Rotation	77.2	75.3	53.69	17.14	19.87	10.59
w/ Sentence Shuffling	85.4	81.5	41.87	10.93	16.67	7.89
w/ Text Infilling + Sentence Shuffling	90.8	83.8	24.17	6.62	11.12	5.41

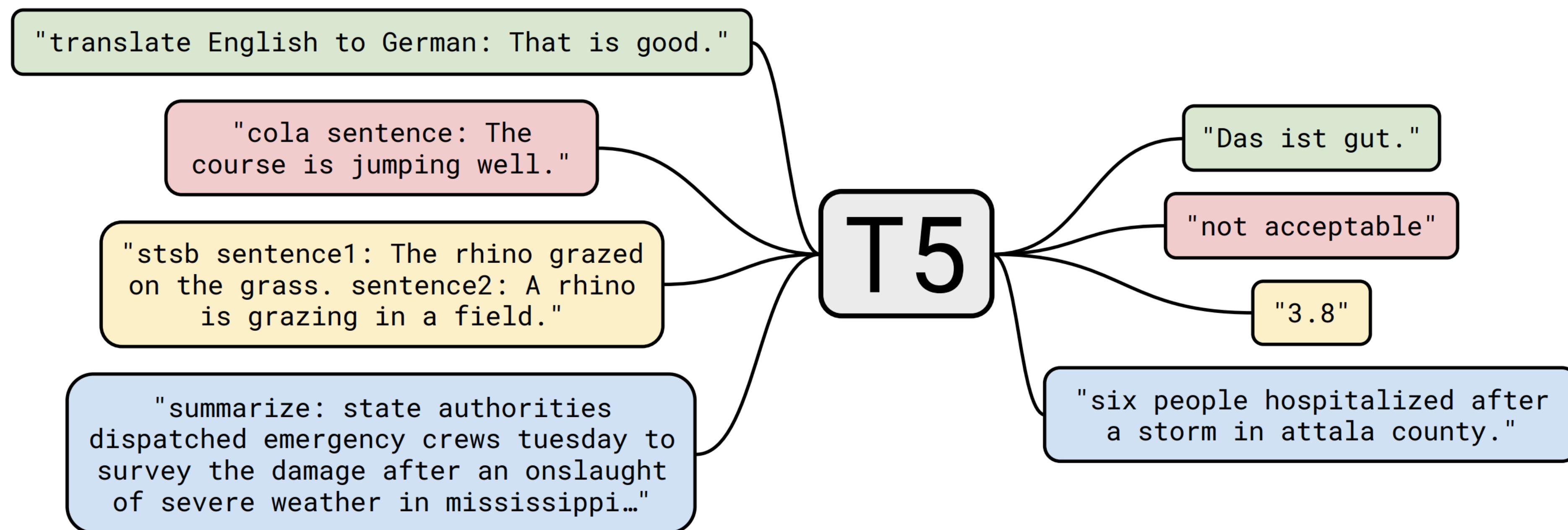


- Different corruption better for transfer to different tasks
- **Use combination of text infilling + sentence permutation**

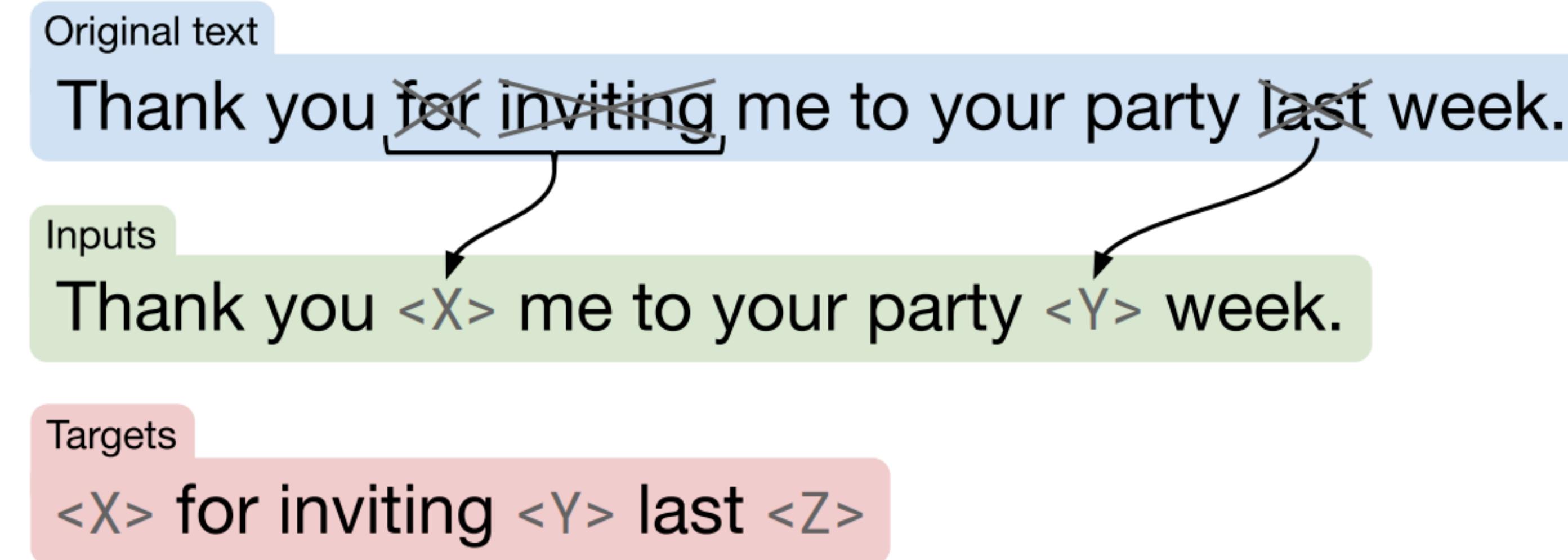
Finetuning = all tasks are converted
into text generat^on problems

T5

- Similar idea as BART: Any problem can be cast as sequence-to-sequence

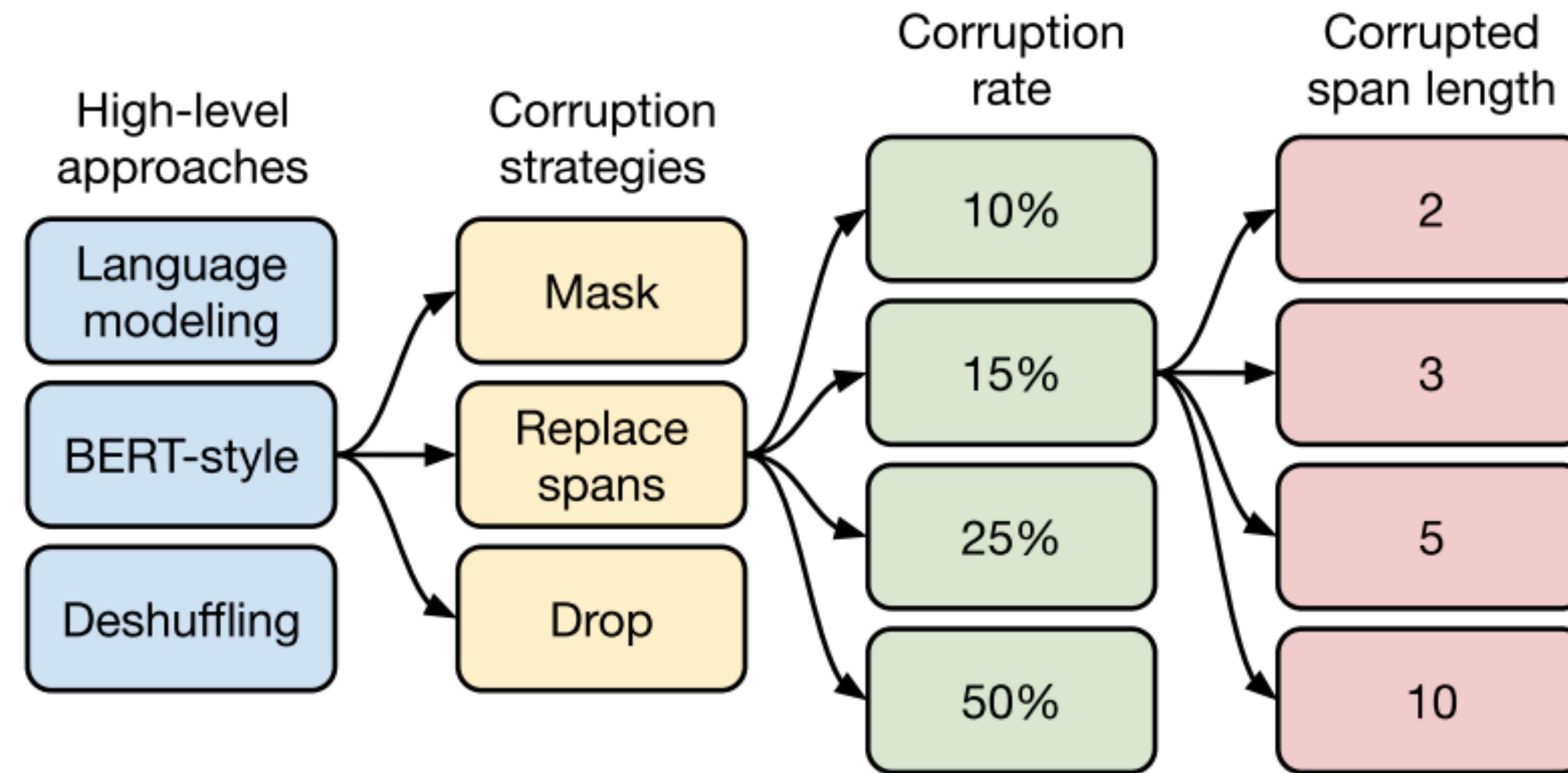


T5 Pretraining



- Similar to BART
- Uses the infilling objective where tokens are reconstructed from underspecified mask corruptions

T5 Pretraining Decisions



- Explored many dimensions of pretraining in seq2seq framework
- Took findings to train much larger model — 11B parameters!

Recap

- **Contextual representations:** Let us model words and sequences conditioned on the context around them
- **ELMo:** Based on bidirectional LSTMs. **Good for pretrained embeddings.**
- **GPT:** Uses a transformer decoder. **Good for generating text as a language model.**
- **BERT:** Uses a transformer encoder. **Good for classification and sequence labelling.**
- **BART + T5:** Pretraining sequence-to-sequence transformer models. **Extendable to all task types!**

References

- Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI blog*, 1(8), 9.
- Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (2018). Deep Contextualized Word Representations. *North American Chapter of the Association for Computational Linguistics*.
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Clark, K., Luong, M. T., Le, Q. V., & Manning, C. D. (2020). Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*.
- Clark, K., Khandelwal, U., Levy, O., & Manning, C. D. (2019). What does bert look at? an analysis of bert's attention. *arXiv preprint arXiv:1906.04341*.