# BERT: Model and Applications
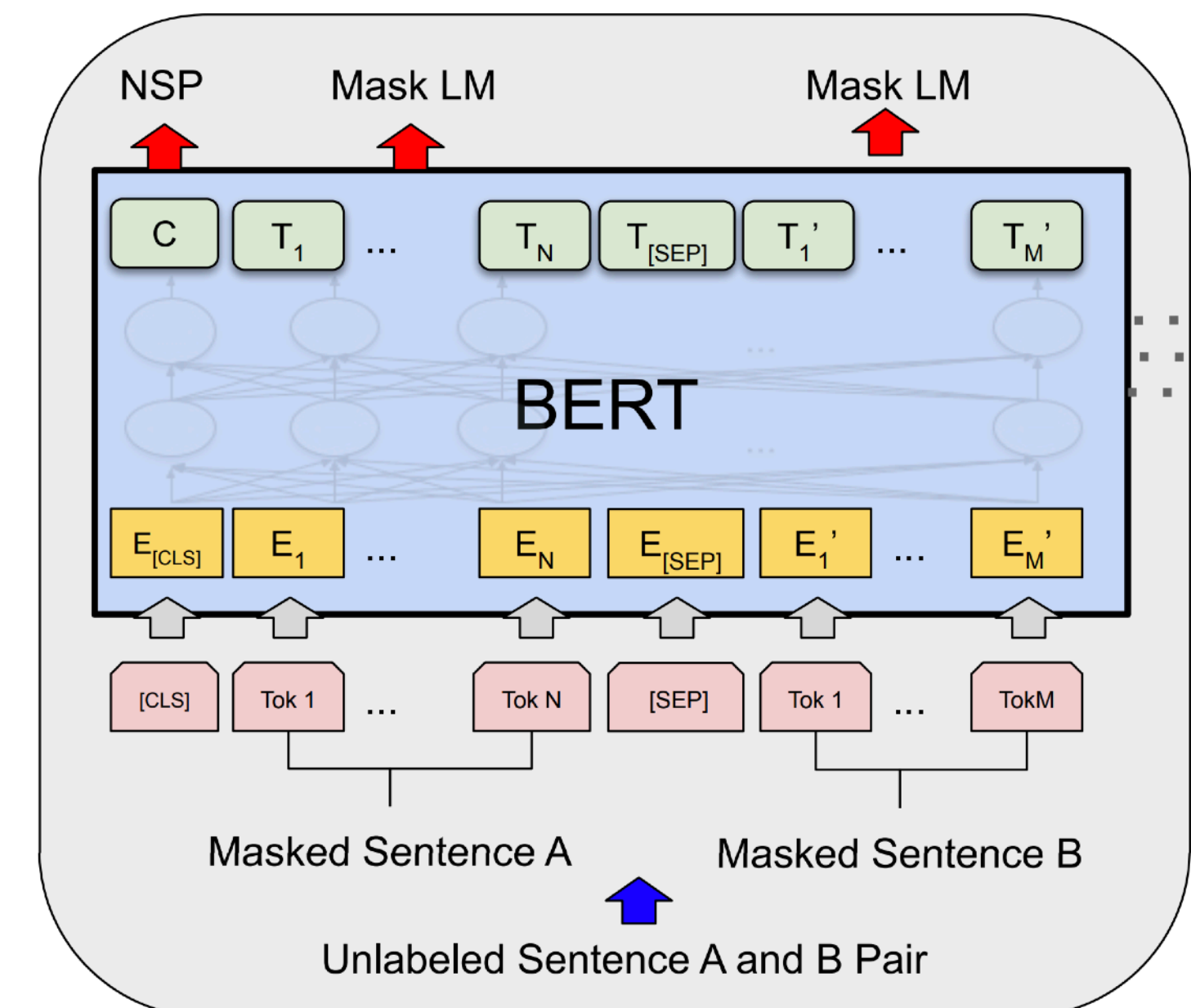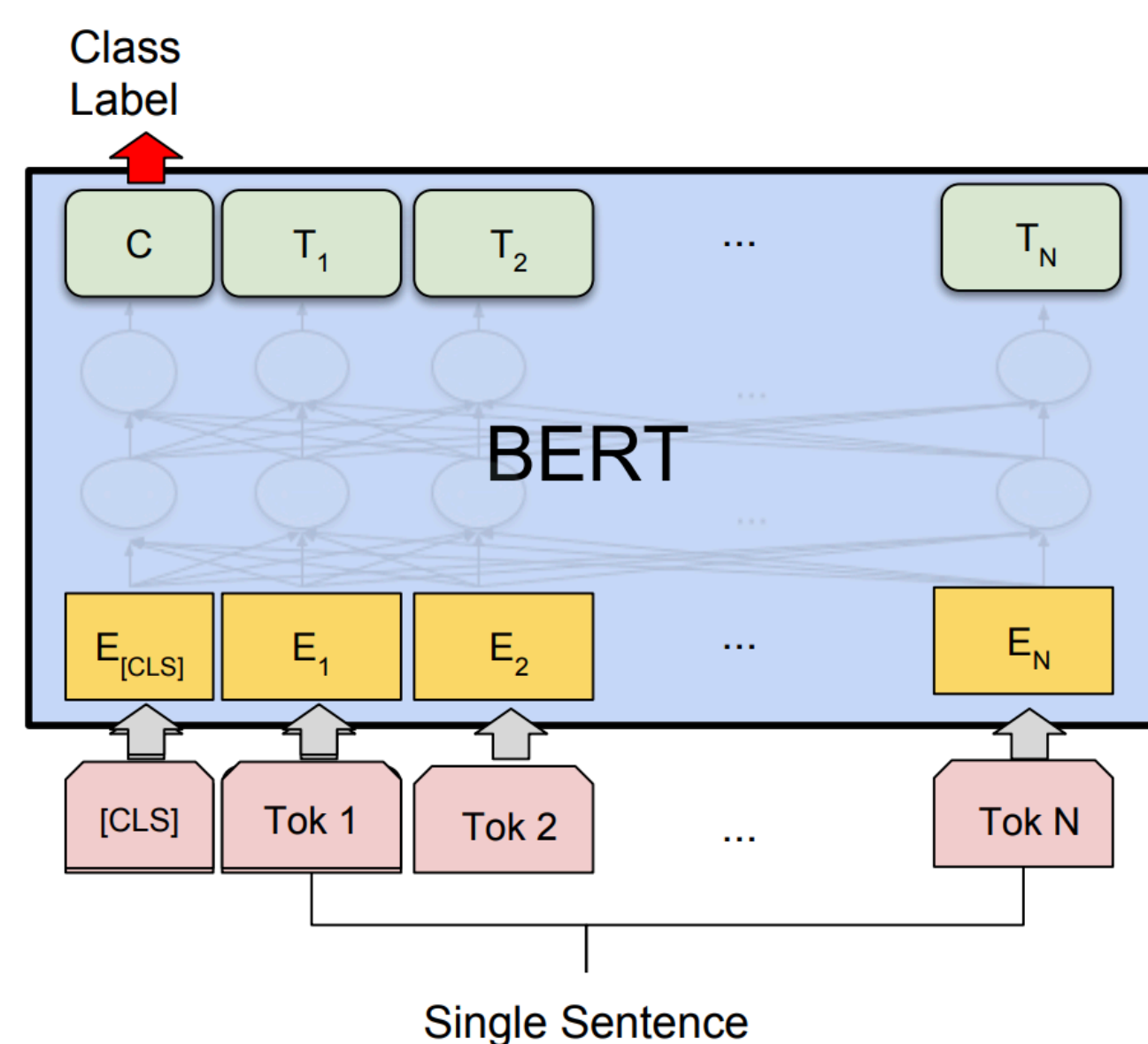
▸ BERT Base: 12 layers, 768-dim per wordpiece token, 12 heads. Total params = 110M

▸ BERT Large: 24 layers, 1024-dim per wordpiece token, 16 heads. Total params = 340M

▸ Positional embeddings and segment embeddings, 30k word pieces

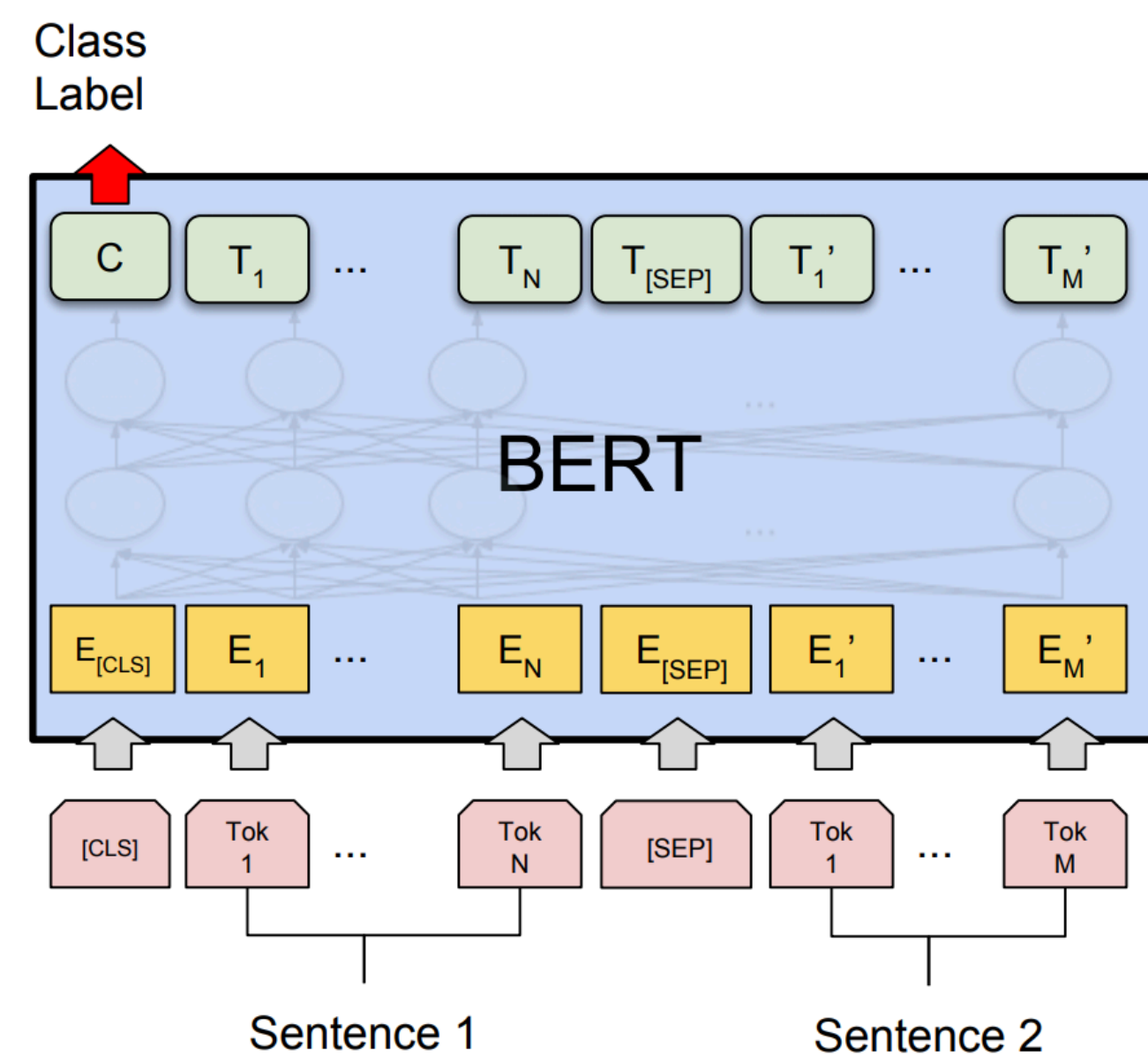▸ This is the model that gets **pre-trained** on a large corpus





Devlin et al. (2019)

# What can BERT do?



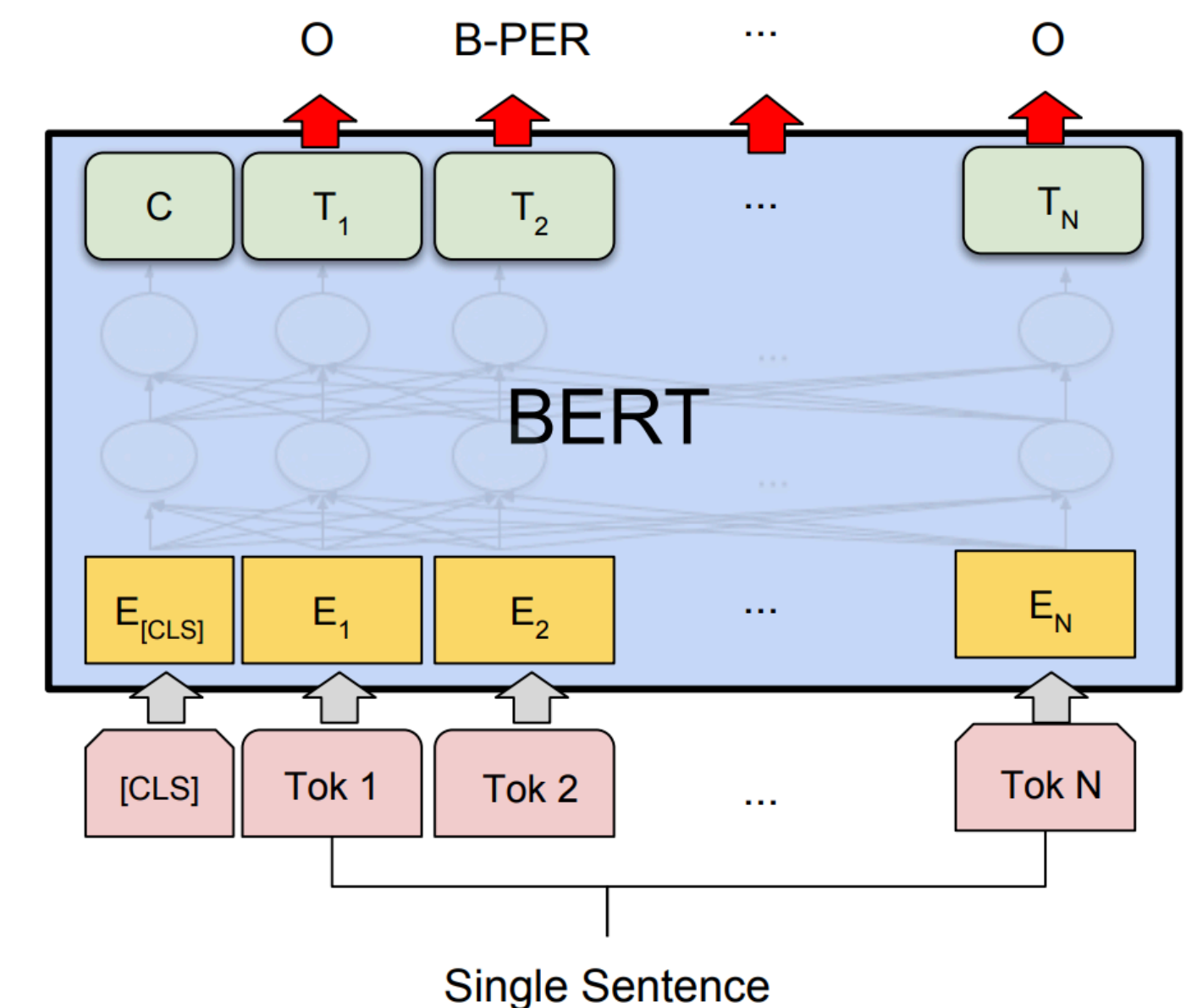(b) Single Sentence Classification Tasks: SST-2, CoLA

(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 sentence pair tasks?

▸ Transformers can capture interactions between the two sentences (even though the NSP objective doesn't really cause this to happen)

Devlin et al. (2019)

# 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 (you can put MASK at the end repeatedly, but this is slow)

▸ Masked language models are intended to be used primarily for "analysis" tasks

Devlin et al. (2019)

# Fine-tuning BERT



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

- ▸ Fine-tune for 1-3 epochs, small learning rate

- ▸ 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 CoNLL 2003 | SA SST-2 | Nat. lang. inference MNLI | SICK-E | Semantic textual similarity 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 is typically better if the whole network is fine-tuned, unlike ELMo

Peters et al. (2019)

# Evaluation

| 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 |

Wang et al. (2019)

# Evaluation

| 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** |

▸ 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. (2019)

# Analysis



Head 1-1
Attends broadly

Head 3-1
Attends to next token
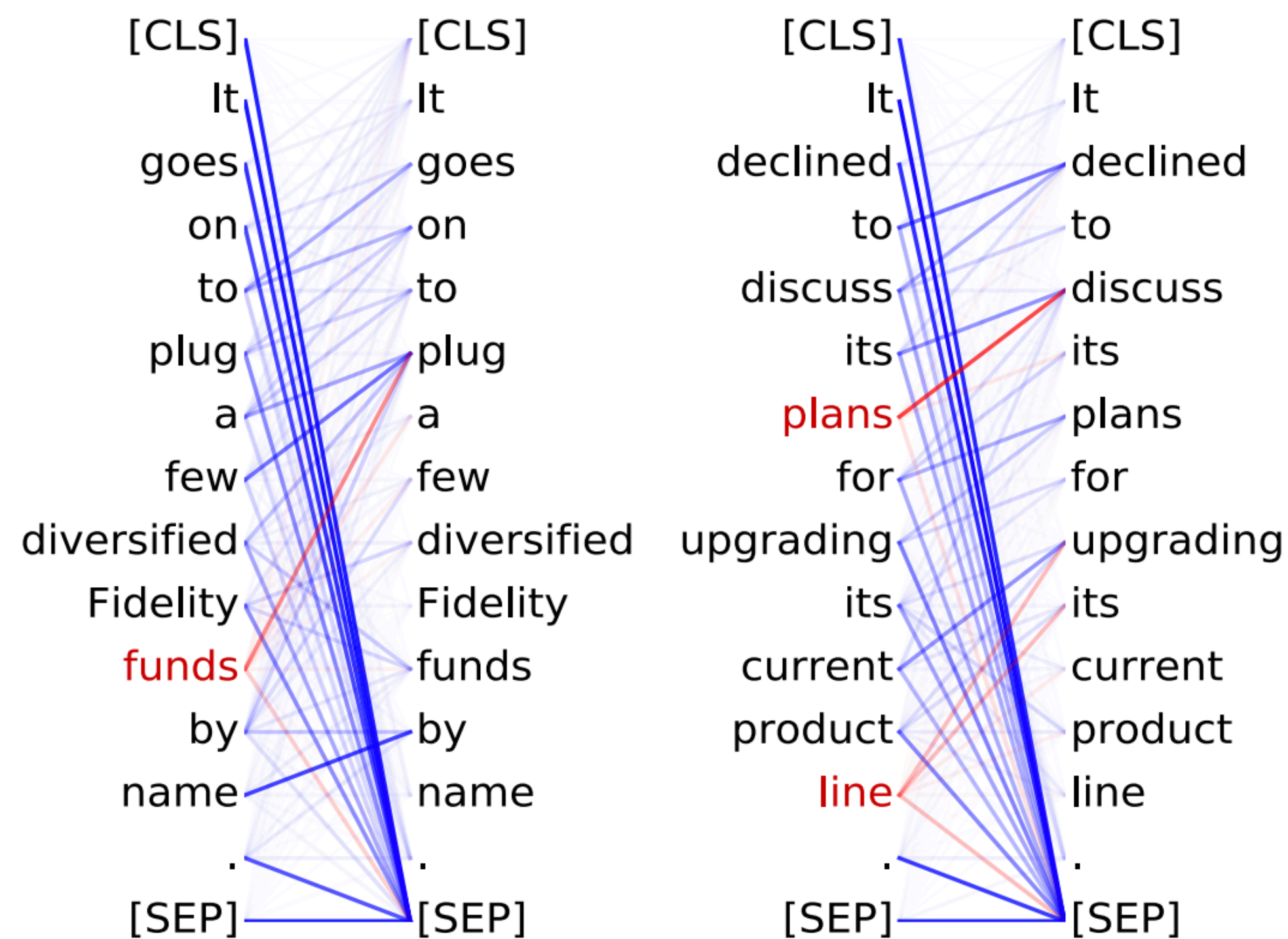
Head 8-7
Attends to [SEP]

Head 11-6
Attends to periods

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

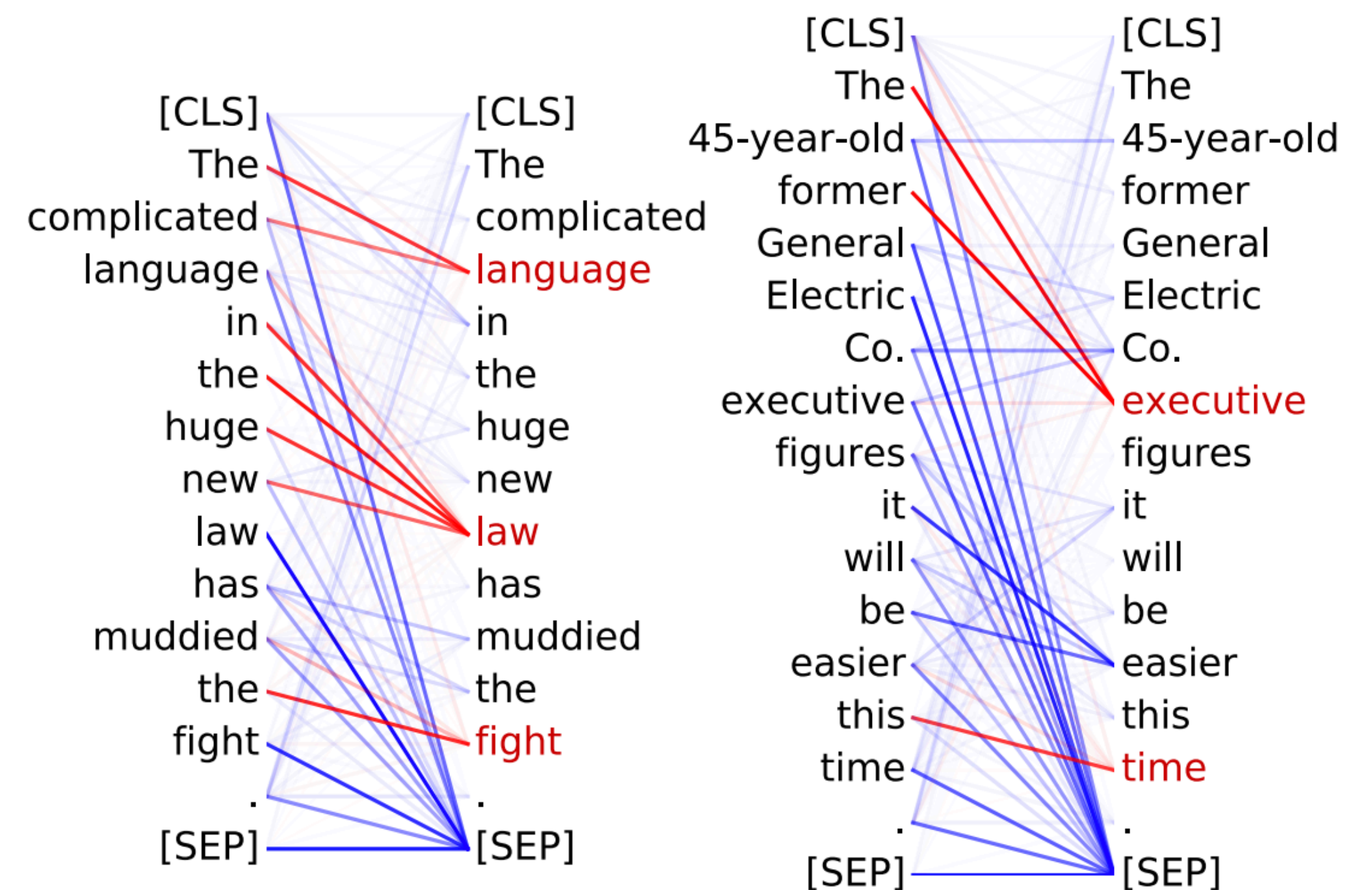Clark et al. (2019)

# Analysis

**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



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

# 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

| 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$_{\text{LARGE}}$ | | | | | | |
|    with BOOKS + WIKI | 13GB | 256 | 1M | 90.9/81.8 | 86.6 | 93.7 |

▸ New training + more data = better performance

▸ For this and more: check out Huggingface Transformers or fairseq

Liu et al. (2019)