

BERT and T5 models



Prof Dr Marko Robnik-Šikonja

Natural Language Processing, Edition 2024

Contents

- subword tokenization
 - BERT models
 - Cross-lingual transfer
 - T5 models
-
- some slides by Jay Alammar, Jacob Devlin and Andrej Miščič

Subword Encoding tokenization

- Learn tokenization based on statistics
- Relevant for modern neural networks
- Use the data to tell us how to tokenize.
- **Subword tokenization** (because tokens are often parts of words)
- Can include common morphemes like *-est* or *-er*.
 - (A morpheme is the smallest meaning-bearing unit of a language; *unlikeliest* has morphemes *un-*, *likely*, and *-est*.)
- Relevant for all languages, but crucial for morphologically-rich languages such as Slovene
- What happens if subword tokenization is inadequate?

Subword tokenization

- Common algorithms:
 - Byte-Pair Encoding (BPE) (Sennrich et al., 2016)
 - WordPiece (Schuster and Nakajima, 2012)
- Both have 2 parts:
 - A token learner that takes a raw training corpus and induces a vocabulary (a set of tokens).
 - A token segmenter that takes a raw test sentence and tokenizes it according to that vocabulary

Byte Pair Encoding (BPE)

Let vocabulary be the set of all individual characters

= {A, B, C, D,...,a, b, c, d....}

- Repeat:
 - choose the two symbols that are most frequently adjacent in training corpus (say 'A', 'B'),
 - adds a new merged symbol 'AB' to the vocabulary
 - replace every adjacent 'A' 'B' in corpus with 'AB'.
- Until k merges have been done.

BPE token learner algorithm

function BYTE-PAIR ENCODING(strings C , number of merges k) **returns** vocab V

$V \leftarrow$ all unique characters in C # initial set of tokens is characters

for $i = 1$ **to** k **do** # merge tokens til k times

$t_L, t_R \leftarrow$ Most frequent pair of adjacent tokens in C

$t_{NEW} \leftarrow t_L + t_R$ # make new token by concatenating

$V \leftarrow V + t_{NEW}$ # update the vocabulary

 Replace each occurrence of t_L, t_R in C with t_{NEW} # and update the corpus

return V

BPE in use

- Most subword algorithms are run inside white-space separated tokens.
- So first add a special end-of-word symbol '___' before whitespace in training corpus
- Next, separate into letters.

BPE token learner

An example corpus :(

low low low low low lowest lowest newer newer newer newer newer newer wider
wider wider new new

Add end-of-word tokens and segment:

corpus

5 l o w _
2 l o w e s t _
6 n e w e r _
3 w i d e r _
2 n e w _

vocabulary

_, d, e, i, l, n, o, r, s, t, w

BPE token learner

corpus

5 l o w _
2 l o w e s t _
6 n e w e r _
3 w i d e r _
2 n e w _

vocabulary

_, d, e, i, l, n, o, r, s, t, w

Merge **e r** to **er**

corpus

5 l o w _
2 l o w e s t _
6 n e w e r _
3 w i d e r _
2 n e w _

vocabulary

_, d, e, i, l, n, o, r, s, t, w, er

BPE

corpus

5 l o w _
2 l o w e s t _
6 n e w e r _
3 w i d e r _
2 n e w _

vocabulary

_, d, e, i, l, n, o, r, s, t, w, e r

Merge **er _** to **er_**

corpus

5 l o w _
2 l o w e s t _
6 n e w e r_
3 w i d e r_
2 n e w _

vocabulary

, d, e, i, l, n, o, r, s, t, w, e r, e r

BPE

corpus

5 l o w _
2 l o w e s t _
6 n e w e r_
3 w i d e r_
2 n e w _

Merge n e to ne

corpus

5 l o w _
2 l o w e s t _
6 n e w e r_
3 w i d e r_
2 n e w _

vocabulary

, d, e, i, l, n, o, r, s, t, w, er, er

vocabulary

, d, e, i, l, n, o, r, s, t, w, er, er, ne

BPE

The next merges are:

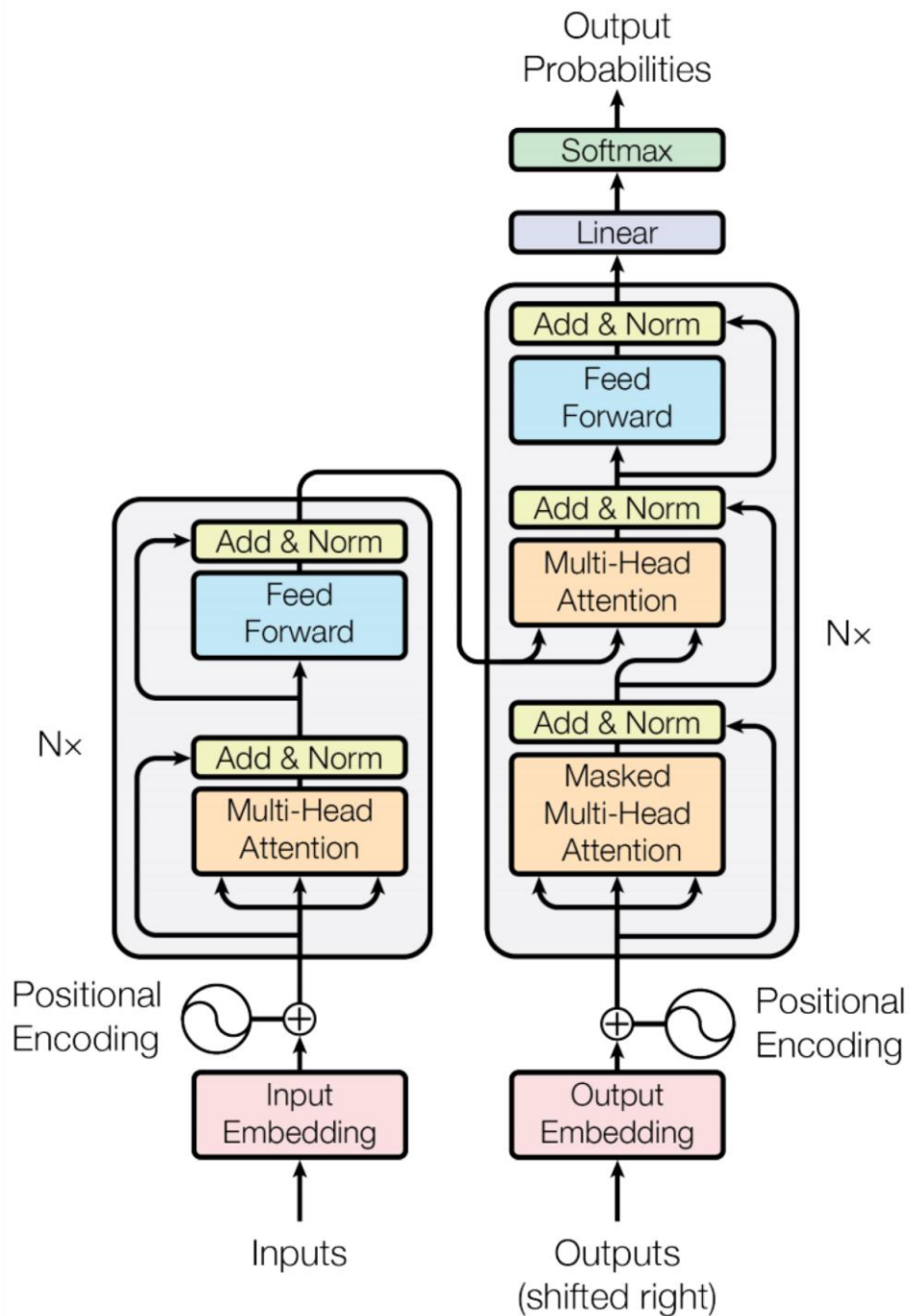
Merge	Current Vocabulary
(ne, w)	—, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new
(l, o)	—, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new, lo
(lo, w)	—, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new, lo, low
(new, er—)	—, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new, lo, low, newer—
(low, —)	—, d, e, i, l, n, o, r, s, t, w, er, er—, ne, new, lo, low, newer—, low—

BPE token learner algorithm

- On the test data, run each merge learned from the training data:
 - Greedily
 - In the order we learned them
 - (test frequencies don't play a role)
- So: merge every `e r` to `er`, then merge `er _` to `er_`, etc.
- Result:
 - Test set "n e w e r _" would be tokenized as a full word
 - Test set "l o w e r _" would be two tokens: "low er_"

Transformer architecture

- typically, the input is first tokenized with subword encoding
- what is the alternative?



BERT

- BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding
- State-of-the-art pretrained LM based on transformer architecture (only the encoder part)
- Idea:
 - use large unlabeled corpora and an auxiliary task to pretrain a model for a general language representation
 - fine-tune the model on a (possibly small) dataset for a specific downstream task (typically classification)
- presentation based on slides from Jacob Devlin and Jay Alamar

Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Volume 1, pp. 4171-4186.

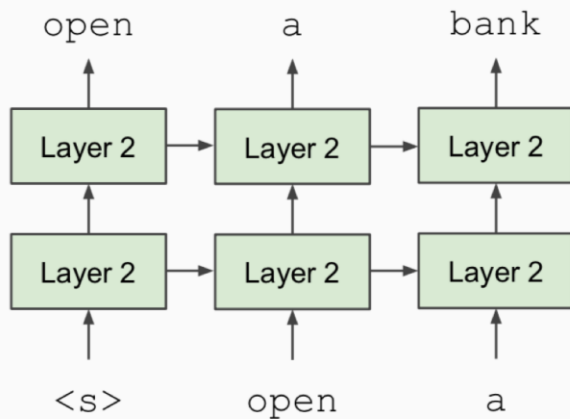
BERT: motivation 1/3

- **Problem:** Classical language models only use the left or right context, but language understanding is bidirectional.
- Why are LMs unidirectional?
 - Reason 1: Directionality is needed to generate a well-formed probability distribution.
 - We don't care about this.
 - Reason 2: Words can “see themselves” in a bidirectional encoder.

BERT: motivation 2/3

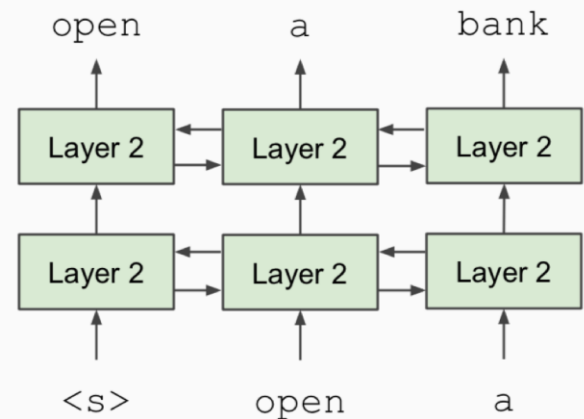
Unidirectional context

Build representation incrementally



Bidirectional context

Words can “see themselves”



BERT: motivation 3/3

- **Solution:** Mask out $k\%$ of the input words, and then predict the masked words
- BERT uses $k = 15\%$

store

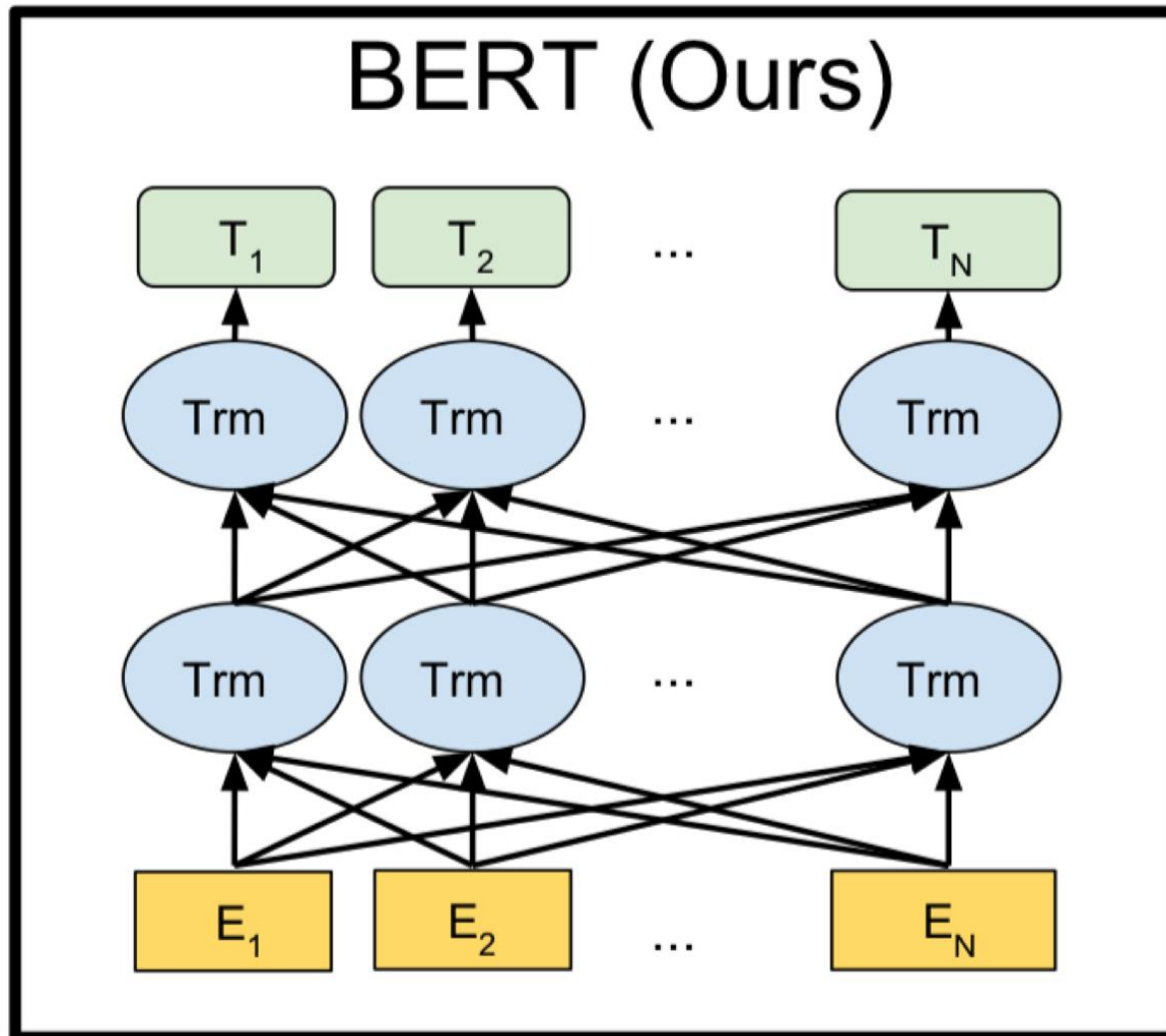
gallon



the man went to the [MASK] to buy a [MASK] of milk

- Too little masking: Too expensive to train (not enough masks)
- Too much masking: Not enough context

BERT architecture



BERT uses several tasks

- besides masked LM, BERT learns relationships between sentences
- predict whether Sentence B is actual sentence that proceeds Sentence A, or a random sentence

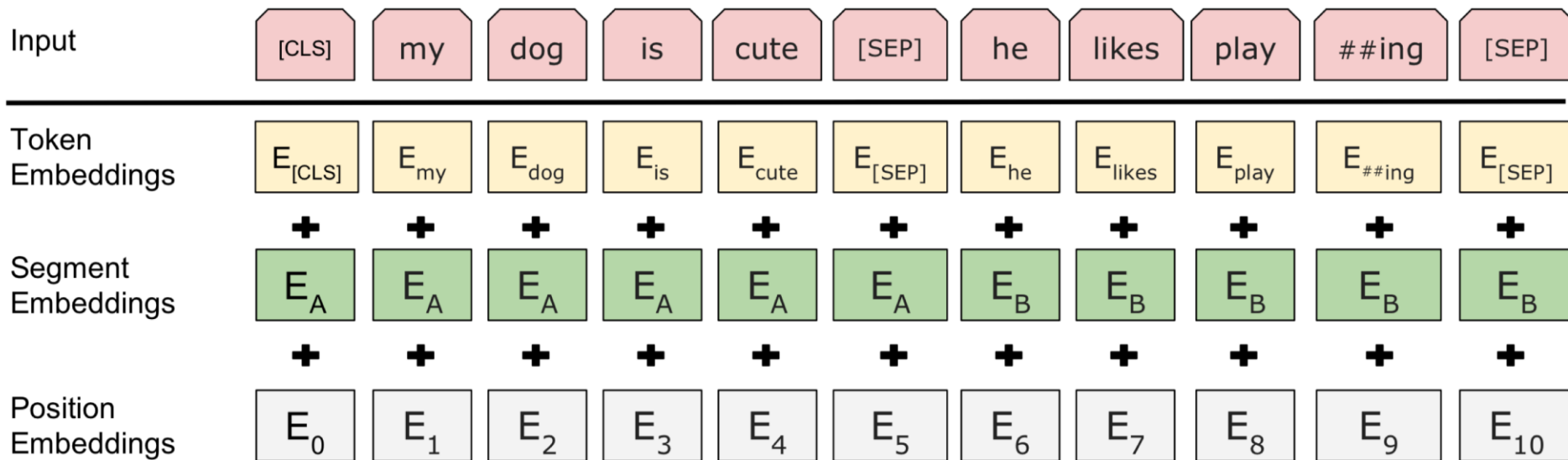
Sentence A = The man went to the store.
Sentence B = He bought a gallon of milk.
Label = IsNextSentence

Sentence A = The man went to the store.
Sentence B = Penguins are flightless.
Label = NotNextSentence

- some follow-up BERT-like models, e.g., RoBERTa, drop this task and claim better performance on downstream tasks

Sentence-pair encoding for BERT

- Token embeddings are word pieces (sub-word encoding)
- (Relatively) common words are in the vocabulary: *at, fairfax, 1910s*
- Other words are built from wordpieces: *hypatia = h ##yp ##ati ##a*
- Learned segmented embeddings represents each sentence
- Positional embedding is the same as for other transformer architectures

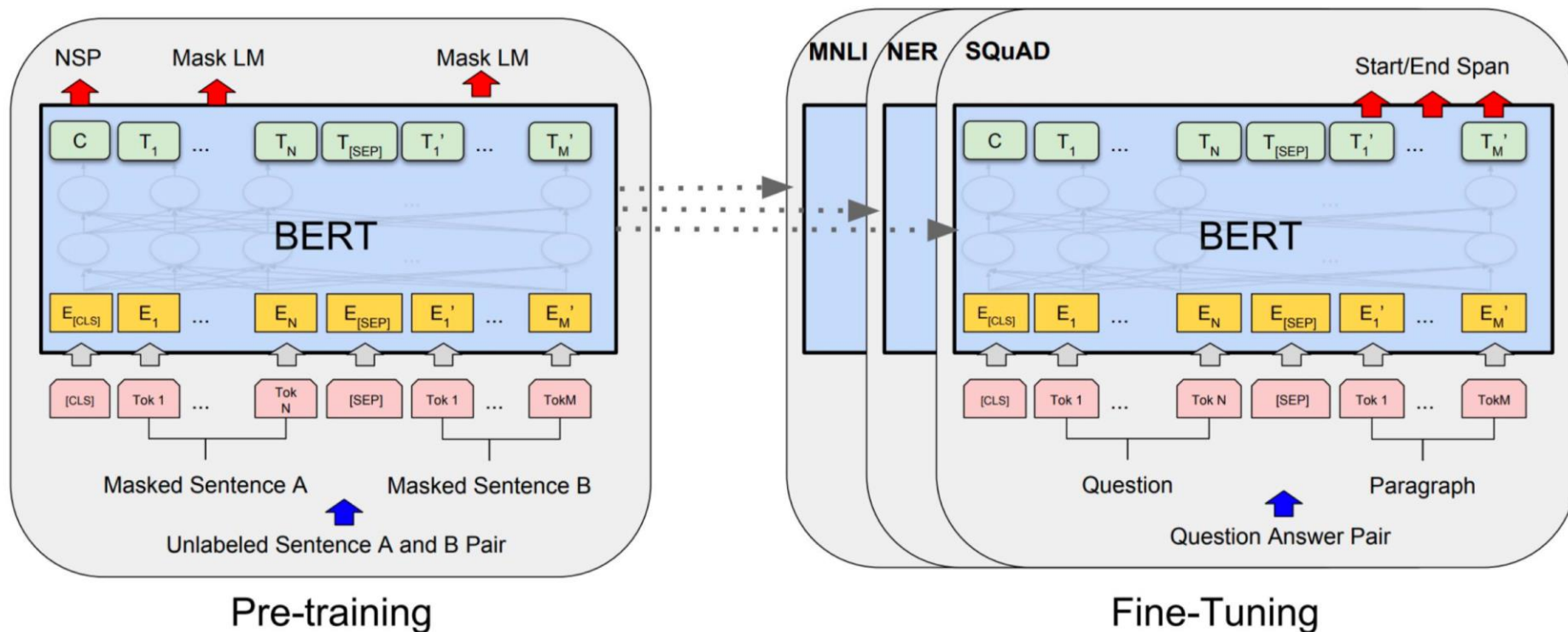


BERT training

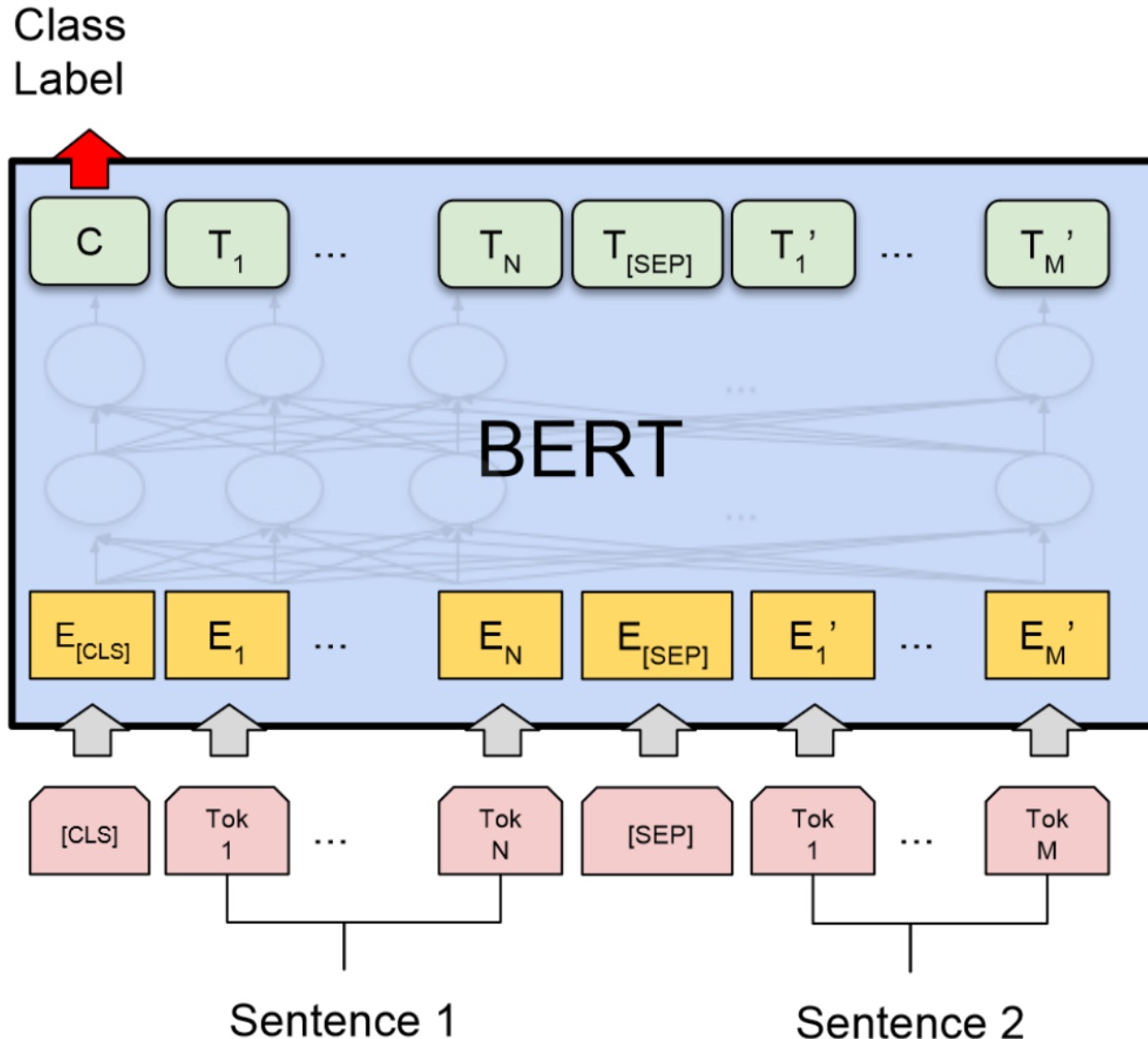
- Transformer encoder
- Self-attention \Rightarrow no locality bias
- Long-distance context has “equal opportunity”
- Single multiplication per layer \Rightarrow efficiency on GPU/TPU
- Trained on Wikipedia + BookCorpus
- English BERT was trained with 2 model sizes:
 - BERT-Base: 12-layer, 768-hidden parameters, 12-head, 110M parameters
 - BERT-Large: 24-layer, 1024-hidden parameters, 16-head, 340M parameters
- Trained on 4x4 or 8x8 TPU slice for 4 days

Use of BERT

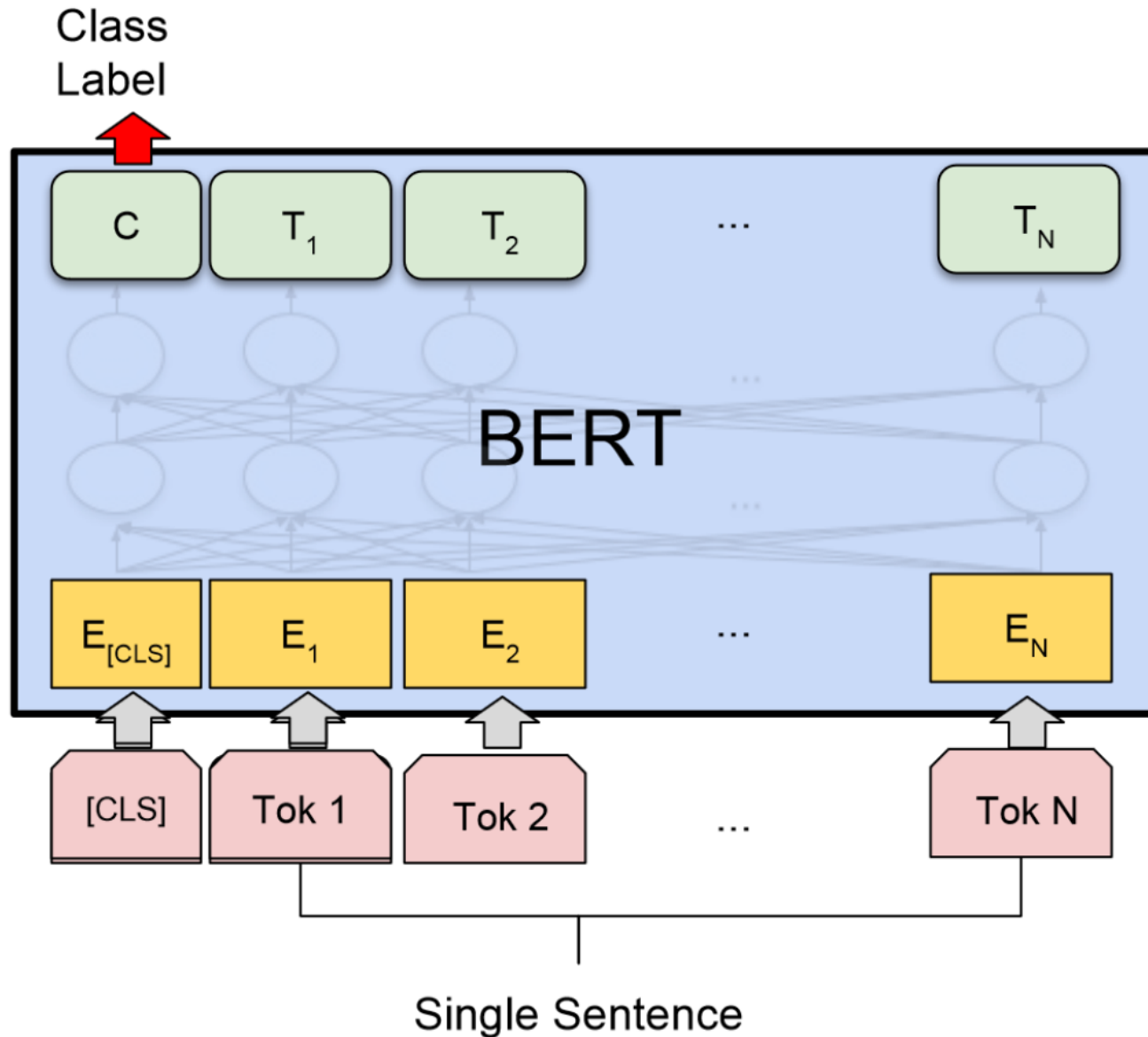
- train a classifier built on the top layer for each task that you fine-tune for, e.g., Q&A, NER, inference
- achieved state-of-the-art results for many tasks
- GLUE and SuperGLUE tasks for natural language understanding



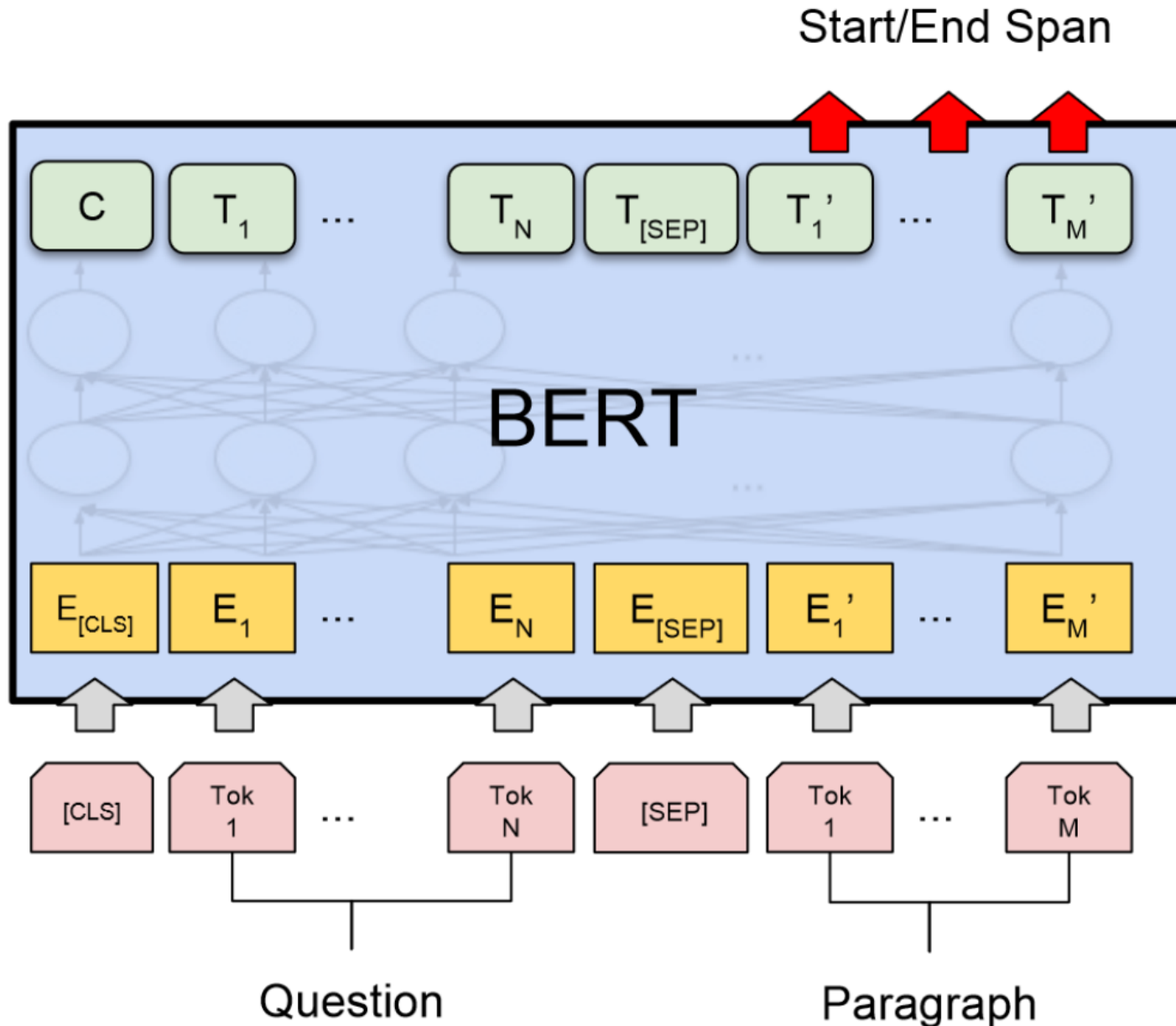
Two sentence classification using BERT- e.g., natural language inference



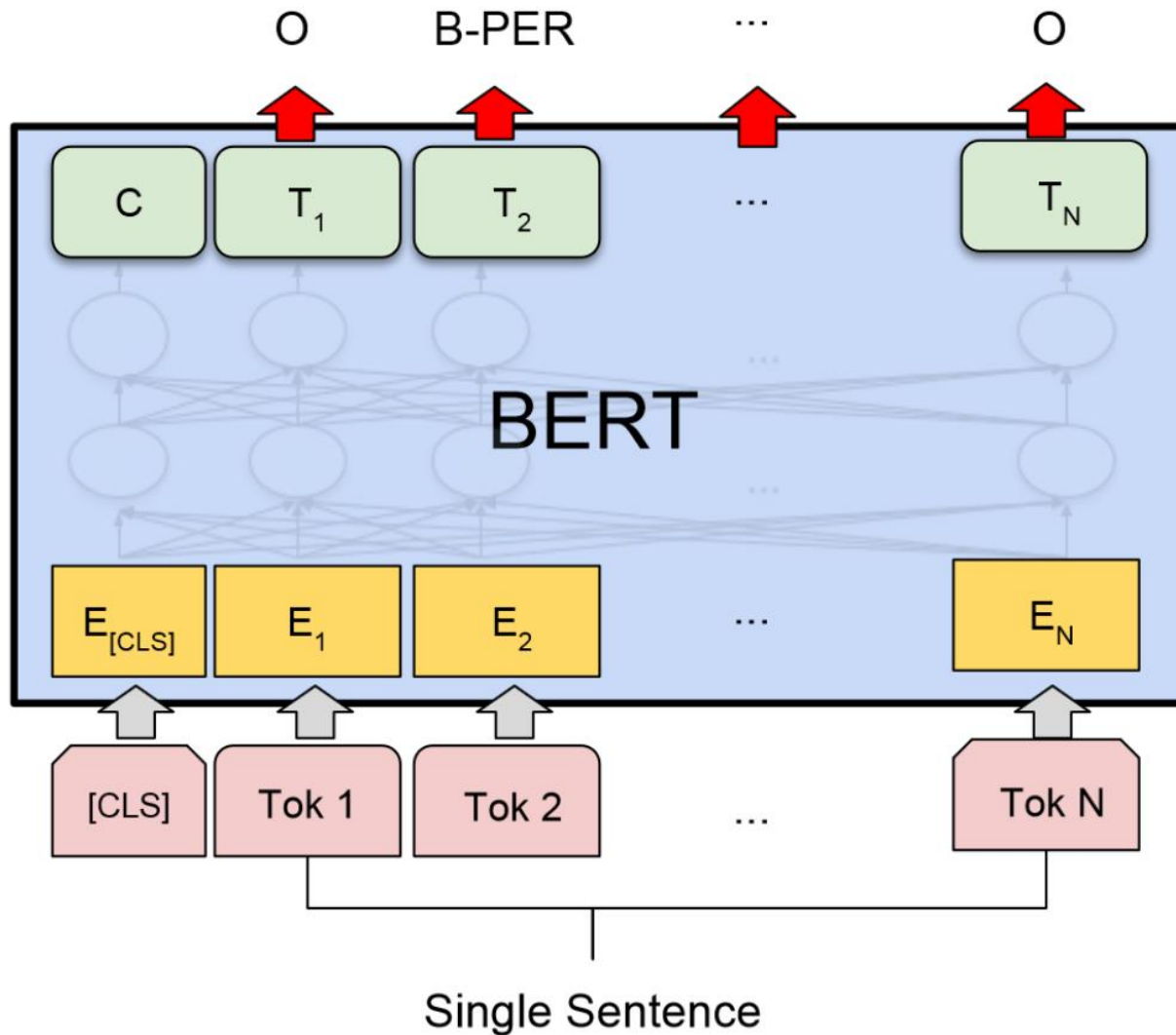
Sentence classification using BERT – sentiment, grammatical correctness



Questions and answers with BERT

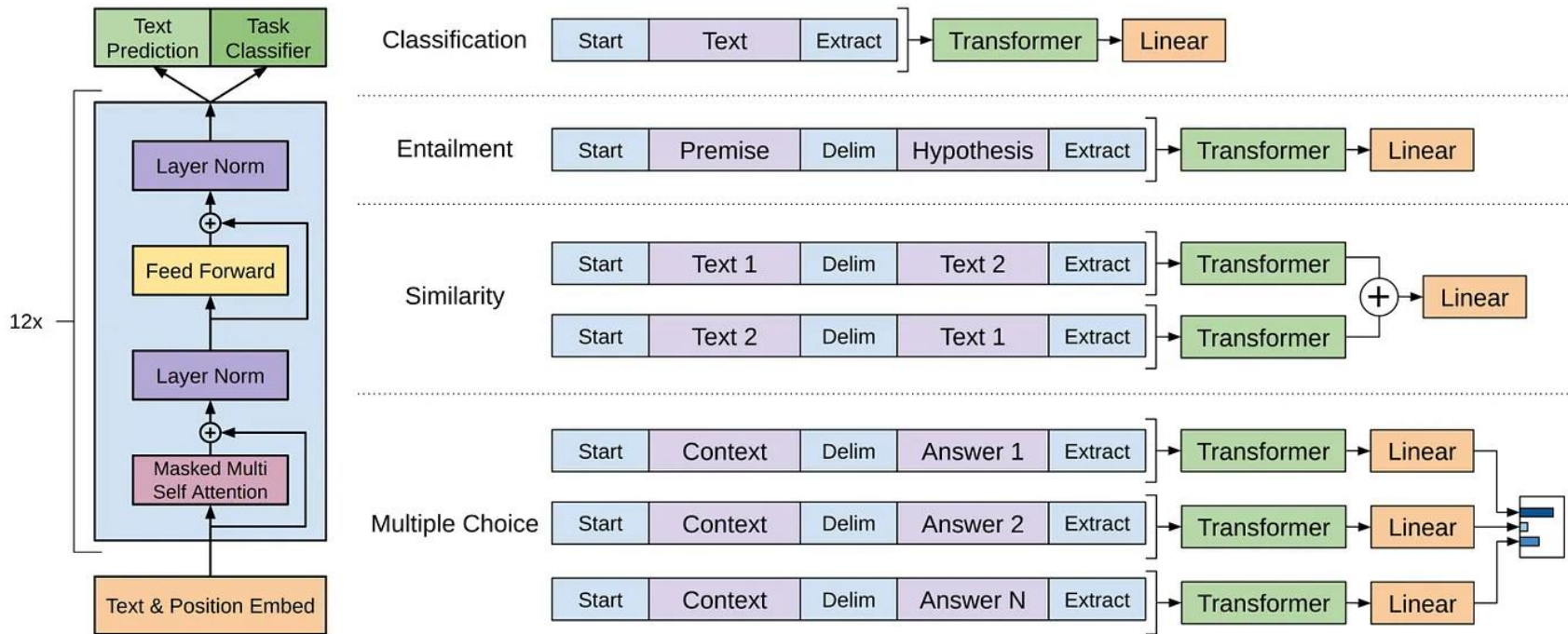


Sentence tagging with BERT- NER, POS tagging, SRL





Finetuning

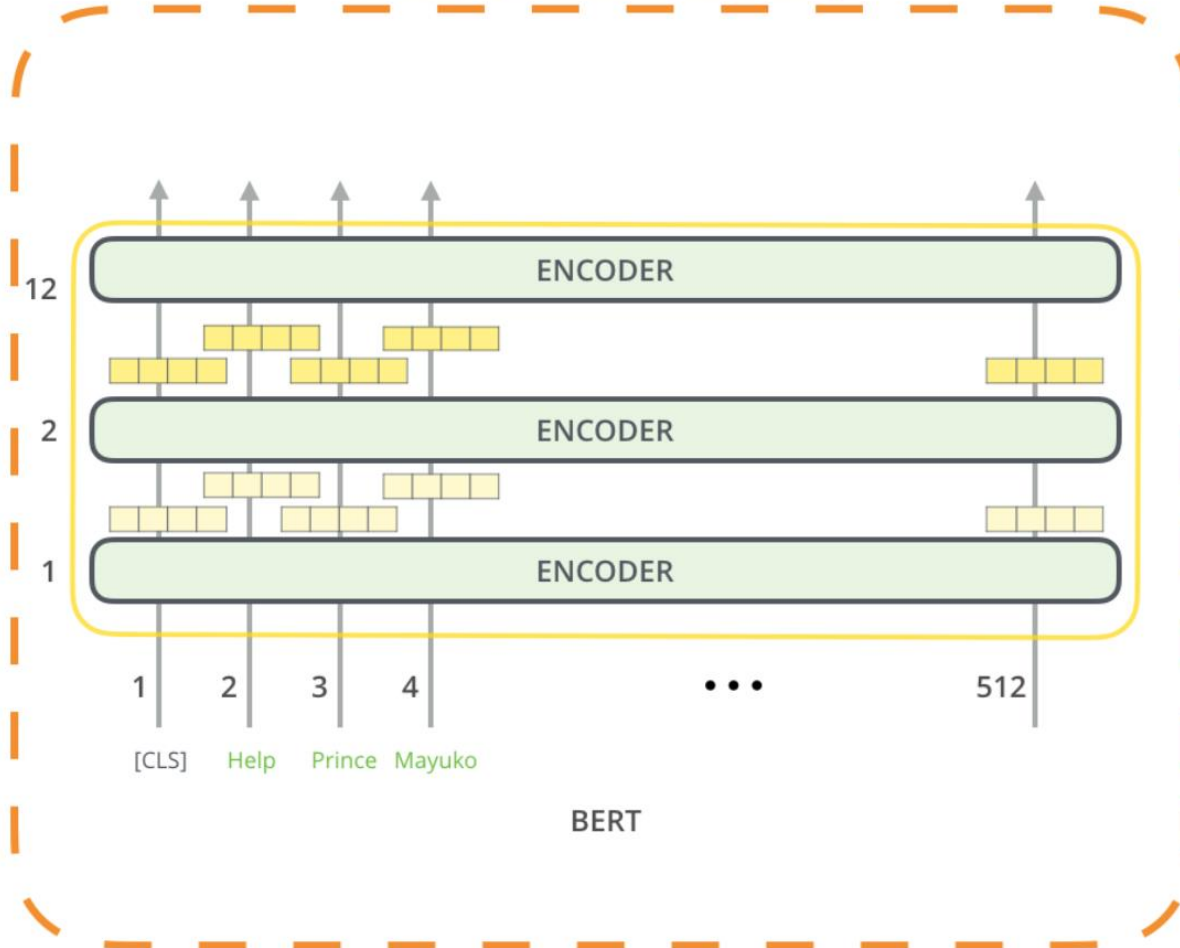


BERT can produce embeddings

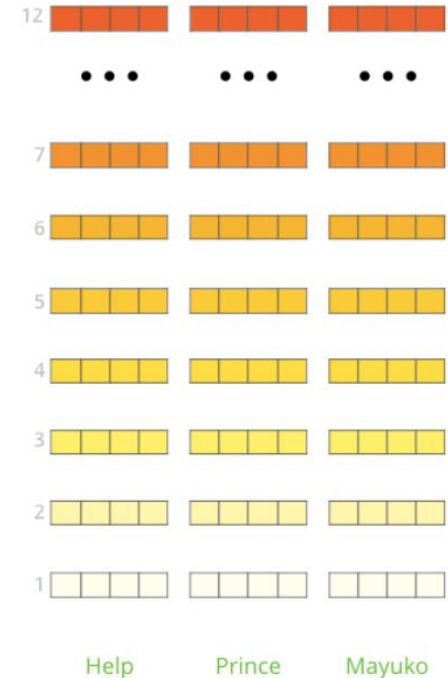
- one can extract fixed size contextual vectors from BERT, achieving slightly lower accuracy than using the whole BERT as the first stage model

Layer-wise embeddings

Generate Contextualized Embeddings



The output of each encoder layer along each token's path can be used as a feature representing that token.



But which one should we use?

Which layer of BERT to use as embeddings?

What is the best contextualized embedding for “Help” in that context?
 For named-entity recognition task CoNLL-2003 NER

		Dev F1 Score	
12	First Layer Embedding	91.0	
...			
7	Last Hidden Layer	94.9	
6	Sum All 12 Layers	95.5	
5			
4			
3			
2	Second-to-Last Hidden Layer	95.6	
1	Sum Last Four Hidden	95.9	
...			
...			
...			
Help	Concat Last Four Hidden	96.1	

Examples of GLUE tasks

- GLUE benchmark is dominated by natural language inference tasks, but also has sentence similarity and sentiment

MultiNLI

Premise: Hills and mountains are especially sanctified in Jainism.

Hypothesis: Jainism hates nature.

Label: Contradiction

CoLA (Corpus of Linguistic Acceptability)

Sentence: The wagon rumbled down the road. Label: Acceptable

Sentence: The car honked down the road. Label: Unacceptable

SuperGLUE tasks

BoolQ - Boolean Questions

CB – Commitment Bank

COPA - Choice of Plausible Alternatives

MultiRC - Multi-Sentence Reading Comprehension

ReCoRD - Reading Comprehension with

Commonsense Reasoning Dataset

RTE - Recognizing Textual Entailment

WiC - Word-in-Context

WSC - Winograd Schema Challeng

Table 2: Development set examples from the tasks in SuperGLUE. **Bold** text represents part of the example format for each task. Text in *italics* is part of the model input. Underlined text is specially marked in the input. Text in a monospaced font represents the expected model output.

BoolQ	Passage: <i>Barq’s – Barq’s is an American soft drink. Its brand of root beer is notable for having caffeine. Barq’s, created by Edward Barq and bottled since the turn of the 20th century, is owned by the Barq family but bottled by the Coca-Cola Company. It was known as Barq’s Famous Olde Tyme Root Beer until 2012.</i>
	Question: <i>is barq’s root beer a pepsi product</i> Answer: No
CB	Text: <i>B: And yet, uh, I we-, I hope to see employer based, you know, helping out. You know, child, uh, care centers at the place of employment and things like that, that will help out. A: Uh-huh. B: What do you think, do you think we are, setting a trend?</i>
	Hypothesis: <i>they are setting a trend</i> Entailment: Unknown
COPA	Premise: <i>My body cast a shadow over the grass.</i> Question: <i>What’s the CAUSE for this?</i>
	Alternative 1: <i>The sun was rising.</i> Alternative 2: <i>The grass was cut.</i>
	Correct Alternative: 1

MultiRC

Paragraph: *Susan wanted to have a birthday party. She called all of her friends. She has five friends. Her mom said that Susan can invite them all to the party. Her first friend could not go to the party because she was sick. Her second friend was going out of town. Her third friend was not so sure if her parents would let her. The fourth friend said maybe. The fifth friend could go to the party for sure. Susan was a little sad. On the day of the party, all five friends showed up. Each friend had a present for Susan. Susan was happy and sent each friend a thank you card the next week*

Question: *Did Susan's sick friend recover?* **Candidate answers:** *Yes, she recovered (T), No (F), Yes (T), No, she didn't recover (F), Yes, she was at Susan's party (T)*

ReCoRD

Paragraph: *(CNN) Puerto Rico on Sunday overwhelmingly voted for statehood. But Congress, the only body that can approve new states, will ultimately decide whether the status of the US commonwealth changes. Ninety-seven percent of the votes in the nonbinding referendum favored statehood, an increase over the results of a 2012 referendum, official results from the State Electoral Commission show. It was the fifth such vote on statehood. "Today, we the people of Puerto Rico are sending a strong and clear message to the US Congress ... and to the world ... claiming our equal rights as American citizens, Puerto Rico Gov. Ricardo Rossello said in a news release. @highlight Puerto Rico voted Sunday in favor of US statehood*

Query For one, they can truthfully say, "Don't blame me, I didn't vote for them," when discussing the <placeholder> presidency **Correct Entities:** US

RTE

Text: *Dana Reeve, the widow of the actor Christopher Reeve, has died of lung cancer at age 44, according to the Christopher Reeve Foundation.*

Hypothesis: *Christopher Reeve had an accident.* **Entailment:** False

WiC

Context 1: *Room and board.* **Context 2:** *He nailed boards across the windows.*

Sense match: False

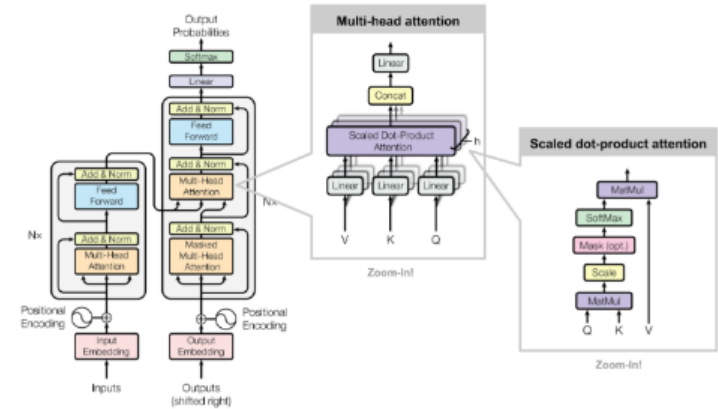
WSC

Text: *Mark told Pete many lies about himself, which Pete included in his book. He should have been more truthful.* **Coreference:** False



Pretrained large language models

- large pretrained neural language models
- trained on large text corpora to capture relations in language
- finetuned to specific tasks
- publicly available LLMs
- for Slovene: fastText, ELMo, SloBERTa, CroSloEngual BERT, SIEng BERT, SloT5, Slo GPT
- for Croatian: fastText, ELMo, BERTić, CroSloEngual BERT
- on Clarin.si and HuggingFace
- hundreds of papers investigating BERT-like models in major ML & NLP conferences



- Ulčar, M., & Robnik-Šikonja, M. (2020). High Quality ELMo Embeddings for Seven Less-Resourced Languages. In *Proceedings of the 12th Language Resources and Evaluation Conference* (pp. 4731-4738).
- Ulčar, M. and Robnik-Šikonja, M., 2021. SloBERTa: Slovene monolingual large pretrained masked language model. *Proceedings of SI-KDD within the Information Society 2021*, pp.17-20.
- Ljubešić, N., & Lauc, D. (2021). BERTić-The Transformer Language Model for Bosnian, Croatian, Montenegrin and Serbian. In *Proceedings of the 8th Workshop on Balto-Slavic Natural Language Processing* (pp. 37-42).
- Ulčar, M., & Robnik-Šikonja, M. (2020). FinEst BERT and CroSloEngual BERT. In *International Conference on Text, Speech, and Dialogue* (pp. 104-111).
- Ulčar, M. & Robnik-Šikonja, M. (2023) Sequence-to-sequence pretraining for a less-resourced Slovenian language. *Frontiers in Artificial Intelligence, Section on Language and Computation, Volume 6 – 2023*, <https://doi.org/10.3389/frai.2023.932519>



SloBERTa

- Currently the best Slovene LLM
- Many thousands downloads from HuggingFace
- Training set: 3.41 B words (corpora Gigafida, KAS, partially Janes, siParl, slWaC)
- Training duration: 4 weeks on Nvidia DGX A100 using 4xGPU
- An example of direct use:
 - <mask> je najlepše mesto na svetu.
 - Odgovori: Ljubljana, Barcelona, London, Madrid, To

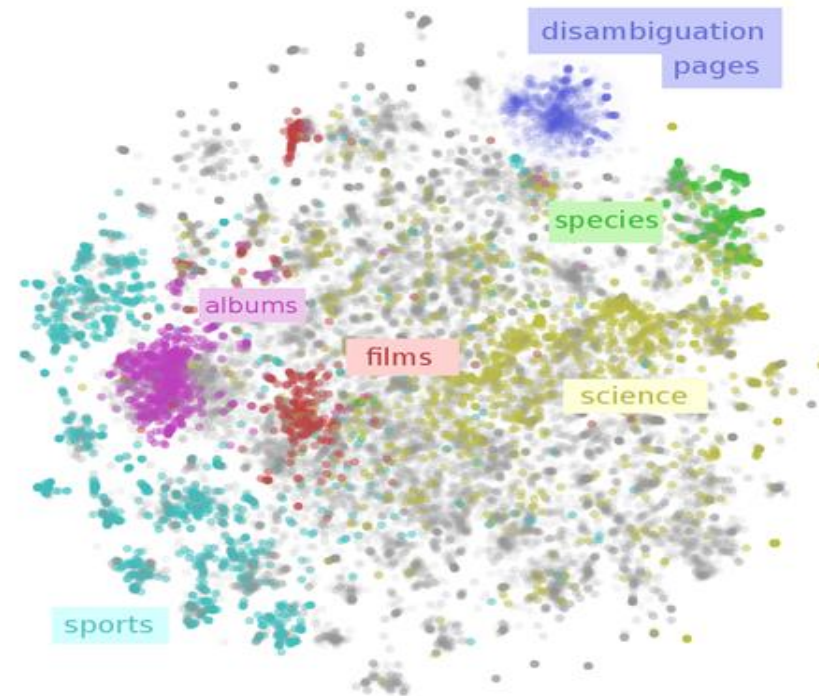


Classical cross-lingual transfer

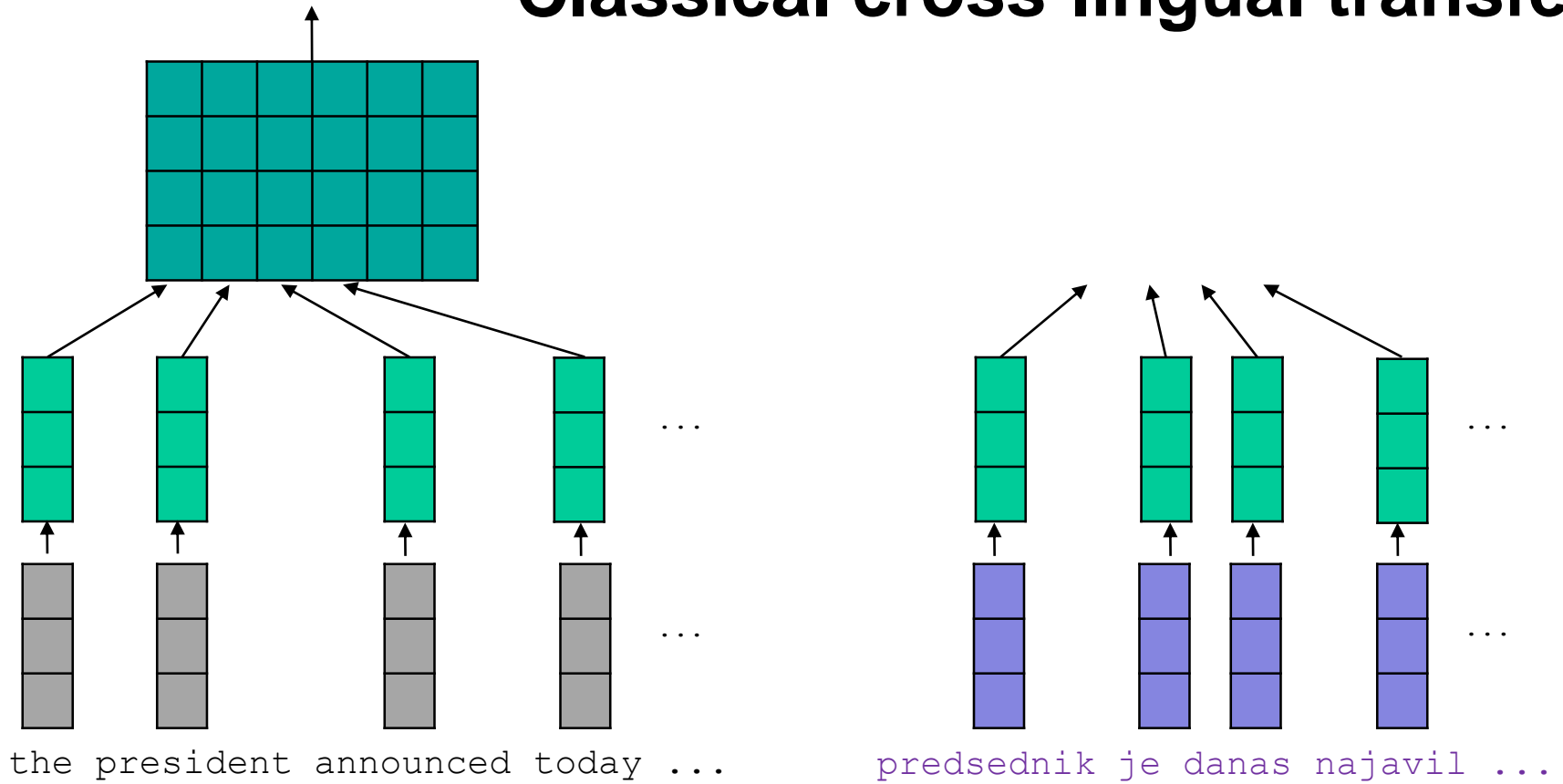
Explicit alignment of vector spaces

$$WS \approx E$$

Nowadays: use multilingual LLMs directly



Classical cross-lingual transfer

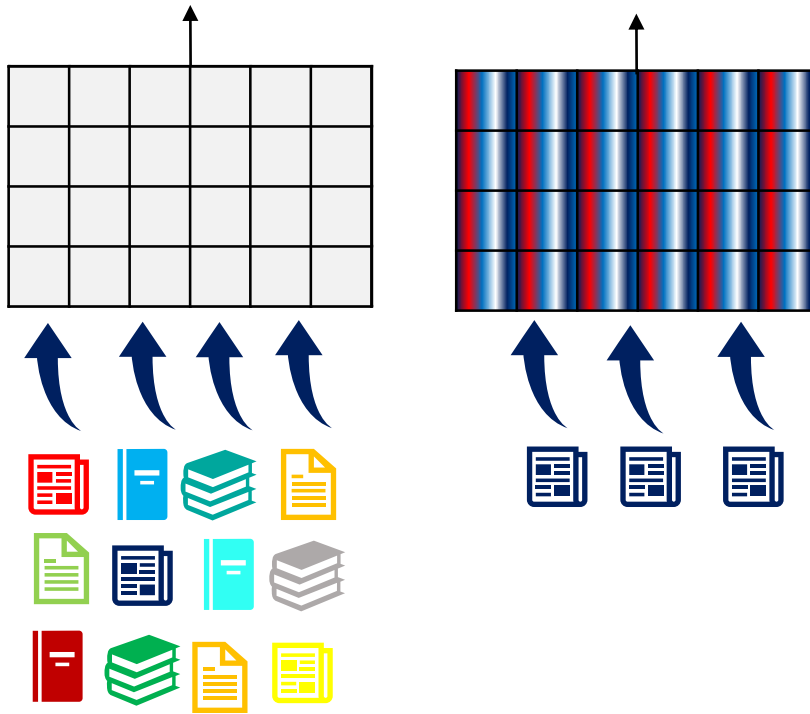




Multilingual PLMs

- Pretrained on multiple languages simultaneously
- multilingual BERT supports 104 languages by training on Wikipedia
- XLM-R was trained on 2.5 TB of texts
- these models allow cross-lingual transfer
- solve problem of insufficient training resources for less-resourced languages
- zero-shot transfer and few-shot transfer

Using multilingual models



Pretraining

Fine-tuning

Classification

predsednik je danas najavil ...

Zero-shot transfer and few-shot transfer





Why not only multilingual?

- performance on many tasks drops with more languages
- results for a few tasks in Slovene (Named Entity Recognition – NER, Part-of-Speech Tagging – POS, Dependency Parsing – DP, Sentiment Analysis – SA, Word Analogy – WA)

Model	NER	POS	DP	SA	WA
mBERT	0.885	0.984	0.681	0.576	0.061
XLM-R	0.912	0.988	0.793	0.604	0.146
SloBERTa	0.933	0.991	0.844	0.623	0.405

Dictionaries in PLMs

- Tokenization depends on the dictionary
- The dictionary is constructed statistically (BPE or SentencePiece algorithm)
- Sentence: “Letenje je bilo predmet precej starodavnih zgodb.”

- SloBERTa:

'_Le', 'tenje', '_je', '_bilo', '_predmet', '_precej', '_staroda', 'vnih', '_zgodb', '.'

- mBERT:

'Let', '##en', '##je', 'je', 'bilo', 'pred', '##met', 'pre', '##cej', 'star', '##oda', '##vnih', 'z', '##go', '##d', '##b', '.'



Trade off: trilingual models

- BERT trained with only a few languages
- more data for training
- more specific dictionary
- good for cross-lingual transfer
- Trilingual models
 - CroSloEngual BERT
 - FinEst BERT
 - LitLat BERT
- SlavBERT (ru, pl, cs, bg; DeepPavlov)

Model	NER	POS	DP	SA	WA
mBERT	0.885	0.984	0.681	0.576	0.061
XLM-R	0.912	0.988	0.793	0.604	0.146
CSE-BERT	0.928	0.990	0.854	0.610	0.195
SloBERTa	0.933	0.991	0.844	0.623	0.405

• Ulčar, M., & Robnik-Šikonja, M. (2020). FinEst BERT and CroSloEngual BERT. In International Conference on Text, Speech, and Dialogue (pp. 104-111).



XL transfer in classification

- Excellent XL transfer between similar languages like Slovene and Croatian
- But: the transfer quality is problem dependent

sentiment analysis

Source	Target	LASER		mBERT		CSE BERT		Both target	
		\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA
Croatian	Slovene	0.53	0.53	0.53	0.54	0.61	0.60	0.60	0.60
Croatian	English	0.63	0.63	0.63	0.66	0.62	0.64	0.62	0.65
English	Slovene	0.54	0.57	0.50	0.53	0.59	0.57	0.60	0.60
English	Croatian	0.62	0.67	0.67	0.63	0.73	0.67	0.73	0.68
Slovene	English	0.63	0.64	0.65	0.67	0.63	0.64	0.62	0.65
Slovene	Croatian	0.70	0.65	0.64	0.63	0.73	0.69	0.73	0.68
Croatian English	Slovene	0.54	0.54	0.53	0.54	0.60	0.58	0.60	0.60
Croatian Slovene	English	0.62	0.61	0.65	0.67	0.63	0.65	0.62	0.65
English Slovene	Croatian	0.64	0.68	0.63	0.63	0.68	0.70	0.73	0.68
Average performance gap		0.04	0.03	0.04	0.03	0.00	0.01		

idiom detection

Language	Slovene ELMo	mBERT	Default F_1
Slovene	0.8163	0.8359	0.667
Croatian	0.9191	0.8970	0.667
Polish	0.2863	0.6987	0.667

- Robnik-Šikonja, M., Reba, K., & Mozetič, I. (2021). Cross-lingual transfer of sentiment classifiers. *Slovenščina 2.0: empirical, applied and interdisciplinary research*, 9(1), 1-25.
- Škvorc, T., Gantar, P., & Robnik-Šikonja, M. (2022). MICE: mining idioms with contextual embeddings. *Knowledge-Based Systems*, 235, 107606.



What PLMs learn?

- We would like to travel to [MASK], ki je najlepši otok v Mediteranu.

SloBERTa: ..., Slovenija, I, Koper, Slovenia

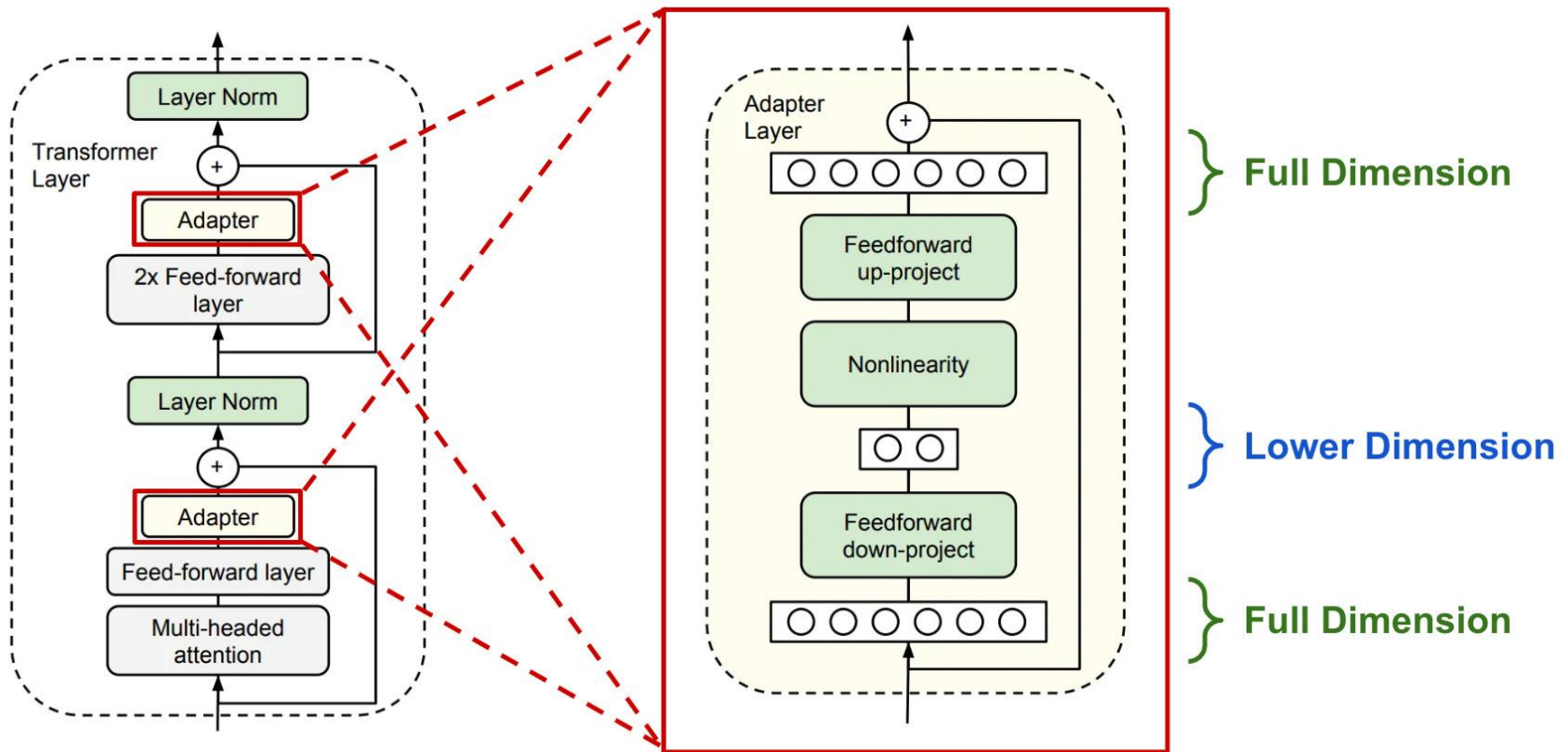
CSE-BERT: Hvar, Rab, Cres, Malta, Brač

XLM-R: Mallorca, Tenerife, otok, Ibiza, Zadar

mBERT: Ibiza, Gibraltar, Tenerife, Mediterranean, Madeira

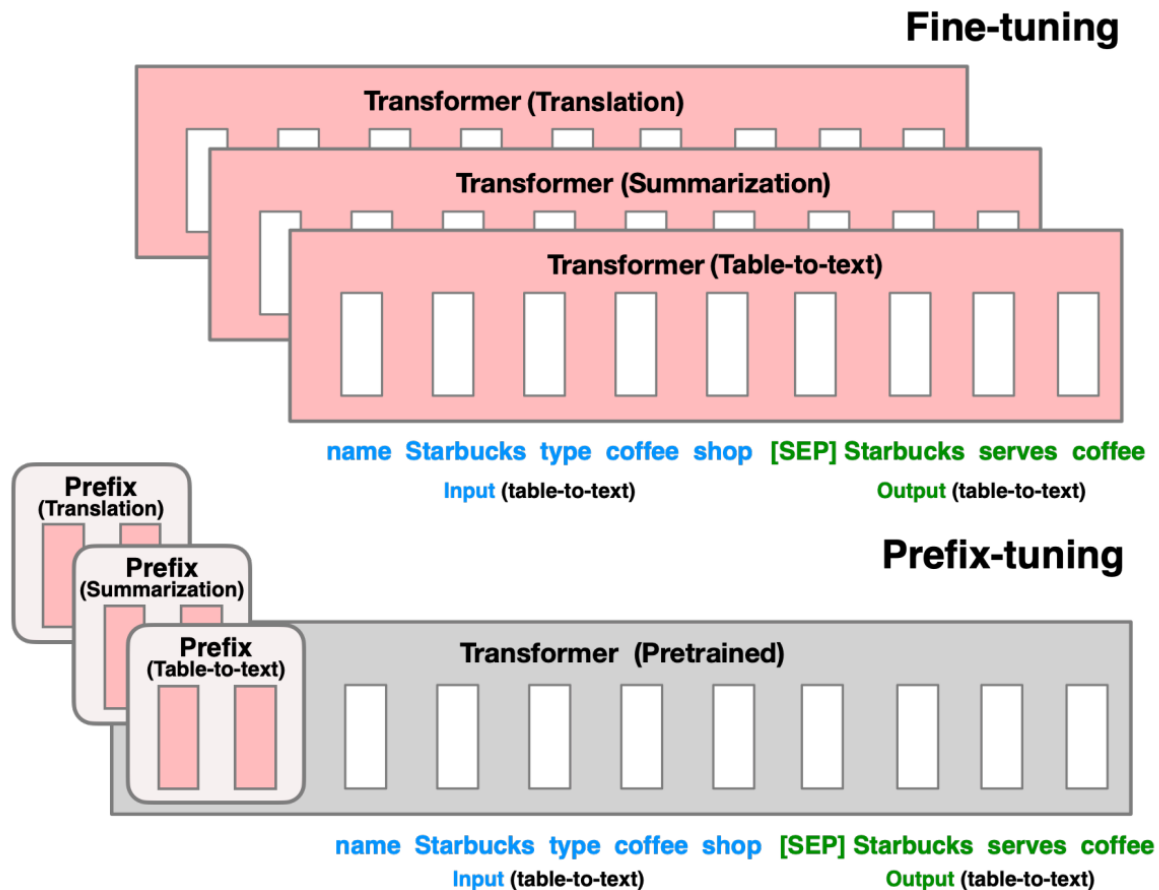
BERT (en): Belgrade, Italy, Serbia, Prague, Sarajevo

Fine-tuning with adapters for transformers



Prefix tuning

- Comparison of classical finetuning (top) with prefix tuning which adds different prefix blocks for each task but other weights are frozen



LoRA – Low Rank Approximation

- See the Substack [post](#)
- injects AB into each layer
- r can be very small, e.g., 1% of original d

Finetuned Weights

Weight Update

$$W_{\text{ft}} = \underbrace{W_{\text{pt}}}_{\text{Pretrained Weights}} + \underbrace{\Delta W}_{\text{Weight Update}}$$

Rank Decomposition Matrix

$$W_{\text{ft}} = W_{\text{pt}} + \underbrace{\Delta W}_{\text{Approximation}} = W_{\text{pt}} + \underbrace{AB}_{\text{Rank Decomposition Matrix}}$$

where $W_{\text{ft}}, W_{\text{pt}}, \Delta W, AB \in \mathbb{R}^{d \times d}$

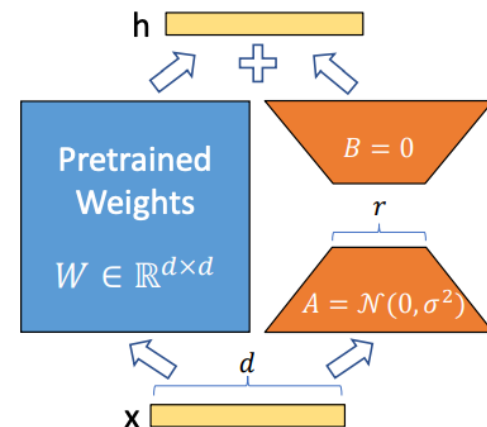
$$\text{and } \underbrace{A \in \mathbb{R}^{d \times r}, B \in \mathbb{R}^{r \times d}}_{\text{Low Rank}}$$

Low Rank

Scaling Factor

$$W_{\text{ft}} = W_{\text{pt}} + \underbrace{\frac{\alpha}{r}}_{\text{Scaling Factor}} \underbrace{AB}_{\text{Rank Decomposition Matrix}}$$

Rank Decomposition Matrix



Attention efficiency

- time and space complexity of self-attention grows quadratically with n (size of input)

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \quad K, V, Q \in \mathbb{R}^{n \times d}$$

- not suitable for very long sequences like
 - documents
 - character-level language models
 - images (as sequences of pixels);
 - protein sequences.

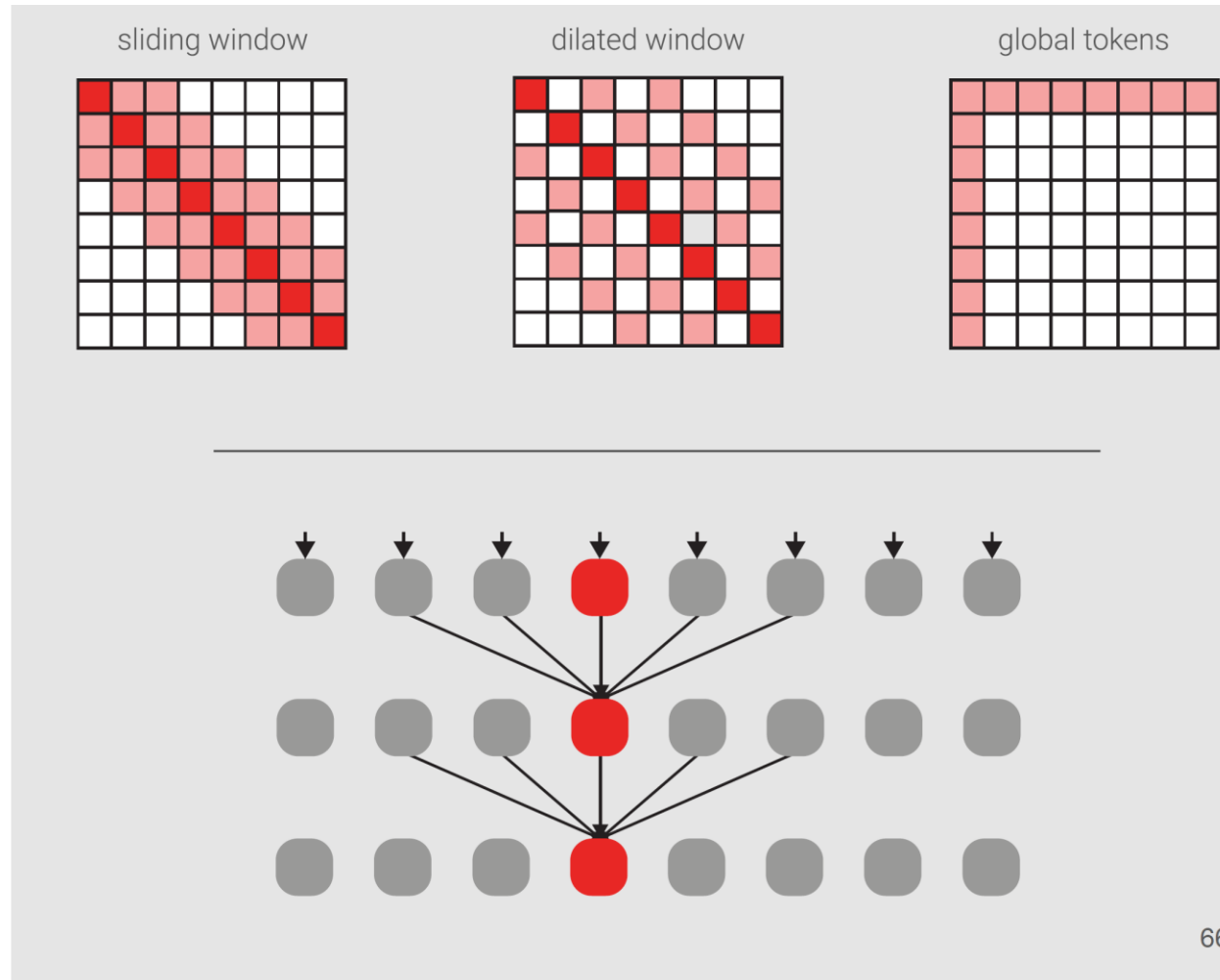
Longformer [1]

- sliding window attention:
each position can attend to $1/2w$
tokens on each side - $O(w \times n)$

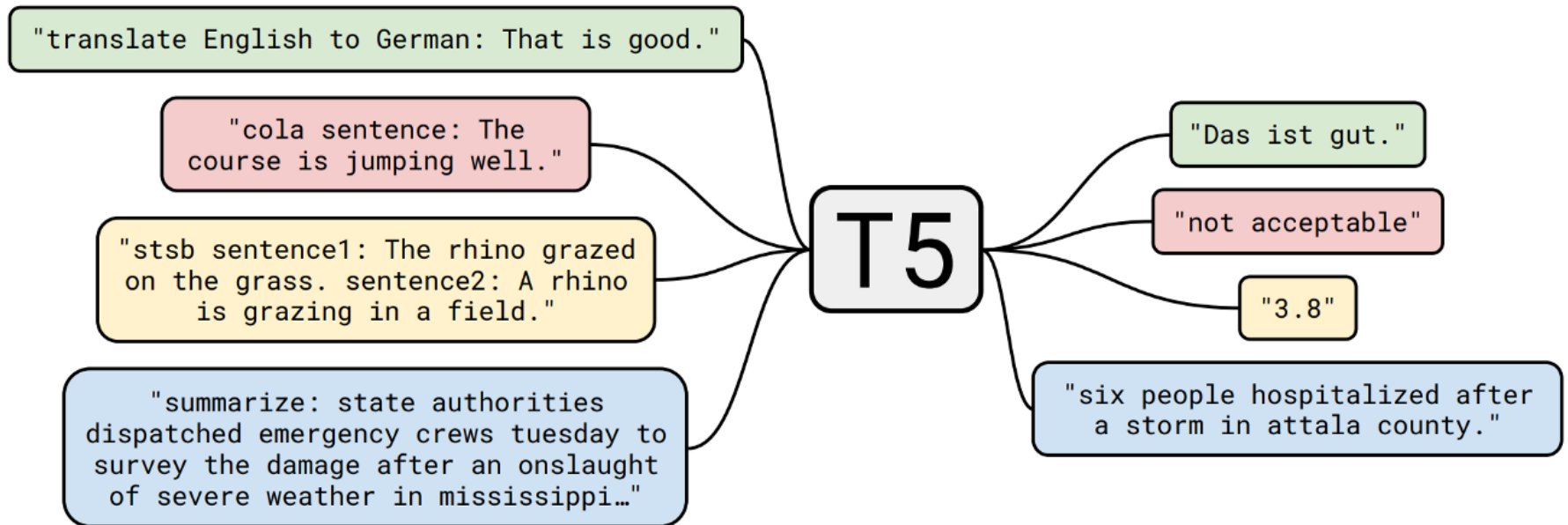
- dilated window attention:
increases the receptive field of the
attention layer - $O(w \times n)$

- global attention:
 k special tokens that aggregate
information from whole sequence
(e.g. [CLS] as in BERT) - $O(k \times n)$

[1] [Beltagy et al.: Longformer: The Long-Document Transformer, 2020.](#)



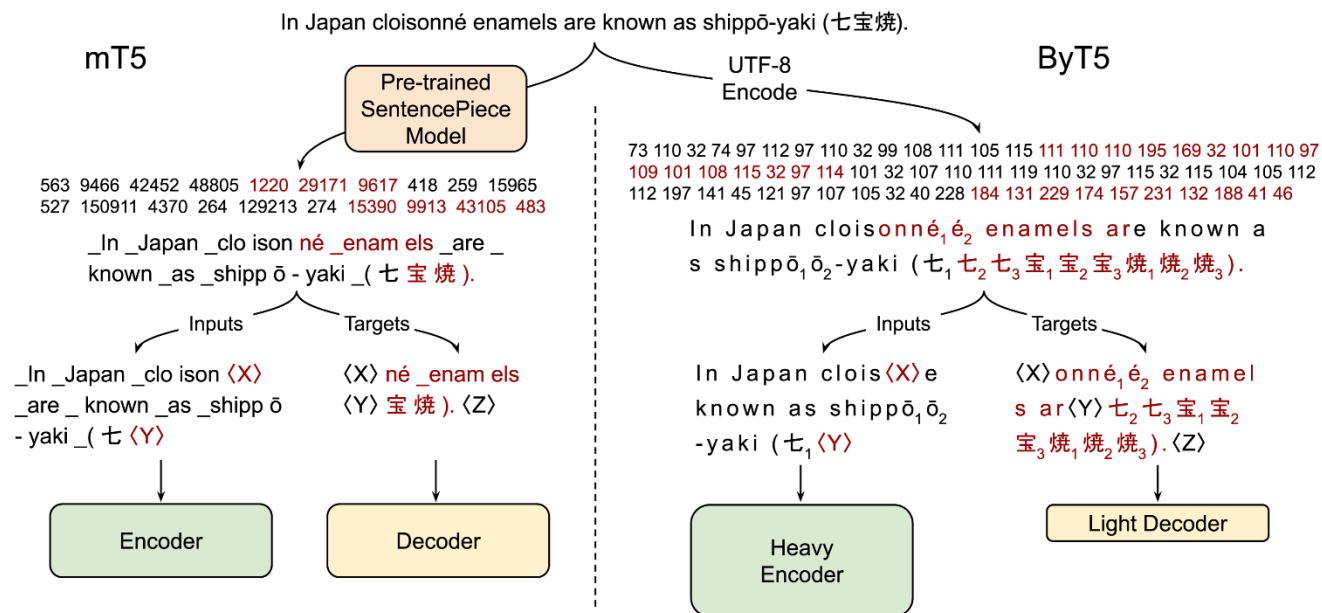
T5 (Text-To-Text Transfer Transformer) models



- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y, Li, W. & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140), 1-67.
- Ulčar, M. & Robnik-Šikonja, M. (2023) Sequence-to-sequence pretraining for a less-resourced Slovenian language. *Frontiers in Artificial Intelligence*, Section on Language and Computation, Volume 6 – 2023, <https://doi.org/10.3389/frai.2023.932519>

Character-based T5

- Character-based input requires more input layers
- But is very suitable for nonstandard spelling, misspelled words, dialects
- Can process text in any language out of the box
- More robust to noise
- Simplifies text preprocessing



Transformers are everywhere

- music: LLM where the vocabulary consists of MIDI pitches, pauses, velocity
- object detection (attention to objects)
- time series, where signals are discretized and values are treated as letters

