

```

Transformers

What does one do with that Transformer?



```

<code>
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="transform.xsl"?>
<!-- This is the XML document -->
<root>
    <node>A</node>
    <node>B</node>
    <node>C</node>
</root>
</code>

```

Transformer has a method transform() which takes an InputSource and an OutputWriter. So we can do something like this:



```

<code>
Transformer transformer = TransformerFactory.newInstance().newTransformer();
InputSource source = new InputSource(new StringReader(xml));
OutputWriter writer = new OutputWriter(System.out);
transformer.transform(source, writer);
writer.close();
</code>

```



And that's it! That's all there is to it. The Transformer class is a very simple class. It's just a thin layer of abstraction over the underlying Transformer interface. The Transformer interface is implemented by the TransformerImpl class, which is part of the com.sun.org.apache.xalan.internal.xsltc.trax package. This class is responsible for doing the actual transformation work. It uses the Transformer interface to interact with the underlying XSLT processor. The Transformer class also provides some convenience methods for working with DOM documents, such as transform(DOMSource, DOMResult) and transform(DOMSource, ResultTreeHandler).
```



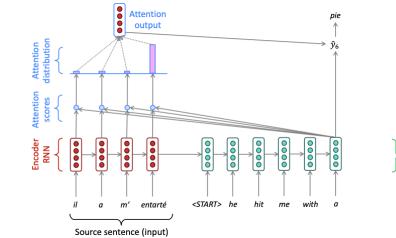
My own latex definitions

- Today's lecture
- Limitations of BERT
- Self and Cross-Attention
- The Transformer Architecture

This lecture material is taken from

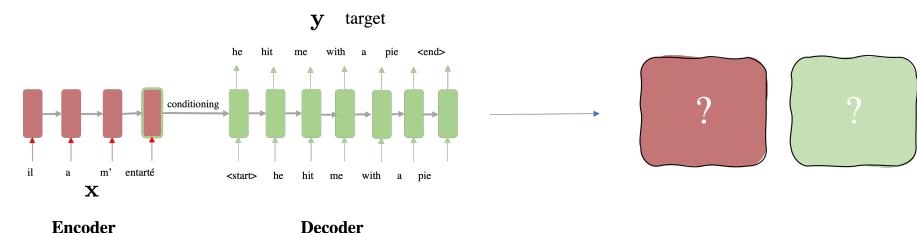
- [Chapter 6, 10, 11, 12 of the book]
- [Chapter 6 of the slides]

A detailed note on the Transformer architecture is available at <https://www.csail.mit.edu/courses/6.S090/Fall2018/transformer.pdf>.

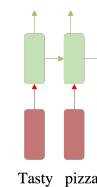


2014-2017

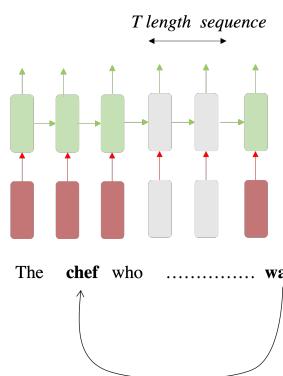
2023



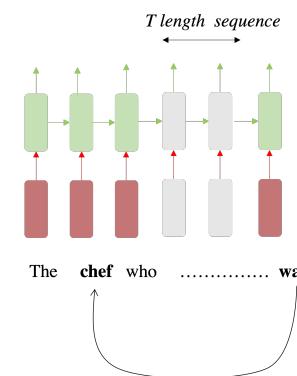
- **Limitations of RNN**
- Linear interaction distance
- Infinite Lack of Parallelizability
- **Limitations of RNN: Linear interaction distance**
  - All states are connected via input
  - All states are connected via hidden-to-hidden
  - Highly sensitive often requires each step's memory
  - Problem: width was  $O(N)$ , for character width  $N$  it's the sequence length.



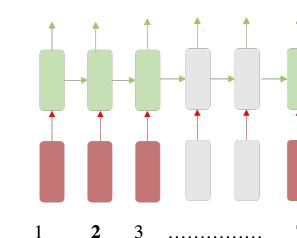
- RNNs are unrolled with right
- The activities based directly on the previous hidden state
- The activities based directly on the input



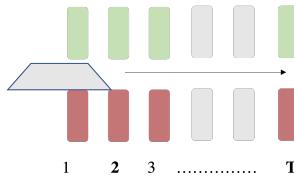
**Limitations of RNN: Linear interaction distance**  
diff. steps for distant word pairs in linear model  
A TINY LINEAR-DISTANCE BIAS: because of vanishing gradient problem.



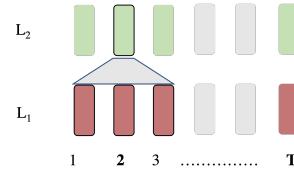
**Limitations of RNN: Lack of Parallelizability**



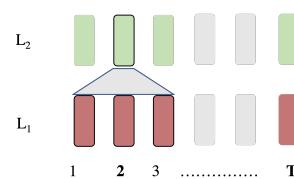
If not recurrence, then what?  
What about window-based classifier?  
What about parallel computation?



What about window-based classifier?  
What about parallel computation?

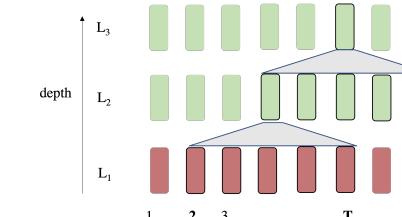


OK, good but how we make long-distance dependencies?

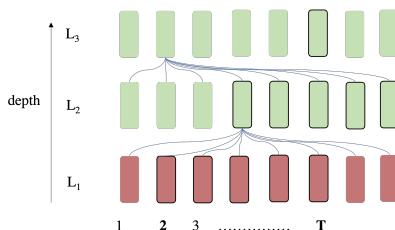
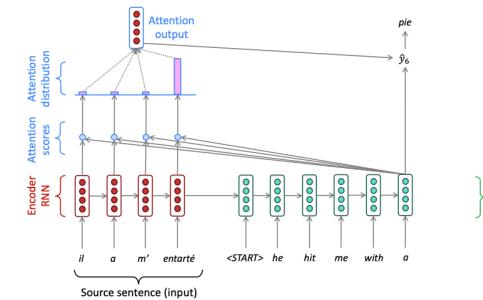


Idea?

Long-distance dependencies with local context [Receptive Field]  
Memory access time is proportional to the receptive field size  
- we can "press on" parallelizability over depth



Good idea, memory access  
parallelizability in time is a must (RNN do not have that)  
- we can "press on" parallelizability over depth  
What about Attention? [\[1\]](#)



## Attention Is All You Need

Ashish Vaswani\* Noam Shazeer\* Niki Parmar\* Jakob Uszkoreit\*  
Google Brain Google Brain Google Research Google Research  
avaswani@google.com noam@google.com nikip@google.com usz@google.com

Llion Jones\* Aidan N. Gomez\* † Lukasz Kaiser\*  
Google Research University of Toronto Google Brain  
llion@google.com aidan@cs.toronto.edu lukaszkaiser@google.com

Illia Polosukhin\* ‡  
illia.polosukhin@gmail.com

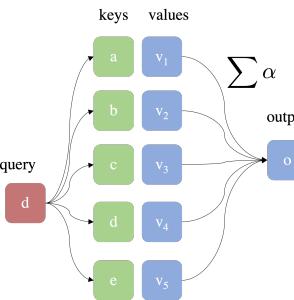
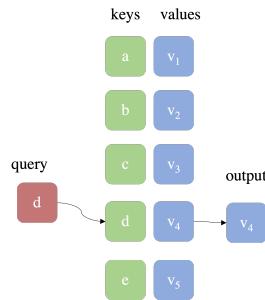
### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

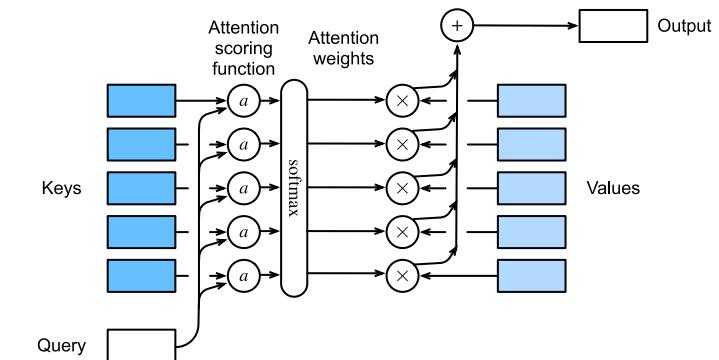
A few words on query, key, values  
Memory access time is proportional to the receptive field size  
- we can "press on" parallelizability over depth  
What about Attention? [\[1\]](#)

Query, key, values terminology is from databases  
 Keys are the columns of the database table, and values are the rows.  
 A query is a row of the database table, and the output is the result of the query.

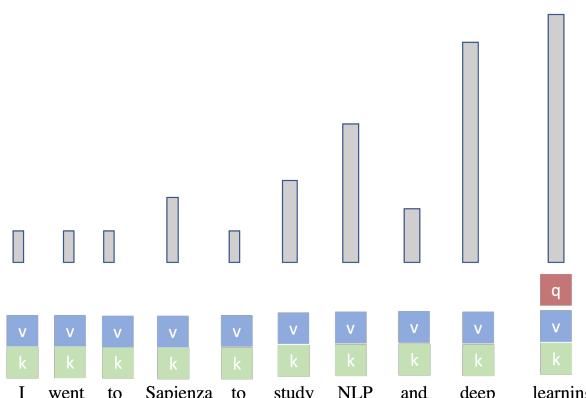
Lookup table vs Soft-average Lookup table



Self-Attention: soft, averaging/lookup table



Self-Attention



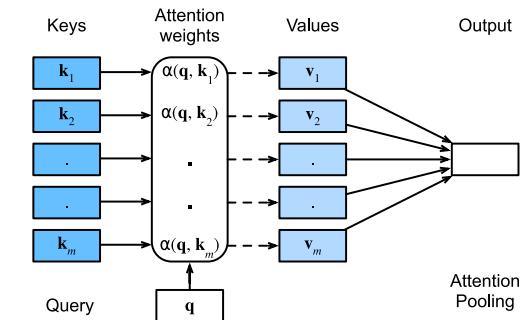
Self-Attention: soft, averaging/lookup table

Properties of  $\alpha(q, k_i)$ :  
 1.  $\alpha(q, k_i) \geq 0$  for all  $i$ . This means that attention weights are non-negative.  
 2.  $\sum_i \alpha(q, k_i) = 1$ . This means that the sum of attention weights is 1.  
 3.  $\alpha(q, k_i) = 0$  if  $k_i$  is orthogonal to  $q$ . This means that attention weights are zero for words that are orthogonal to the query.  
 4.  $\alpha(q, k_i) = 1$  if  $k_i$  is parallel to  $q$ . This means that attention weights are 1 for words that are parallel to the query.  
 5.  $\alpha(q, k_i) = 0$  if  $k_i$  is orthogonal to  $q$ . This means that attention weights are zero for words that are orthogonal to the query.

Self-Attention Properties

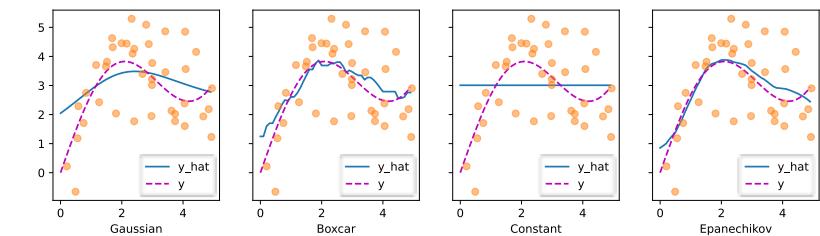
Softmax normalization

Self-Attention

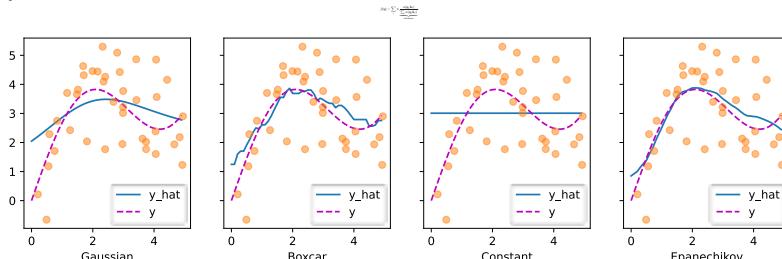


So far, nothing is learned just linear combination  
 Before going into the learning part, let us make a connection with classic machine learning  
 Attention Pooling via Nadarya-Watson Regression [1964]

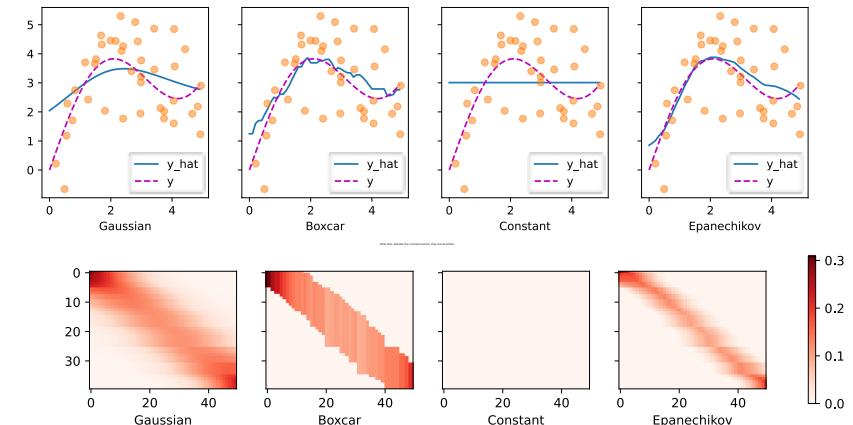
$$\hat{y} = \sum \frac{\alpha_i y_i}{\sum \alpha_i}$$



Attention Pooling via Nadaraya-Watson Regression



Attention Weights

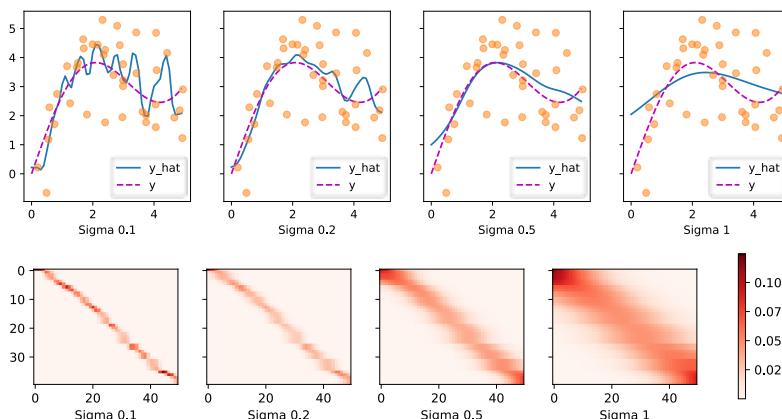


Attention Pooling via Nadaraya-Watson Regression

Attention Pooling via Nadaraya-Watson Regression

Attention Pooling via Nadaraya-Watson Regression

Let's tune the bandwidth parameter in the [Gaussian] kernel



Instead of tuning the bandwidth why not learning?

Learnable self-attention

Attention weights are learned by the model

Learnable self-attention

The learnable self-attention weights are computed as:

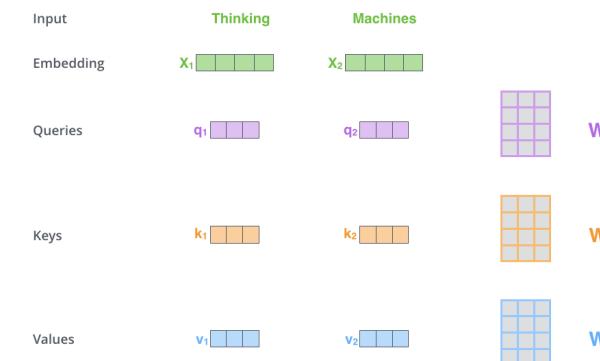
$\alpha_{ij} = \text{softmax}\left(\frac{\langle \mathbf{q}_i, \mathbf{k}_j \rangle}{\sqrt{d_k}}\right)$

where  $\mathbf{q}_i$  is the query vector and  $\mathbf{k}_j$  is the key vector.

Attention weights are computed as:

$\mathbf{v}_i = \sum_j \alpha_{ij} \mathbf{v}_j$

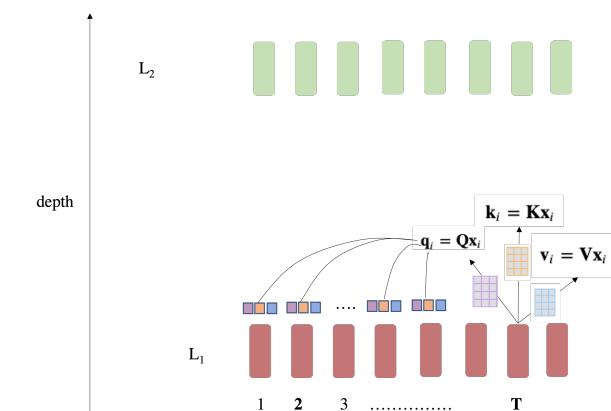
Computation with two word tokens

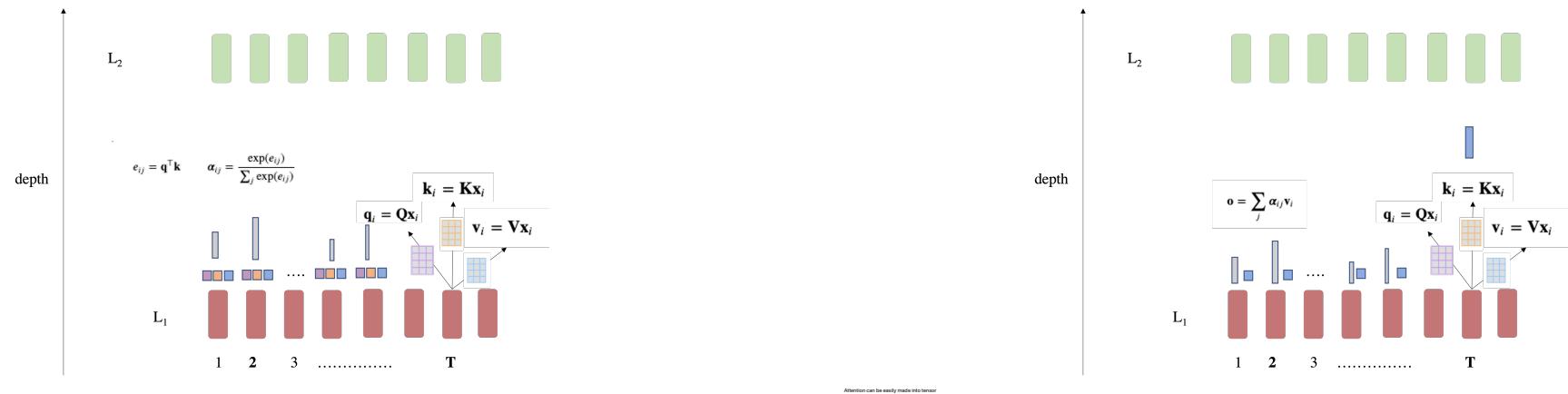




Computation with two word tokens

Input	<b>Thinking</b>	<b>Machines</b>
Embedding	$x_1$ [green green green green]	$x_2$ [green green green green]
Queries	$q_1$ [purple purple purple]	$q_2$ [purple purple purple]
Keys	$k_1$ [orange orange orange]	$k_2$ [orange orange orange]
Values	$v_1$ [blue blue blue]	$v_2$ [blue blue blue]
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$
Divide by 8 ( $\sqrt{d_k}$ )	14	12
Softmax	0.88	0.12
Softmax X Value	$v_1$ [blue blue blue]	$v_2$ [blue blue blue]
Sum	$z_1$ [pink pink pink]	$z_2$ [pink pink pink]





$X$

$$\begin{matrix} \begin{matrix} \begin{matrix} \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \end{matrix} \end{matrix} \times \begin{matrix} \begin{matrix} \text{purple} & \text{purple} & \text{purple} \\ \text{purple} & \text{purple} & \text{purple} \\ \text{purple} & \text{purple} & \text{purple} \end{matrix} \end{matrix} = \begin{matrix} \begin{matrix} \text{purple} & \text{purple} & \text{purple} \\ \text{purple} & \text{purple} & \text{purple} \\ \text{purple} & \text{purple} & \text{purple} \end{matrix} \end{matrix} \quad Q$$

$X$

$$\begin{matrix} \begin{matrix} \begin{matrix} \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \end{matrix} \end{matrix} \times \begin{matrix} \begin{matrix} \text{orange} & \text{orange} & \text{orange} \\ \text{orange} & \text{orange} & \text{orange} \\ \text{orange} & \text{orange} & \text{orange} \end{matrix} \end{matrix} = \begin{matrix} \begin{matrix} \text{orange} & \text{orange} & \text{orange} \\ \text{orange} & \text{orange} & \text{orange} \\ \text{orange} & \text{orange} & \text{orange} \end{matrix} \end{matrix} \quad K$$

$X$

$$\begin{matrix} \begin{matrix} \begin{matrix} \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} \end{matrix} \end{matrix} \times \begin{matrix} \begin{matrix} \text{blue} & \text{blue} & \text{blue} \\ \text{blue} & \text{blue} & \text{blue} \\ \text{blue} & \text{blue} & \text{blue} \end{matrix} \end{matrix} = \begin{matrix} \begin{matrix} \text{blue} & \text{blue} & \text{blue} \\ \text{blue} & \text{blue} & \text{blue} \\ \text{blue} & \text{blue} & \text{blue} \end{matrix} \end{matrix} \quad V$$

Attention with tensor

$$\text{softmax}\left(\frac{\mathbf{Q}^T \mathbf{K}}{\sqrt{d_k}}\right) \mathbf{V} = \mathbf{Z}$$

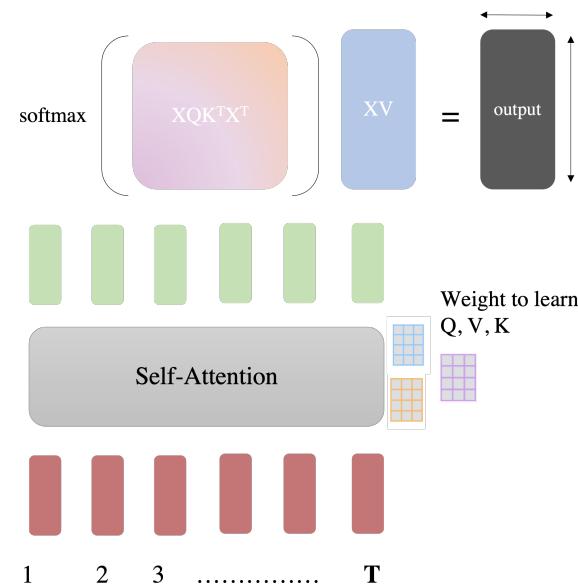
$XQ$

$K^T X^T$

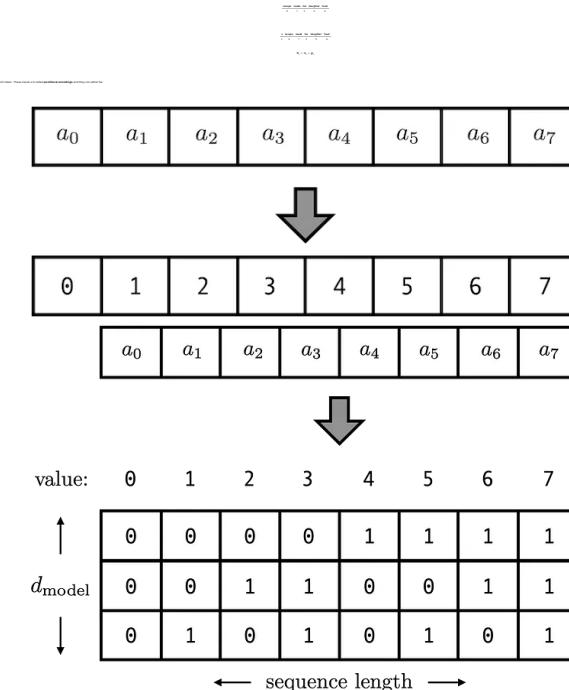
$= X Q K^T X^T$

n n

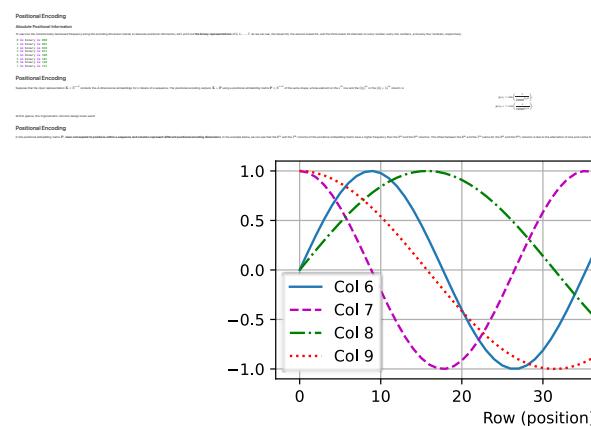
Self-Attention with tensor



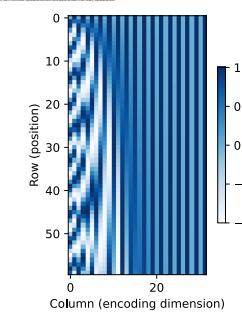
Self-attention as a NLP block: (File 1) No notion of order!



No Skill many things to file  
Self-attention as a NLP block: (File 1) No notion of order!  
Positional Encoding - Absolute Positional Info

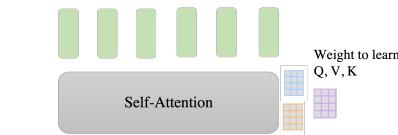


Positional Encoding - Absolute Positional Info



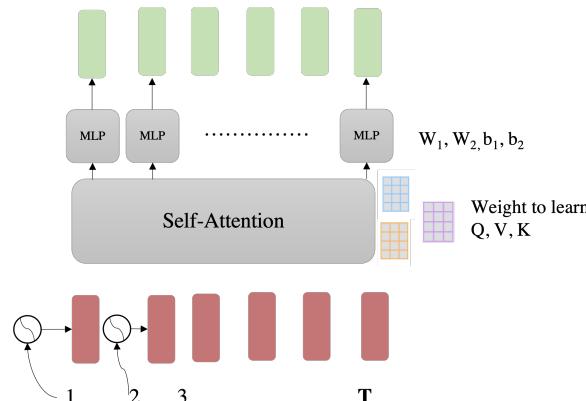
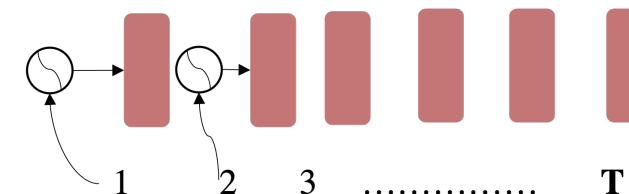
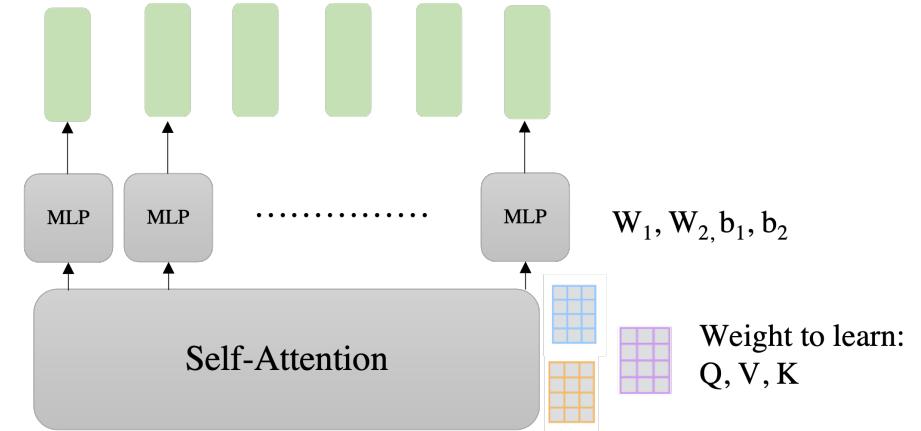
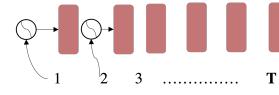
Positional Encoding - Relative Positional Info

$$\begin{bmatrix} \sin(\frac{\pi}{N}r_1) & \cos(\frac{\pi}{N}r_1) & \dots & \cos(\frac{\pi}{N}r_N) \\ \vdots & \vdots & \ddots & \vdots \\ \sin(\frac{\pi}{N}r_1) & \cos(\frac{\pi}{N}r_1) & \dots & \cos(\frac{\pi}{N}r_N) \end{bmatrix}$$



Self-attention as a NLP block (Fix 2) No nonlinearity! Just weighted average.

Fix 2) Adding non-linearity



Self-attention as a NLP block (Fix 2) We can look in the future!

Only used in the Decoder

Masking Self-attention

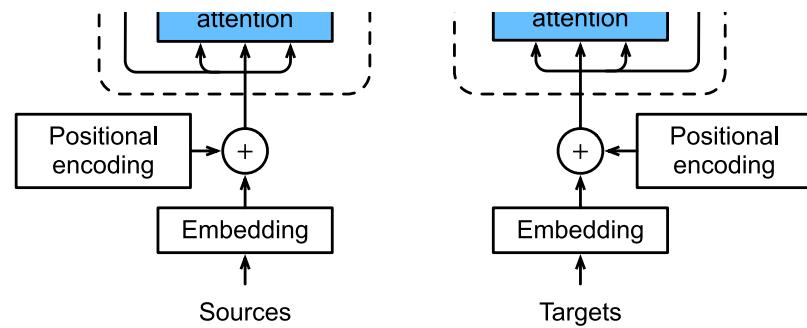
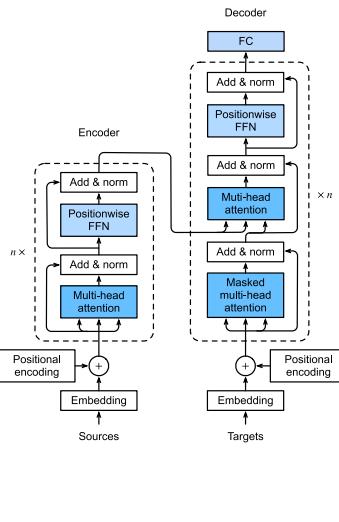
For encoding  
this word

We attend to those non-gray

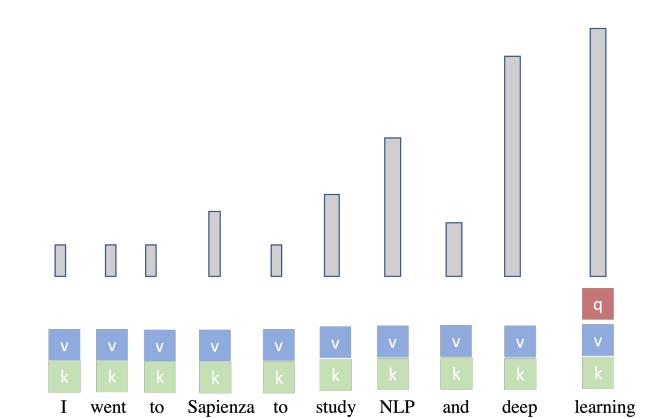
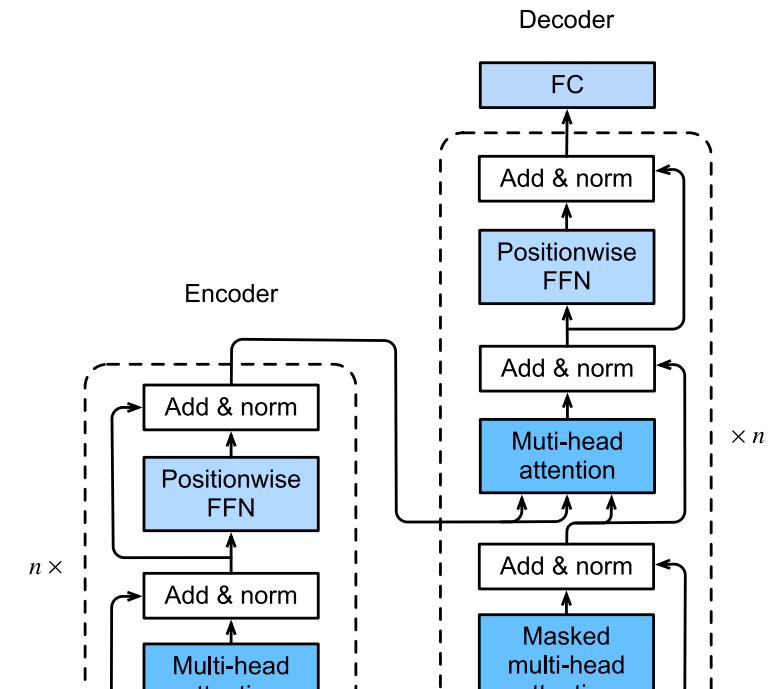
	<start>	Iacopo	made	his
<start>		-∞	-∞	-∞
Iacopo			-∞	-∞
made				-∞
his				



The Transformer



The world is "multi-modal"

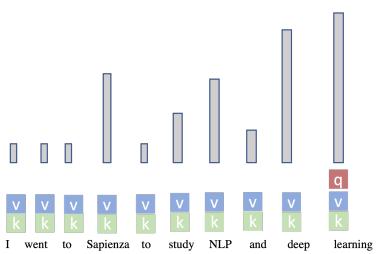


Multi-head self-attention

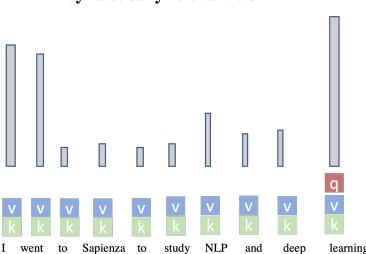


Multi-head self-attention

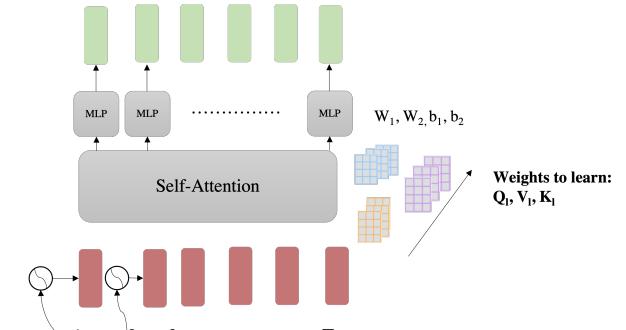
Attention mechanism #1 attends to “entities”



Attention mechanism #2 attends to “syntactically relevant words”

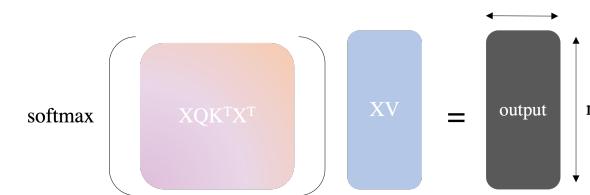


Multi-head self-attention

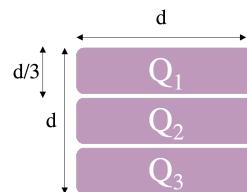
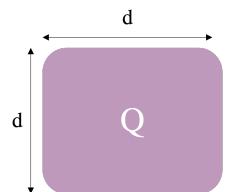


Weights to learn:  
 $Q_1, V_1, K_1$

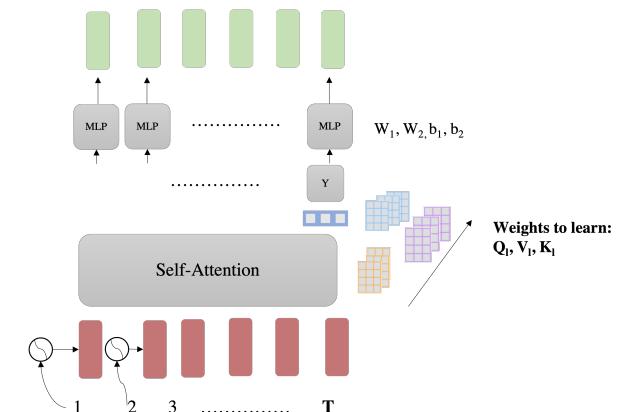
Self-Attention with tensor



Multi Head Self-Attention with tensor



Multi Head Self-Attention with tensor



Weights to learn:  
 $Q_1, V_1, K_1$

2) Residual Connection [He et al., 2016]

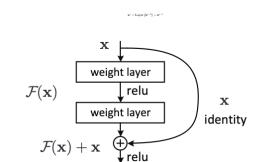


Figure 2. Residual learning: a building block.

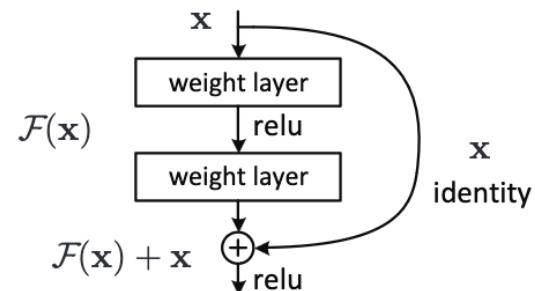
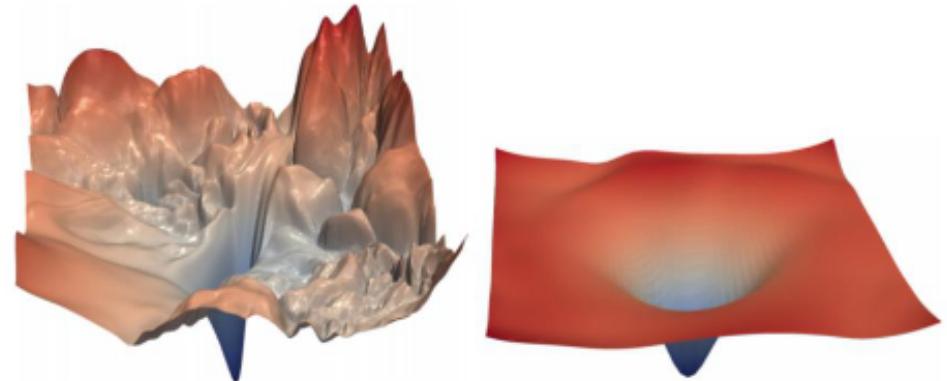
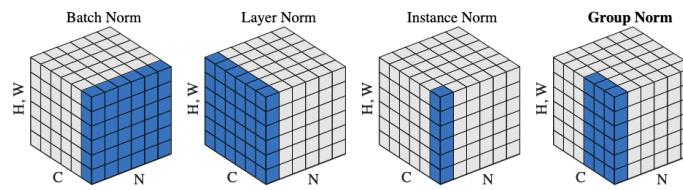


Figure 2. Residual learning: a building block.

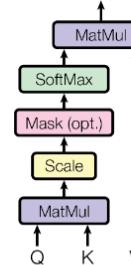


[no residuals]

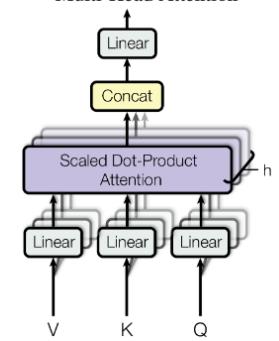
[residuals]



### Scaled Dot-Product Attention



### Multi-Head Attention



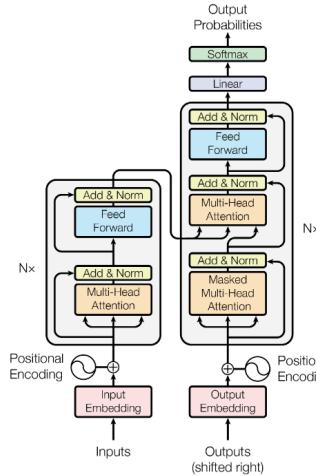
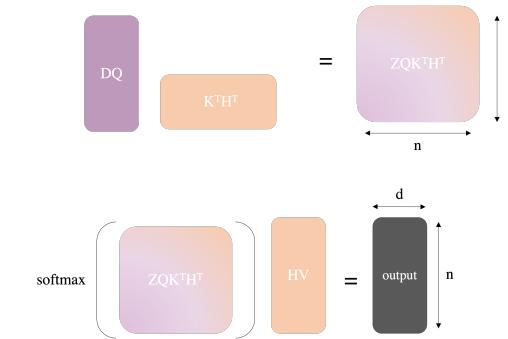


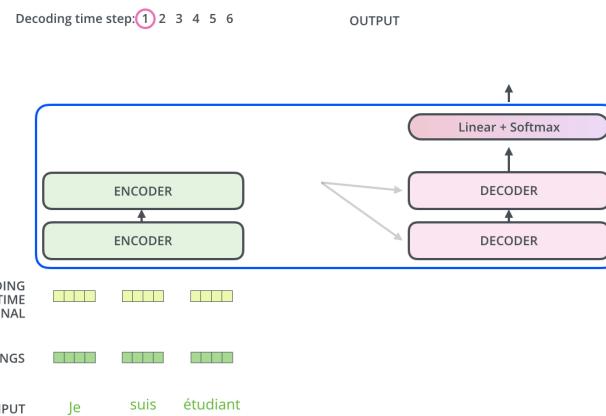
Figure 1: The Transformer - model architecture.

Cross-Attention with Tensor  
 $\mathbf{Q} = \mathbf{A}_1 \cdots \mathbf{A}_n \mathbf{U}^T$  query matrix of cross-attention  
 $\mathbf{K} = \mathbf{A}_1 \cdots \mathbf{A}_n \mathbf{V}^T$  key matrix of cross-attention  
 $\mathbf{H} = \mathbf{A}_1 \cdots \mathbf{A}_n \mathbf{U}^T \mathbf{V}$  value matrix of cross-attention  
 $\mathbf{DQ} = \mathbf{Q}^T \mathbf{Q}$  dot product of query and query

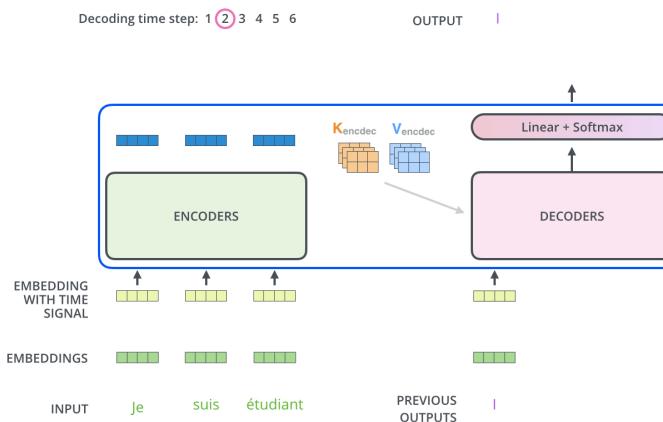


Cross-Attention  
 $\mathbf{Q} = \mathbf{A}_1 \cdots \mathbf{A}_n \mathbf{U}^T$  query matrix of cross-attention  
 $\mathbf{K} = \mathbf{A}_1 \cdots \mathbf{A}_n \mathbf{V}^T$  key matrix of cross-attention  
 $\mathbf{H} = \mathbf{A}_1 \cdots \mathbf{A}_n \mathbf{U}^T \mathbf{V}$  value matrix of cross-attention  
 $\mathbf{DQ} = \mathbf{Q}^T \mathbf{Q}$  dot product of query and query

Transformers Decoding



Transformers Decoding



Which word in our vocabulary  
is associated with this index?

am

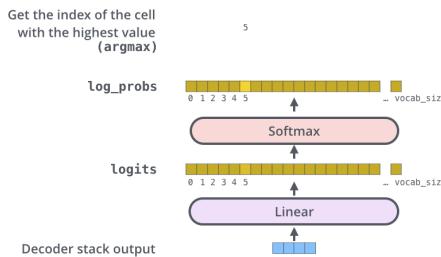


Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

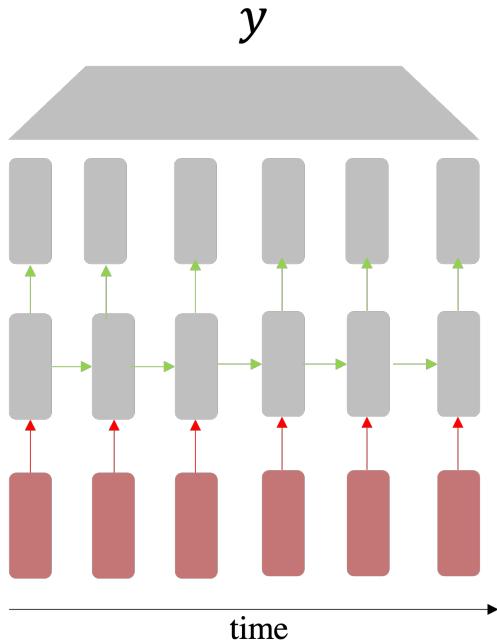
Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	<b><math>3.3 \cdot 10^{18}</math></b>	
Transformer (big)	<b>28.4</b>	<b>41.8</b>	$2.3 \cdot 10^{19}$	

Table 4: Performance of best models of each model architecture using the combined corpus and tf-idf extractor.

Model	Test perplexity	ROUGE-L
seq2seq-attention, $L = 500$	5.04952	12.7
Transformer-ED, $L = 500$	2.46645	34.2
Transformer-D, $L = 4000$	2.22216	33.6
Transformer-DMCA, no MoE-layer, $L = 11000$	2.05159	36.2
Transformer-DMCA, MoE-128, $L = 11000$	1.92871	37.9
Transformer-DMCA, MoE-256, $L = 7500$	1.90325	38.8



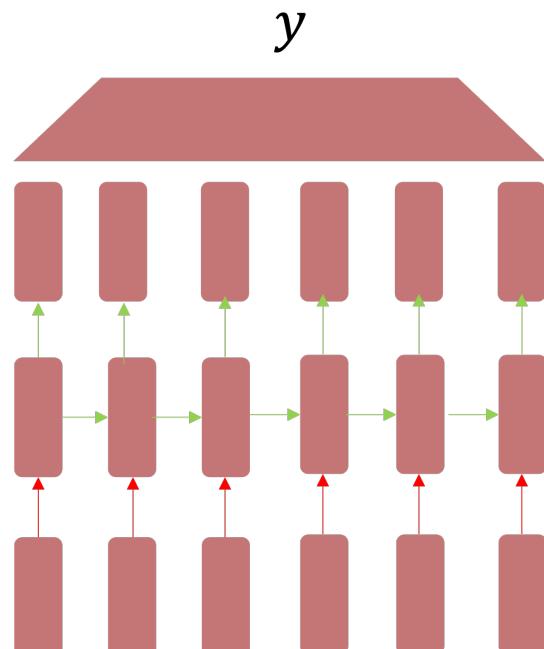
Pretrained word embeddings



Not-pretrained (random weights)

Pre-trained (word-embedding)

Pretrained model



Pre-trained jointly

Where we are: pre-trained whole model



Pre-training (Supervised Learning)

What can we learn from reconstructing the input?

Answer: There are 6 words in \_\_\_\_\_.

Q: Are \_\_\_\_\_ words in the table?

The answer method is the one that checks for words and \_\_\_\_\_.

It asks the user to read the table and \_\_\_\_\_.

Answer: The value of  $\pi$  is 3.141592653589793.

The user needs to calculate the area of the circle and the table.

Q: What is the value of  $\pi$ ?

The answer method is the one that checks for words and \_\_\_\_\_.

It asks the user to calculate the area of the circle and the table.

Q: What is the value of  $\pi$ ?

The answer method is the one that checks for words and \_\_\_\_\_.

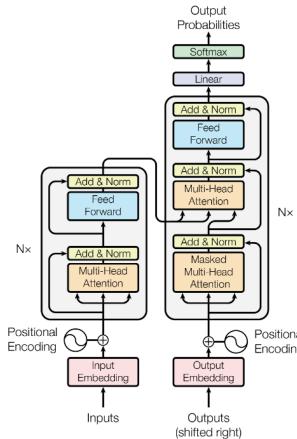
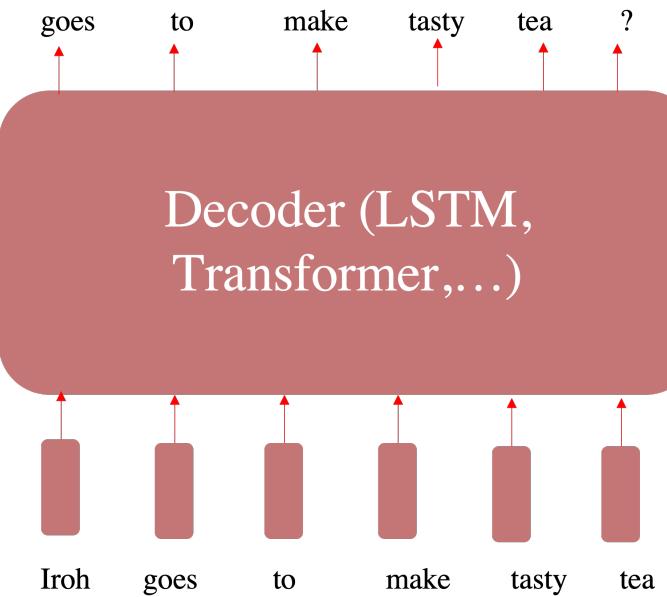
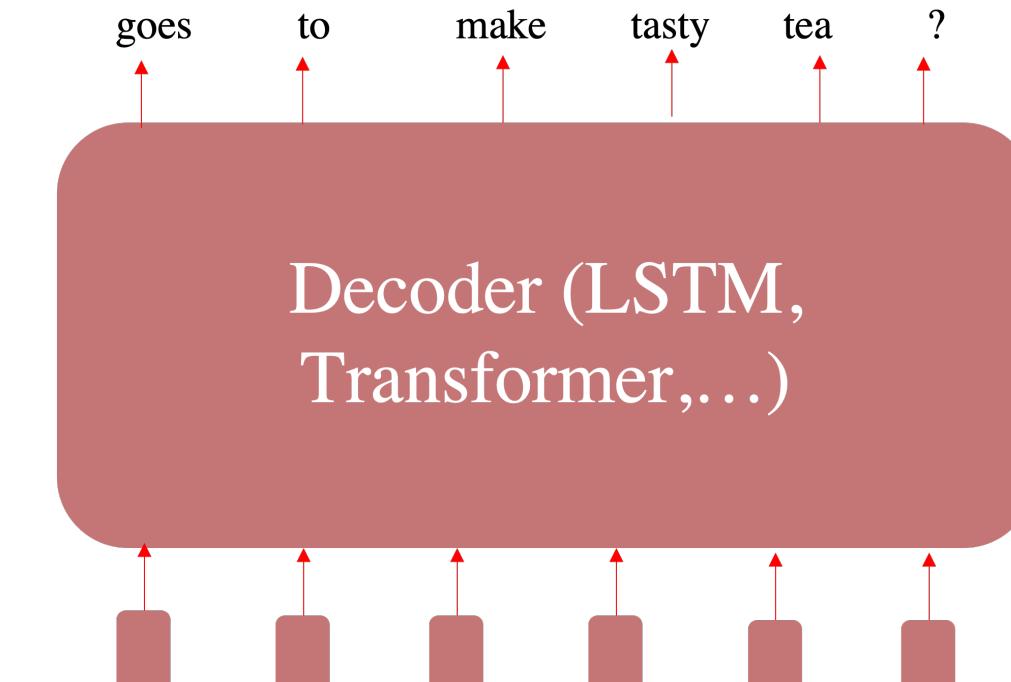
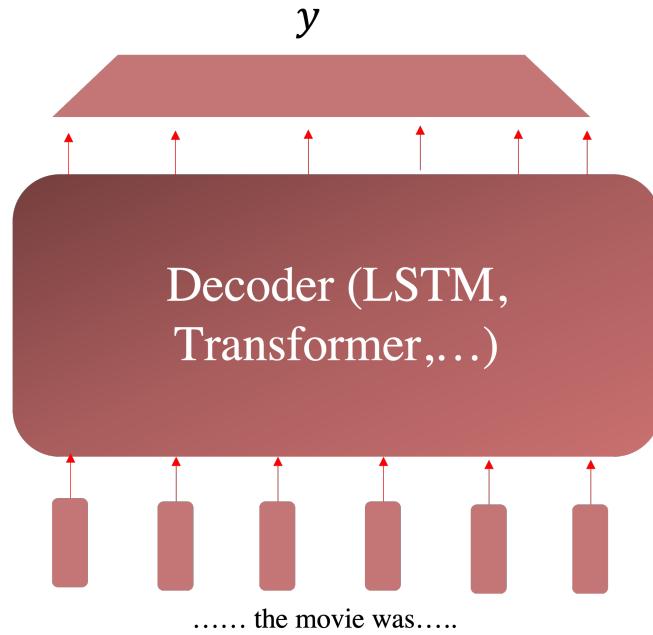


Figure 1: The Transformer - model architecture.

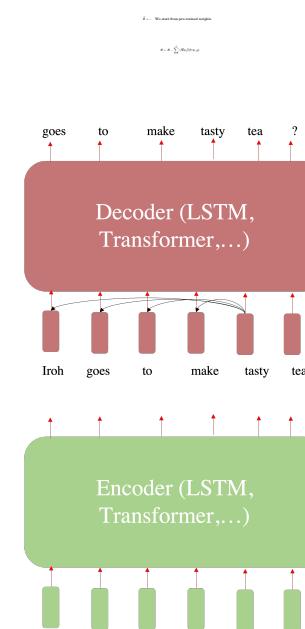
Iroh      goes      to      make      tasty      tea



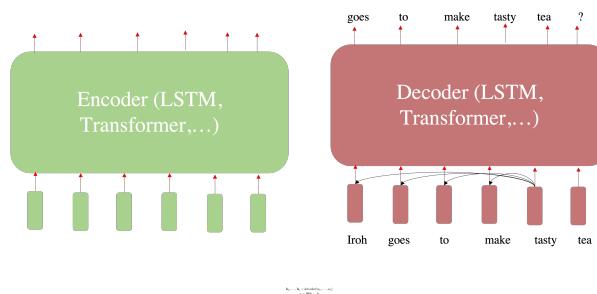
Pretuning Paragin  
Data & Pretraining: fine-tuning pre-trained models in the wild



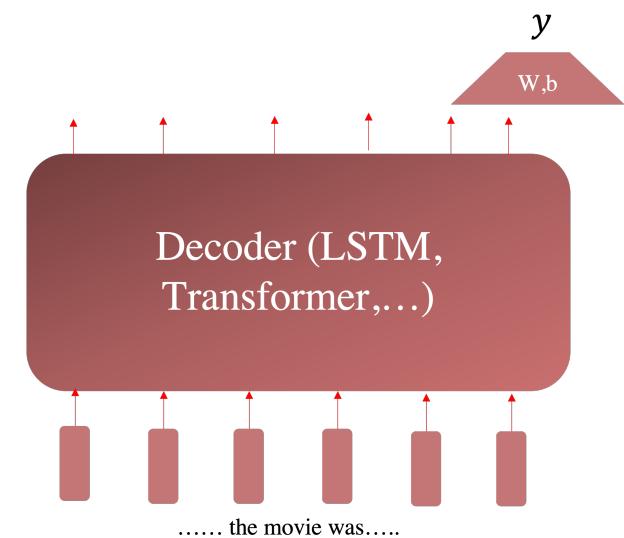
Pretuning/fine-tuning: optimization perspective  
Decoder architectures & training perspectives  
• The “standard” perspective:  $\arg\max_{y_t} p_{\theta}(y_t | \mathcal{X}_t)$   
• Aiming for “good” pre-training:  $\arg\min_{\theta} \mathbb{E}_{\mathcal{D}}[\mathcal{L}(\theta, \mathcal{D})]$   
• Aiming for “good” fine-tuning:  $\arg\min_{\theta} \mathbb{E}_{\mathcal{D}}[\mathcal{L}(\theta, \mathcal{D}) + \lambda \mathcal{L}_{\text{pre}}(\theta)]$   
Pretuning/fine-tuning: optimization perspective  
Three ways to pre-train a model  
Pretuning to three types of architectures  
1) Decoder Only



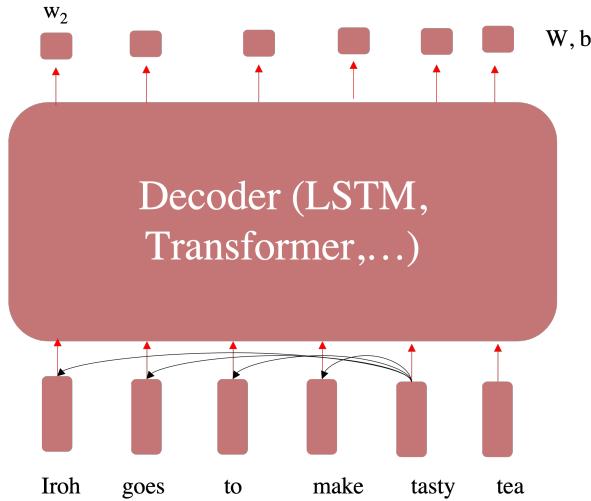
B) Encoder-Decoder  
Data & Pretraining: fine-tuning pre-trained models in the wild



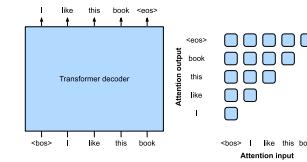
T) Decoder only - Option 1  
Data & Pretraining: fine-tuning pre-trained models in the wild  
Model & Optimization: backpropagation through time



T) Decoder only - Option 2  
Data & Pretraining: fine-tuning pre-trained models in the wild  
Model & Optimization: backpropagation through time



Generative Pretrained Transformer (GPT) [Radford et al., 2018]

<bos> I like this book  
Attention input

Attention output

Generative Pretrained Transformer (GPT) [Radford et al., 2018]

## Improving Language Understanding by Generative Pre-Training

Alec Radford      Karthik Narasimhan      Tim Salimans      Ilya Sutskever  
 OpenAI      OpenAI      OpenAI      OpenAI  
 alec@openai.com      karthikn@openai.com      tim@openai.com      ilyasu@openai.com

### Abstract

Natural language understanding comprises a wide range of diverse tasks such as textual entailment, question answering, semantic similarity assessment, and document classification. Although large unlabeled text corpora are abundant, labeled data for learning these specific tasks is scarce, making it challenging for discriminatively trained models to perform adequately. We demonstrate that large gains on these tasks can be realized by *generative pre-training* of a language model on a diverse corpus of unlabeled text, followed by *discriminative fine-tuning* on each specific task. In contrast to previous approaches, we make use of task-aware input transformations during fine-tuning to achieve effective transfer while requiring minimal changes to the model architecture. We demonstrate the effectiveness of our approach on a wide range of benchmarks for natural language understanding. Our general task-agnostic model outperforms discriminatively trained models that use architectures specifically crafted for each task, significantly improving upon the state of the art in 9 out of the 12 tasks studied. For instance, we achieve absolute improvements of 8.9% on commonsense reasoning (Stories Cloze Test), 5.7% on question answering (RACE), and 1.5% on textual entailment (MultiNLI).

### 1 Introduction

The ability to learn effectively from raw text is crucial to alleviating the dependence on supervised learning in natural language processing (NLP). Most deep learning methods require substantial

amounts of manually labeled data, which restricts their applicability in many domains that suffer from a dearth of annotated resources [61]. In these situations, models that can leverage linguistic information from unlabeled data provide a valuable alternative to gathering more annotation, which

Generative Pretrained Transformer (GPT) [Radford et al., 2018]

Table 2: Experimental results on natural language inference tasks, comparing our model with current state-of-the-art methods. 5x indicates an ensemble of 5 models. All datasets use accuracy as the evaluation metric.

Method	MNLI-m	MNLI-mm	SNLI	SciTail	QNLI	RTE
ESIM + ELMo [44] (5x)	-	-	<u>89.3</u>	-	-	-
CAFE [58] (5x)	80.2	79.0	<u>89.3</u>	-	-	-
Stochastic Answer Network [35] (3x)	<u>80.6</u>	<u>80.1</u>	-	-	-	-
CAFE [58]	78.7	77.9	<u>88.5</u>	<u>83.3</u>	-	-
GenSen [64]	71.4	71.3	-	-	<u>82.3</u>	59.2
Multi-task BiLSTM + Attn [64]	72.2	72.1	-	-	<u>82.1</u>	<b>61.7</b>
Finetuned Transformer LM (ours)	<b>82.1</b>	<b>81.4</b>	<b>89.9</b>	<b>88.3</b>	<b>88.1</b>	56.0

**Context (human-written):** In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**GPT-2:** The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

## Abstract

Recent work has demonstrated substantial gains on many NLP tasks and benchmarks by pre-training on a large corpus of text followed by fine-tuning on a specific task. While typically task-agnostic in architecture, this method still requires task-specific fine-tuning datasets of thousands or tens of thousands of examples. By contrast, humans can generally perform a new language task from only a few examples or from simple instructions – something which current NLP systems still largely struggle to do. Here we show that scaling up language models greatly improves task-agnostic, few-shot performance, sometimes even reaching competitiveness with prior state-of-the-art fine-tuning approaches. Specifically, we train GPT-3, an autoregressive language model with 175 billion parameters, 10x more than any previous non-sparse language model, and test its performance in the few-shot setting. For all tasks, GPT-3 is applied without any gradient updates or fine-tuning, with tasks and few-shot demonstrations specified purely via text interaction with the model. GPT-3 achieves strong performance on many NLP datasets, including translation, question-answering, and cloze tasks, as well as several tasks that require on-the-fly reasoning or domain adaptation, such as unscrambling words, using a novel word in a sentence, or performing 3-digit arithmetic. At the same time, we also identify some datasets where GPT-3's few-shot learning still struggles, as well as some datasets where GPT-3 faces methodological issues related to training on large web corpora. Finally, we find that GPT-3 can generate samples of news articles which human evaluators have difficulty distinguishing from articles written by humans. We discuss broader societal impacts of this finding and of GPT-3 in general.

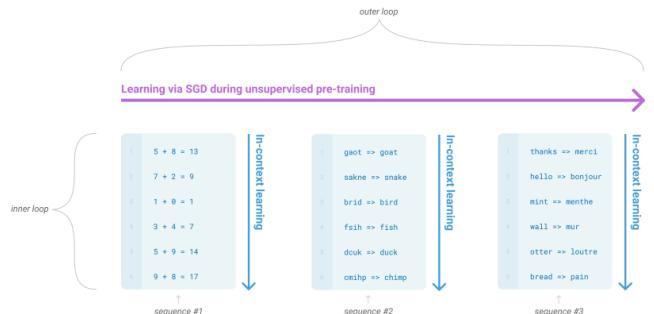
## Language Models are Few-Shot Learners

Tom B. Brown*	Benjamin Mann*	Nick Ryder*	Melanie Subbiah*	
Jared Kaplan†	Prafulla Dhariwal	Arvind Neelakantan	Pranav Shyam	Girish Sastry
Amanda Askell	Sandhini Agarwal	Ariel Herbert-Voss	Gretchen Krueger	Tom Henighan
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher Hesse	Mark Chen	Eric Sigler	Mateusz Litwin	Scott Gray
Benjamin Chess	Jack Clark	Christopher Berner		
Sam McCandlish	Alec Radford	Ilya Sutskever	Dario Amodei	

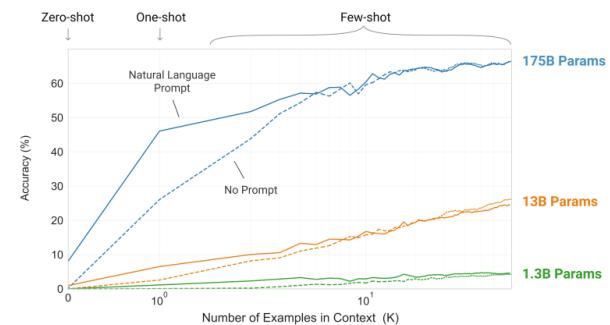
OpenAI

Dataset	Quantity (tokens)	Weight in training mix	Epochs elapsed when training for 300B tokens
Common Crawl (filtered)	410 billion	60%	0.44
WebText2	19 billion	22%	2.9
Books1	12 billion	8%	1.9
Books2	55 billion	8%	0.43
Wikipedia	3 billion	3%	3.4

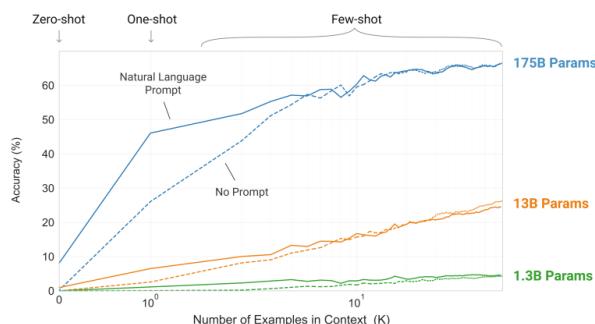
**Table 2.2: Datasets used to train GPT-3.** “Weight in training mix” refers to the fraction of examples during training that are drawn from a given dataset, which we intentionally do not make proportional to the size of the dataset. As a result, when we train for 300 billion tokens, some datasets are seen up to 3.4 times during training while other datasets are seen less than once.



**Figure 1.1: Language model meta-learning.** During unsupervised pre-training, a language model develops a broad set of skills and pattern recognition abilities. It then uses these abilities at inference time to rapidly adapt to or recognize the desired task. We use the term “in-context learning” to describe the inner loop of this process, which occurs within the forward-pass upon each sequence. The sequences in this diagram are not intended to be representative of the data a model would see during pre-training, but are intended to show that there are sometimes repeated sub-tasks embedded within a single sequence.



**Figure 1.2: Larger models make increasingly efficient use of in-context information.** We show in-context learning performance on a simple task requiring the model to remove random symbols from a word, both with and without a natural language task description (see Sec. 3.9.2). The steeper “in-context learning curves” for large models demonstrate improved ability to learn a task from contextual information. We see qualitatively similar behavior across a wide range of tasks.



**Figure 1.2: Larger models make increasingly efficient use of in-context information.** We show in-context learning performance on a simple task requiring the model to remove random symbols from a word, both with and without a natural language task description (see Sec. 3.9.2). The steeper “in-context learning curves” for large models demonstrate improved ability to learn a task from contextual information. We see qualitatively similar behavior across a wide range of tasks.

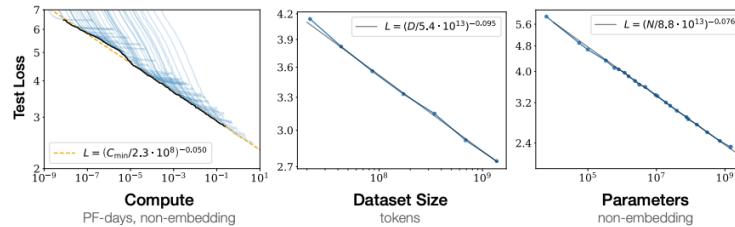
## Scaling Laws for Neural Language Models

Jared Kaplan \*  
Johns Hopkins University, OpenAI  
[jaredk@jhu.edu](mailto:jaredk@jhu.edu)

Sam McCandlish \*  
OpenAI  
[sam@openai.com](mailto:sam@openai.com)

Tom Henighan  
OpenAI  
[henighan@openai.com](mailto:henighan@openai.com) Tom B. Brown  
OpenAI  
[tom@openai.com](mailto:tom@openai.com) Benjamin Chess  
OpenAI  
[bchess@openai.com](mailto:bchess@openai.com) Rewon Child  
OpenAI  
[rewon@openai.com](mailto:rewon@openai.com)

Scott Gray  
OpenAI  
[scott@openai.com](mailto:scott@openai.com) Alec Radford  
OpenAI  
[alec@openai.com](mailto:alec@openai.com) Jeffrey Wu  
OpenAI  
[jeffwu@openai.com](mailto:jeffwu@openai.com) Dario Amodei  
OpenAI  
[damodei@openai.com](mailto:damodei@openai.com)



**Figure 1** Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute<sup>2</sup> used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

**Performance depends strongly on scale, weakly on model shape:** Model performance depends most strongly on scale, which consists of three factors: the number of model parameters  $N$  (excluding embeddings), the size of the dataset  $D$ , and the amount of compute  $C$  used for training. Within reasonable limits, performance depends very weakly on other architectural hyperparameters such as depth vs. width. (Section 3)

**Smooth power laws:** Performance has a power-law relationship with each of the three scale factors  $N, D, C$  when not bottlenecked by the other two, with trends spanning more than six orders of magnitude (see Figure 1). We observe no signs of deviation from these trends on the upper end, though performance must flatten out eventually before reaching zero loss. (Section 3)

**Universality of overfitting:** Performance improves predictably as long as we scale up  $N$  and  $D$  in tandem, but enters a regime of diminishing returns if either  $N$  or  $D$  is held fixed while the other increases. The performance penalty depends predictably on the ratio  $N^{0.74}/D$ , meaning that every time we increase the model size 8x, we only need to increase the data by roughly 5x to avoid a penalty. (Section 4)

**Universality of training:** Training curves follow predictable power-laws whose parameters are roughly independent of the model size. By extrapolating the early part of a training curve, we can roughly predict the loss that would be achieved if we trained for much longer. (Section 5)

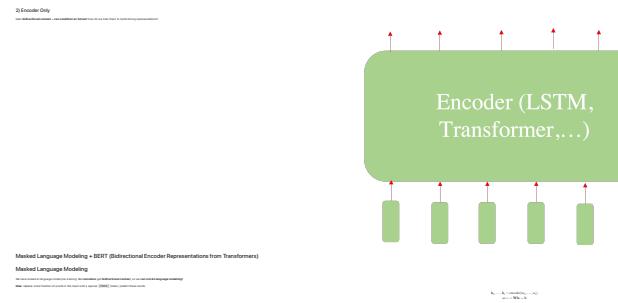
**Transfer improves with test performance:** When we evaluate models on text with a different distribution than they were trained on, the results are strongly correlated to those on the training validation set with a roughly constant offset in the loss – in other words, transfer to a different distribution incurs a constant penalty but otherwise improves roughly in line with performance on the training set. (Section 3.2.2)

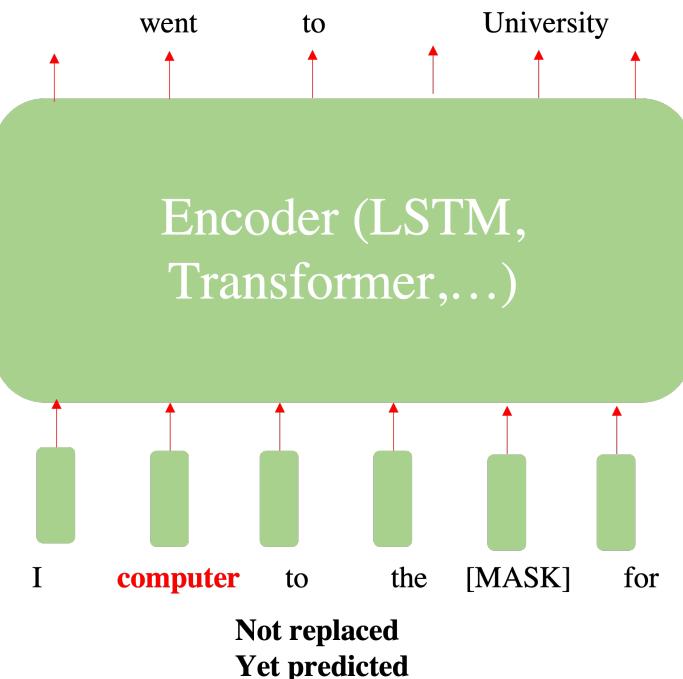
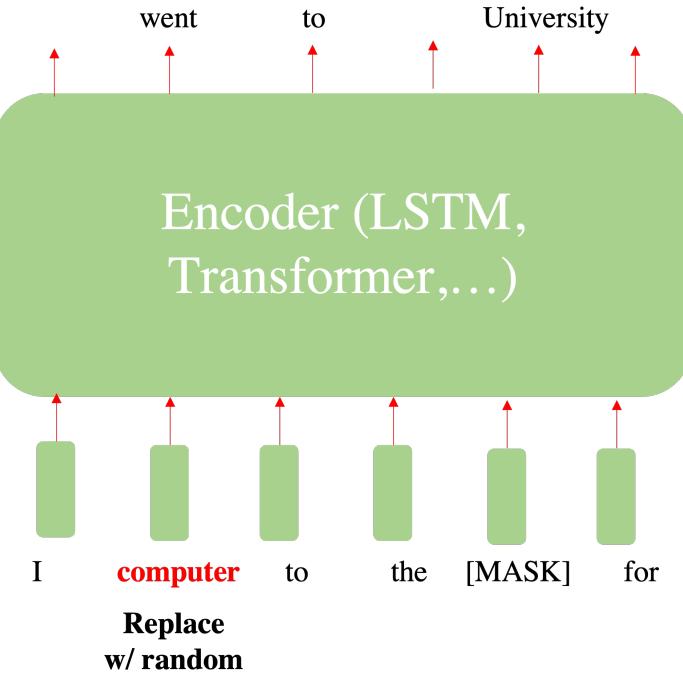
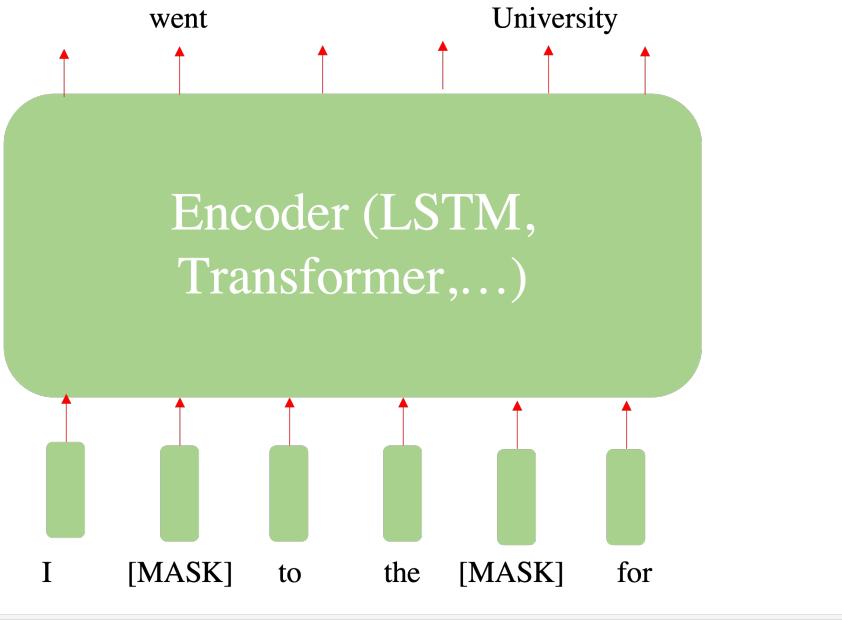
**Sample efficiency:** Large models are more sample-efficient than small models, reaching the same level of performance with fewer optimization steps (Figure 2) and using fewer data points (Figure 4).

**Convergence is inefficient:** When working within a fixed compute budget  $C$  but without any other restrictions on the model size  $N$  or available data  $D$ , we attain optimal performance by training *very large models* and stopping *significantly short of convergence* (see Figure 3). Maximally compute-efficient training would therefore be far more sample efficient than one might expect based on training small models to convergence, with data requirements growing very slowly as  $D \sim C^{0.27}$  with training compute. (Section 6)

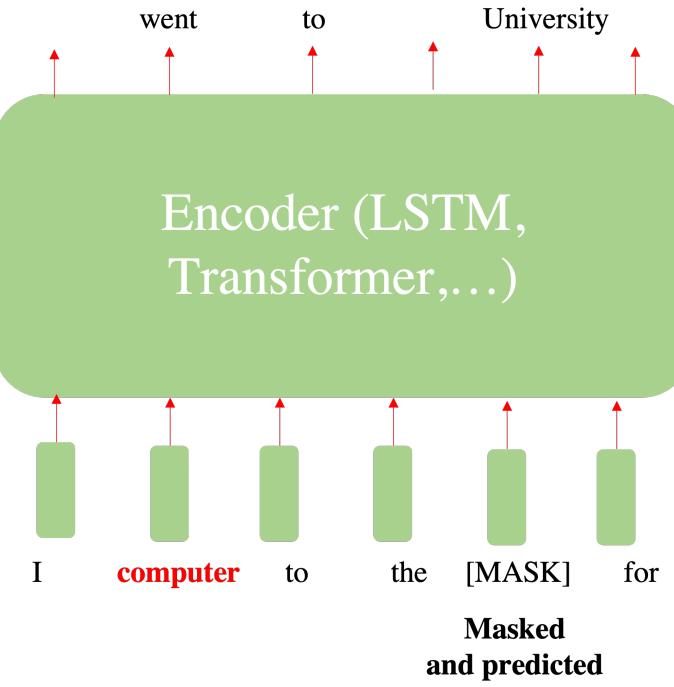
**Optimal batch size:** The ideal batch size for training these models is roughly a power of the loss only, and continues to be determinable by measuring the gradient noise scale [MKAT18]; it is roughly 1-2 million tokens at convergence for the largest models we can train. (Section 5.1)

Taken together, these results show that language modeling performance improves smoothly and predictably as we appropriately scale up model size, data, and compute. We expect that larger language models will perform better and be more sample efficient than current models.





BERT: Masked LM  
BERT is a pre-training task that predicts the original word in a masked sequence. A mask is randomly selected from the following categories:  
 1. Mask words: words that are not part of the input sequence.  
 2. Masked words: words that are part of the input sequence.  
 3. Masked words: words that are part of the input sequence.



BERT: Next Sentence Prediction  
  
Input: <cls> this movie is great <sep> i like it <sep>  
  
Token Embeddings:  $e_{\text{cls}}$ ,  $e_{\text{this}}$ ,  $e_{\text{movie}}$ ,  $e_{\text{is}}$ ,  $e_{\text{great}}$ ,  $e_{<\text{sep}>}$ ,  $e_i$ ,  $e_{\text{like}}$ ,  $e_{\text{it}}$ ,  $e_{<\text{sep}>}$   
  
Segment Embeddings:  $e_A$ ,  $e_A$ ,  $e_A$ ,  $e_A$ ,  $e_A$ ,  $e_A$ ,  $e_B$ ,  $e_B$ ,  $e_B$ ,  $e_B$   
  
Positional Embeddings:  $e_0$ ,  $e_1$ ,  $e_2$ ,  $e_3$ ,  $e_4$ ,  $e_5$ ,  $e_6$ ,  $e_7$ ,  $e_8$ ,  $e_9$

BERT Details  
BERT uses a multi-layered architecture, with each layer learning more complex representations of the input sequence. The first layer processes the input sequence, and subsequent layers build upon the previous ones. This allows BERT to capture long-range dependencies and context information across the entire sequence.

BERT: One model many tasks

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	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	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>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

Table 1: GLUE Test results, scored by the evaluation server (<https://gluebenchmark.com/leaderboard>). The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set.<sup>8</sup> BERT and OpenAI GPT are single-model, single-task. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks. We exclude entries that use BERT as one of their components.

Rank Name	Model	URL	Score	CoLA	SST-2	MRPC	STS-B	QQP MNLI-m	MNLI-mm	CoLA	RTE	WNLI	
1	Microsoft Alexander-v-team	Turing-NLP v6	91.3	73.3	97.5	94.9/92.6	93.9/91.7	76.4/90.9	92.5	91.9	96.7	93.6	97.9
2	JDDExploit-d-team	Vega v1	91.3	73.8	97.9	94.9/92.6	93.5/91.7	76.7/91.1	92.1	91.9	96.7	92.4	97.9
3	Microsoft Alexander-v-Team	Turing-NLP v5	91.2	72.6	97.6	93.9/91.7	93.7/93.3	76.4/91.1	92.6	92.4	97.9	94.1	96.9
4	DILY Team	DeBERTa + CLEVER	91.1	74.7	97.8	93.9/91.7	93.4/93.1	76.5/91.0	92.1	91.8	96.7	93.2	96.6
5	ERNIE Team - Baidu	ERNIE	91.1	75.6	97.8	93.9/91.0	93.6/92.6	75.7/90.9	92.3	91.7	97.3	92.6	96.9
6	AlceMild + DILY	StreudBERT + CLEVER	91.0	75.5	97.7	93.9/91.0	93.5/91.7	75.6/90.8	91.7	91.5	97.4	92.5	95.2
7	DeBERTa Team - Microsoft	DeBERTa / Turing-NLP	90.8	71.6	97.5	94.9/92.0	92.9/92.6	76.2/90.8	91.9	91.6	93.2	93.2	94.5
8	HFL-IVY-TICK	MacALBERT + DILY	90.7	74.8	97.0	94.9/92.0	92.6/92.6	74.7/90.5	91.3	91.1	97.8	92.0	94.5
9	PINN-JN-Omni-Slot5	ALBERT + DAF + NAS	90.6	73.9	97.2	94.9/92.0	93.6/92.4	76.1/91.0	91.6	91.3	97.5	91.7	94.5
10	T5 Team - Google	T5	90.3	71.6	97.5	92.8/90.4	93.9/92.8	75.1/90.6	92.2	91.9	96.9	92.8	94.5
11	Microsoft DIBS AI & MSR AI & DATECH	MT-DNN+SMART	89.9	69.5	97.5	93.7/91.6	92.9/92.5	73.9/90.3	91.0	90.8	99.2	99.7	94.5
12	Huawei Noah's Ark Lab	MEZHA-Large	89.8	71.7	97.3	93.2/91.0	92.4/91.0	75.2/90.7	91.5	91.3	96.2	93.3	94.5
13	LG AI Research	ANNA	89.6	68.7	97.0	92.7/90.1	93.6/92.6	75.3/90.6	91.8	91.6	96.0	91.8	95.9

GLUE Benchmarks  
  
BERT: Next Sentence Prediction  
  
Input: <cls> this movie is great <sep> i like it <sep>  
  
Token Embeddings:  $e_{\text{cls}}$ ,  $e_{\text{this}}$ ,  $e_{\text{movie}}$ ,  $e_{\text{is}}$ ,  $e_{\text{great}}$ ,  $e_{<\text{sep}>}$ ,  $e_i$ ,  $e_{\text{like}}$ ,  $e_{\text{it}}$ ,  $e_{<\text{sep}>}$   
  
Segment Embeddings:  $e_A$ ,  $e_A$ ,  $e_A$ ,  $e_A$ ,  $e_A$ ,  $e_A$ ,  $e_B$ ,  $e_B$ ,  $e_B$ ,  $e_B$   
  
Positional Embeddings:  $e_0$ ,  $e_1$ ,  $e_2$ ,  $e_3$ ,  $e_4$ ,  $e_5$ ,  $e_6$ ,  $e_7$ ,  $e_8$ ,  $e_9$

BERT Details  
BERT uses a multi-layered architecture, with each layer processing the input sequence. The first layer captures local context, while subsequent layers capture increasingly global context and dependencies. This allows BERT to handle a wide variety of NLP tasks with a single, pre-trained model.

BERT: One model many tasks

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
XLNet <sub>LARGE</sub>						
with BOOKS + WIKI	13GB	256	1M	94.0/87.8	88.4	94.4
+ additional data	126GB	2K	500K	94.5/88.8	89.8	95.6

Table 4: Development set results for RoBERTa as we pretrain over more data (16GB → 160GB of text) and pretrain for longer (100K → 300K → 500K steps). Each row accumulates improvements from the rows above. RoBERTa matches the architecture and training objective of BERT<sub>LARGE</sub>. Results for BERT<sub>LARGE</sub> and XLNet<sub>LARGE</sub> are from Devlin et al. (2019) and Yang et al. (2019), respectively. Complete results on all GLUE tasks can be found in the Appendix.

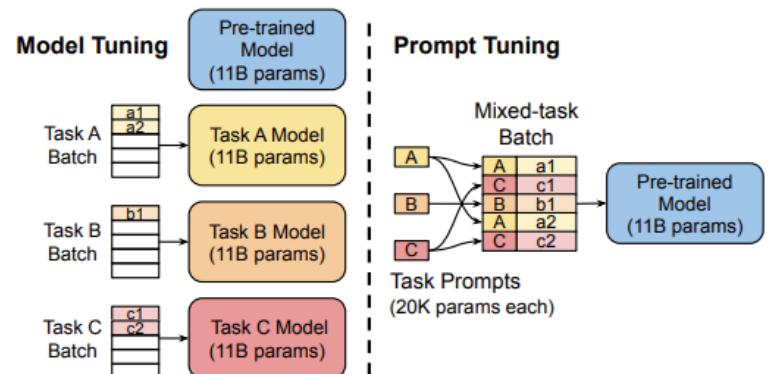
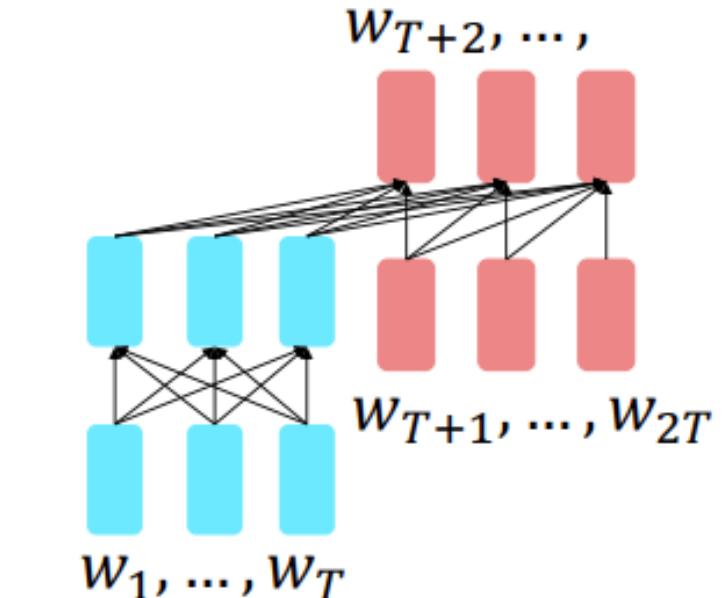
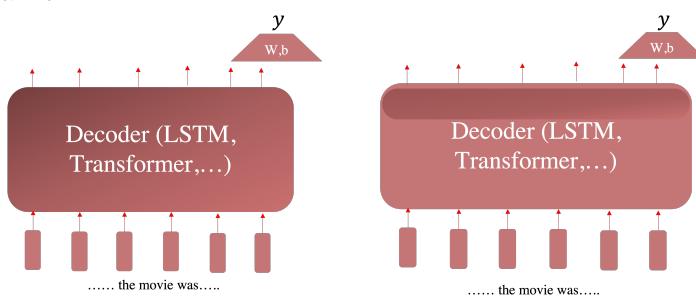
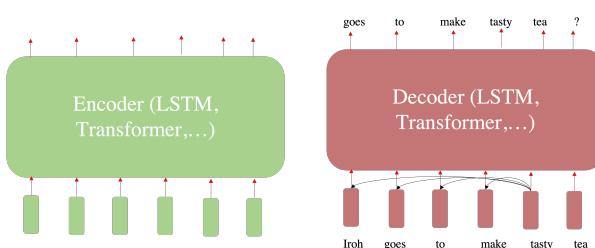


Figure 2: **Model tuning** requires making a task-specific copy of the entire pre-trained model for each downstream task and inference must be performed in separate batches. **Prompt tuning** only requires storing a small task-specific prompt for each task, and enables mixed-task inference using the original pre-trained model. With a T5 “XXT” model, each copy

trained model. With a 12 ALB model, each copy of the tuned model requires 11 billion parameters. By contrast, our tuned prompts would only require 20,480 parameters per task—a reduction of *over five orders of magnitude*—assuming a prompt length of 5 tokens.

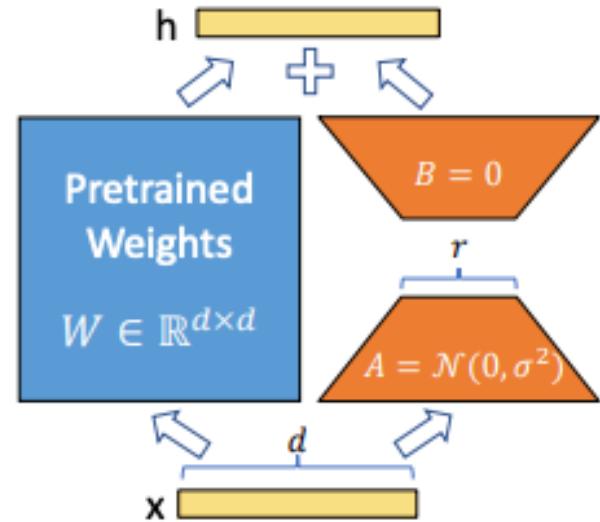


Figure 1: Our reparametrization. We only train  $A$  and  $R$ .