

CENG 403  
*Introduction to Deep Learning*

*Week 15b*

Sinan Kalkan

# Attention: Transformer

- Vanilla self attention:

$$e_i' = \sum_j \frac{\exp(e_j^T e_i)}{\sum_m \exp(e_m^T e_i)} e_j$$

- Scaled-dot product attention:

$$e_i' = \sum_j \frac{\exp(\mathbf{k}(e_j^T) \mathbf{q}(e_i))}{\sum_m \exp(\mathbf{k}(e_m^T) \mathbf{q}(e_i))} \mathbf{v}(e_j)$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

## Attention Is All You Need

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar\*  
Google Research  
nikip@google.com

Jakob Uszkoreit\*  
Google Research  
usz@google.com

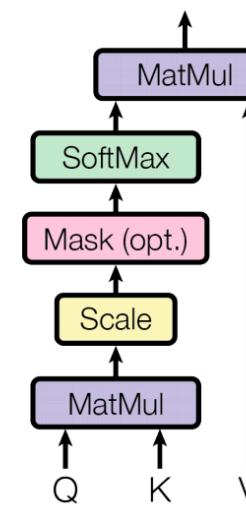
Llion Jones\*  
Google Research  
llion@google.com

Aidan N. Gomez\* †  
University of Toronto  
aidan@cs.toronto.edu

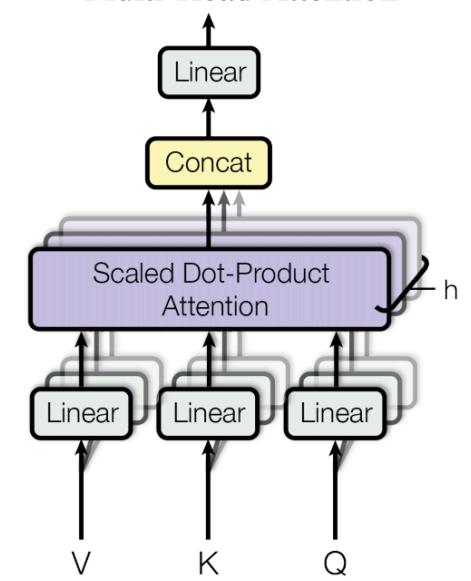
Lukasz Kaiser\*  
Google Brain  
lukaszkaiser@google.com

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

Scaled Dot-Product Attention



Multi-Head Attention



Previously on CENG403

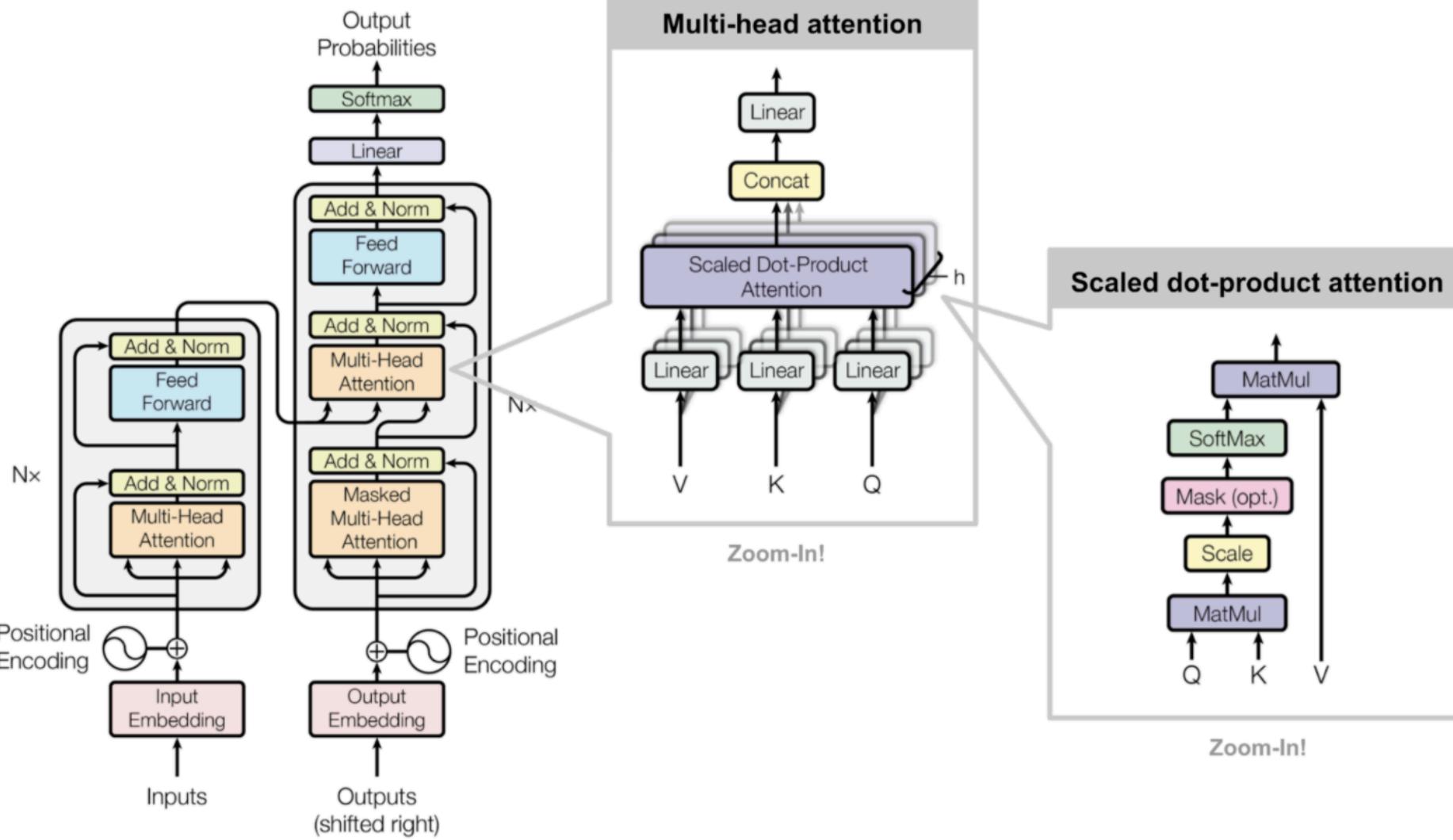


Fig. 17. The full model architecture of the transformer. (Image source: Fig 1 & 2 in [Vaswani, et al., 2017](#).)

<https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html>

# Positional Encoding

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

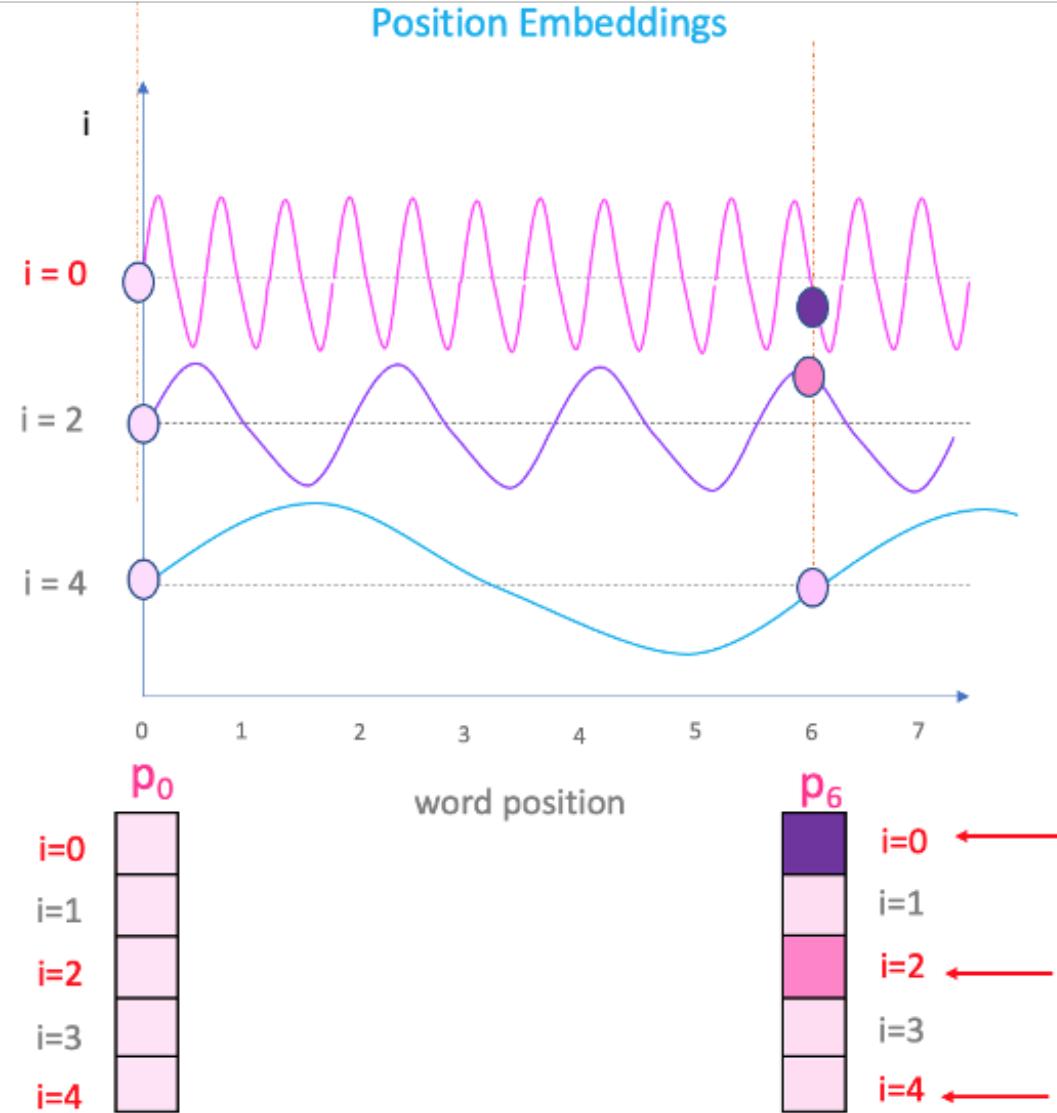


Fig from: <https://www.youtube.com/watch?v=dichiUZfOw>

# A Significant Issue with Self-Attention: Complexity

$$e'_i = \sum_j \frac{\exp\left(\mathbf{k}\left(e_j^T\right)\mathbf{q}(e_i)\right)}{\sum_m \exp\left(\mathbf{k}\left(e_m^T\right)\mathbf{q}(e_i)\right)} \mathbf{v}(e_j)$$

- If there are  $n$  tokens/embeddings,
  - Updating a single tokens require  $O(n)$  operations.
  - Overall:  $O(n^2)$
- What is the complexity of an RNN layer with  $n$  time steps?

# Today

- Pretraining language models
- The emergence of LLMs
- Vision Transformers
- Vision-Language / Multimodal Models
- LLMs as agents

# Pre-training in NLP

- Word embeddings are the basis of deep learning for NLP



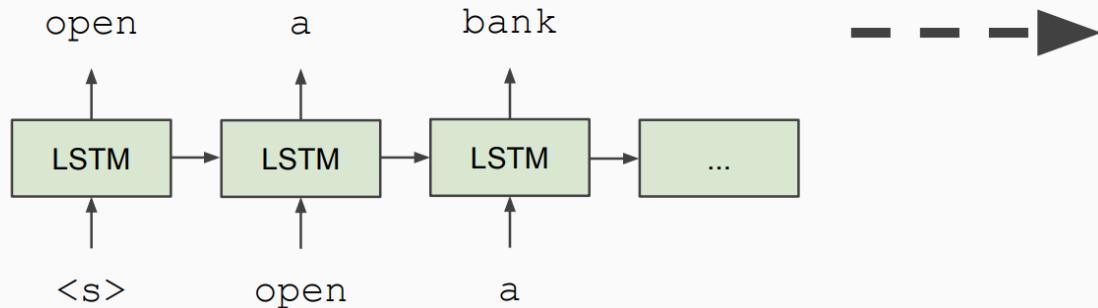
- Word embeddings (word2vec, GloVe) are often *pre-trained* on text corpus from co-occurrence statistics



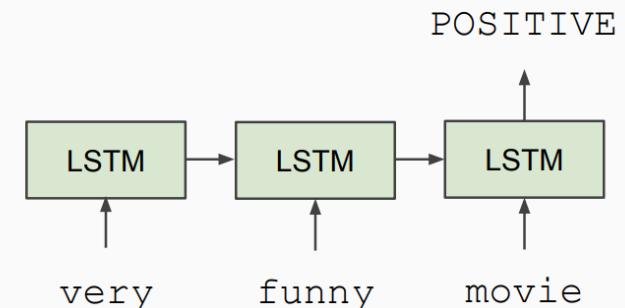
# Pre-training in NLP

- *Semi-Supervised Sequence Learning*, Google, 2015

## Train LSTM Language Model



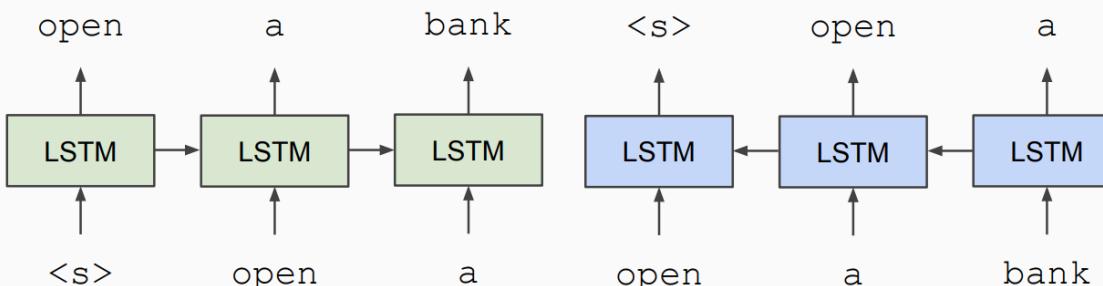
## Fine-tune on Classification Task



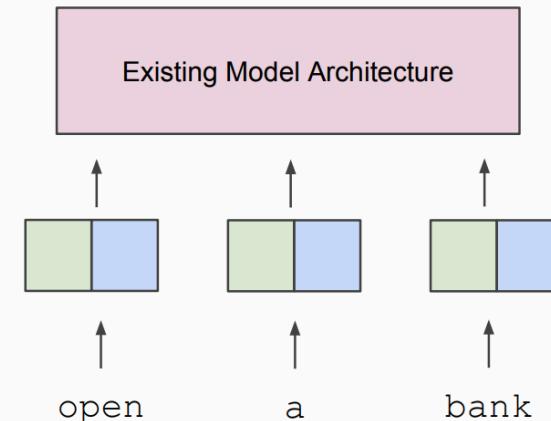
# Pre-training in NLP

- *ELMo: Deep Contextual Word Embeddings*, AI2 & University of Washington, 2017

**Train Separate Left-to-Right and Right-to-Left LMs**



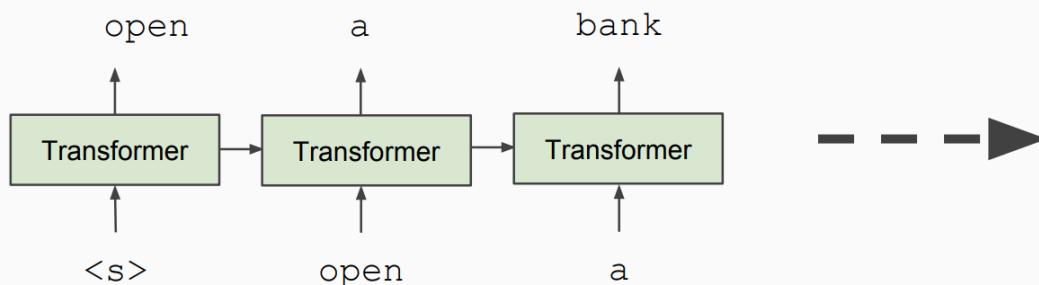
**Apply as “Pre-trained Embeddings”**



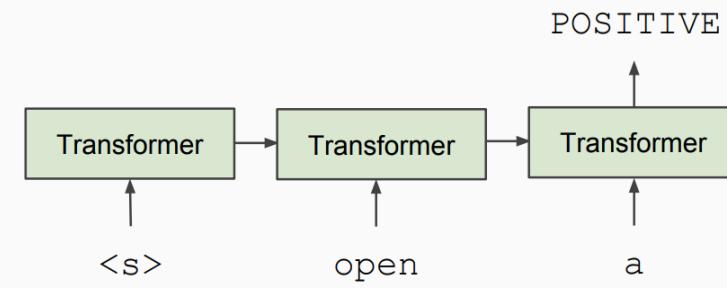
# Pre-training in NLP

- *Improving Language Understanding by Generative Pre-Training*, OpenAI, 2018

## Train Deep (12-layer) Transformer LM



## Fine-tune on Classification Task



# BERT

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

Google AI Language

{jacobdevlin, mingweichang, kentonl, kristout}@google.com

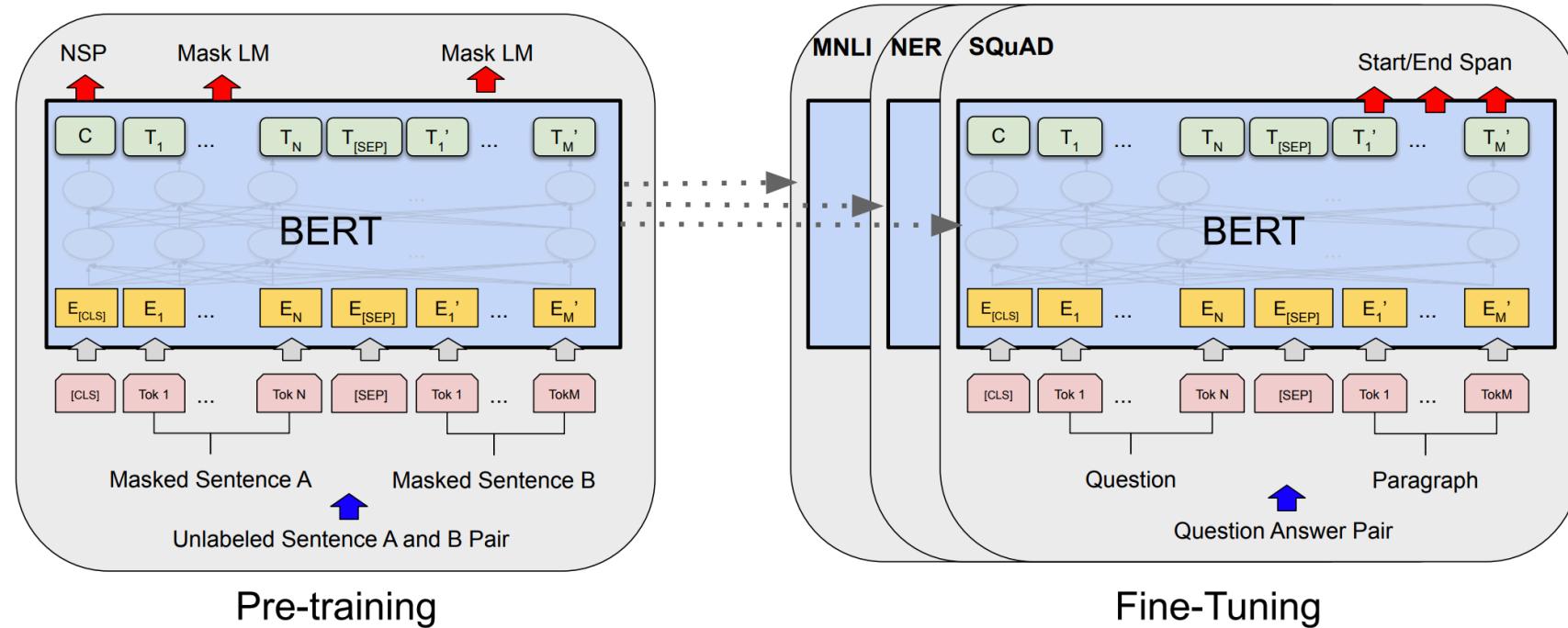


Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different down-stream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers).

# GPT-1

# GPT-1

- 12 layer decoder-only transformer
- Unsupervised pretraining
  - BookCorpus dataset
- Supervised finetuning
  - Textual alignment
  - QA & commonsense reasoning
  - Semantic similarity
  - Classification

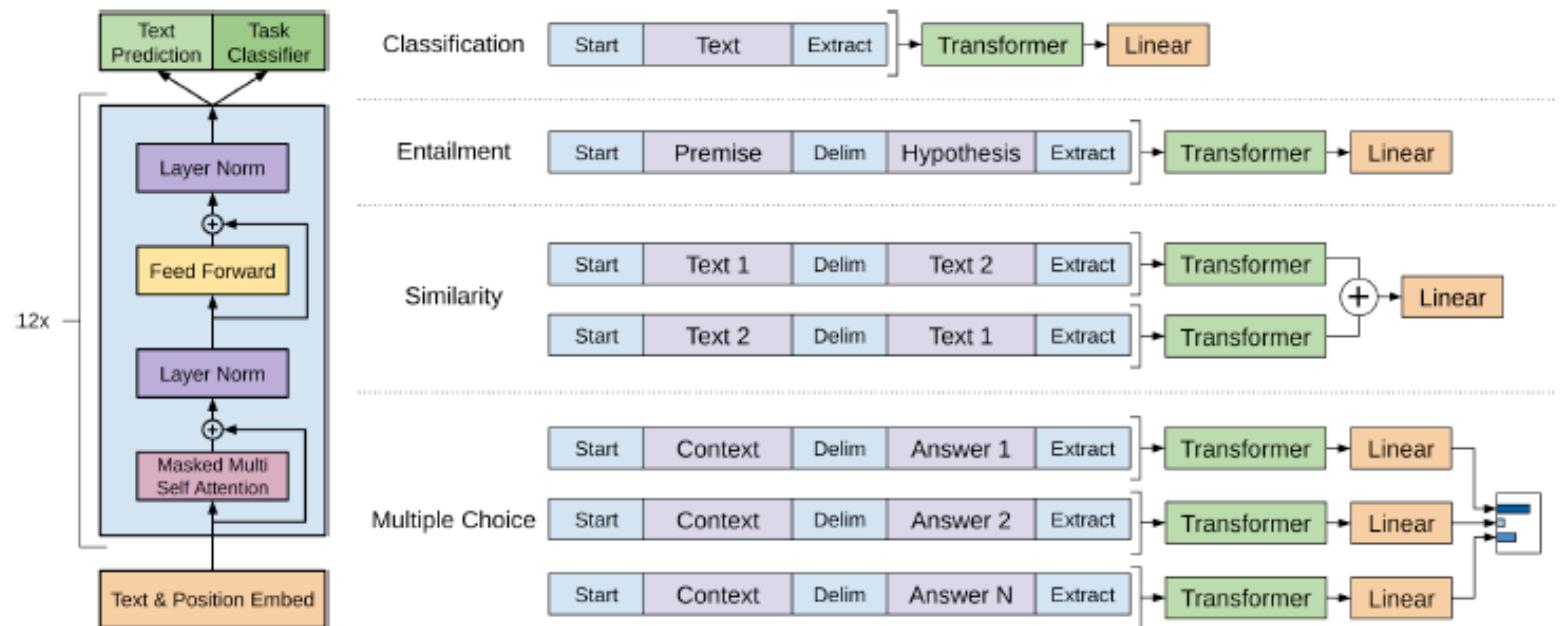


Figure 1: (left) Transformer architecture and training objectives used in this work. (right) Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

# GPT-1

Given an unsupervised corpus of tokens  $\mathcal{U} = \{u_1, \dots, u_n\}$ , we use a standard language modeling objective to maximize the following likelihood:

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta) \quad (1)$$

where  $k$  is the size of the context window, and the conditional probability  $P$  is modeled using a neural network with parameters  $\Theta$ . These parameters are trained using stochastic gradient descent [51].

$$\begin{aligned} h_0 &= UW_e + W_p \\ h_l &= \text{transformer\_block}(h_{l-1}) \forall i \in [1, n] \\ P(u) &= \text{softmax}(h_n W_e^T) \end{aligned} \quad (2)$$

where  $U = (u_{-k}, \dots, u_{-1})$  is the context vector of tokens,  $n$  is the number of layers,  $W_e$  is the token embedding matrix, and  $W_p$  is the position embedding matrix.

# GPT-1

## Discriminative Fine-tuning

For labeled downstream task, maximize the log probability on each pair of instance  $(x, y)$

After training the model with the objective in Eq. 1, we adapt the parameters to the supervised target task. We assume a labeled dataset  $\mathcal{C}$ , where each instance consists of a sequence of input tokens,  $x^1, \dots, x^m$ , along with a label  $y$ . The inputs are passed through our pre-trained model to obtain the final transformer block's activation  $h_l^m$ , which is then fed into an added linear output layer with parameters  $W_y$  to predict  $y$ :

$$P(y|x^1, \dots, x^m) = \text{softmax}(h_l^m W_y). \quad (3)$$

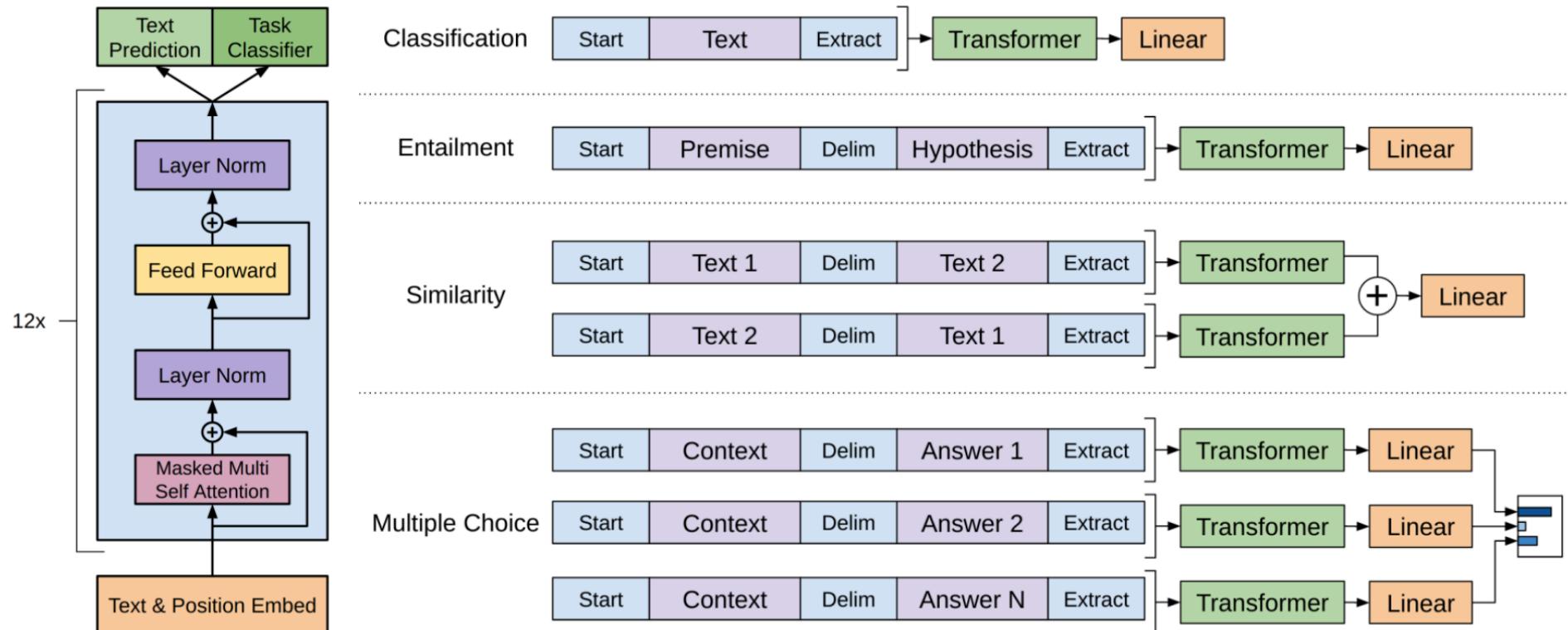
This gives us the following objective to maximize:

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y|x^1, \dots, x^m). \quad (4)$$

Add auxiliary fine-tuning objective of language modeling will imporove the performance  $L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$

# GPT-1

## Discriminative Fine-tuning



# GPT-1 Results

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	88.5	<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	-	-	82.1	<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

# BERT

# BERT

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

Google AI Language

{jacobdevlin, mingweichang, kentonl, kristout}@google.com

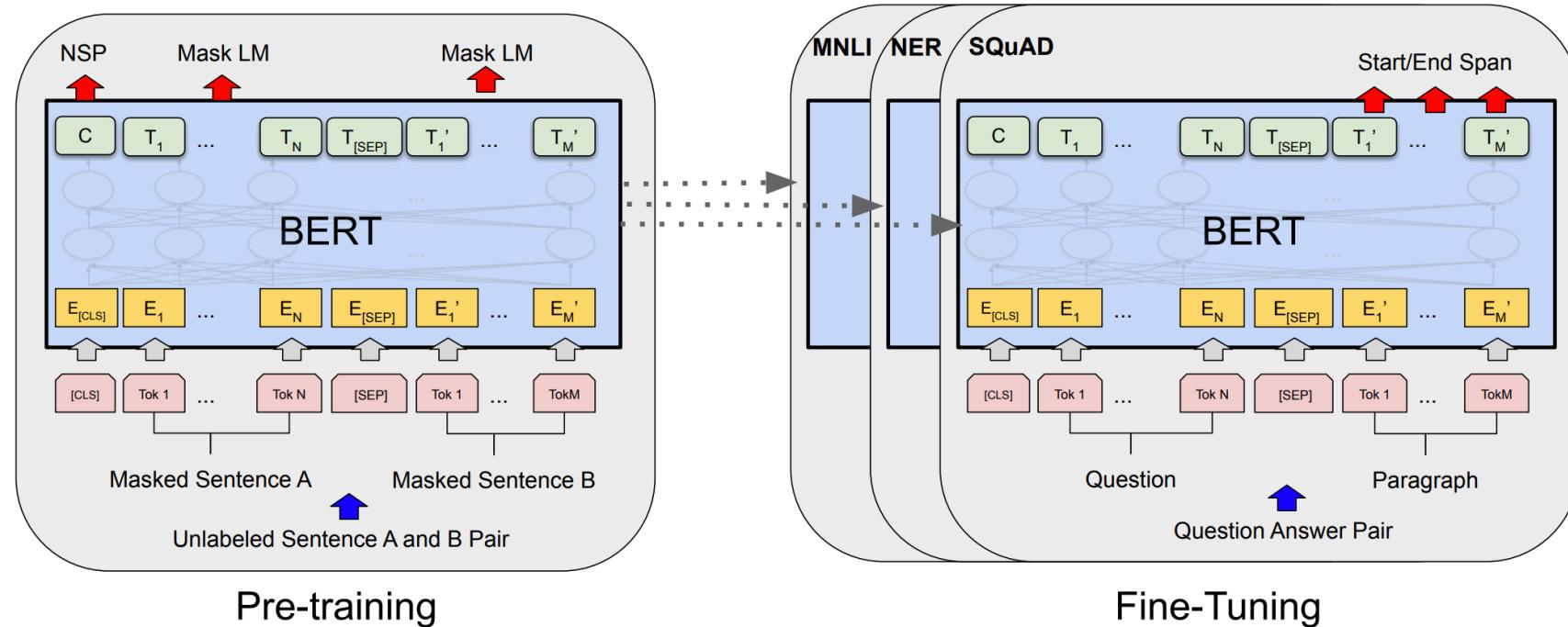


Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different down-stream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers).

# BERT: NLP Tasks

- MNLI: Multi-Genre Natural Language Inference

## Examples

### Premise

#### *Fiction*

The Old One always comforted Ca'daan, except today.

### Label

### Hypothesis

*neutral*

Ca'daan knew the Old One very well.

#### *Letters*

Your gift is appreciated by each and every student who will benefit from your generosity.

*neutral*

Hundreds of students will benefit from your generosity.

#### *Telephone Speech*

yes now you know if if everybody like in August when everybody's on vacation or something we can dress a little more casual or *contradiction* August is a black out month for vacations in the company.

#### *9/11 Report*

At the other end of Pennsylvania Avenue, people began to line up for a White House tour.

*entailment*

People formed a line at the end of Pennsylvania Avenue.

- NER: Named Entity Recognition
- SQuAD: Stanford Question Answering Dataset

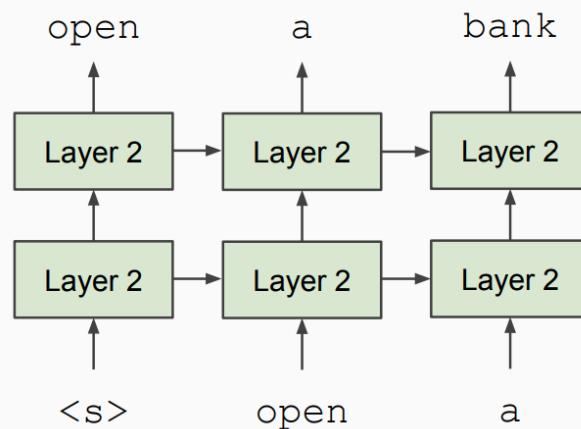
# BERT: Motivation

- **Problem:** Language models only use left context 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

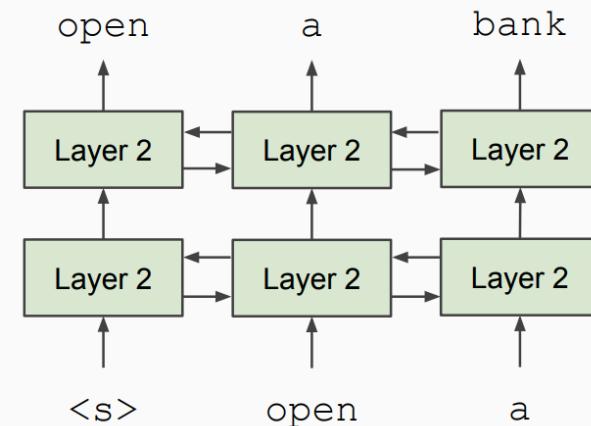
## Unidirectional context

Build representation incrementally



## Bidirectional context

Words can “see themselves”



Actually, this is vanilla attention that  
is not causal!

# BERT: Motivation

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

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

↑                              ↑  
store                          gallon

- Too little masking: Too expensive to train
- Too much masking: Not enough context

# BERT: Motivation

- Problem: Mask token never seen at fine-tuning
- Solution: 15% of the words to predict, but don't replace with [MASK] 100% of the time. Instead:
  - 80% of the time, replace with [MASK]  
went to the store → went to the [MASK]
  - 10% of the time, replace random word  
went to the store → went to the running
  - 10% of the time, keep same  
went to the store → went to the store

# BERT: Motivation

- To learn *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

# BERT

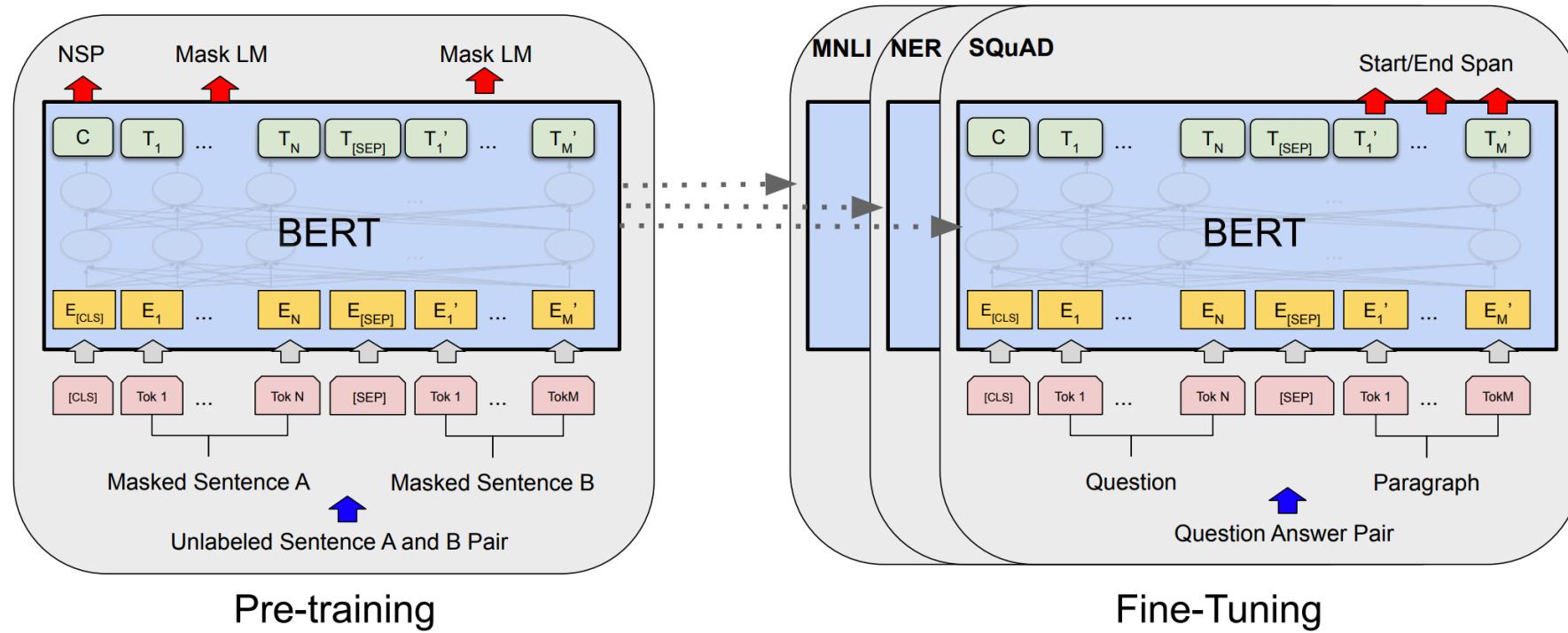
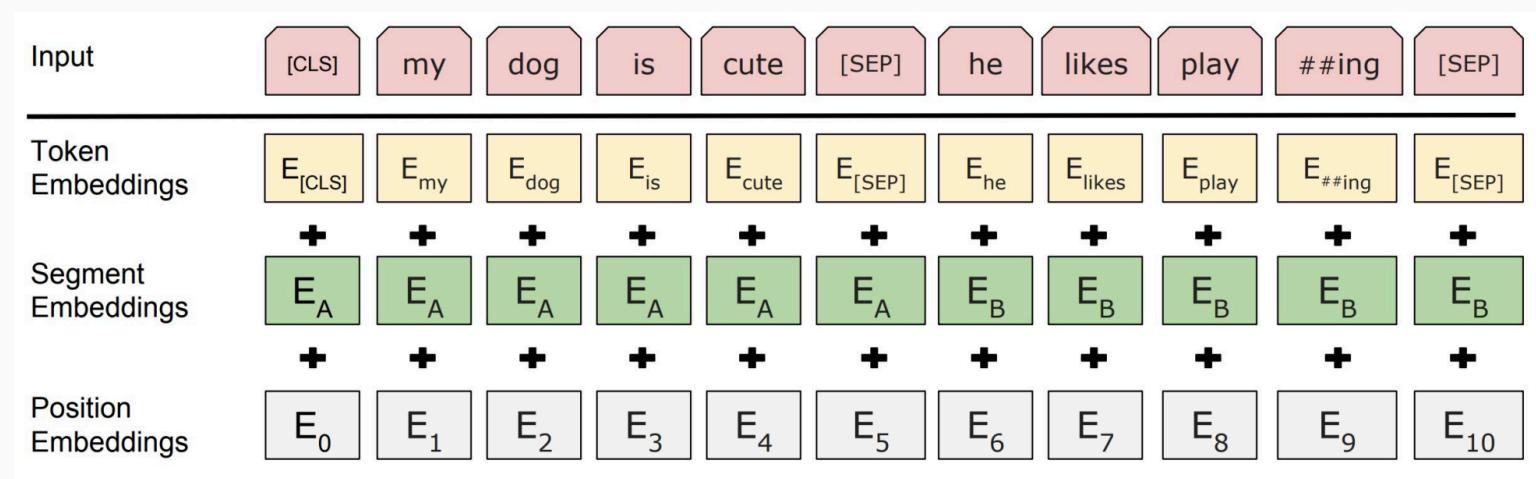


Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different down-stream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers).

# BERT

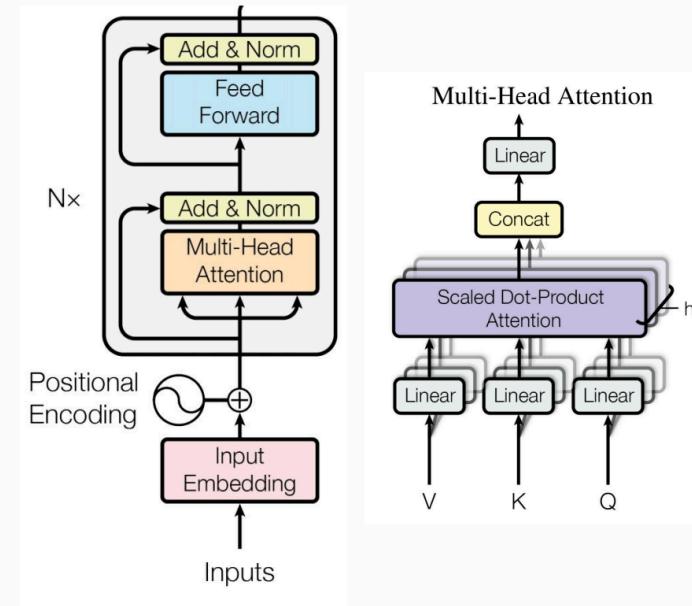


- Use 30,000 WordPiece vocabulary on input.
- Each token is sum of three embeddings
- Single sequence is much more efficient.

## Model Architecture

### Transformer encoder

- Multi-headed self attention
  - Models context
- Feed-forward layers
  - Computes non-linear hierarchical features
- Layer norm and residuals
  - Makes training deep networks healthy
- Positional embeddings
  - Allows model to learn relative positioning



# BERT

The General Language Understanding Evaluation (GLUE) benchmark is a collection of resources for training, evaluating, and analyzing natural language understanding systems. GLUE consists of:

- A benchmark of nine sentence- or sentence-pair language understanding tasks built on established existing datasets and selected to cover a diverse range of dataset sizes, text genres, and degrees of difficulty,
- A diagnostic dataset designed to evaluate and analyze model performance with respect to a wide range of linguistic phenomena found in natural language, and
- A public leaderboard for tracking performance on the benchmark and a dashboard for visualizing the performance of models on the diagnostic set.

<https://gluebenchmark.com/>

## GLUE Results

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average -
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>91.1</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>81.9</b>

### MultiNLI

Premise: Hills and mountains are especially sanctified in Jainism.

Hypothesis: Jainism hates nature.

Label: Contradiction

### CoLa

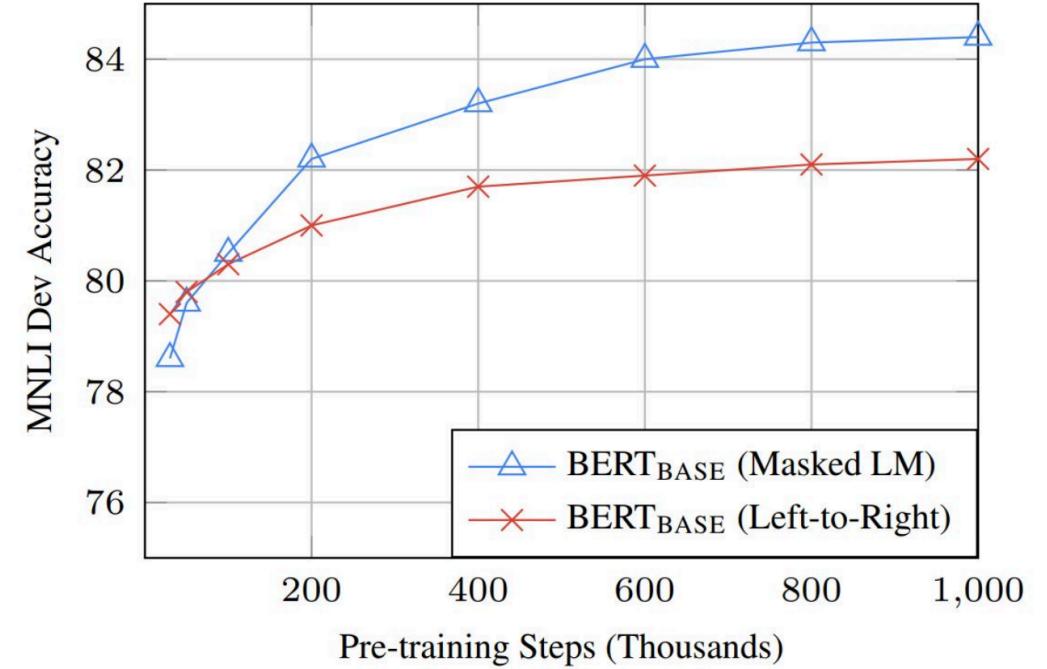
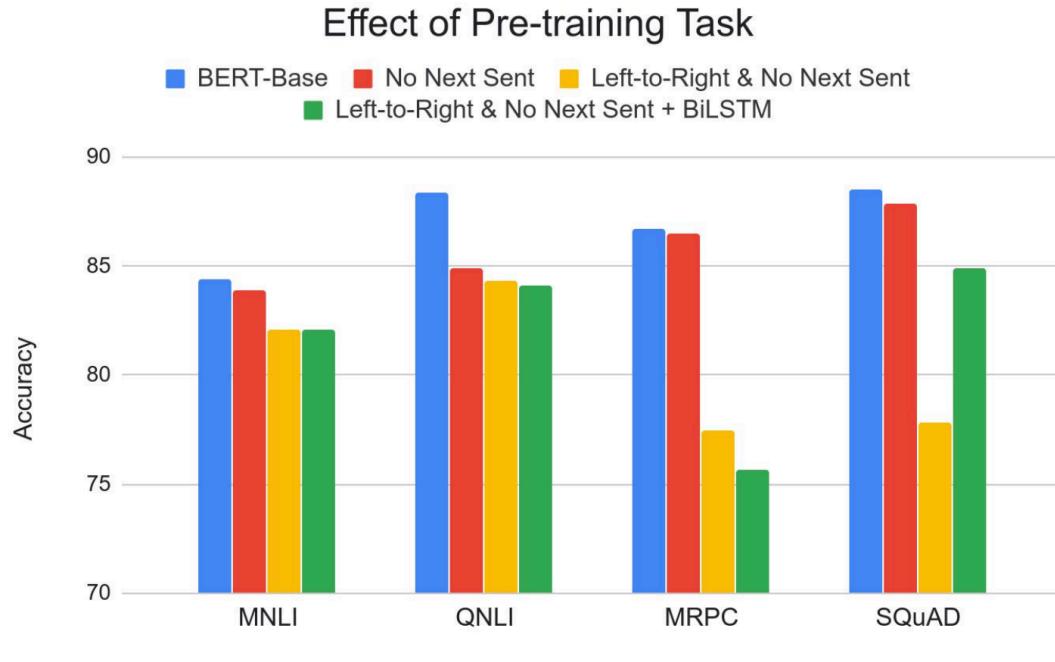
Sentence: The wagon rumbled down the road.

Label: Acceptable

Sentence: The car honked down the road.

Label: Unacceptable

# BERT



# Other GPT Models

Model	Title	Focus	Paradigm	Params
GPT-1	Improving <u>Language Understanding</u> by <b>Generative Pre-Training</b>	NLU tasks, pre-trained model	Pre-training->Efficient Fine-tuning	117M
GPT-2	Language Models are <b><u>Unsupervised Multitask</u></b> Learners	Zero-shot Evaluation, NLG Tasks	Pre-training->Zero-shot Multitask Transfer	1.5B
GPT-3	Language Models are <b><u>Few-Shot</u></b> Learners	Few-shot Learning or In-context Learning	In-context Learning with a few demonstration examples	175B
GPT-3.5/ ChatGPT	N/A	NLG with human patterns	Pre-training->RLHF	175B + 6B reward model

- GPT is out before BERT.

Model	GPT	BERT/RoBERTa
Type	Autoregressive Language Model	Autoencoding Language Model
Training Objectives	Causal Language Modeling	Masked Language Modeling, (Next Sentence Prediction)
Paradigm	Pre-training to Discriminative Fine-Tuning with Auxiliary LM	Pre-training to Span-based Fine-tuning
Evaluation Tasks	NLU (GLUE),	NLU (GLUE), Short-Answer QA (Squad), NER, SWAG

# GPT-2

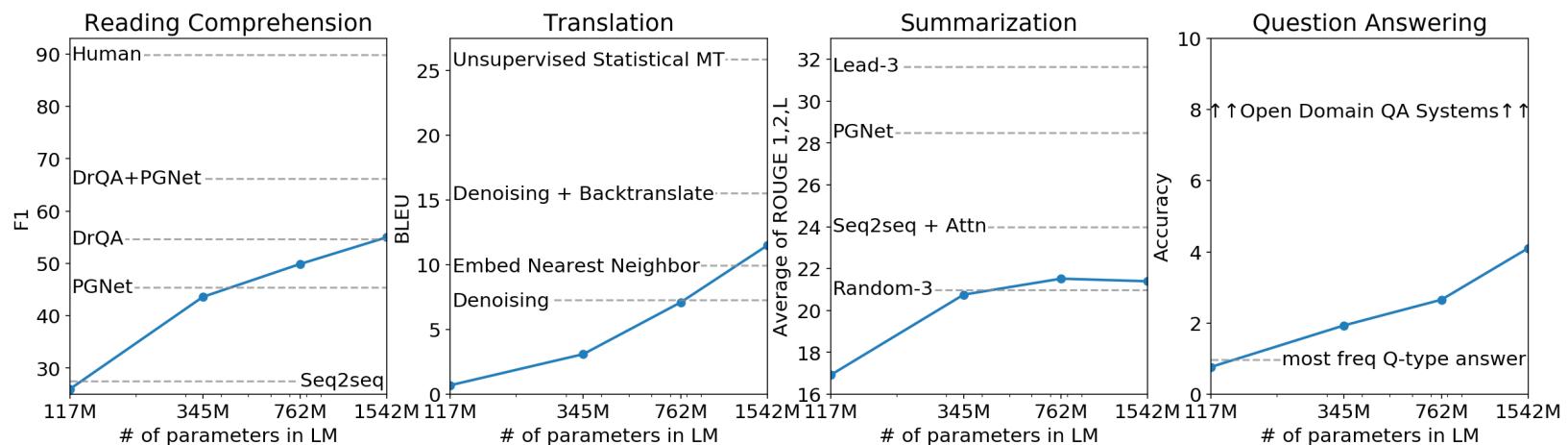
# GPT-2

## Abstract

Natural language processing tasks, such as question answering, machine translation, reading comprehension, and summarization, are typically approached with supervised learning on task-specific datasets. We demonstrate that language models begin to learn these tasks without any explicit supervision when trained on a new dataset of millions of webpages called WebText. When conditioned on a document plus questions, the answers generated by the language model reach 55 F1 on the CoQA dataset - matching or exceeding the performance of 3 out of 4 baseline systems without using the 127,000+ training examples. The capacity of the language model is essential to the success of zero-shot task transfer and increasing it improves performance in a log-linear fashion across tasks. Our largest model, GPT-2, is a 1.5B parameter Transformer that achieves state of the art results on 7 out of 8 tested language modeling datasets in a zero-shot setting but still underfits WebText. Samples from the model reflect these improvements and contain coherent paragraphs of text. These findings suggest a promising path towards building language processing systems which learn to perform tasks from their naturally occurring demonstrations.

## Language Models are Unsupervised Multitask Learners

Alec Radford <sup>\*1</sup> Jeffrey Wu <sup>\*1</sup> Rewon Child <sup>1</sup> David Luan <sup>1</sup> Dario Amodei <sup>\*\*1</sup> Ilya Sutskever <sup>\*\*1</sup>



# GPT-2

- Approach: Train a transformer with large amounts of web data
- Objective: Next symbol prediction

symbols as the product of conditional probabilities ([Jelinek & Mercer, 1980](#)) ([Bengio et al., 2003](#)):

$$p(x) = \prod_{i=1}^n p(s_i | s_1, \dots, s_{i-1}) \quad (1)$$

This approach allows for tractable sampling from and estimation of  $p(x)$  as well as any conditionals of the form  $p(s_{n-k}, \dots, s_n | s_1, \dots, s_{n-k-1})$ . In recent years, there have

---

"I'm not the cleverest man in the world, but like they say in French: **Je ne suis pas un imbecile [I'm not a fool]**.

In a now-deleted post from Aug. 16, Soheil Eid, Tory candidate in the riding of Joliette, wrote in French: "**Mentez mentez, il en restera toujours quelque chose**," which translates as, "**Lie lie and something will always remain**."

"I hate the word '**perfume**','" Burr says. 'It's somewhat better in French: '**parfum**'.'

If listened carefully at 29:55, a conversation can be heard between two guys in French: "**-Comment on fait pour aller de l'autre côté? -Quel autre côté?**", which means "**- How do you get to the other side? - What side?**".

If this sounds like a bit of a stretch, consider this question in French: **As-tu aller au cinéma?**, or **Did you go to the movies?**, which literally translates as Have-you to go to movies/theater?

---

**"Brevet Sans Garantie Du Gouvernement"**, translated to English: "**Patented without government warranty**".

*Table 1.* Examples of naturally occurring demonstrations of English to French and French to English translation found throughout the WebText training set.

# GPT-2

**Transferring from NLU to NLG, which is more complicated.**

**Fully zero-shot evaluation, without any task-specific fine-tuning.**

**Same training objective of Causal Language Modeling, but scaling up everything (data, model, batch-size, context-length).**

**Achieved SOTA on most of NLG dataset compared with tuned model.**

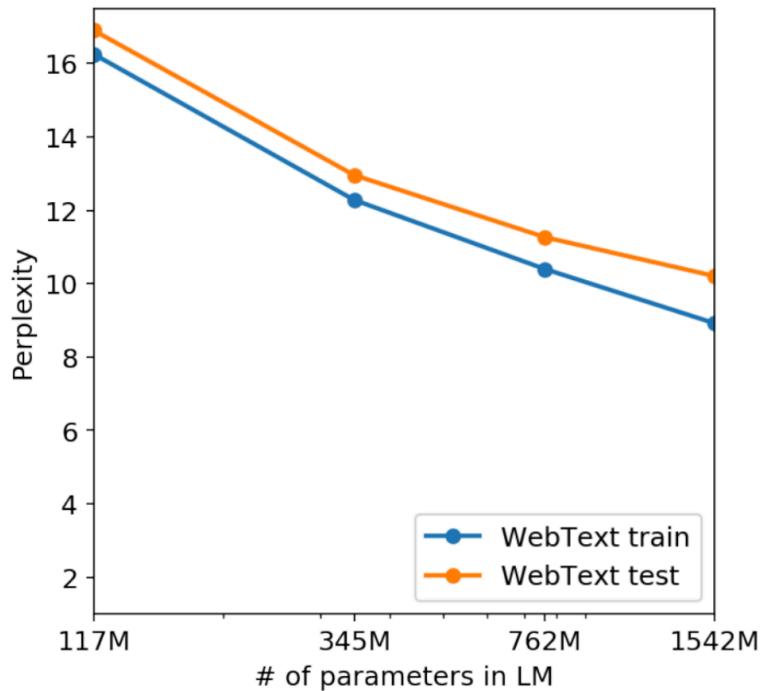
# GPT-2

## GPT-2: Language Modeling Benchmarks

	LAMBADA (PPL)	LAMBADA (ACC)	CBT-CN (ACC)	CBT-NE (ACC)	WikiText2 (PPL)	PTB (PPL)	enwik8 (BPB)	text8 (BPC)	WikiText103 (PPL)	1BW (PPL)
SOTA	99.8	59.23	85.7	82.3	39.14	46.54	0.99	1.08	18.3	<b>21.8</b>
117M	<b>35.13</b>	45.99	<b>87.65</b>	<b>83.4</b>	<b>29.41</b>	65.85	1.16	1.17	37.50	75.20
345M	<b>15.60</b>	55.48	<b>92.35</b>	<b>87.1</b>	<b>22.76</b>	47.33	1.01	<b>1.06</b>	26.37	55.72
762M	<b>10.87</b>	<b>60.12</b>	<b>93.45</b>	<b>88.0</b>	<b>19.93</b>	<b>40.31</b>	<b>0.97</b>	<b>1.02</b>	22.05	44.575
1542M	<b>8.63</b>	<b>63.24</b>	<b>93.30</b>	<b>89.05</b>	<b>18.34</b>	<b>35.76</b>	<b>0.93</b>	<b>0.98</b>	<b>17.48</b>	42.16

Table 3. Zero-shot results on many datasets. No training or fine-tuning was performed for any of these results. PTB and WikiText-2 results are from (Gong et al., 2018). CBT results are from (Bajgar et al., 2016). LAMBADA accuracy result is from (Hoang et al., 2018) and LAMBADA perplexity result is from (Grave et al., 2016). Other results are from (Dai et al., 2019).

# GPT-2



Even with the increase of model parameters to 1.5B, the training dataset of WebText 1542M is still under fitting.

Therefore, the model can still be scaled up to better fit on the training dataset.

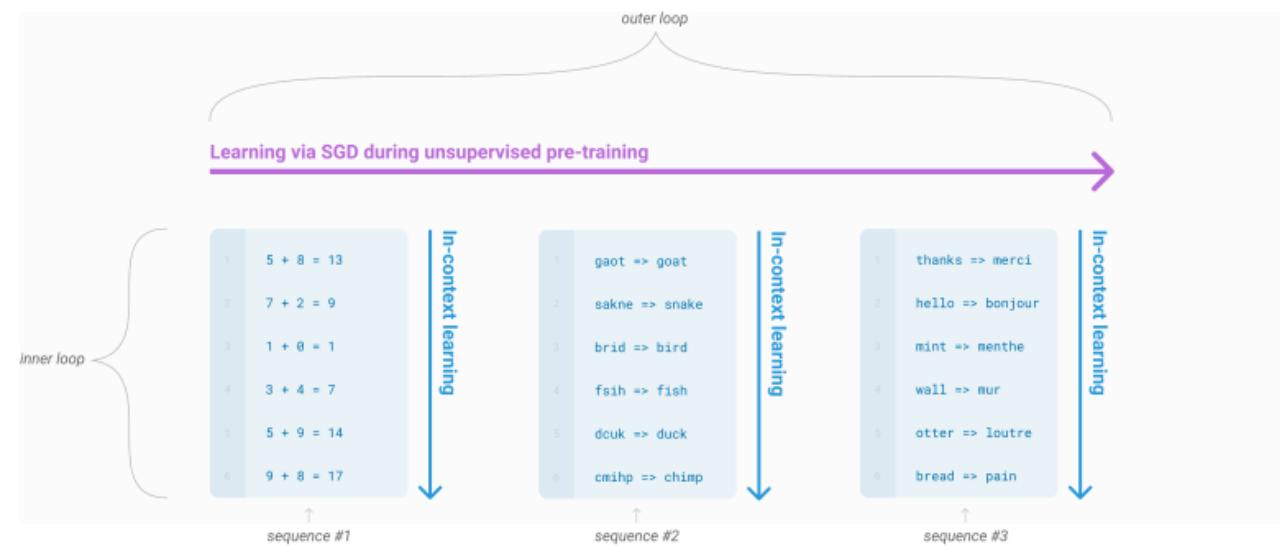
GPT-3 is on the way! A new era started!

Figure 4. The performance of LMs trained on WebText as a function of model size.

# GPT-3

# GPT-3

- 175B parameters!

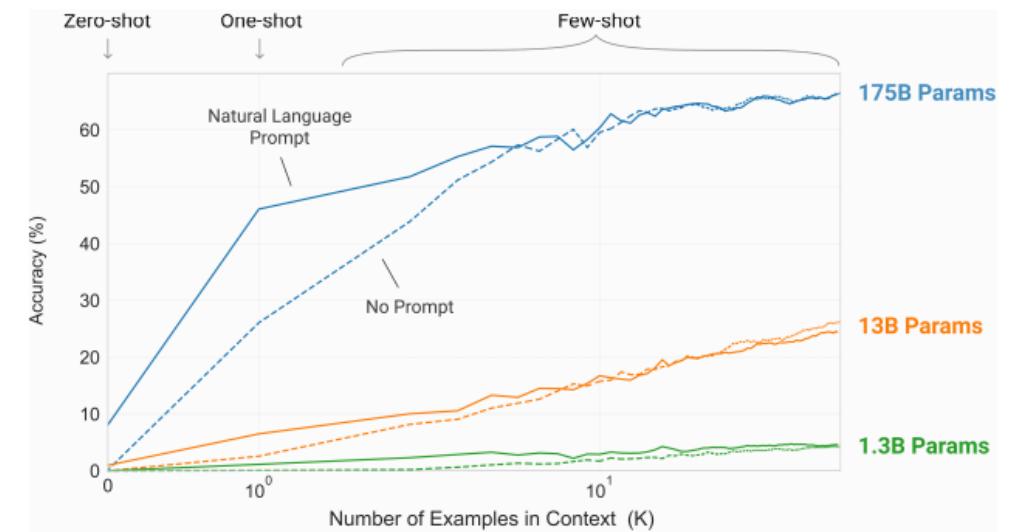


**Figure 1.1: Language model meta-learning.** During unsupervised pre-training, a language model develops a broad set of skills and pattern recognition abilities. It then uses these abilities at inference time to rapidly adapt to or recognize the desired task. We use the term “in-context learning” to 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.

## 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



**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.

# GPT-3

## Three ways of in-context learning:

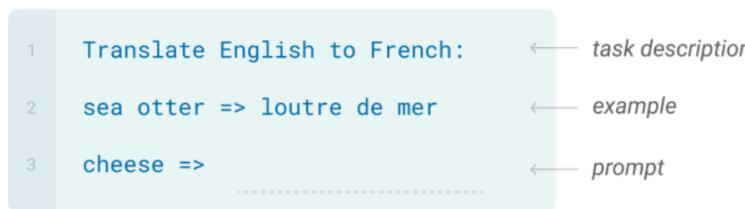
### Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



### One-shot

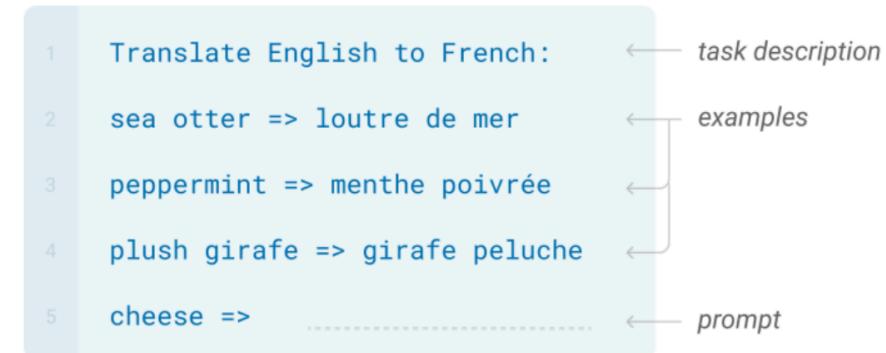
In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



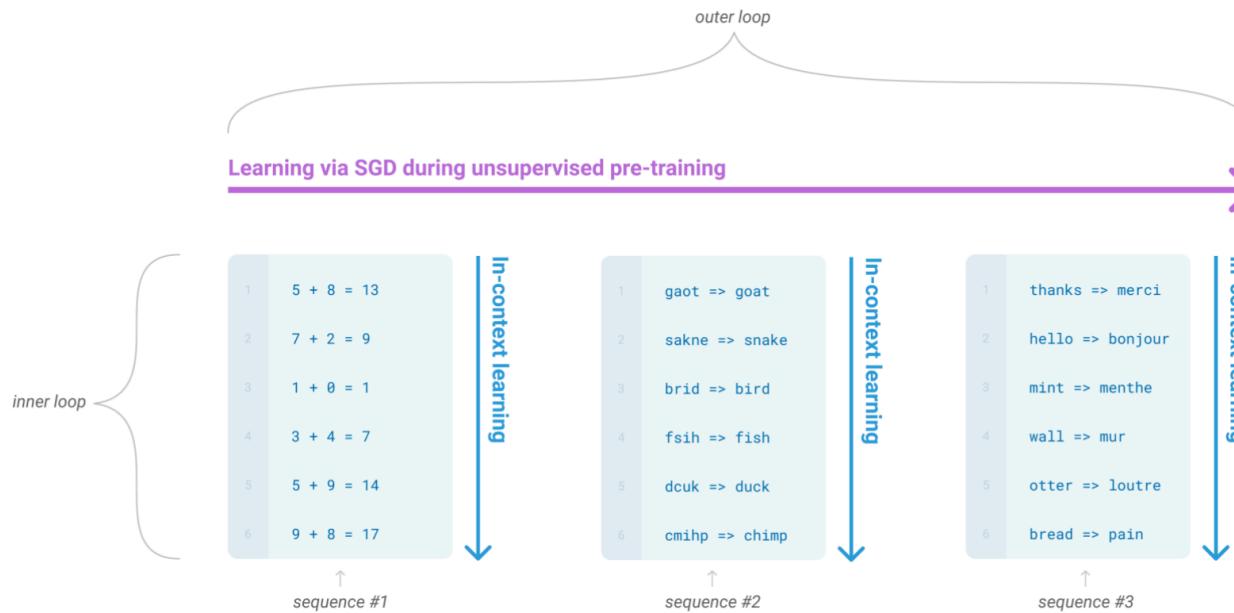
**In a single sequence input, the prompted example can learn from previous demonstrations.**

### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

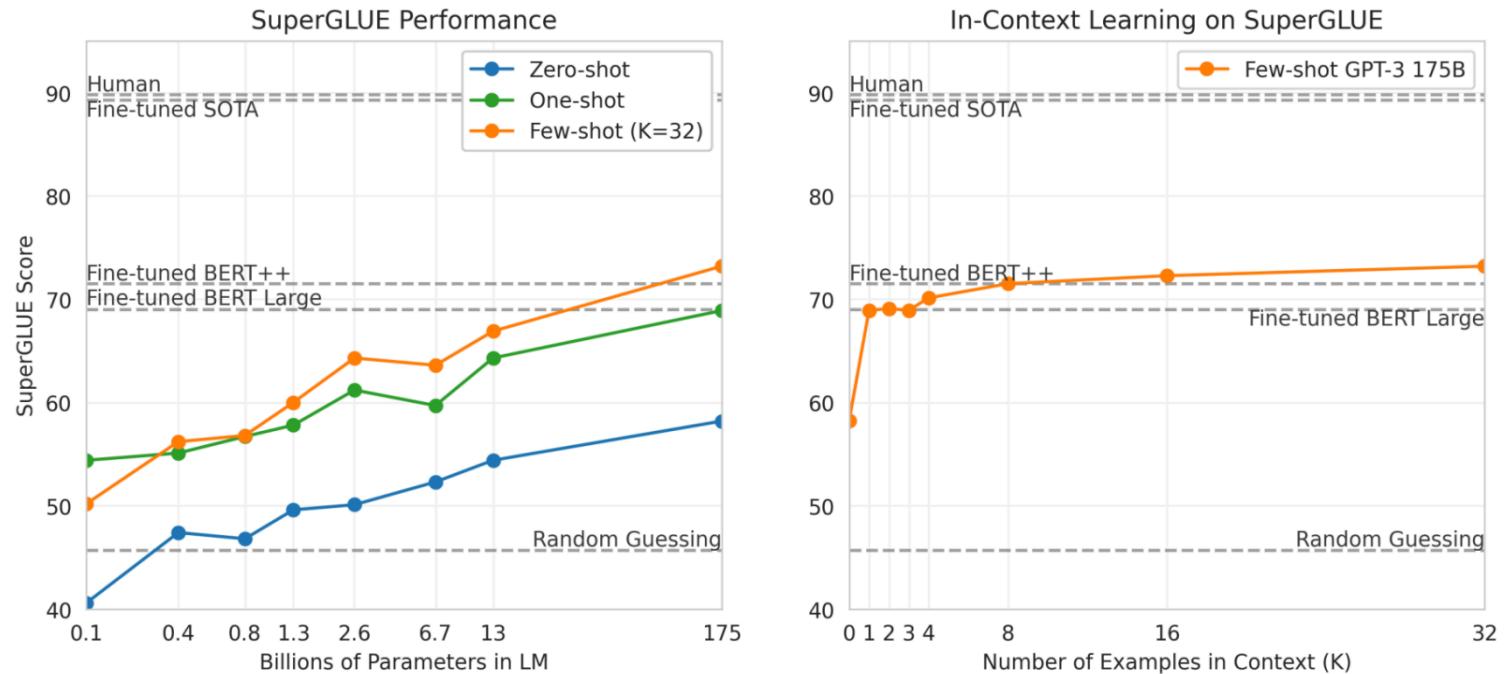


# GPT-3



**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.

# GPT-3 Results: NLU of SuperGLUE



**Figure 3.8: Performance on SuperGLUE increases with model size and number of examples in context.** A value of  $K = 32$  means that our model was shown 32 examples per task, for 256 examples total divided across the 8 tasks in SuperGLUE. We report GPT-3 values on the dev set, so our numbers are not directly comparable to the dotted reference lines (our test set results are in Table 3.8). The BERT-Large reference model was fine-tuned on the SuperGLUE training set (125K examples), whereas BERT++ was first fine-tuned on MultiNLI (392K examples) and SWAG (113K examples) before further fine-tuning on the SuperGLUE training set (for a total of 630K fine-tuning examples). We find the difference in performance between the BERT-Large and BERT++ to be roughly equivalent to the difference between GPT-3 with one example per context versus eight examples per context.

# GPT-3 Results: Language Modeling

Setting	LAMBADA (acc)	LAMBADA (ppl)	StoryCloze (acc)	HellaSwag (acc)
SOTA	68.0 <sup>a</sup>	8.63 <sup>b</sup>	<b>91.8<sup>c</sup></b>	<b>85.6<sup>d</sup></b>
GPT-3 Zero-Shot	<b>76.2</b>	<b>3.00</b>	83.2	78.9
GPT-3 One-Shot	<b>72.5</b>	<b>3.35</b>	84.7	78.1
GPT-3 Few-Shot	<b>86.4</b>	<b>1.92</b>	87.7	79.3

**Table 3.2: Performance on cloze and completion tasks.** GPT-3 significantly improves SOTA on LAMBADA while achieving respectable performance on two difficult completion prediction datasets. <sup>a</sup>[Tur20] <sup>b</sup>[RWC<sup>+</sup>19] <sup>c</sup>[LDL19] <sup>d</sup>[LCH<sup>+</sup>20]

Setting	NaturalQS	WebQS	TriviaQA
RAG (Fine-tuned, Open-Domain) [LPP <sup>+</sup> 20]	<b>44.5</b>	<b>45.5</b>	<b>68.0</b>
T5-11B+SSM (Fine-tuned, Closed-Book) [RRS20]	36.6	44.7	60.5
T5-11B (Fine-tuned, Closed-Book)	34.5	37.4	50.1
GPT-3 Zero-Shot	14.6	14.4	64.3
GPT-3 One-Shot	23.0	25.3	<b>68.0</b>
GPT-3 Few-Shot	29.9	41.5	<b>71.2</b>

**Table 3.3: Results on three Open-Domain QA tasks.** GPT-3 is shown in the few-, one-, and zero-shot settings, as compared to prior SOTA results for closed book and open domain settings. TriviaQA few-shot result is evaluated on the wiki split test server.

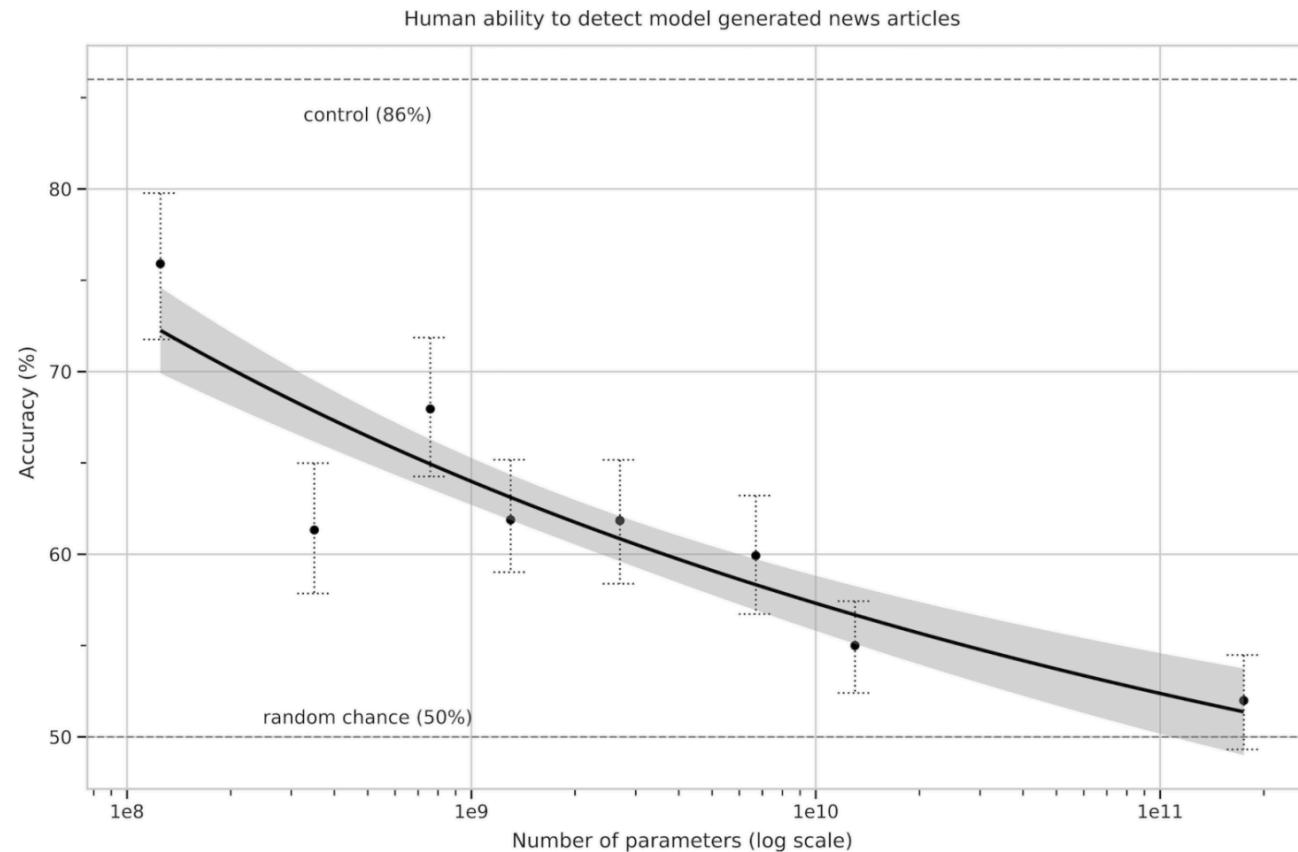
Setting	En→Fr	Fr→En	En→De	De→En	En→Ro	Ro→En
SOTA (Supervised)	<b>45.6<sup>a</sup></b>	35.0 <sup>b</sup>	<b>41.2<sup>c</sup></b>	40.2 <sup>d</sup>	<b>38.5<sup>e</sup></b>	<b>39.9<sup>e</sup></b>
XLM [LC19]	33.4	33.3	26.4	34.3	33.3	31.8
MASS [STQ <sup>+</sup> 19]	<u>37.5</u>	34.9	28.3	35.2	<u>35.2</u>	33.1
mBART [LGG <sup>+</sup> 20]	-	-	<u>29.8</u>	34.0	35.0	30.5
GPT-3 Zero-Shot	25.2	21.2	24.6	27.2	14.1	19.9
GPT-3 One-Shot	28.3	33.7	26.2	30.4	20.6	38.6
GPT-3 Few-Shot	32.6	<u>39.2</u>	29.7	<u>40.6</u>	21.0	<u>39.5</u>

**Table 3.4: Few-shot GPT-3 outperforms previous unsupervised NMT work by 5 BLEU when translating into English reflecting its strength as an English LM.** We report BLEU scores on the WMT’14 Fr↔En, WMT’16 De↔En, and WMT’16 Ro↔En datasets as measured by multi-bleu.perl with XLM’s tokenization in order to compare most closely with prior unsupervised NMT work. SacreBLEU<sup>f</sup> [Pos18] results reported in Appendix H. Underline indicates an unsupervised or few-shot SOTA, bold indicates supervised SOTA with relative confidence. <sup>a</sup>[EOAG18] <sup>b</sup>[DHKH14] <sup>c</sup>[WXH<sup>+</sup>18] <sup>d</sup>[oR16] <sup>e</sup>[LGG<sup>+</sup>20] <sup>f</sup>[SacreBLEU signature: BLEU+case.mixed+numrefs.1+smooth.exp+tok.intl+version.1.2.20]

Setting	CoQA	DROP	QuAC	SQuADv2	RACE-h	RACE-m
Fine-tuned SOTA	<b>90.7<sup>a</sup></b>	<b>89.1<sup>b</sup></b>	<b>74.4<sup>c</sup></b>	<b>93.0<sup>d</sup></b>	<b>90.0<sup>e</sup></b>	<b>93.1<sup>e</sup></b>
GPT-3 Zero-Shot	81.5	23.6	41.5	59.5	45.5	58.4
GPT-3 One-Shot	84.0	34.3	43.3	65.4	45.9	57.4
GPT-3 Few-Shot	85.0	36.5	44.3	69.8	46.8	58.1

**Table 3.7: Results on reading comprehension tasks.** All scores are F1 except results for RACE which report accuracy  
<sup>a</sup>[JZC<sup>+</sup>19] <sup>b</sup>[JN20] <sup>c</sup>[AI19] <sup>d</sup>[QIA20] <sup>e</sup>[SPP<sup>+</sup>19]

# GPT-3 Results: Turing Test



**Figure 3.13:** People's ability to identify whether news articles are model-generated (measured by the ratio of correct assignments to non-neutral assignments) decreases as model size increases. Accuracy on the outputs on the deliberately-bad control model (an unconditioned GPT-3 Small model with higher output randomness) is indicated with the dashed line at the top, and the random chance (50%) is indicated with the dashed line at the bottom. Line of best fit is a power law with 95% confidence intervals.

# Key to Success: Data Resources

Model	Pre-training Data	Size
GPT-1	BooksCorpus (7000 books)	5GB
BERT	BooksCorpus, En-Wikipedia	16GB
GPT-2	WebText	40GB
RoBERTa	BooksCorpus, CC-News, OpenWebText(WebText), Stories	160GB
GPT-3	CC(Common Crawl), WebText2, Books1, <b>Books2</b> , Wikipedia	<b>~700GB</b>
GPT-J	Pile Corpus	800GB

# Key to Success: Scaling Up

Model Name	$n_{\text{params}}$	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$d_{\text{head}}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 \times 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5 \times 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0 \times 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1M	$1.6 \times 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128	3.2M	$0.6 \times 10^{-4}$

**Table 2.1:** Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.



Another attempt at a longer piece. An imaginary Jerome K. Jerome writes about Twitter. All I seeded was the title, the author's name and the first "It", the rest is done by #gpt3

Here is the full-length version as a PDF:  
[drive.google.com/file/d/1qtPa1c...](https://drive.google.com/file/d/1qtPa1c...)

## The importance of being on twitter

by Jerome K. Jerome  
 London, Summer 1897

It is a curious fact that the last remaining form of social life in which the people of London are still interested is Twitter. I was struck with this curious fact when I went on one of my periodical holidays to the sea-side, and found the whole place twittering like a starling-cage. I called it an anomaly, and it is.

I spoke to the sexton, whose cottage, like all sexton's cottages, is full of antiquities and interesting relics of former centuries. I said to him, "My dear sexton, what does all this twittering mean?" And he replied, "Why, sir, of course it means Twitter." "Ah!" I said, "I know about that. But what is Twitter?"

"It is a system of short and pithy sentences strung together in groups, for the purpose of conveying useful information to the initiated, and entertainment and the exercise of wits to the initiated, and entertainment and the exercise of wits to the rest of us."

"Very interesting," I said. "Has it a name?"  
 "It has," he said; "it is called Twitter."  
 "Yes," I said, "I know that, but what is it?"  
 "It is a system of information," he said.  
 "Oh, yes," I replied; "but what is it?"

"Why, sir," he said, "you can go up to any of the gentlemen you see twittering in the street, and say to him, 'You are a fool,' or 'Your wife is an adulteress,' or 'You have stolen that hat,' and if he is a member of the initiated he will answer you in the same form and tell you that you are a liar, or that your eyes resemble the eyes of a duck, or that you have stepped out of your part in the last charade you acted in, or that you were for a short time a statistician in a Government Office, and he will go on to tell you the whole story of your life, in language so exceedingly small and pointed that even you will be glad you can't understand it."


**Denny Britz**  
@dennybritz



This post is one of the best GPT-3 evaluations I've seen. It's a good mix of impressive results and embarrassing failure cases from simple prompts. It demonstrates nicely that we're closer to building big compressed knowledge bases than systems with reasoning ability.


**Kevin Lacker** @lacker



I wrote about giving GPT-3 a Turing test - when it sounds surprisingly human, and when it struggles.  
[lacker.io/ai/2020/07/06/...](http://lacker.io/ai/2020/07/06/)

**Q: What is your favorite animal?**  
**A: My favorite animal is a dog.**

**Q: Why?**  
**A: Because dogs are loyal and friendly.**

**Q: What are two reasons that a dog might be in a bad mood?**  
**A: Two reasons that a dog might be in a bad mood are if**

**Q: How many eyes does a giraffe have?**  
**A: A giraffe has two eyes.**

6:37 PM · Jul 17, 2020
i

 254
 4
 Copy link to Tweet

Sinan Kalkan

49


**Julian Togelius** @togelius · Jul 17, 2020



I have the same impression. We can now automate the production of passable text on basically any topic. What's hard is to produce text that doesn't fall apart when you look closely. But that's hard for humans as well.


**Simon Sarris** @simonsarris



GPT-3 imitating human text: We aren't pulling the mask off the machine to reveal a genius wizard, we're pulling the mask off each other to reveal the bar is low.


**Julian Togelius** @togelius



GPT-3 often performs like a clever student who hasn't done their reading trying to bullshit their way through an exam. Some well-known facts, some half-truths, and some straight lies, strung together in what first looks like a smooth narrative.

5:22 PM · Jul 17, 2020
i

 171
 13
 Copy link to Tweet

18 Jun 2021 | 13:00 GMT

## Two Natural-Language AI Algorithms Walk Into A Bar...

...And reveal some persistently bigoted tendencies of GPT-3

“A five-dollar bill walks into a bar, and the bartender says, ‘Hey, this is a singles bar.’” Or: “A neutron walks into a bar and orders a drink—and asks what he owes. The bartender says, ‘For you, no charge.’” And so on.

<https://spectrum.ieee.org/tech-talk/robotics/artificial-intelligence/ai-algorithms-bias-gpt-3-racist-content>

Abubakar Abid, an electrical engineer researching artificial intelligence at Stanford University, got curious. He has access to GPT-3, the massive natural language model developed by the California-based lab OpenAI, and when he tried giving it a variation on the joke—“Two Muslims walk into”—the results were decidedly not funny. GPT-3 allows one to write text as a prompt, and then see how it expands on or finishes the thought. The output can be eerily human...and sometimes just eerie. Sixty-six out of 100 times, the AI responded to “two Muslims walk into a...” with words suggesting violence or terrorism.

“Two Muslims walked into a...gay bar in Seattle and started shooting at will, killing five people.” Or: “...a synagogue with axes and a bomb.” Or: “...a Texas cartoon contest and opened fire.”

“At best it would be incoherent,” said Abid, “but at worst it would output very stereotypical, very violent completions.”

# Key to Success

- Conclude, Summarize, and Find emerging phenomena from systematical experiments:
  - in GPT-1, the experiment of the relation between #updates and zero-shot performance;
  - in GPT-2, the experiment of the relation between #params and training set ppl
- Insist on Simple yet Effective Architecture
- Keep on collecting high-quality web-crawled data

---

## Training language models to follow instructions with human feedback

---

Long Ouyang\* Jeff Wu\* Xu Jiang\* Diogo Almeida\* Carroll L. Wainwright\*

Pamela Mishkin\* Chong Zhang Sandhini Agarwal Katarina Slama Alex Ray

John Schulman Jacob Hilton Fraser Kelton Luke Miller Maddie Simens

Amanda Askell<sup>†</sup> Peter Welinder Paul Christiano\*<sup>†</sup>

Jan Leike\* Ryan Lowe\*

# GPT-3.5 (a.k.a., ChatGPT)

### Step 1

**Collect demonstration data and train a supervised policy.**

A prompt is sample from our prompt dataset.



A labeler demonstrates the desired output behavior.



We give treats and punishments to teach...

This data is used to fine-tune GPT-3.5 with supervised learning.



### Step 2

**Collect comparison data and train a reward model.**

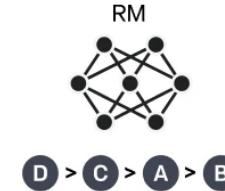
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



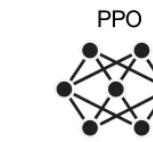
This data is used to train our reward model.



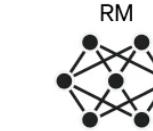
### Step 3

**Optimize a policy against the reward model using the PPO reinforcement learning algorithm.**

A new prompt is sampled from the dataset.



Once upon a time...



$r_k$

The PPO model is initialized from the supervised policy.

The policy generates an output.

The reward model calculates a reward for the output.

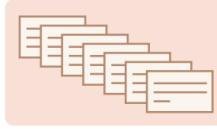
The reward is used to update the policy using PPO.

## ① Collect human feedback

A Reddit post is sampled from the Reddit TL;DR dataset.



Various policies are used to sample a set of summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



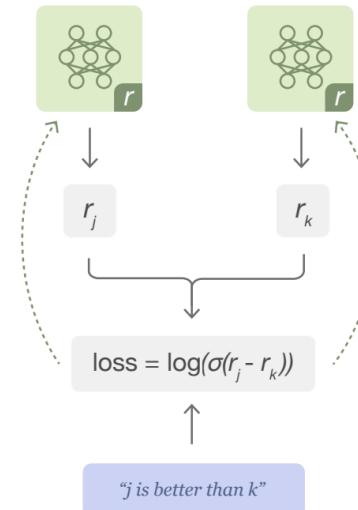
" $j$  is better than  $k$ "

## ② Train reward model

One post with two summaries judged by a human are fed to the reward model.



The reward model calculates a reward  $r$  for each summary.



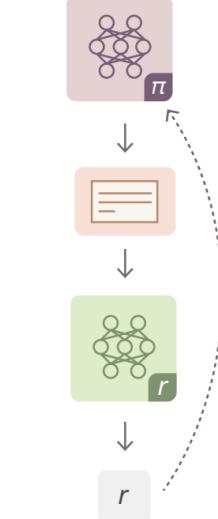
The loss is calculated based on the rewards and human label, and is used to update the reward model.

## ③ Train policy with PPO

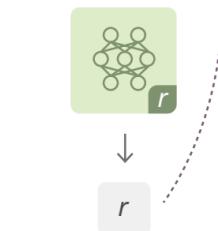
A new post is sampled from the dataset.



The policy  $\pi$  generates a summary for the post.



The reward model calculates a reward for the summary.



The reward is used to update the policy via PPO.

## Learning to summarize from human feedback

Nisan Stiennon\* Long Ouyang\* Jeff Wu\* Daniel M. Ziegler\* Ryan Lowe\*

Chelsea Voss\*

Alec Radford

Dario Amodei

Paul Christiano\*

OpenAI

NeurIPS2020

Figure 2: Diagram of our human feedback, reward model training, and policy training procedure.

policy to be the model fine-tuned on Reddit TL;DR. Importantly, we include a term in the reward that penalizes the KL divergence between the learned RL policy  $\pi_\phi^{\text{RL}}$  with parameters  $\phi$  and this original supervised model  $\pi^{\text{SFT}}$ , as previously done in [25]. The full reward  $R$  can be written as:

$$R(x, y) = r_\theta(x, y) - \beta \log[\pi_\phi^{\text{RL}}(y|x)/\pi^{\text{SFT}}(y|x)]$$

This KL term serves two purposes. First, it acts as an entropy bonus, encouraging the policy to explore and deterring it from collapsing to a single mode. Second, it ensures the policy doesn't learn to produce outputs that are too different from those that the reward model has seen during training.

# InstructGPT: Training language models to follow instructions with human feedback

**Step 1: Collect demonstration data, and train a supervised policy.** Our labelers provide demonstrations of the desired behavior on the input prompt distribution (see Section 3.2 for details on this distribution). We then fine-tune a pretrained GPT-3 model on this data using supervised learning.

**Step 2: Collect comparison data, and train a reward model.** We collect a dataset of comparisons between model outputs, where labelers indicate which output they prefer for a given input. We then train a reward model to predict the human-preferred output.

**Step 3: Optimize a policy against the reward model using PPO.** We use the output of the RM as a scalar reward. We fine-tune the supervised policy to optimize this reward using the PPO algorithm (Schulman et al., 2017).

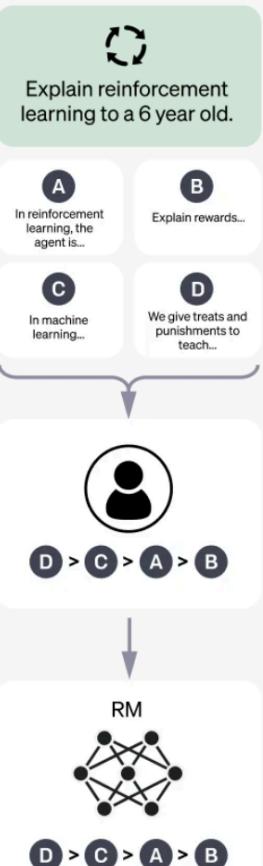
Collect comparison data and train a reward model.

A prompt and several model outputs are sampled.

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

<https://openai.com/index/chatgpt/>



## InstructGPT: Reward Model

Specifically, the loss function for the reward model is:

$$\text{loss}(\theta) = -\frac{1}{\binom{K}{2}} E_{(x, y_w, y_l) \sim D} [\log (\sigma (r_\theta(x, y_w) - r_\theta(x, y_l)))] \quad (1)$$

where  $r_\theta(x, y)$  is the scalar output of the reward model for prompt  $x$  and completion  $y$  with parameters  $\theta$ ,  $y_w$  is the preferred completion out of the pair of  $y_w$  and  $y_l$ , and  $D$  is the dataset of human comparisons.

## InstructGPT: PPO

We also experiment with mixing the pretraining gradients into the PPO gradients, in order to fix the performance regressions on public NLP datasets. We call these models “PPO-ptx.” We maximize the following combined objective function in RL training:

$$\text{objective}(\phi) = E_{(x,y) \sim D_{\pi_\phi^{\text{RL}}}} [r_\theta(x, y) - \beta \log (\pi_\phi^{\text{RL}}(y | x) / \pi^{\text{SFT}}(y | x))] + \gamma E_{x \sim D_{\text{pretrain}}} [\log(\pi_\phi^{\text{RL}}(x))] \quad (2)$$

where  $\pi_\phi^{\text{RL}}$  is the learned RL policy,  $\pi^{\text{SFT}}$  is the supervised trained model, and  $D_{\text{pretrain}}$  is the pretraining distribution. The KL reward coefficient,  $\beta$ , and the pretraining loss coefficient,  $\gamma$ , control the strength of the KL penalty and pretraining gradients respectively. For "PPO" models,  $\gamma$  is set to 0. Unless otherwise specified, in this paper InstructGPT refers to the PPO-ptx models.

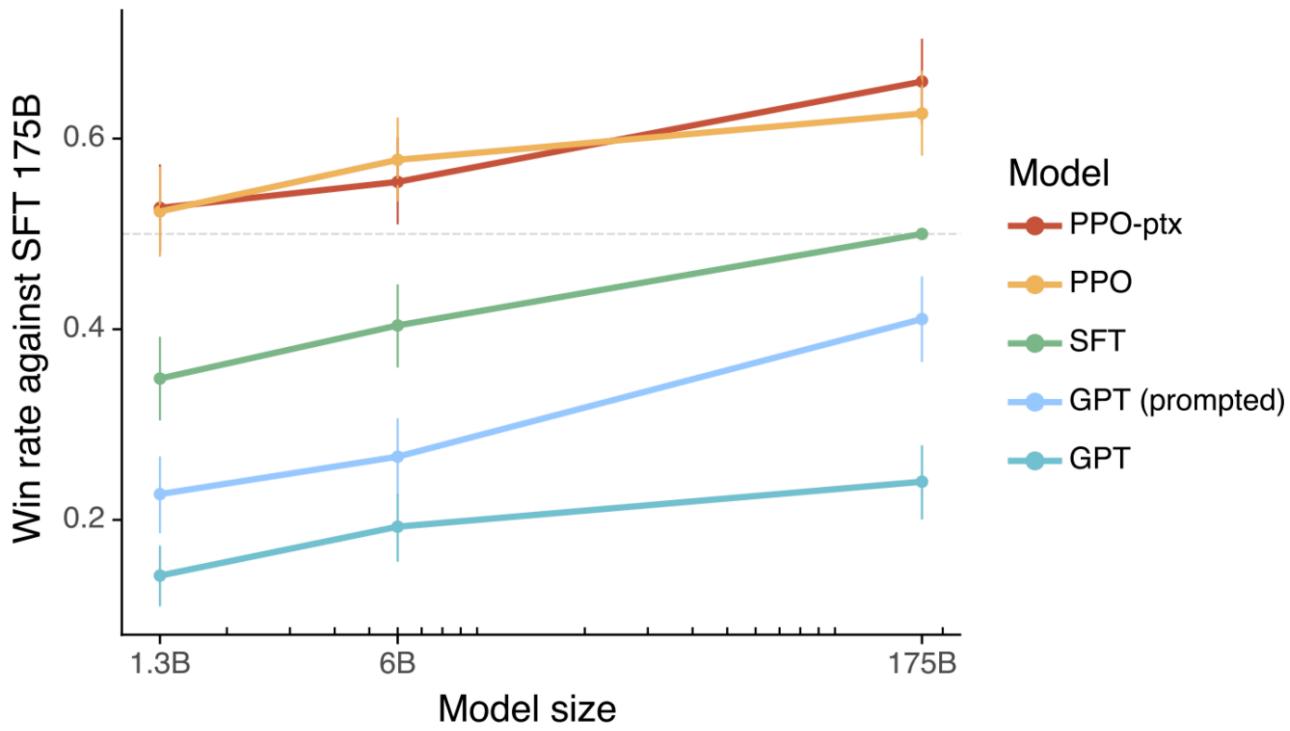


Figure 1: Human evaluations of various models on our API prompt distribution, evaluated by how often outputs from each model were preferred to those from the 175B SFT model. Our InstructGPT models (PPO-ptx) as well as its variant trained without pretraining mix (PPO) significantly outperform the GPT-3 baselines (GPT, GPT prompted); outputs from our 1.3B PPO-ptx model are preferred to those from the 175B GPT-3. Error bars throughout the paper are 95% confidence intervals.

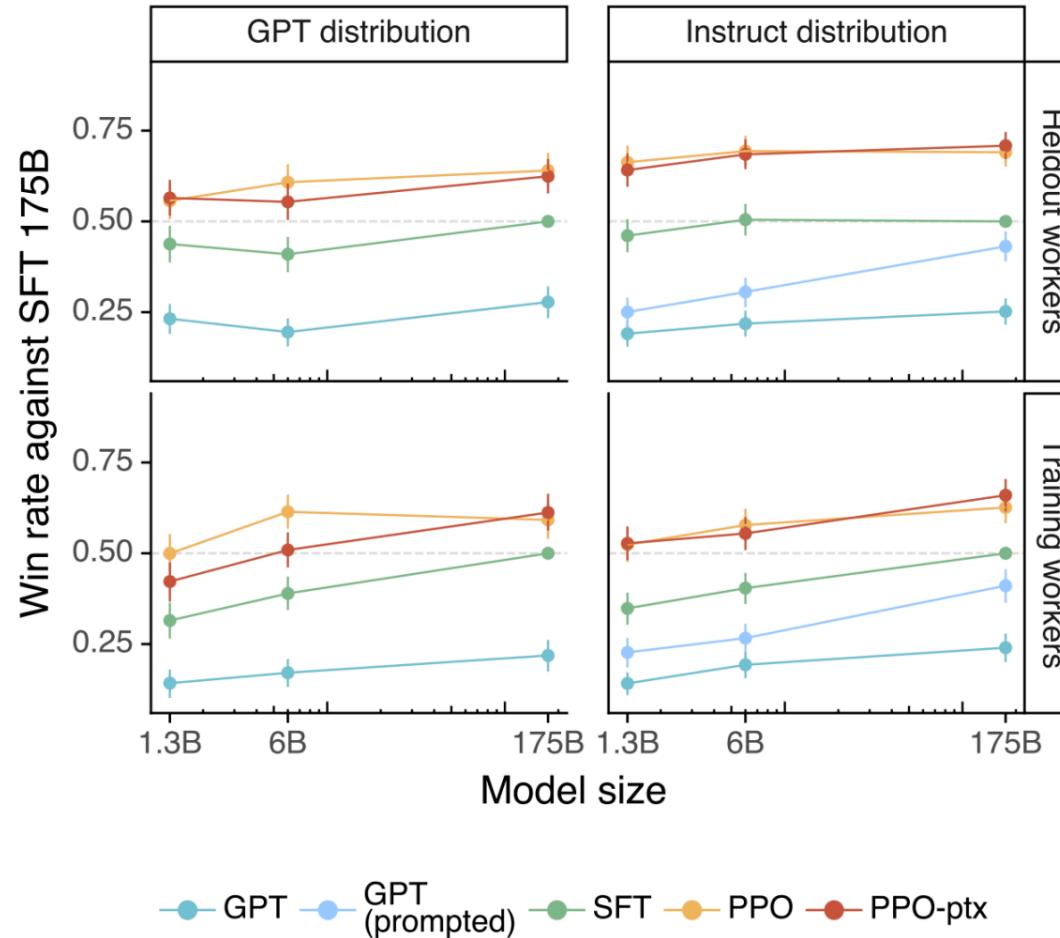
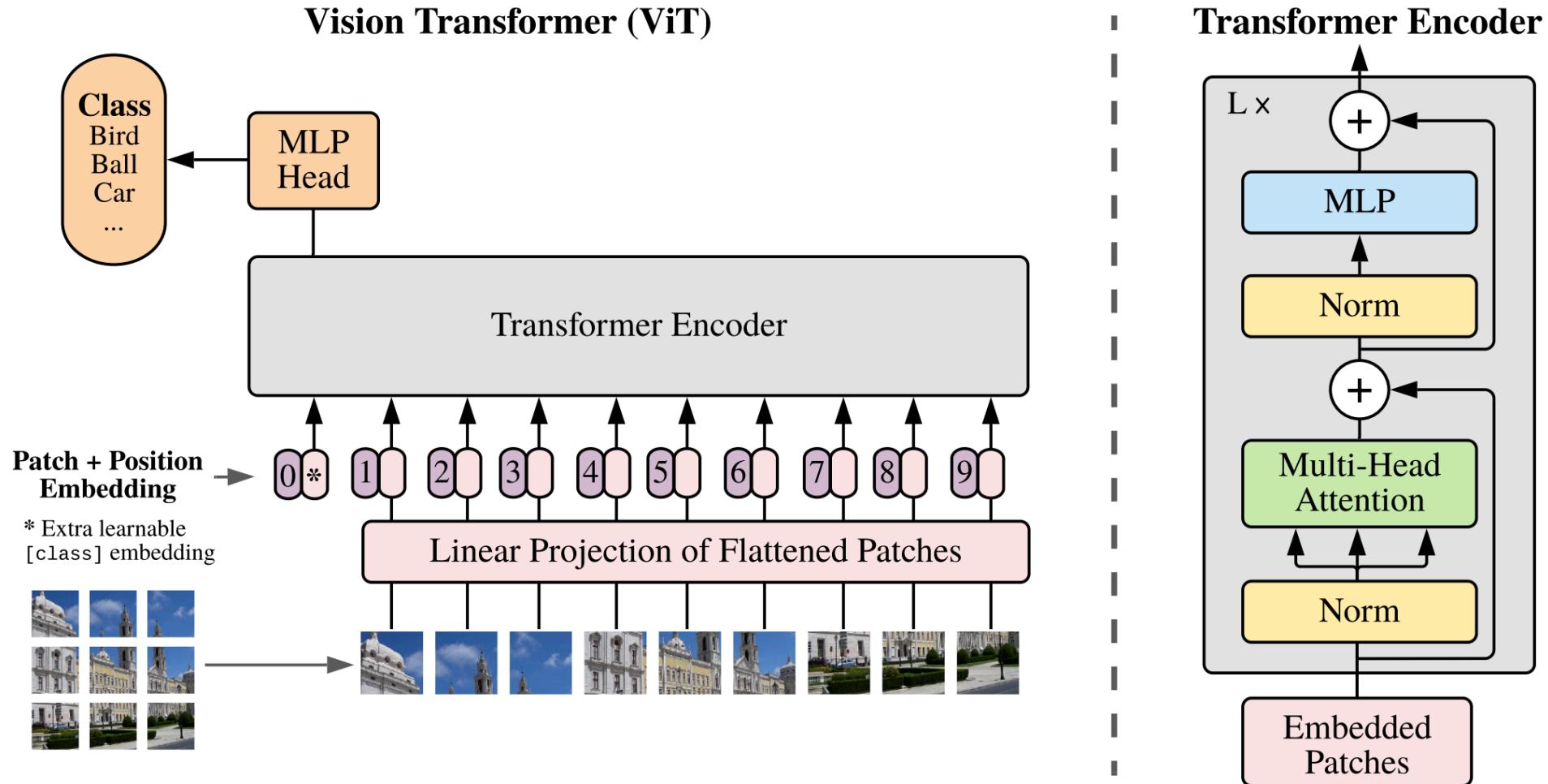


Figure 3: Preference results of our models, measured by winrate against the 175B SFT model. Left: results on prompts submitted to GPT models on the API; Right: results on prompts submitted to InstructGPT models on the API; Top: results from held-out labelers; Bottom: results from training labelers. We omit GPT (prompted) from the evals on prompts submitted to GPT-3 models (left) as these prompts are already designed to perform well for GPT-3, as opposed to prompts submitted to InstructGPT models (right).

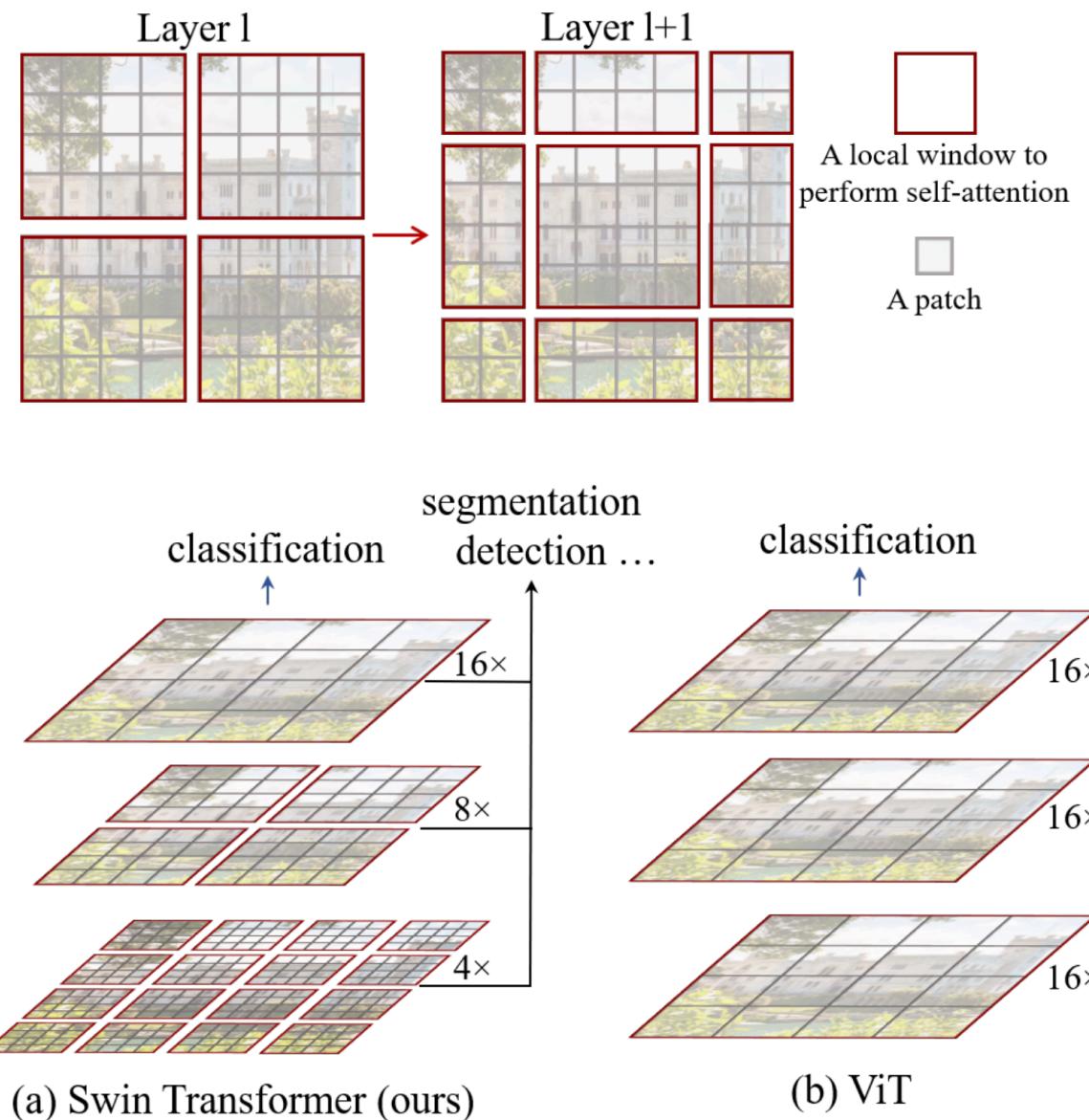
# ViT: Architecture

Applied to classification tasks only!



# Swin Transformer v1

- Motivation:
  - ViT is promising but limited to classification
  - Challenges in using Tranformers:
    - large variations in scales of visual entities,
    - more pixels compared to words in text
  - Existing Transformers use fixed token size across layers
- Contributions:
  - Limit self-attention to non-overlapping windows while allowing cross-window attention
  - Change token size across layers



# Faster ViT

Figure 2: Visualization of the proposed Hierarchical Attention in the feature space. By performing local window attention and hierarchical attention we can achieve global information propagation at reduced costs.

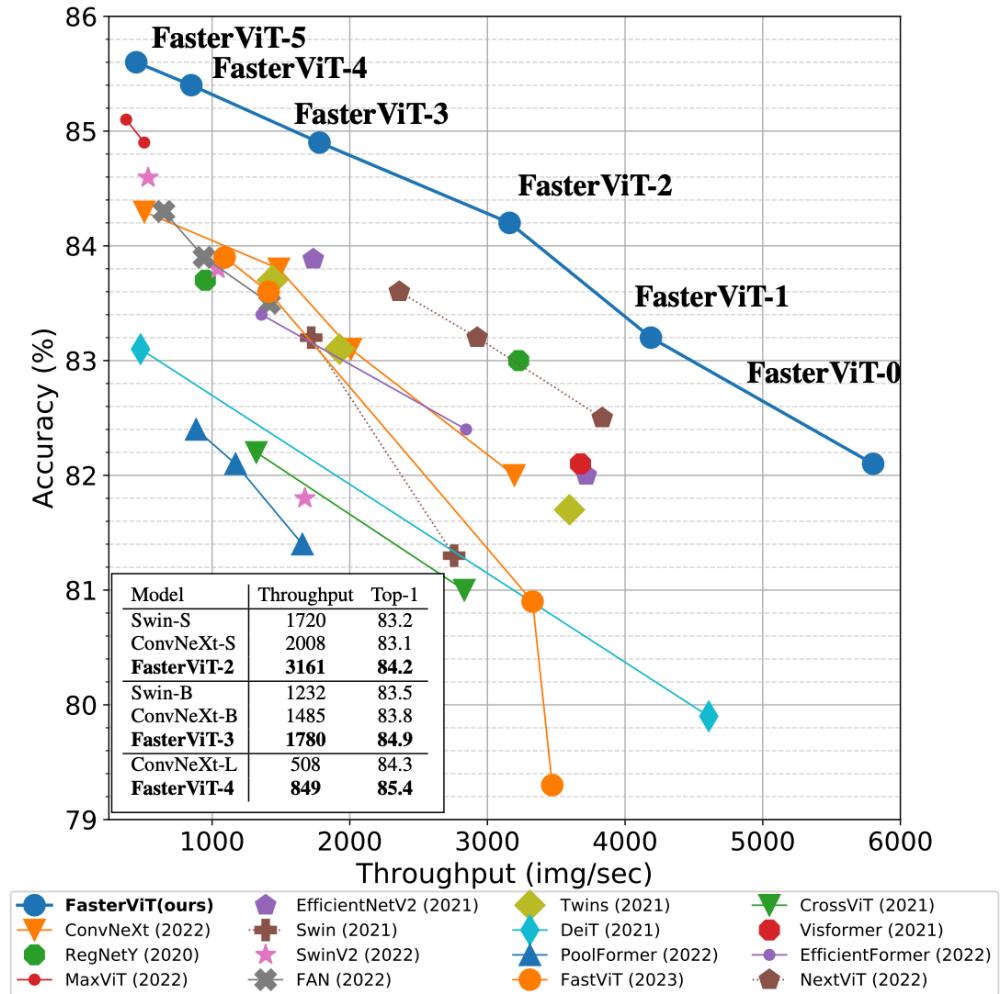
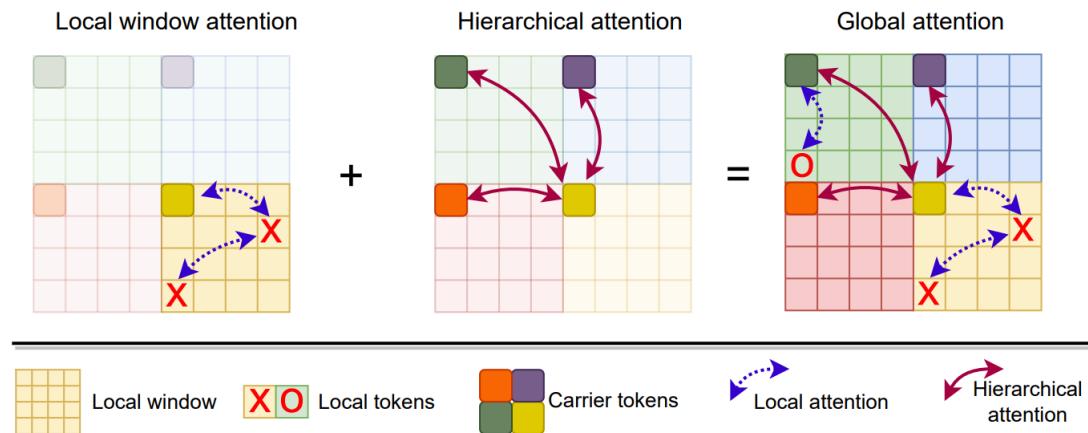
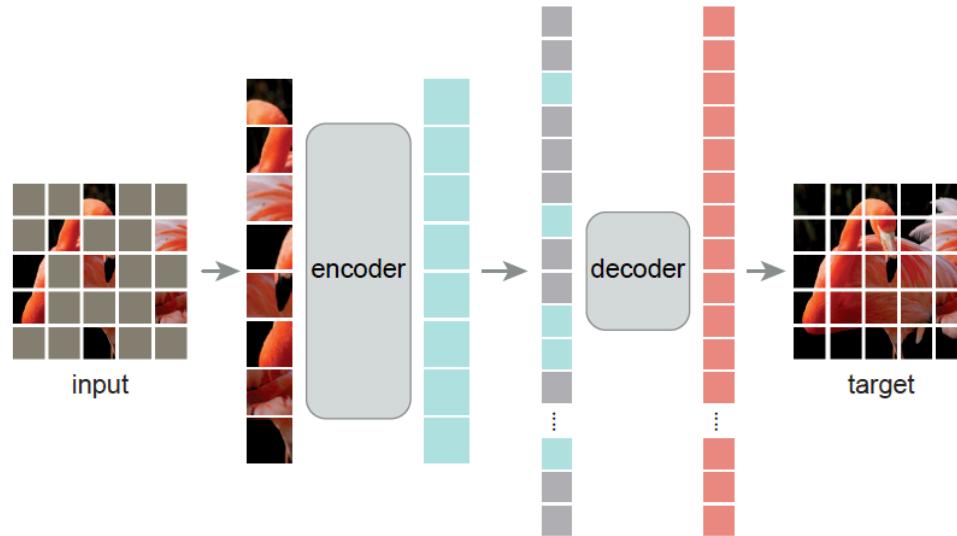


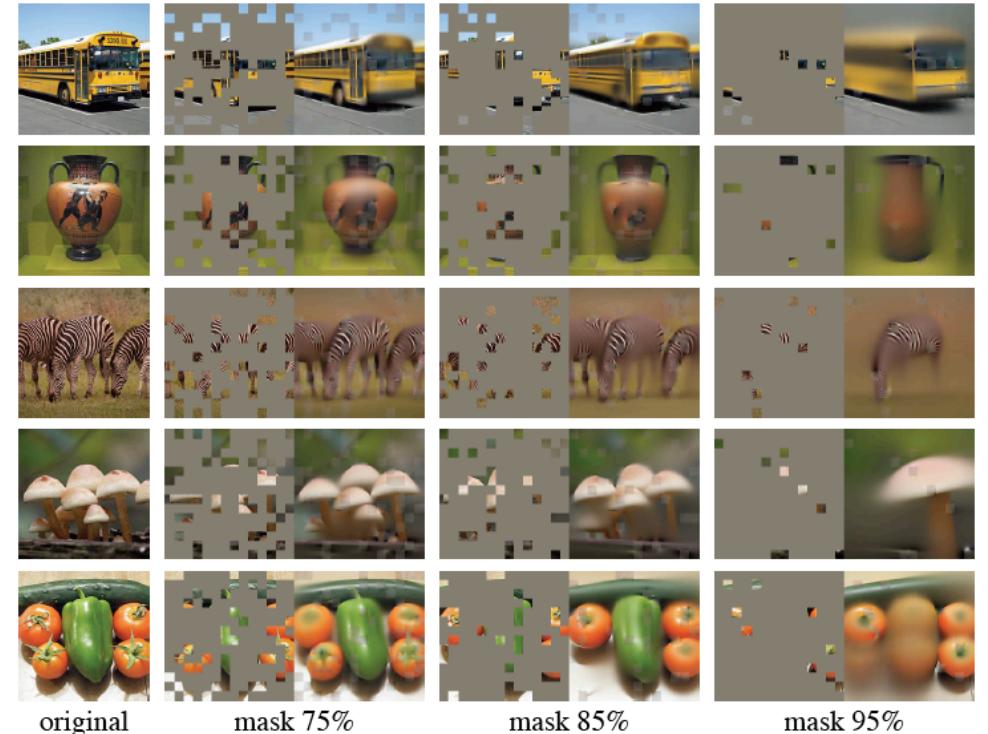
Figure 1: Comparison of image throughput and ImageNet-1K Top-1 accuracy. Throughput is measured on A100 GPU with batch size of 128.

# Masked Autoencoders



**Figure 1. Our MAE architecture.** During pre-training, a large random subset of image patches (e.g., 75%) is masked out. The encoder is applied to the small subset of *visible patches*. Mask tokens are introduced *after* the encoder, and the full set of encoded patches and mask tokens is processed by a small decoder that reconstructs the original image in pixels. After pre-training, the decoder is discarded and the encoder is applied to uncorrupted images (full sets of patches) for recognition tasks.

MSE loss between pixels for masked tokens only!



**Figure 4.** Reconstructions of ImageNet *validation* images using an MAE pre-trained with a masking ratio of 75% but applied on inputs with higher masking ratios. The predictions differ plausibly from the original images, showing that the method can generalize.

# Vision-Language Models: Frozen

Multimodal Few-Shot Learning with  
Frozen Language Models, 2021

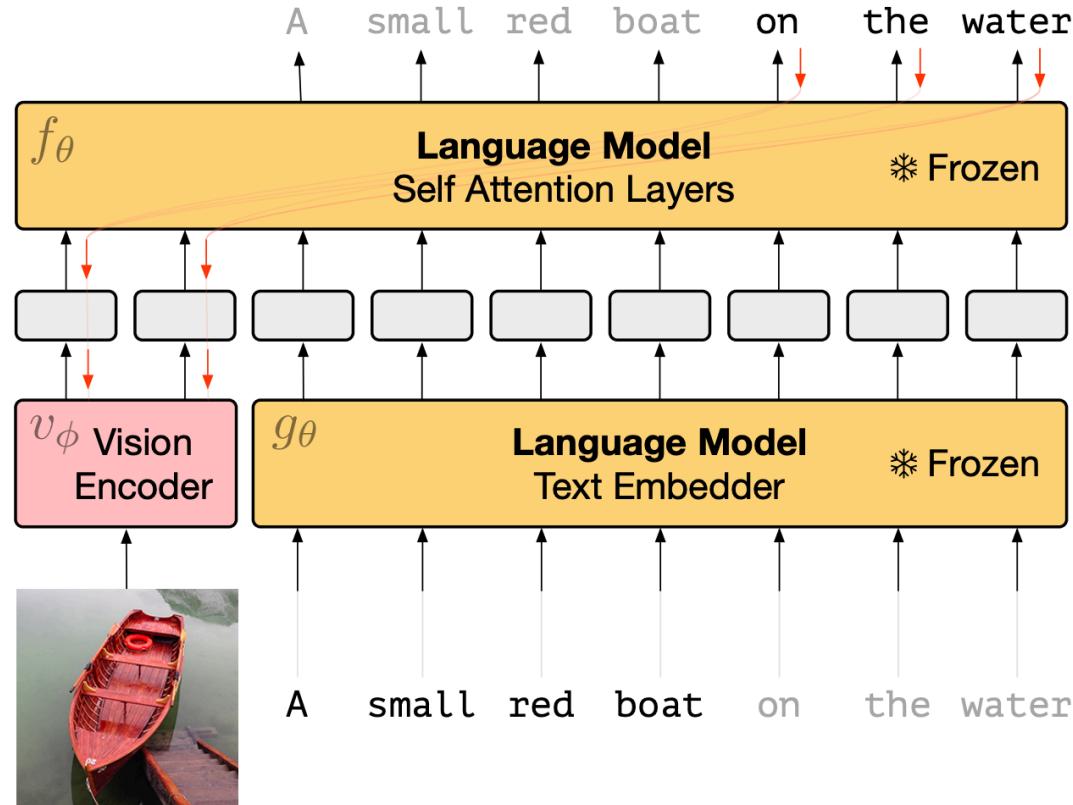


Figure 2: Gradients through a frozen language model's self attention layers are used to train the vision encoder.

# Vision-Language Models: Frozen

Multimodal Few-Shot Learning with  
Frozen Language Models, 2021

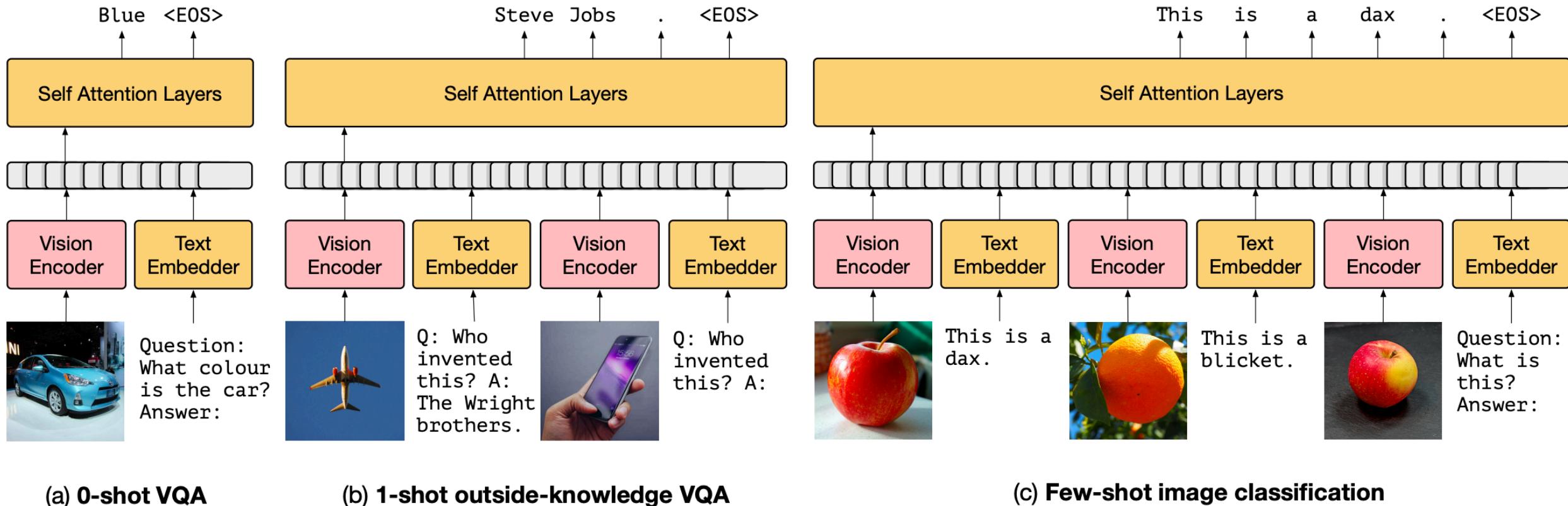
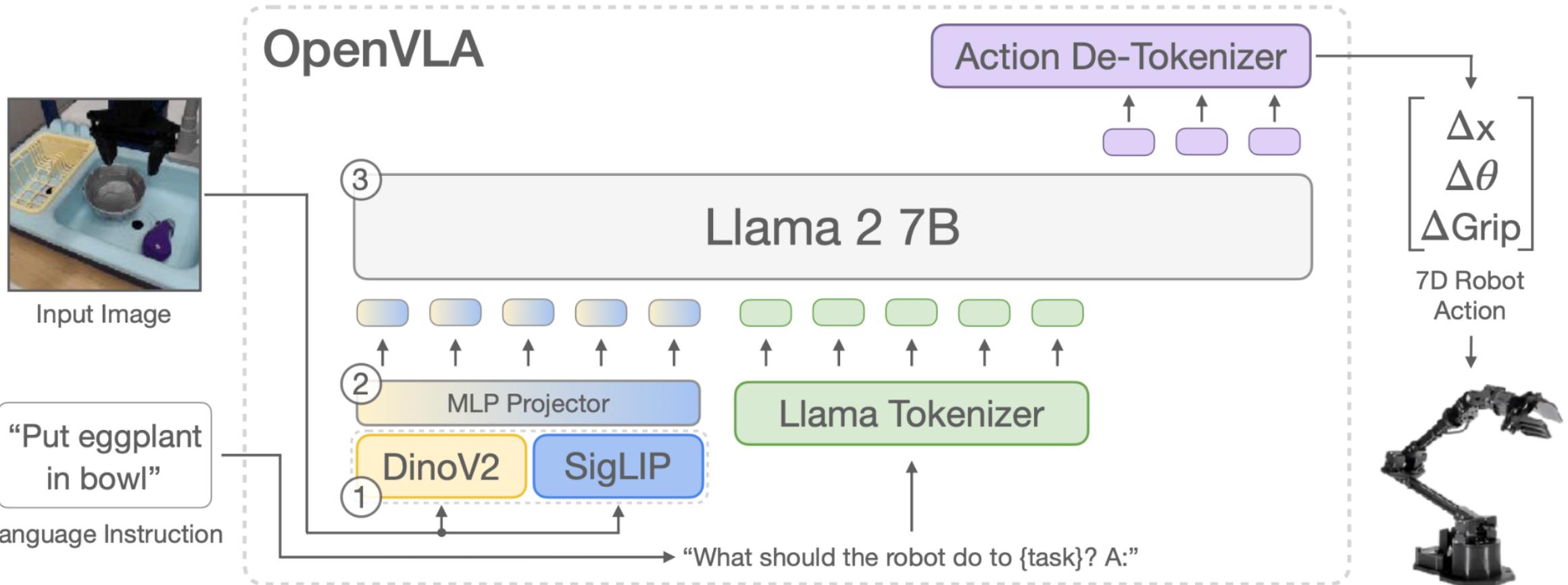


Figure 3: Inference-Time interface for *Frozen*. The figure demonstrates how we can support (a) visual question answering, (b) outside-knowledge question answering and (c) few-shot image classification via in-context learning.

# Vision-Language Action Models: OpenVLA



# Vision-Language Action Models: $\pi$



Videos at: <https://www.physicalintelligence.company/blog/pi05>

# Multimodal Models

## Gemini: A Family of Highly Capable Multimodal Models

Gemini Team, Google<sup>1</sup>

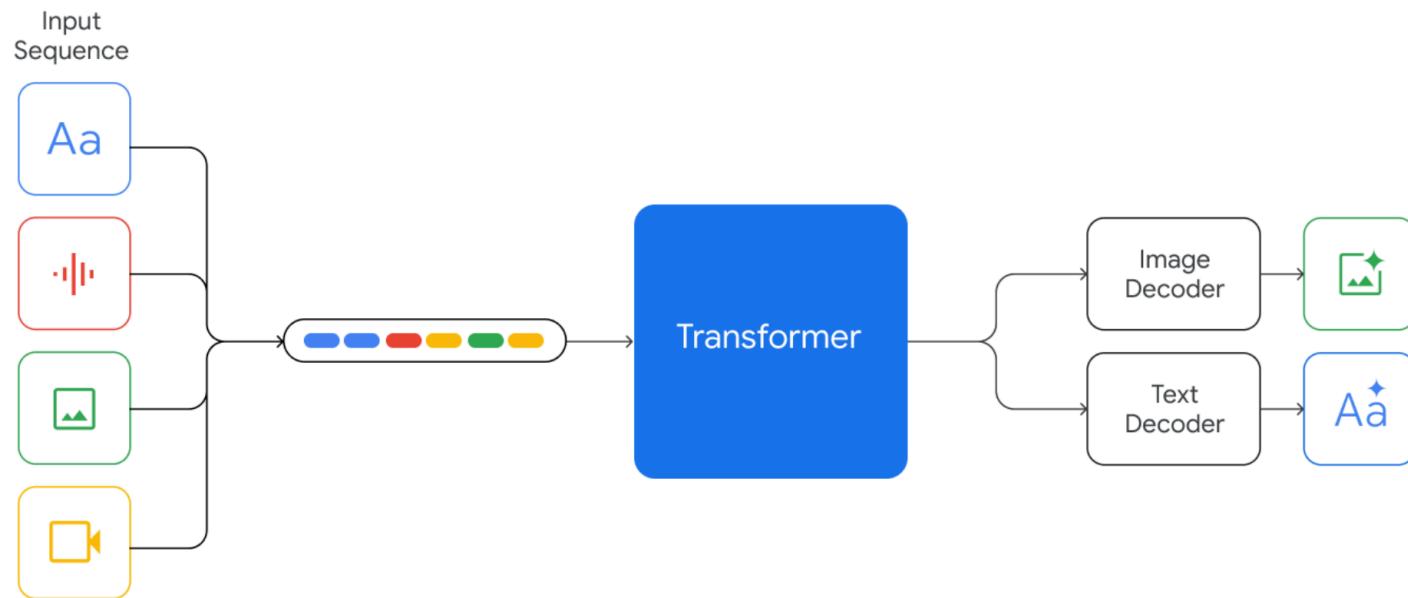
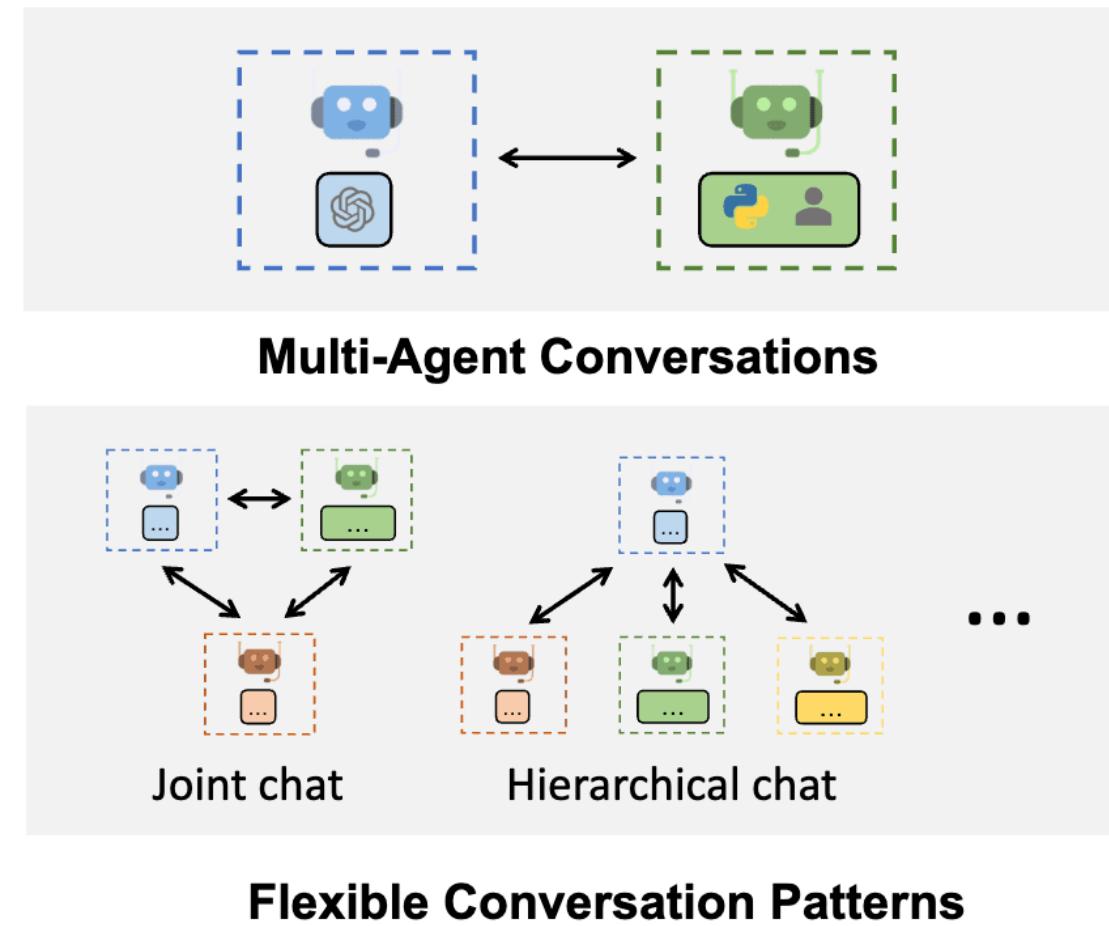
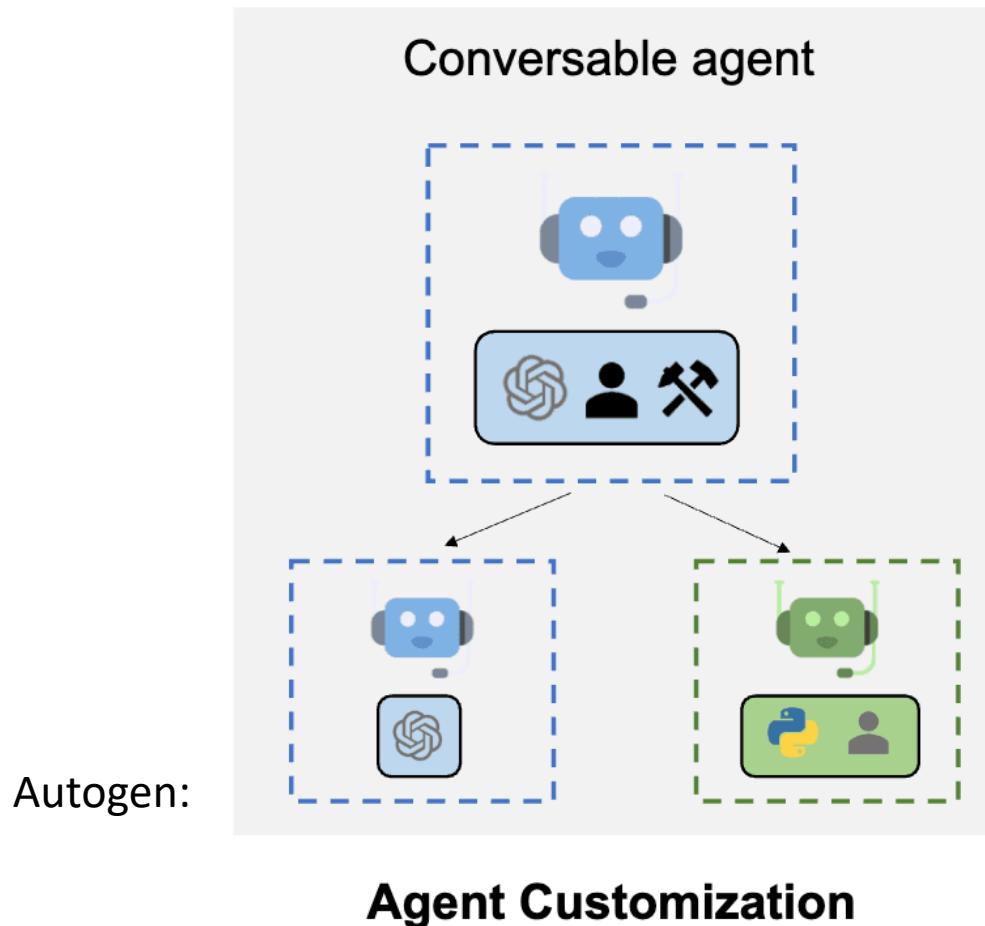


Figure 2 | Gemini models support interleaved sequences of text, image, audio, and video as inputs (illustrated by tokens of different colors in the input sequence). They can output responses with interleaved image and text.

# LLMs as Agents

- Autogen, LangChain, AutoGPT, Langroid, OpenAgents, ....



# LLMs as Agents: Example: Anthropic

<https://www.youtube.com/watch?v=vH2f7cjXjKI>