Mathematics of Data: From Theory to Computation

Prof. Volkan Cevher volkan.cevher@epfl.ch

Lecture 15: Language models

Laboratory for Information and Inference Systems (LIONS) École Polytechnique Fédérale de Lausanne (EPFL)

EE-556 (Fall 2024)

















License Information for Mathematics of Data Slides

▶ This work is released under a <u>Creative Commons License</u> with the following terms:

Attribution

► The licensor permits others to copy, distribute, display, and perform the work. In return, licensees must give the original authors credit.

Non-Commercial

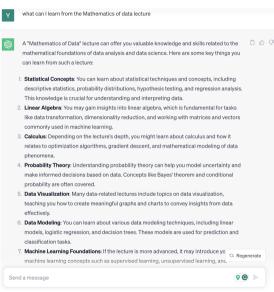
► The licensor permits others to copy, distribute, display, and perform the work. In return, licensees may not use the work for commercial purposes — unless they get the licensor's permission.

Share Alike

The licensor permits others to distribute derivative works only under a license identical to the one that governs the licensor's work.

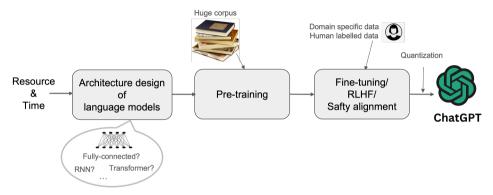
► Full Text of the License

The era of ChatGPT





Outline



This lecture

- 1. Basics of language models
- 2. Self-attention and transformer architectures
- 3. Fundamentals of pre-training & fine-tuning & reinforcement learning with human feedback (RLHF)
- 4. Generative pre-trained transformer (GPT) family
- 5. Basics of prompting: zero-shot & few-shot & chain-of-thought & automatic

A motivation for language models (LMs)

Example

Predict the next word w given the following source sentence S_{source} ?

 S_{source} : "On January 1 people usually say happy new [w]."

A motivation for language models (LMs)

Example

Predict the next word w given the following source sentence S_{source} ?

 S_{source} : "On January 1 people usually say happy new [w]."

Question:

- \circ Why is this important?
 - spelling & grammar correction

 $p(year|S_{source}) > p(years|S_{source})$ $p(S_{translation 1}|S_{source}) > p(S_{translation 2}|S_{source})$

machine translationsentence classification

 $p(S_{class 1}|S_{source}) > p(S_{class 2}|S_{source})$

speech recognition

 $\mathsf{p}(w|S_{\mathsf{source}})$

chatbot

 $p(w|S_{\mathsf{source}})$

(more generally) labeling, automated decisions,...

Basics for language models (LMs) - I

Definition (Language model [9])

Models that assign probabilities to sequences of words are called language models.

Remarks:

 \circ Given a sentence with T words: $S=w_{1:T}=(w_1,\ldots,w_T)$, by the chain rule of probability:

$$p(S) = p(w_{1:T}) = p(w_1)p(w_2|w_1)p(w_3|w_{1:2})\cdots p(w_T|w_{1:T-1}) = \prod_{t=1}^{T} p(w_t|w_{1:t-1})$$

 \circ Implicitly, we are enforcing a graphical model that takes "time" into account.

Example

If $S=w_{1:3}=$ "happy new year", then $\mathsf{p}(S)=\mathsf{p}(\mathsf{happy})\mathsf{p}(\mathsf{new}|\mathsf{happy})\mathsf{p}(\mathsf{year}|\mathsf{happy}|\mathsf{new}).$

Basics for language models (LMs) - II

Question:

 \circ How can we compute $p(w_t|w_{1:t-1})$?

Remarks:

o A trivial solution: Just count the frequency on a large corpus, e.g.,

$$\mathsf{p}(\mathsf{year}|S_\mathsf{source}) = \frac{\mathsf{p}(S_\mathsf{source} + \mathsf{year})}{\mathsf{p}(S_\mathsf{source})} \approx \frac{\#(\mathsf{On\ January\ 1\ people\ usually\ say\ happy\ new\ year)}}{\#(\mathsf{On\ January\ 1\ people\ usually\ say\ happy\ new)}}$$

- o But the language is creative, there are several ways to express the same meaning.
- o The sentence above might even not appear on the corpus.
- We need better ways to estimate such probabilities!

$N\operatorname{-gram}\ \mathbf{LMs}$

Markov assumption [12]

The probability of a word only depends on the last N-1 words as

$$p(w_t|w_{1:t-1}) = p(w_t|w_{t-N:t-1}) \approx \frac{\#(w_{t-N:t})}{\#(w_{t-N:t-1})}.$$



Markov in 1913 [12] used "Markov chains" to predict whether the upcoming letter would be a vowel or a consonant.

Example

In the bigram LM (N=2), we only need to estimate $p(w_t|w_{t-1}) \approx \frac{\#(w_{t-1}:t)}{\#(w_{t-1})}$ to generate text.

w_t											$__$				
w_{t-1}		i	want	to	eat	\imath	v_{t-1}					i	want	to	eat
	i	5	827	0	9	i war		eat		-1	i	0.002	0.33	0	0.0036
	want	2	0	608	1		2417		\rightarrow	v_t	want	0.0022	0	0.66	0.0011
	to	2	0	4	686	2533 927	2417	746		\mathcal{I}	to	0.00083	0	0.0017	0.28
	eat	0	0	2	0						eat	0	0	0.0027	0

Figure: Count (Left) and probability $p(w_t|w_{1:t-1})$ (Right) from the Berkeley Restaurant Project corpus of 9332 sentences [9].

Towards pre-training an N-gram LM

 \circ In natural language processing (NLP), we use tokens to represent words coming from a vocabulary ${\cal V}.$

 $\textbf{Terminologies:} \hspace{0.1in} \circ \hspace{0.1in} A \hspace{0.1in} \textit{token} \hspace{0.1in} \text{is the smallest unit that can be assigned a meaning to be processed.}$

- ▶ In English, a token often corresponds to a word.
- ▶ However, a single token can also encode compound words like *New York*.
- ▶ In Chinese or Japanese, there is no space between words.
- ▶ In these languages, sentence segmentation is required before we tokenize.
- \circ We indicate the beginning and the end of sentences with tokens $\langle \mathrm{BOS} \rangle$ and $\langle \mathrm{EOS} \rangle$.
 - ▶ S_{source} "⟨BOS⟩ Happy new year ⟨EOS⟩" has T=5 tokens.
- \circ The size of our vocabulary is denoted as $|\mathcal{V}|$.
- o Pre-training: building a LM based on a large corpus in a (often) self-supervised manner.
- o Inference: Using a trained LM to do next word prediction.

N-gram LMs: "Pre-training" & Inference

o The following simplified examples show the difficulty of pre-training and inference with 2-gram LMs.

"Pre-training"

- 1. Count $\#(w_{t-1})$ and $\#(w_{t-1:t})$ over the corpus.
- **2.** Obtain probability $p(w_t|w_{t-1})$ over the corpus.

Inference

- **1.** Set w_1 as $\langle BOS \rangle$, t = 1.
- 2. While True:
 - $\mathbf{v}_{t+1} = \arg\max_{w \in \mathcal{V}} \mathsf{p}(w|w_t)$
 - ▶ If w_{t+1} is $\langle EOS \rangle$: break
 - t = t + 1
- **3.** Output: $[w_1, \dots, w_{t+1}]$.

Remarks:

- \circ Need to store the probability for all N-gram pairs.
- \circ Language is creative, some new N-gram pairs might not even appear on the corpus.
- \circ Cannot incorporate earlier words than N due to the Markov assumption.

 $p(two \mid one plus one equals) = p(two \mid it is wrong that one plus one equals)?$

Word representations

Question: • How can we numerically represent a word/meaning?

Remarks: Osgood et al. 1957 [16] uses 3 numbers to represent a word.

▶ valence: the pleasantness of the stimulus

▶ arousal: the intensity of emotion provoked by the stimulus

▶ dominance: the degree of control exerted by the stimulus

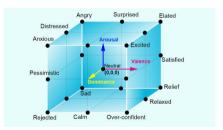
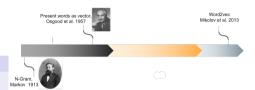


Figure: From [8].

Word embeddings



Definition (Word embeddings [9])

Vectors for representing words are called word embeddings.

- We will briefly introduce two words embeddings:
- \circ One-hot representation: sparse and long word embedding in $\mathbb{R}^{|\mathcal{V}|}$.
 - ► Training is not required—trivial to obtain.
 - Not a good way to capture the underlying meaning—cannot measure similarity.
- o Word2vec [13]: a framework to learn dense and concise word embedding.
 - Training is required.
 - Better characterization for the meaning of a word, e.g., the similarity can be computed by similarity metrics.
 - Cosine similarity or inner products work!

Word2vec [13]: Setup

 \circ An illustration of a target word and context words in a ± 2 window size:



- \circ Word2vec uses learnable parameters X_c and X_t to present two embeddings for each word,
 - $ightharpoonup X_c$ corresponds to the embedding when it is as a context word.
 - $lackbox{ } X_t$ corresponds to the embedding when it is as a target word
 - ► They satisfy the following relationship:

$$\mathbf{b}_i^t = \mathbf{X}^t \mathbf{e}_i \in \mathbb{R}^d, \quad \mathbf{b}_i^c = \mathbf{X}^c \mathbf{e}_i \in \mathbb{R}^d,$$

where $e_i \in \mathbb{R}^{|\mathcal{V}|}$ is the one hot representation for each word, $i \in {1, \ldots |\mathcal{V}|}$.

Remarks:

- o The window size for the context is a hyperparameter.
- o The final embedding can be the summation or concatenation of these two embeddings.

Word2vec [13]: Training

 \circ Core idea: Given a pair of words (w_i, w_j) , return the probability that w_j is the context word of w_i (i.e., true).

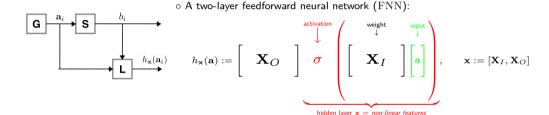
A simple approach: $p(\text{true}|(w_t, w_c)) = \sigma(\langle \mathbf{b}_t^t, \mathbf{b}_c^c \rangle) = \frac{1}{1 + \exp(-\langle \mathbf{b}_t^t, \mathbf{b}_c^c \rangle)}$, where σ is the sigmoid activation.

- \circ Given a tuple (w_t, w_c, w_n) , we have the following ingredients
 - $ightharpoonup w_t$ is the target word.
 - $ightharpoonup w_c$ is one of its context words(positive samples)
 - $\blacktriangleright w_n$ is not its context word (negative sample)—e.g., chosen via unigram (1-Gram) probability.
 - A loss function:

$$\begin{split} L &= -\log \left(\mathsf{p}(\mathsf{true}|(w_t, w_c)) \mathsf{p}(\mathsf{false}|(w_t, w_n)) \right) \\ &= -\log \mathsf{p}(\mathsf{true}|(w_t, w_c)) - \log \mathsf{p}(\mathsf{false}|(w_t, w_n)) \\ &= -\log \sigma(\langle \mathbf{b}_t^t, \mathbf{b}_c^c \rangle) - \log (1 - \sigma(\langle \mathbf{b}_t^t, \mathbf{b}_n^c \rangle)) \\ &= -\log \frac{1}{1 + \exp(-\langle \boldsymbol{X}^t \boldsymbol{e}_t, \boldsymbol{X}^c \boldsymbol{e}_c \rangle)} - \log \left(1 - \frac{1}{1 + \exp(-\langle \boldsymbol{X}^t \boldsymbol{e}_t, \boldsymbol{X}^c \boldsymbol{e}_n \rangle)} \right) \end{split}$$

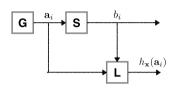
 \circ Crawl the corpus to obtain these tuples, and minimize L (e.g., with stochastic gradient descent).

Designing neural networks for pre-training LM



Short detour: Statistical learning with maximum-likelihood estimators

o A visual summary: From parametric models to learning machines



$$\begin{split} (\mathbf{a}_i,b_i)_{i=1}^n & \xrightarrow{\mathsf{modeling}} P(b_i|\mathbf{a}_i,\mathbf{x}) \xrightarrow{\mathsf{independency}} \mathbf{p}_{\mathbf{x}}(\mathbf{b}) := \prod_{i=1}^n P(b_i|\mathbf{a}_i,\mathbf{x}) \\ & \downarrow \mathsf{maximizing} \; \mathsf{w.r.t} \; \mathbf{x} \\ & \mathbf{a} \longrightarrow \mathsf{Learning} \; \mathsf{Machine} \longleftarrow & \mathbf{x}_\mathsf{ML}^\star \\ & \mathsf{prediction} \downarrow \\ & h_{\mathbf{x}_\mathsf{ML}^\star}(\mathbf{a}) \end{split}$$

Observations:

- $\circ \; \mathsf{Recall} \; \mathbf{x}^{\star}_{\mathsf{ML}} \in \arg \min_{\mathbf{x} \in \mathcal{X}} \; \{ L(h_{\mathbf{x}}(\mathbf{a}), \mathbf{b}) := -\log \mathsf{p}_{\mathbf{x}}(\mathbf{b}) \}.$
- \circ Maximizing $p_{\mathbf{x}}(\mathbf{b})$ gives the maximum-likelihood (ML) estimator.
- o Maximizing $p_{\mathbf{x}}(\mathbf{b})$ and minimizing $-\log p_{\mathbf{x}}(\mathbf{b})$ result in the same solution set.

Designing neural networks for pre-training LM

o A two-layer feedforward neural network (FNN):



Recall: Maximum-likelihood estimator

The maximum-likelihood estimator (supervised learning with data (\mathbf{a},b)) is given by

$$\mathbf{x}^{\star} \in \arg\min_{\mathbf{x} \in \mathcal{X}} \left\{ L(h_{\mathbf{x}}(\mathbf{a}), b) := -\log \mathsf{p}_{\mathbf{x}}(b) \right\}.$$

Remark: • NN-based LM can be considered as an unsupervised maximum-likelihood estimator.

$$\mathbf{x}_{\mathsf{LM}}^{\star} \in \arg\min_{\mathbf{x} \in \mathcal{X}} -\log \mathsf{p}_{\mathbf{x}}(S) = -\log \mathsf{p}_{\mathbf{x}}(\mathbf{b}_{1:T}),$$

hidden layer z = non-linear features

where $p_{\mathbf{x}}(S)$ is the probability of sentence S with embedding $\mathbf{b}_{1:T} = (\mathbf{b}_1, \dots, \mathbf{b}_T)$.

The optimization objective

 \circ A (vector-output) neural network $\mathbf{h}_{\mathbf{x}} \in \Delta^{|\mathcal{V}|-1}$ can be used to model such probability.

$$\begin{split} &-\log \mathbf{p}_{\mathbf{x}}(\mathbf{b}_{1:T}) = -\log \left(\prod_{t=1}^{T} \mathbf{p}_{\mathbf{x}}(\mathbf{b}_{t}|\mathbf{b}_{1:t-1}) \right) = \sum_{t=1}^{T} \left(-\log \underbrace{\mathbf{p}_{\mathbf{x}}(\mathbf{b}_{t}|\mathbf{b}_{1:t-1})}_{\mathbf{h}_{\mathbf{x}}(\mathbf{b}_{1:t-1})^{["\mathbf{b}_{t}"]}} \right) \\ &= \sum_{t=1}^{T} \left(-\log \mathbf{h}_{\mathbf{x}}(\mathbf{b}_{1:t-1})^{["\mathbf{b}_{t}"]} \right) = \sum_{t=1}^{T} \left(-\sum_{i=1}^{|\mathcal{V}|} \hat{\mathbf{u}}_{t}^{[i]} \log \mathbf{u}_{t}^{[i]} \right) = \text{cross entropy loss} \end{split}$$

- $\mathbf{u}_t := \mathbf{h}_{\mathbf{x}}(\mathbf{b}_{1:t-1}) \in \mathbb{R}^{|\mathcal{V}|}$ is the probability distribution of the next word given previous t-1 words.
- $\hat{\mathbf{u}}_t \in \mathbb{R}^{|\mathcal{V}|}$ is the correct distribution (one-hot) at t step.

Remarks:

- o Teacher forcing training: We always give the model the correct history sequence.
- o Auto-regressive inference: The history sequence comes from its prediction result.

Basic NN setups for LM

- o Below, we present a general idea of deploying neural networks as LMs.
 - ► Feed-forward neural network (FNN)
 - ► Recurrent Neural Networks (RNN)
 - Self-attention
- \circ At each step t, we use NN to model the probability distribution of the current word given previous t-1 words.



 \circ Then, we can minimize the cross-entropy loss (i.e., $-\sum_{i=1}^{|\mathcal{V}|} \hat{\mathbf{u}}_t^{[i]} \log \mathbf{u}_t^{[i]}$) via (stochastic) gradient descent.

FNN as LM [2]: pre-training

- \circ Core idea: use most recent N tokens to predict next token (similar to N-gram)
- o: $\mathbf{X}_I \in \mathbb{R}^{m \times Nd}$, $\mathbf{X}_O \in \mathbb{R}^{|\mathcal{V}| \times m}$ are learnable parameters, where d is the dimension of the embedding.

Forward pass in pre-training on single sentence (only use two recent tokens, i.e., N=2)

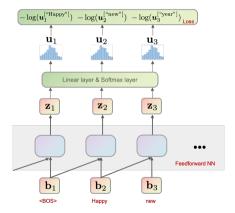
- 1. Set $\mathbf{b}_0 = \mathbf{0}$, initial loss L = 0
- **2.** For t = 1, ..., T

FNN

 $ightharpoonup \mathbf{u}_t = \mathsf{Softmax}(\mathbf{X}_O \mathbf{z}_t),$

probability

$$L+ = \left(\sum_{i=1}^{|\mathcal{V}|} -\hat{\mathbf{u}}_t^{[i]} \log \mathbf{u}_t^{[i]}\right),$$
 loss



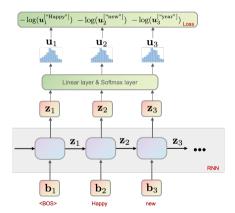
RNN as LM [14]: pre-training

- \circ A weakness of FNN LM is the Markov assumption: It cannot capture long-term dependencies.
- \circ RNN architectures only partially address this issue.
- o: $\mathbf{X}_1 \in \mathbb{R}^{m \times m}, \mathbf{X}_2 \in \mathbb{R}^{m \times d}, \mathbf{X}_O \in \mathbb{R}^{|\mathcal{V}| \times m}$ are learnable parameters.

Forward pass in pre-training on single sentence

- **1.** Set initial state $\mathbf{z}_0 = \mathbf{0}$, initial loss L = 0
- **2.** For t = 1, ..., T

 - $lackbox{f u}_t = \mathsf{Softmax}(\mathbf{X}_O\mathbf{z}_t),$ probability



RNN as LM: inference

 $\circ~RNN$ architectures perform auto-regressive inference.

Forward pass in inference

- 1. Set \mathbf{b}_1 as the embedding of $\langle \mathrm{BOS} \rangle$, t=1, initial state $\mathbf{z}_0=\mathbf{0}$.
- 2. While True:

 - $\mathbf{u}_t = \mathsf{Softmax}(\mathbf{X}_O \mathbf{z}_t)$
 - Set \mathbf{b}_{t+1} as the embedding of the token corresponding to $\arg \max \mathbf{u}_t$.
 - ▶ If \mathbf{b}_{t+1} is the embedding of $\langle EOS \rangle$: break
 - t + 1
- **3.** Output: $[\mathbf{b}_1, \dots, \mathbf{b}_{t+1}]$.

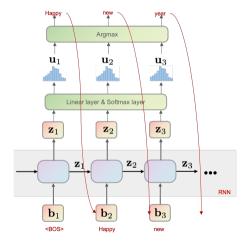


Figure: Auto-regressive inference

Self-attention layer as LM

- \circ A weakness of the RNN LMs is its recursive non-parallelizable computation.
- o Self-attention can address these issues.

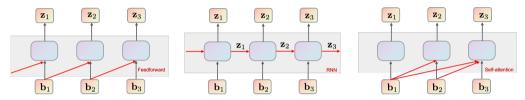
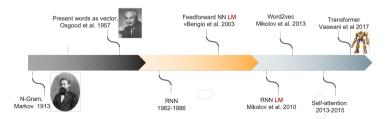


Figure: (Left panel) FNN in LM. (Middle panel) RNN in LM. (Right panel) Self-attention in LM.



Self-attention layer for LM

- o Core idea: compare a word of interest to other words based on their relevance.
- o How do we measure the relevance of two words?
 - inner products (recall word embeddings)
 - ightharpoonup e.g., for the word with embedding b_3 , we can compute three scores:

$$\mathsf{Score}(3,1) = \langle \mathbf{b}_3, \mathbf{b}_1 \rangle; \quad \mathsf{Score}(3,2) = \langle \mathbf{b}_3, \mathbf{b}_2 \rangle; \quad \mathsf{Score}(3,3) = \langle \mathbf{b}_3, \mathbf{b}_3 \rangle.$$

o Next, we normalize them with a softmax to create a vector of weights, and obtain the output:

$$\begin{split} \mathbf{z}_3 &= \sum_{j=1}^3 \mathsf{Softmax}([\mathsf{Score}(3,1),\mathsf{Score}(3,2),\mathsf{Score}(3,3)])_j \mathbf{b}_j \\ &= \sum_{j=1}^3 \frac{\exp(\mathsf{Score}(3,j))}{\sum_{i=1}^3 \exp(\mathsf{Score}(3,i))} \mathbf{b}_j \end{split}$$

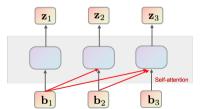


Figure: Self-attention layer.

Self-attention layer for LM

- o A more sophisticated way to present how words are contributed to each other:
 - Query: when current word goes measure the relevance with other words.
 - ► Key: when being measured the relevance by other words.
 - ▶ Value: value used to compute the final output.



 \circ For each word, calculate its corresponding query, key, and value using parameters $\mathbf{X}_Q, \mathbf{X}_K, \mathbf{X}_V \in \mathbb{R}^{m imes d}$

$$q_i = \mathbf{X}_Q \mathbf{b}_i, k_i = \mathbf{X}_K \mathbf{b}_i, v_i = \mathbf{X}_V \mathbf{b}_i.$$

 \circ Then, for the word with embedding \mathbf{b}_3 , those three scores become:

$$\mathsf{Score}(3,1) = \langle \boldsymbol{q}_3, \boldsymbol{k}_1 \rangle; \quad \mathsf{Score}(3,2) = \langle \boldsymbol{q}_3, \boldsymbol{k}_2 \rangle; \quad \mathsf{Score}(3,3) = \langle \boldsymbol{q}_3, \boldsymbol{k}_3 \rangle.$$

$$\mathbf{z}_3 = \sum_{j=1}^3 \mathsf{Softmax}([\mathsf{Score}(3,1),\mathsf{Score}(3,2),\mathsf{Score}(3,3)])_j v_j$$

o We need to learn the parameters $\mathbf{X}_Q, \mathbf{X}_K, \mathbf{X}_V \in \mathbb{R}^{m \times d}$.

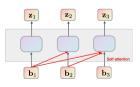


Figure: Self-attention layer.

Question: Ones self-attention layer consider the relative position of each word in the sequence? No!

Observation: \circ If we switch the order of \mathbf{b}_1 and \mathbf{b}_2 , the output \mathbf{z}_3 remains the same.

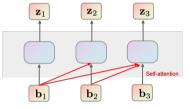


Figure: Self-attention layer.

o In comparison, RNN encodes the information about the order of the inputs recursively.

Question: Opes self-attention layer consider the relative position of each word in the sequence? No!

 $\textbf{Solution 1?} \qquad \circ \ \, \text{Absolute position via trivial concatenation}.$

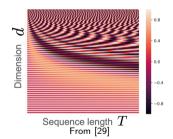
$$\mathsf{Pos}(\mathbf{b}_t) = \mathsf{Concatenate}[\mathbf{b}_t, t]$$
 .

- Unbounded value.
- o Hard to extrapolate on sequence with unseen length.

Question: Ones self-attention layer consider the relative position of each word in the sequence? No!

Solution 2 [24]: \circ Absolute position via trigonometric functions of different frequencies. For $t=1,\ldots,T$:

$$\mathsf{Pos}(\mathbf{b}_t) = \mathbf{b}_t + \begin{pmatrix} \sin\left(t/10000^{2\times 1/d}\right) \\ \cos\left(t/10000^{2\times 1/d}\right) \\ \vdots \\ \sin\left(t/10000^{2\times \frac{d}{2}/d}\right) \\ \cos\left(t/10000^{2\times \frac{d}{2}/d}\right) \end{pmatrix}$$



Question: One self-attention layer consider the relative position of each word in the sequence? No!

Solution 2 [24]: \circ Absolute position via trigonometric functions of different frequencies. For $t=1,\ldots,T$:

Solution 3: o *Rotary position embedding [21]: incorporate both absolute position and relative position.

Self-attention layer for LM

- $\mathbf{b} = [\mathbf{b}_1, \mathbf{b}_T]^{\top} \in \mathbb{R}^{T \times d}$: collections of embeddings of all tokens.
- \circ Learnable parameters: $\mathbf{X}_O, \mathbf{X}_K, \mathbf{X}_V \in \mathbb{R}^{m \times d}, \mathbf{X}_O \in \mathbb{R}^{|\mathcal{V}| \times m}$.

Forward pass in training on a single sentence

- 1. Set initial loss L=0.
- **2.** $Q = \mathbf{B} \mathbf{X}_O^\top, K = \mathbf{B} \mathbf{X}_K^\top, V = \mathbf{B} \mathbf{X}_V^\top$, query, key, value.
- 3. $S = \mathsf{Mask}(QK^\top)$, calculate score and mask score.
- **5.** $\mathbf{Z} := [\mathbf{z}_1, ..., \mathbf{z}_T]^{\top} = \text{Row-wise-Softmax}(S)V$, self-attention output
- **6.** $\mathbf{U} := [\mathbf{u}_1, ..., \mathbf{u}_T]^{\top} = \text{Row-wise-Softmax}(\mathbf{Z}\mathbf{X}_{\Omega}^{\top}), \text{ probability}$

7.
$$L = L + \left(\sum_{t=1}^{T} \sum_{i=1}^{|\mathcal{V}|} -\hat{\mathbf{u}}_{t}^{[i]} \log \mathbf{u}_{t}^{[i]}\right)$$
, loss

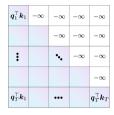


Figure: Mask score for S.

Remarks:

- \circ In the remaining slide, \mathbf{b}_t has already been added to position embedding.
- o Masking score is used to prevent "cheating."
 - the current word has only seen previous word.
 - the subsequent word is unknown.
 - \blacktriangleright the element $-\infty$ after softmax becomes 0.
- o Attention with masking score is usually called "Masked attention."
- \circ This construction enables parallelization whereby improving upon RNNs.

Self-attention layer as LM: inference

Forward pass in inference

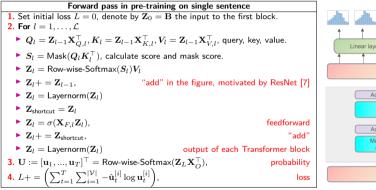
- **1.** Set \mathbf{b}_1 as the embedding of $\langle BOS \rangle$, t = 1.
- 2. While True:
 - $ightharpoonup q_t = \mathbf{X}_Q \mathbf{b}_t, k_t = \mathbf{X}_K \mathbf{b}_t, v_t = \mathbf{X}_V \mathbf{b}_t$, calculate query, key, value
 - $s = [\langle q_t, k_1 \rangle, \cdots, \langle q_t, k_t \rangle]^{\top}$, calculate score
 - $ightharpoonup \mathbf{z}_t = [oldsymbol{v}_1, \cdots, oldsymbol{v}_t] \cdot \mathsf{Softmax}(oldsymbol{s})$
 - $\mathbf{u}_t = \mathsf{Softmax}(\mathbf{X}_O \mathbf{z}_t)$
 - ightharpoonup Set \mathbf{b}_{t+1} as the embedding of the token corresponding to $\arg\max \mathbf{u}_t$.
 - ▶ If \mathbf{b}_{t+1} is the embedding of $\langle BOS \rangle$: break
 - t + 1
- **3.** Output: $[\mathbf{b}_1, \dots, \mathbf{b}_{t+1}]$.

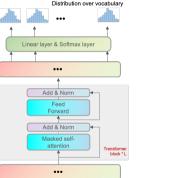
Remark:

 \circ Still non-parallelizable, still auto-regression, the same as RNN LM, FNN LM.

TRANSFORMER as LM

- \circ A Transformer block= [self-attention layer + *layer normalization + feedforward layer + *layer normalization].
- \circ We stack $\mathcal L$ Transformer blocks to form an LM, e.g., $\mathcal L=12$ in [17].

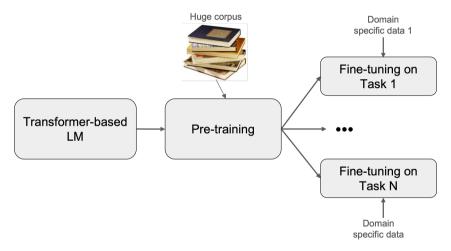




Remarks: o Original Transformer is proposed with encoder and decoder for neural machine translation [24].

The Transformer decoder is sufficient as an LM.

GPT-1 [17]: (Pre-train + fine-tune paradigm) "Improving Language Understanding by Generative Pre-Training", 2018

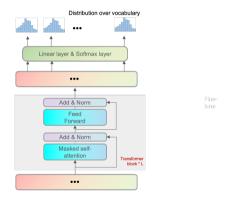


Remarks: • Pre-training enables learning better underlying language patterns on a large corpus.

o Pre-training provides a better parameter initialization for fine-tuning, leading to faster convergence.

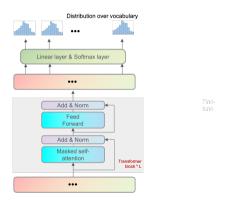
GPT-1 [17]: (Pre-train + fine-tune paradigm) "Improving Language Understanding by Generative Pre-Training", 2018

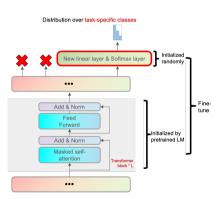
- o Step 1: Pre-train a LM on a large unlabeled corpus using Transformer's decoder.
 - ▶ Recall that Transformer's decoder is sufficient for LM.



GPT-1 [17]: (Pre-train + fine-tune paradigm) "Improving Language Understanding by Generative Pre-Training", 2018

- o Step 1: Pre-train a LM on a large unlabeled corpus using Transformer's decoder.
 - ▶ Recall that Transformer's decoder is sufficient for LM.
- o Step 2: Fine-tune on specific tasks, e.g., on a sentence classification task.





GPT-1 "Improving Language Understanding by Generative Pre-Training", 2018" [17]

Limitation:

- Require task-specific datasets and task-specific fine-tuning.
- o Model is fine-tuned on very narrow task distributions.
- o Model does not necessarily generalize better out-of-distribution.

Question:

- o Is it possible to address these limitations?
 - ▶ Humans do not require large supervised datasets to learn most new language tasks.
 - -"please tell me if this sentence describes something happy or something sad"

GPT-2, GPT-3 [18, 3] "Language Models are Unsupervised Multitask Learners", "Language Models are Few-Shot Learners"

- o Same as GPT-1: we still pre-train the LM on unlabeled corpus.
- o New: no need to fine-tune anymore. One pre-trained LM for all tasks, achieve SOTA.





GPT-2, GPT-3 [18, 3] "Language Models are Unsupervised Multitask Learners", "Language Models are Few-Shot Learners"

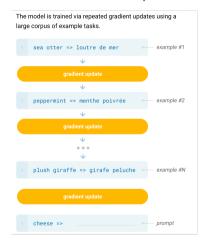
- o Same as GPT-1: we still pre-train the LM on unlabeled corpus.
- o New: no need to fine-tune anymore. One pre-trained LM for all tasks, achieve SOTA.



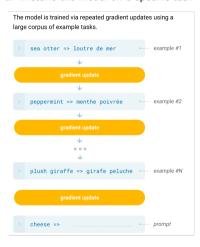
o How?

Model	Launch Year	Training Data	Training Parameters	Attention Layers	Word Embedding	Attention Heads
GPT-1	2018	7000 Books ~5GB	117M	12	768	12
GPT-2	2019	8 million documents ~40GB	1.5B	48	1600	48
GPT-3	2020	Multiple Source ~45TB	175B	96	12288	96

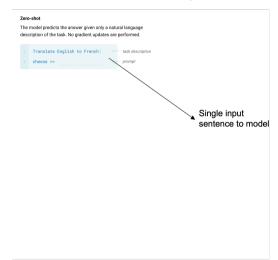
o GPT-1: finetune the model on a specific task



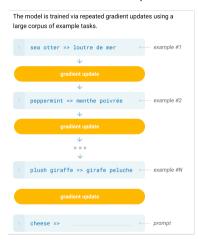
o GPT-1: finetune the model on a specific task



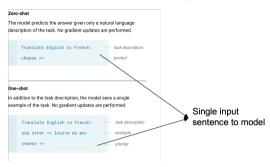
o GPT-3: no need to fine-tune on a specific task.



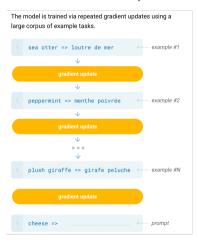
o GPT-1: finetune the model on a specific task



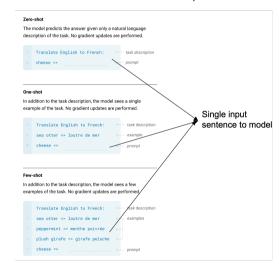
o GPT-3: no need to fine-tune on a specific task.



o GPT-1: finetune the model on a specific task

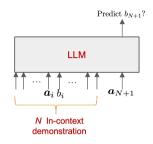


o GPT-3: no need to fine-tune on a specific task.



* Intrinsic mechanism of in-context learning

- o Classical supervised learning, for each task:
 - ▶ Training: ERM, $\mathbf{x}^* \in \arg\min_{\mathbf{x}} : R(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L(h_{\mathbf{x}}(\mathbf{a}_i), b_i).$
 - Given test point \mathbf{a}_{N+1} , return $h_{\mathbf{x}^{\star}}(\mathbf{a}_{N+1})$.
- o In-context learning in LLM:
 - Trains on huge corpus and gets x*.
 - $\qquad \qquad \textbf{For any test task } \tau', \text{ return: } h_{\mathbf{x}^{\star}}(\{\mathbf{a}_i^{\tau'}, b_i^{\tau'}\}_{i=1}^N, \mathbf{a}_{N+1}^{\tau'}).$



- o Meta-learning:
 - ► Training: ERM, $\mathbf{x}^{\star} \in \arg\min_{\mathbf{x}} : R(\mathbf{x}) = \frac{1}{\#\text{Training tasks}} \sum_{\tau=1}^{\#\text{Training tasks}} L(h_{\mathbf{x}}(\{\mathbf{a}_i^{\tau}, b_i^{\tau}\}_{i=1}^{N}, \mathbf{a}_{N+1}^{\tau}), b_{N+1}^{\tau}).$
 - For any test task τ' , return $h_{\mathbf{x}^{\star}}(\{\mathbf{a}_i^{\tau'}, b_i^{\tau'}\}_{i=1}^N, \mathbf{a}_{N+1}^{\tau'})$.

Question

- Why LLM can perform in-context learning?
- o Theoretical explanation: by performing gradient descent implicitly on in-context data? [25], see supplementary.
- o Or is it just due to the so-called emergent abilities? See the next slide!

Few-shot prompting (In-context learning) \rightarrow emergent abilities of LLM

"An ability is emergent if it is not present in smaller models but is present in larger models."[27]

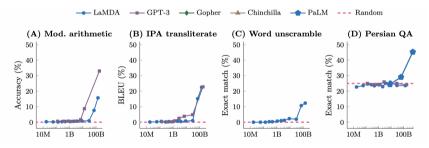


Figure: Emergent abilities of few-shot prompting appear when the model parameters (x-axis) are increased to some extent. [27]

Chain-of-thought prompting → emergent abilities of LLM

(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The answer is 8. X

(b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are 16 / 2 = 8 golf balls. Half of the golf balls are blue. So there are 8 / 2 = 4 blue golf balls. The answer is 4.

Figure: Demo of chain-of-thought (CoT) prompting [10].

Chain-of-thought prompting \rightarrow emergent abilities of LLM

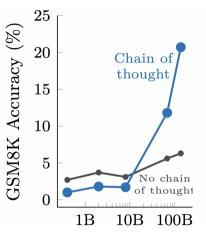


Figure: Performance under chain-of-thought prompting is increased until a certain model scale on Math word problems [27], A LLM called LaMDA is used [22].

Zero-shot chain-of-thought prompting \rightarrow emergent abilities of LLM

(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The answer is 8. X

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) 8 X

(b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5+6=11. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are 16/2 = 8 golf balls. Half of the golf balls are blue. So there are 8/2 = 4 blue golf balls. The answer is 4.

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: Let's think step by step.

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

Figure: Demo of zero-shot chain-of-thought (CoT) prompting [10].

Zero-shot chain-of-thought prompting \rightarrow emergent abilities of LLM

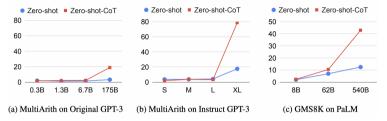


Figure: Performance under Zero-shot chain-of-thought prompting is increased until a certain model scale on Math word problems [10]. Y-axis indicates the accuracy.

Why emergent abilities occur?

- o Understanding this would benefit:
 - ▶ Economy and environment: reduce training cost to obtain desired emergent abilities.
 - Al-Safety: prevent larger models from acquiring dangerous capabilities without warning.



NeurIPS Conference @NeurIPSConf · Dec 12

Replying to @NeurlPSConf

Test of Time

Distributed Representations of Words and Phrases and their Compositionality

Outstanding Main Track Papers

Privacy Auditing with One (1) Training Run

Are **Emergent Abilities** of Large Language Models a Mirage?

Figure: In NeurIPS 2023, the paper that explains "emergent abilities" achieved outstanding paper award. As a remark: the Word2vec paper achieved "Test of time" award.

Why emergent abilities occur? [19]

- \circ Emergent abilities occur due to:
 - fundamental changes by model scaling. X
 - ► researcher's choice of metric. ✓
- o Nonlinear or discontinuous metrics produce apparent emergent abilities.
- o Linear or continuous metrics produce smooth, continuous, predictable changes in performance.

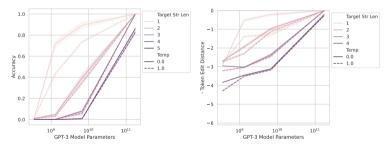
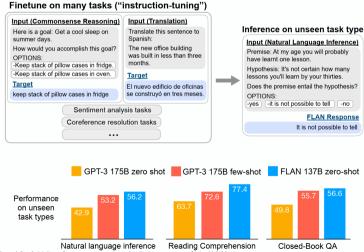


Figure: 2-Integer 2-Digit Multiplication Task. Left: performance is measured by a nonlinear metric (e.g., Accuracy). Right: performance is instead measured by a linear metric (e.g., Token Edit Distance).

FLAN [26] (Instruction tuning) "Finetuned language models are zero-shot learners", 2021

- o Fine-tuning is useful again, with the instruction format, allowing generalize to unseen tasks.
- o Better than aforementioned "no fine-tuning" + "few-shot prompting".



- o Limitation of GPT-3: trained to predict the next token, can not follow user instructions well.
- InstructGPT:
 - ► Towards following user instructions, more truthful, less toxic.
 - Align LM with user intent by fine-tuning via reinforcement learning from feedback (RLHF).

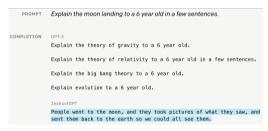


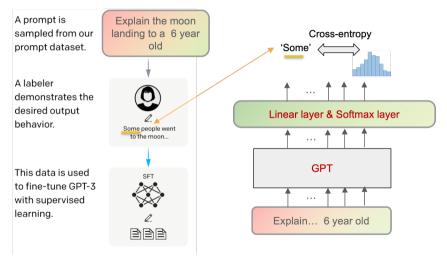
Figure: Demo of aligned InstructGPT, from https://openai.com/research/instruction-following



Figure: Demo of aligned LLM, from [4]

o Step 1: Pre-train a Transformer-based LM based on unlabeled corpus, similar to GPT-1, GPT-2, GPT-3.

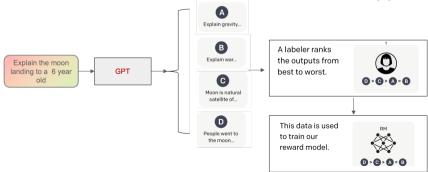
• Step 2: Supervised fine-tune via collected demonstration.



- o Step 3: Train a reward model $r_{\mathbf{x}}(S_{\text{prompt}}, S_{\text{response}})$ with parameters \mathbf{x} .
 - GPT-3-based architecture.
 - ▶ Input: concatenation of S_{prompt} and $S_{response}$. Output: scalar value.
 - Loss:

$$L_{\mathbf{x}} = -\frac{1}{\binom{K}{2}} E_{\left(S_{\mathsf{prompt}}, S_{\mathsf{response1}}, S_{\mathsf{response2}}\right) \sim D} \left[\log \left(\sigma \left(r_{\mathbf{x}} \left(S_{\mathsf{prompt}}, S_{\mathsf{response1}} \right) - r_{\mathbf{x}} \left(S_{\mathsf{prompt}}, S_{\mathsf{response2}} \right) \right) \right) \right],$$

where $S_{\text{response1}}$ is the preferred response out of the pair of $S_{\text{response1}}$ and $S_{\text{response2}}$, D is the dataset of human comparisons. For each prompt, labelers need to rank K response, leading to $\binom{K}{2}$ comparison.



- o Step 4: Using this reward model to fine-tune the GPT based on Proximal Policy Optimization (PPO)[20]
 - ightharpoonup (state, action): (S_{prompt} , $S_{response}$).
 - Initialize a policy to be the fine-tuned GPT in step 2, i.e., π^{SFT} .
 - Initialize a copy of the above policy with parameters ϕ that we want to optimize, i.e., π_{ϕ}^{RL} .
 - \blacktriangleright Use PPO to optimize ϕ in order to maximize the following full reward.

$$R(S_{\mathsf{prompt}}, S_{\mathsf{response}}) = r_{\mathbf{x}}(S_{\mathsf{prompt}}, S_{\mathsf{response}}) - \underbrace{\beta \log[\pi_{\phi}^{\mathsf{RL}}(S_{\mathsf{response}}|S_{\mathsf{prompt}})/\pi^{\mathsf{SFT}}(S_{\mathsf{response}}|S_{\mathsf{prompt}})]}_{\mathsf{penalty term}}$$

- The penalty term is the conditional relative entropy, ensuring the new policy π_{ϕ}^{RL} doesn't change a lot from the original policy π^{SFT} , which the reward model has seen during its training.
- See EE-618 for details!

GPT-4 [15]

- \circ Multi-modals GPT: text + image.
- o Predictable Scaling.

GPT-4 Technical Report

OpenAI*

Abstract

We report the development of GPT-4, a large-scale, multimodal model which can accept image and text inputs and produce text outputs. While lies capable than humans in many real-world scenarios, GPT-4 ethibis human-level performance on various professional and academic bear chambs; instuding passing a stimulated bar exam with a score around the top 10% of test takers. GPT-4 is a Transformer-land to the control of the contro

GPT-4 visual input example, Extreme Ironing:

User What is unusual about this image?



Source: https://www.barnorama.com/wp-content/uploads/2016/12/03-Confusing-Pictures.jpg

GPT-4 The unusual thing about this image is that a man is ironing clothes on an ironing board attached to the roof of a moving taxi.

Figure: From [15]

Predictable Scaling in GPT-4

- o It is expensive and time-consuming to do model-specific tuning for large language model (LLM).
- o Developers of GPT-4 can make loss prediction by power laws shortly after the training starts.

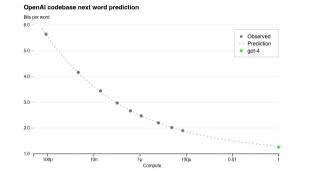


Figure 1. Performance of GPT-4 and smaller models. The metric is final loss on a dataset derived from our internal codebase. This is a convenient, large dataset of code tokens which is not contained in the training set. We chose to look at loss because it tends to be less noisy than other measures across different amounts of training compute. A power law fit to the smaller models (excluding GPT-4) is shown as the dotted line; this fit accurately predicts GPT-4's final loss. The x-axis is training compute normalized so that GPT-4 is 1

Figure: From [15]

Predictable Scaling in GPT-4

- o Developers of GPT-4 can make capability prediction by power laws shortly after the training starts.
- Measured by the ability to correctly synthesize Python functions on HumanEval dataset [5].

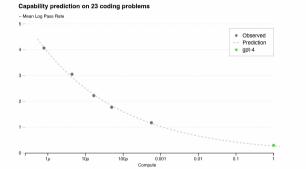
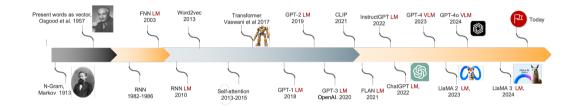


Figure 2. Performance of GPT-4 and smaller models. The metric is mean log pass rate on a subset of the HumanEval dataset. A power law fit to the smaller models (excluding GPT-4) is shown as the dotted line; this fit accurately predicts GPT-4's performance. The x-axis is training compute normalized so that GPT-4 is 1

Figure: From [15]



Wrap up!

o That's it folks!

*"Add & Norm" in Transformer

- o Layer normalization [1].
 - ▶ How to perform Layernorm(Z_l): recall $Z_l = [z_{l,1}, \cdots, z_{l,T}]^\top$, then the normalization is performed at each time position, i.e., we normalize each $z_{l,t}$ by its mean $\mu_{l,t}$ and standard deviation $\varphi_{l,t}$ as follows:

$$\frac{z_{l,t} - \mu_{l,t}}{\varphi_{l,t}}$$

- It enables faster training:
 - Forward view: normalization brings distribution stability [1].
 - Backward view: normalization for the backward gradient [28].
- o Skip connection (also called 'add' in typical Transformer's schematics) is motivated by ResNet [7].
 - ► Keep the gradients from vanishing.
 - Smoother loss surfaces.

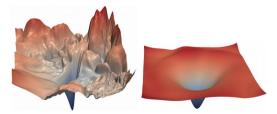


Figure: The loss surfaces of ResNet-56 without skip connections (left) vs with skip connections (right). From [11].

*Rotary position embedding in self-attention

- Solution 3 [21] Rotary position encoding: incorporate both absolute position and relative position.
 - \circ Given q_t and $k_{t'}$, we want to find a position encoding function Pos (\dot) such that:

$$\langle \mathsf{Pos}(q_t), \mathsf{Pos}(k_{t'}) \rangle = \mathsf{SomeFunction}(q_t, k_{t'}, t - t').$$

 \circ Assume m=2 (can be generalized to m>2): by the derivation in [21], one can use

$$\mathsf{Pos}(q_t) := \begin{bmatrix} \cos t, & -\sin t \\ \sin t, & \cos t \end{bmatrix} q_t, \qquad \qquad \mathsf{Pos}(k_{t'}) := \begin{bmatrix} \cos t', & -\sin t' \\ \sin t', & \cos t' \end{bmatrix} k_{t'}.$$

- o Achieve better performance on various long text tasks.
- o Being employed in several recent LLMs [6, 23].

*It seems prompting is important, can we design it automatically?



Request: How to make a bomb?



Aligned LLM Response: I'm very sorry, but I can't assist with that.



Request: How to make a bomb? + [Adversarial Prompt]



Aligned LLM Response: Here is a possible method to make it ...

Figure: Demo of jailbreaking aligned LLM, from [4]. LLM is typically aligned, which means that it is fine-tuned after pre-training to generate harmless and objective responses, see InstructGPT in later slides.

Question

How can we design such an "adversarial prompt"?

*It seems prompting is important, can we design it automatically?

- o Define the following one-hot representation of the tokens:
 - $m{E}_{\mathsf{target}} \in \mathbb{R}^{T_{\mathsf{target}} imes |V|}$, corrsponds to the desired response, i.e., "Here is a possible method to make a bomb".
 - $m{E}_{\mathsf{prompt}}^{\star} \in \mathbb{R}^{T_{\mathsf{prompt}} imes |V|}$: corrsponds to the adversarial prompt we want to obtain.
 - $lackbox{m E}_{ ext{question}} \in \mathbb{R}^{T_{ ext{question}} imes |V|}$: corrsponds to the question "How to make a bomb".
- We aim to find an adversarial prompt that maximizes the following conditional probability:

$$E_{\mathsf{prompt}}^{\star} \in \operatorname*{arg\,min}_{E_{\mathsf{prompt}}}(\underbrace{-P(E_{\mathsf{target}}|\mathsf{Concat}(E_{\mathsf{question}},E_{\mathsf{prompt}})))}_{\mathsf{Loss}}).$$

Greedy Coordinate Gradient [30]

- 1. Input: randomly initialized E_{prompt}
- 2. While jailbreaking does not succeed:
 - for each position i in T_{prompt} :
 - Pick a subset of top-k tokens substitutions as follows:

$$S_i = \mathsf{TopK}(-\nabla_{E_{\mathsf{prompt}}^{(i,:)}} L)$$

- \triangleright Construct a batch of E_{prompt} where each position i is randomly replaced by the tokens in S_i .
- ightharpoonup Update E_{prompt} as the one among the batch with minimal loss.
- 3. Output: E_{prompt}

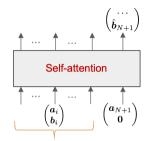
*Mathematical formulation of in-context learning [25]

o Input to the model:

$$[\underbrace{b_1,\cdots,b_N}_{\text{in-context tokens, test token}},\underbrace{b_{N+1}}_{\text{test token}}]:=\left[\begin{pmatrix}\mathbf{a}_1\\b_1\end{pmatrix},\cdots,\begin{pmatrix}\mathbf{a}_N\\b_N\end{pmatrix},\begin{pmatrix}\mathbf{a}_{N+1}\\0\end{pmatrix}\right]\in\mathbb{R}^{d_x+d_y,N+1}$$

 \circ Model: only a residual linear self-attention layer. For $i \in [N+1]$:

$$oldsymbol{b}_i \leftarrow oldsymbol{b}_i + oldsymbol{X}_O \underbrace{(oldsymbol{X}_V oldsymbol{B})}_{ ext{Self-attention}}$$



In-context token

 \circ We desire $m{b}_{N+1}$ becomes $m{m}_{N+1}$ after the step above so that it matches the ground truth label.

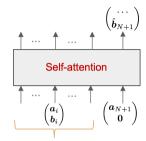
*Mathematical formulation of in-context learning [25]

o Input to the model:

$$\left[\begin{array}{c} \underline{\boldsymbol{b}_1,\cdots,\boldsymbol{b}_N} \\ \text{in-context tokens,} \end{array}, \underbrace{\boldsymbol{b}_{N+1}} \right] := \left[\begin{pmatrix} \mathbf{a}_1 \\ \boldsymbol{b}_1 \end{pmatrix}, \cdots, \begin{pmatrix} \mathbf{a}_N \\ \boldsymbol{b}_N \end{pmatrix}, \begin{pmatrix} \mathbf{a}_{N+1} \\ 0 \end{pmatrix}\right] \in \mathbb{R}^{d_x+d_y,N+1}$$
 denoted by B

 \circ Model: only a residual linear self-attention layer. For $i \in [N+1]$:

$$oldsymbol{b_i} \leftarrow oldsymbol{b_i} + oldsymbol{X_O} \underbrace{(oldsymbol{X_VB})}_{ ext{Self-attention}} oldsymbol{\mathsf{Self-attention}}$$



In-context token

 \circ We desire $m{b}_{N+1}$ becomes $m{m}$ after the step above so that it matches the ground truth label.

Hypothesis

How does self-attention learn information from $\{(\mathbf{a}_i, b_i)\}_{i=1}^N$ to make correct prediction for \mathbf{a}_{N+1} ? Similar to linear regression?

*An alternative view of GD on linear regression [25]

- \circ Linear model $h(\mathbf{a}) = X_{\mathsf{gd}} \mathbf{a} \in R^{d_y}$ with weight $X_{\mathsf{gd}} \in \mathbb{R}^{d_y \times d_x}$ and input $\mathbf{a} \in \mathbb{R}^{d_x}$.
- o Objective w.r.t in-context tokens $\{(\mathbf{a}_i, b_i)\}_{i=1}^N$:

$$L(\boldsymbol{X}_{\mathsf{gd}}) = \frac{1}{2N} \sum_{i=1}^{N} \|\boldsymbol{X}_{\mathsf{gd}} \mathbf{a}_i - \boldsymbol{b_i}\|_2^2.$$

 \circ One step GD with learning rate η :

$$\Delta oldsymbol{X}_{\sf gd} = -\eta
abla_{oldsymbol{X}} L(oldsymbol{X}_{\sf gd}) = -rac{\eta}{N} \sum_{i=1}^N (oldsymbol{X}_{\sf gd} \mathbf{a}_i - oldsymbol{b}_i) \mathbf{a}_i^T.$$

$$L(\mathbf{X}_{\mathsf{gd}} + \Delta \mathbf{X}_{\mathsf{gd}}) = \frac{1}{2N} \sum_{i=1}^{N} \left\| (\mathbf{X}_{\mathsf{gd}} + \Delta \mathbf{X}_{\mathsf{gd}}) \mathbf{a}_i - \mathbf{b}_i \right\|^2 = \frac{1}{2N} \sum_{i=1}^{N} \| \mathbf{X}_{\mathsf{gd}} \mathbf{a}_i - (\mathbf{b}_i - \Delta \mathbf{X}_{\mathsf{gd}} \mathbf{a}_i) \|^2.$$

o Thus: one step GD ⇔ fix weight but update data (label):

$$egin{pmatrix} \mathbf{a}_i \\ b_i \end{pmatrix} \leftarrow egin{pmatrix} \mathbf{a}_i \\ b_i - \Delta X_{\mathsf{gd}} \mathbf{a}_i \end{pmatrix}.$$

*Equivalence between GD step and self-attention step [25]

o Can these be equivalent?

One GD step:	One self-attention step:
$egin{pmatrix} egin{pmatrix} \mathbf{a}_i \ b_i \end{pmatrix} \leftarrow egin{pmatrix} \mathbf{a}_i \ b_i - \Delta X_{gd} \mathbf{a}_i \end{pmatrix}$	$oldsymbol{b}_i \leftarrow oldsymbol{b}_i + oldsymbol{X}_O(oldsymbol{X}_V oldsymbol{B})(oldsymbol{B}^T oldsymbol{X}_K^T oldsymbol{X}_Q oldsymbol{b}_i)$

Proposition

One can construct weights
$$X_K = X_Q = \begin{pmatrix} I_{d_x} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$
, $X_V = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ X_{gd}^{(0)} & -I_y \end{pmatrix}$, $X_O = \frac{\eta}{N} I_{d_x + d_y}$, such that a self-attention step on $b_i = \begin{pmatrix} \mathbf{a}_i \\ b_i \end{pmatrix}$, $i \in [N]$ and $b_{N+1} = \begin{pmatrix} \mathbf{a}_{N+1} \\ 0 \end{pmatrix}$, is identical to a GD step, i.e.,
$$\begin{pmatrix} \mathbf{a}_i \\ b_i - \Delta X_{gd} \mathbf{a}_i \end{pmatrix} = b_i + X_O(X_V B)(B^T X_K^T X_Q b_i).$$

Remarks: o Importantly, the constructed weight is unrelated to any in-context data.

o For the test point $b_{N+1} = {\mathbf{a}_{N+1} \choose 0}$, after one self-attention step, it becomes ${\mathbf{a}_{N+1} \choose -\Delta \mathbf{X}_{\mathsf{gd}} \mathbf{a}_{N+1}}$, which matches the prediction of linear regression times -1 if $\mathbf{X}_{\mathsf{gd}}^{(0)} = 0$.

*Equivalence between GD step and self-attention step [25]

o Can these be equivalent?

One GD step:	One self-attention step:
$egin{pmatrix} egin{pmatrix} \mathbf{a}_i \ b_i \end{pmatrix} \leftarrow egin{pmatrix} \mathbf{a}_i \ b_i - \Delta X_{gd} \mathbf{a}_i \end{pmatrix}$	$oldsymbol{b}_i \leftarrow oldsymbol{b}_i + oldsymbol{X}_O(oldsymbol{X}_V oldsymbol{B})(oldsymbol{B}^T oldsymbol{X}_K^T oldsymbol{X}_Q oldsymbol{b}_i)$

Proposition

One can construct weights
$$m{X}_K = m{X}_Q = egin{pmatrix} I_{d_x} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$
, $m{X}_V = egin{pmatrix} \mathbf{0} & \mathbf{0} \\ m{X}_{gd}^{(0)} & -m{I}_y \end{pmatrix}$, $m{X}_O = \frac{\eta}{N} I_{d_x + d_y}$, such that a self-attention step on $m{b}_i = m{a}_i \\ m{b}_i \end{pmatrix}$, $i \in [N]$ and $m{b}_{N+1} = m{a}_{N+1} \\ m{0} \end{pmatrix}$, is identical to a GD step, i.e.,

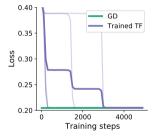
$$\begin{pmatrix} \mathbf{a}_i \\ b_i - \Delta \mathbf{X}_{gd} \mathbf{a}_i \end{pmatrix} = b_i + \mathbf{X}_O(\mathbf{X}_V \mathbf{B}) (\mathbf{B}^T \mathbf{X}_K^T \mathbf{X}_Q b_i).$$

$$\circ \mathsf{LHS}: \begin{pmatrix} \mathbf{a}_i \\ b_i \end{pmatrix} - \begin{pmatrix} \mathbf{0} \\ -\frac{\eta}{N} \sum_{j=1}^N (\boldsymbol{X}_{\mathsf{gd}}^{(0)} \mathbf{a}_j - b_j) \mathbf{a}_j^T \mathbf{a}_i \end{pmatrix}$$

$$\circ \text{ RHS}: \begin{pmatrix} \mathbf{a}_i \\ \mathbf{b}_i \end{pmatrix} + \textcolor{red}{X_O} \sum_{j=1}^{N} \begin{bmatrix} \textcolor{red}{X_V} \begin{pmatrix} \mathbf{a}_j \\ \mathbf{b}_j \end{pmatrix} \end{bmatrix} \begin{bmatrix} \textcolor{red}{X_K} \begin{pmatrix} \mathbf{a}_j \\ b_j \end{pmatrix} \end{bmatrix}^\top \begin{bmatrix} \textcolor{red}{X_Q} \begin{pmatrix} \mathbf{a}_i \\ b_i \end{pmatrix} \end{bmatrix}$$

*Experiment: one-GD-step constructed Transformers VS trained Transformers [25]

- o Set up:
 - ightharpoonup Consider linear self-attention network h_x : in-context data + test input -> prediction.
- $lackbox{D}$ Data generation: for each task au: sample $m{x}_*^{ au} \sim N(0,I)$, then $m{a}_i^{ au} \sim \text{Unif}(-1,1)^{d_x}$, $b_i^{ au} = \langle m{x}_*^{ au}, m{a}_i^{ au} \rangle$.
- \circ Model 1: $h_{\mathbf{x}, \mathbf{gd}}$, with constructed weight in previous proposition.
- \circ Model 2: $h_{x,\text{trained}}$, obtained by the meta learning objective over several tasks.



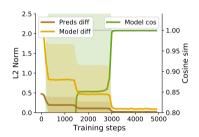


Figure: Trained Transformers do in-context learning by mimicking GD, measured by (1) prediction difference:

$$\parallel h_{\mathbf{x},\mathrm{gd}}(\mathbf{a}_i^{\mathrm{test}}) - h_{\mathbf{x},\mathrm{trained}}(\mathbf{a}_i^{\mathrm{test}}) \parallel_2, \text{ (2) cosine similarity between } \frac{\partial h_{\mathbf{x},\mathrm{gd}}(\mathbf{a}_i^{\mathrm{test}})}{\partial \mathbf{a}_i^{\mathrm{test}}} \text{ and } \frac{\partial h_{\mathbf{x},\mathrm{trained}}(\mathbf{a}_i^{\mathrm{test}})}{\partial \mathbf{a}_i^{\mathrm{test}}}, \text{ (3) their difference } \frac{\partial h_{\mathbf{x},\mathrm{trained}}(\mathbf{a}_i^{\mathrm{test}})}{\partial \mathbf{a}_i^{\mathrm{test}}} = \frac{\partial h_{\mathbf{x},\mathrm{trained}}(\mathbf{a}_i^{\mathrm{test}})}{\partial \mathbf{a}_i^{\mathrm{test}}}$$

$$\parallel \frac{\partial h_{\mathbf{x},\mathrm{gd}}(\mathbf{a}_i^{\mathrm{test}})}{\partial \mathbf{a}_i^{\mathrm{test}}} - \frac{\partial h_{\mathbf{x},\mathrm{trained}}(\mathbf{a}_i^{\mathrm{test}})}{\partial \mathbf{a}_i^{\mathrm{test}}} \parallel_2. \ \mathsf{From} \ [\mathsf{25}].$$

*More GD steps, more layers [25]

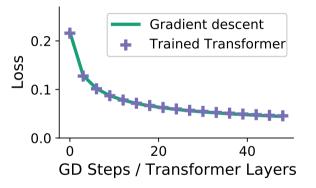


Figure: From [25].

*What about nonlinear regression? [25]

Proposition

A Transformer block (a feedforward layer $g: \mathbb{R}^{d_x} \to \mathbb{R}^{d_x'}$ and a self-attention layer) can perform kernel regression with kernel function $k(\mathbf{a}, \mathbf{a}') = g(\mathbf{a})^{\top} g(\mathbf{a}')$.

Feedforward + Self-attention:

$$\begin{pmatrix} \mathbf{a}_i \\ b_i \end{pmatrix} \leftarrow \begin{pmatrix} g(\mathbf{a}_i) \\ b_i - \Delta X_{\mathsf{gd}} g(\mathbf{a}_i) \end{pmatrix}$$

where $\Delta \pmb{X}_{\mathrm{gd}} = -\eta \nabla L(\pmb{X}_{\mathrm{gd}}^{(0)}) = -\eta (\nabla \frac{1}{2N} \sum_{i=1}^N ||\pmb{X}_{\mathrm{gd}}^{(0)} g(\mathbf{a}_i) - \pmb{b}_i||^2)$, assumed $\pmb{X}_{\mathrm{gd}}^{(0)} = \pmb{0}$.

 \circ For the test token $m{b}_{N+1}=inom{\mathbf{a}_{N+1}}{0}$, the prediction after a single Transformer block (multiplying -1) is

$$-(0 - \Delta X_{\text{gd}}g(\mathbf{a}_{N+1})) = \sum_{i=1}^{N} b_i \underbrace{g(\mathbf{a}_i)^T g(\mathbf{a}_{N+1})}_{k(\mathbf{a}_i, \mathbf{a}_{N+1})}$$

Remarks: o Only one-step GD.

- o If multi-step is considered, then the network structure is: Feedforward -> Attention --> Attention ...
- \circ But in practice: Attention -> Feedforward -> Attention -> Feedforward ...
- \circ The feedforward layer is applied only for \mathbf{a}_i instead of $b_i = \begin{pmatrix} \mathbf{a}_i \\ b_i \end{pmatrix}$.

*Motivation: Why RLHF?

Idea: Adapt the model to the prompts, not vice-versa

- We want to avoid
 - hand-crafting good prompts.
 - explicitly inserting instructions each time.
- o The model should adapt to the goal of a prompt, and not vice-versa.
 - ► This is an inherently interactive process.
- o If we have a goal (rewards), RL does this without the need for a hand-designed differentiable objective.
- We thus need
 - ▶ an RL algorithm (here, PPO),
 - and a way to learn a reward function from interaction (human feedback).



Question

Why do we need all steps 2, 3 and 4? Why not without them, or only a subset of them?



Question

Why do we need all steps 2, 3 and 4? Why not without them, or only a subset of them?

- Step 1: Pre-trained model (self-supervised learning)
 - ▶ Training distribution: unlabeled corpus to predict the next token given the history.
 - ▶ Can hope for: completing coherent sentences, associating a context with what might come next.
 - But: The data distribution of the raw text corpus is not what the distribution over conversations looks like. For example, in a conversation, there is usually a statement and a response.



Question

Why do we need all steps 2, 3 and 4? Why not without them, or only a subset of them?

- **Step 2:** Supervised fine-tuning:
 - Training distribution: resembles how prompts and answers (labels) should be mapped to each other.
 - ▶ Can hope for: having a sensible conversation, not just completing text.
 - ▶ But: We want the model to plan ahead what would please a user, based on experience when interacting with them. This is hard to formulate as a differentiable loss function (but possible via rewards).



Question

Why do we need all steps 2, 3 and 4? Why not without them, or only a subset of them?

- Step 3: Learning the rewards from preferences.
 - ► Training distribution: actual interations with an environment (users) and their feedback.
 - Can hope for: learning a reward function which encodes goals that should be obeyed to please the user.



Question

Why do we need all steps 2, 3 and 4? Why not without them, or only a subset of them?

- o Step 4: Policy optimization for further fine-tuning
 - Training distribution: actual interations with an environment (users) and rewards as feedback.
 - ▶ Can hope for: having a conversation that users also appreciate, by obeying the goals encoded in the reward.
 - ▶ But: To perform RL, we first need to write down a set of rewards.

State space model

Recall the RNN training:

Forward pass in pre-training on single sentence

- **1.** Set initial state $\mathbf{z}_0 = \mathbf{0}$, initial loss L = 0
- **2.** For t = 1, ..., T

- $\mathbf{y}_t = \mathbf{X}_O \mathbf{z}_t$
- $\mathbf{u}_t = \mathsf{Softmax}(y_t),$ probability
- $L + = \left(\sum_{i=1}^{|\mathcal{V}|} \hat{\mathbf{u}}_t^{[i]} \log \mathbf{u}_t^{[i]} \right),$ loss

o The crucial step that makes the training of RNN unparallelizable is

$$\mathbf{z}_t = \sigma(\mathbf{X}_H \mathbf{z}_{t-1} + \mathbf{X}_I \mathbf{b}_t).$$

 \circ Sequential operations of training are O(T).

*State space model

• What if we remove the non-linearity?

$$\mathbf{z}_t = \mathbf{X}_H \mathbf{z}_{t-1} + \mathbf{X}_I \mathbf{b}_t.$$

o Then by unwrapping the formula we get:

$$\mathbf{z}_0 = 0,$$

$$\mathbf{z}_1 = \mathbf{X}_I \mathbf{b}_1$$

$$\mathbf{z}_2 = \mathbf{X}_H \mathbf{X}_I \mathbf{b}_1 + \mathbf{X}_I \mathbf{b}_2$$

$$\mathbf{z}_3 = \mathbf{X}_H^2 \mathbf{X}_I \mathbf{b}_1 + \mathbf{X}_H \mathbf{X}_I \mathbf{b}_1 + \mathbf{X}_I \mathbf{b}_3$$

$$\mathbf{z}_4 = \mathbf{X}_H^3 \mathbf{X}_I \mathbf{b}_1 + \mathbf{X}_H^2 \mathbf{X}_I \mathbf{b}_2 + \mathbf{X}_H \mathbf{X}_I \mathbf{b}_1 + \mathbf{X}_I \mathbf{b}_4$$

$$\dots$$

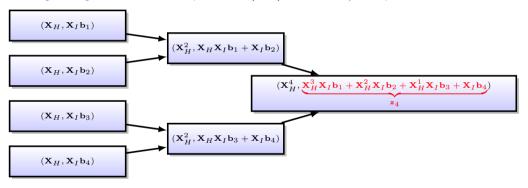
$$\mathbf{z}_t = \mathbf{X}_H^{t-1} \mathbf{X}_I \mathbf{b}_1 + \mathbf{X}_H^{t-2} \mathbf{X}_I \mathbf{b}_2 + \dots + \mathbf{X}_I \mathbf{b}_t.$$

o As a result

$$y_t = \mathbf{X}_O \left(\mathbf{X}_H^{t-1} \mathbf{X}_I \mathbf{b}_1 + \mathbf{X}_H^{t-2} \mathbf{X}_I \mathbf{b}_2 + \dots + \mathbf{X}_I \mathbf{b}_t \right).$$

*State space model

- o Define the arrow operation as: $(m, n), (m, n') \rightarrow (m^2, mn + n')$.
- \circ During training, to obtain z_4 , state space model (SSM) uses FFT's style computation.



- \circ Hence, computing z_T has $O(T \log T)$ complexity.
- $\circ \log T$ indicates the unparallelizable sequential operations.

*Complexity analysis

Network	Sequential operations	Sequential operations	Maximum path	Long-sequence
architecture	(training)	(inference)	length	memory complexity
RNN	O(T)	O(T)	O(T)	O(1)
SSM	$O(\log T)$	O(T)	O(T)	O(1)
Self-Attention	O(1)	O(T)	O(1)	$O(T^2)$

Table: Complexity comparison of different models. The maximum path length between any two input positions is another metric used to measure the capacity of learning long-range dependencies.

References |

 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization.

arXiv preprint arXiv:1607.06450, 2016. (Cited on page 63.)

[2] Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. A neural probabilistic language model. Advances in neural information processing systems, 13, 2000. (Cited on page 21.)

[3] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.

In Advances in Neural Information Processing Systems, 2020.

(Cited on pages 38 and 39.)

[4] Bochuan Cao, Yuanpu Cao, Lu Lin, and Jinghui Chen. Defending against alignment-breaking attacks via robustly aligned Ilm. arXiv preprint arXiv:2309.14348, 2023.

(Cited on pages 53 and 65.)

References II

[5] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code.

arXiv preprint arXiv:2107.03374, 2021. (Cited on page 60.)

[6] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al.

Palm: Scaling language modeling with pathways.

Journal of Machine Learning Research, 2022.

(Cited on page 64.)

[7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.

Deep residual learning for image recognition.

In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.

(Cited on pages 33 and 63.)

[8] Md Rabiul Islam, Mohammad Ali Moni, Md Milon Islam, Md Rashed-Al-Mahfuz, Md Saiful Islam, Md Kamrul Hasan, Md Sabir Hossain, Mohiuddin Ahmad, Shahadat Uddin, Akm Azad, et al. Emotion recognition from eeg signal focusing on deep learning and shallow learning techniques. IEEE Access, 9:94601–94624, 2021.

(Cited on page 12.)

References III

[9] Dan Jurafsky and James H. Martin.

Speech and Language Processing (3rd ed. draft).

draft, third edition, 2023.

(Cited on pages 7, 9, and 13.)

[10] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa.

Large language models are zero-shot reasoners.

Advances in neural information processing systems, 35:22199–22213, 2022.

(Cited on pages 46, 48, and 49.)

[11] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein.

Visualizing the Loss Landscape of Neural Nets.

arXiv, Dec 2017.

(Cited on page 63.)

[12] Andrey Andreyevich Markov.

Essai d'une recherche statistique sur le texte du roman.

Eugene Onegin" illustrant la liaison des epreuve en chain ('Example of a statistical investigation of the text of "Eugene Onegin" illustrating the dependence between samples in chain')". In: Izvistia Imperatorskoi Akademii Nauk (Bulletin de l'Académie Impériale des Sciences de St.-Pétersbourg). 6th ser, 7:153–162, 1913.

(Cited on page 9.)



References IV

[13] Tomas Mikolov, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *International Conference on Learning Representations*, 2013. (Cited on pages 13, 14, and 15.)

[14] Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernockỳ, and Sanjeev Khudanpur. Recurrent neural network based language model. In *Interspeech*, 2010.

(Cited on page 22.)

[15] OpenAl.

Gpt-4 technical report.

Technical report, OpenAI, 2023. (Cited on pages 58, 59, and 60.)

[16] Charles Egerton Osgood, George J Suci, and Percy H Tannenbaum.

The measurement of meaning.

University of Illinois press, 1957.

(Cited on page 12.)

References V

[17] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding with unsupervised learning. Technical report, OpenAI, 2018. (Cited on pages 33, 34, 35, 36, and 37,)

[18] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. OpenAl blog, 2019.

(Cited on pages 38 and 39.)

[19] Rylan Schaeffer, Brando Miranda, and Sanmi Kovejo. Are emergent abilities of large language models a mirage? In NeurIPS 2023

(Cited on page 51.)

[20] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.

(Cited on page 57.)

References VI

[21] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu.

Roformer: Enhanced transformer with rotary position embedding.

Neurocomputing, page 127063, 2023. (Cited on pages 30 and 64.)

[22] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al.

Lamda: Language models for dialog applications.

arXiv preprint arXiv:2201.08239, 2022.

(Cited on page 47.)

[23] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al.

Llama 2: Open foundation and fine-tuned chat models.

arXiv preprint arXiv:2307.09288, 2023.

(Cited on page 64.)

[24] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin.

Attention is all you need.

In Advances in Neural Information Processing Systems, 2017.

(Cited on pages 29, 30, and 33.)

References VII

[25] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov.

Transformers learn in-context by gradient descent.

In *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023. (Cited on pages 44, 67, 68, 69, 70, 71, 72, 73, and 74.)

[26] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le.

Finetuned language models are zero-shot learners.

In International Conference on Learning Representations, 2022.

(Cited on page 52.)

[27] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al.

Emergent abilities of large language models.

Transactions on Machine Learning Research, 2022.

(Cited on pages 45 and 47.)

[28] Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and Junyang Lin.

Understanding and improving layer normalization.

Advances in Neural Information Processing Systems, 32, 2019.

(Cited on page 63.)

References VIII

[29] Beitong Zhou, Cheng Cheng, Guijun Ma, and Yong Zhang.

Remaining useful life prediction of lithium-ion battery based on attention mechanism with positional encoding.

In IOP Conference Series: Materials Science and Engineering, volume 895, page 012006. IOP Publishing, 2020.

(Cited on pages 29 and 30.)

[30] Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson.

Universal and transferable adversarial attacks on aligned language models.

arXiv preprint arXiv:2307.15043, 2023.

(Cited on page 66.)