

# GPTs for Those Who Know and Love OLS: The Statistics of Large Language Models

BY ALYSSA BILINSKI, PhD\*

*Departments of Health Services, Policy, and Practice & Biostatistics,  
Brown University, Providence, RI*  
alyssa\_bilinski@brown.edu

5

JEREMY GOLDWASSER

*Department of Statistics,  
University of California Berkeley, Berkeley, CA*  
jeremy\_goldwasser@berkeley.edu

10

S. OZAI ALI

*Department of Health Services, Policy, and Practice  
Brown University, Providence, RI*

## SUMMARY

Large language models (LLMs) have become increasingly ubiquitous, but most researchers interact with them through chat interfaces rather than as statistical models. This paper provides an overview of how generative pre-trained transformers (GPTs) work for researchers with backgrounds in biostatistics, epidemiology, or health economics. We frame GPTs as an extension of familiar statistical methods: ordinary least squares, generalized linear models, and neural networks. We then describe the specific features that enable GPTs to generate text at scale across applications, including tokenization, embeddings, and the attention mechanism that allows models to weigh the relevance of different parts of an input sequence. Throughout, we emphasize that the mathematical operations underlying these models (e.g., matrix multiplication, gradient descent, softmax transformations) are conceptually accessible to researchers with quantitative training, even as optimal architectures and training procedures remain areas of active research. We conclude by discussing factors that have driven recent improvements in model performance, including innovations in preference learning, increased scale, expanded context windows, chain-of-thought reasoning, and supporting infrastructure.

15

20

25

*Some key words:* large language models; generative pre-trained transformers; ordinary least squares; generalized linear models; neural networks; tutorial; overview; review

30

---

\* We are grateful for inspiration, background and feedback from Alex D'Amour, Natalia Emanuel, Jeffrey Imai-Eaton, Adam Jermyn, Alex Reinhart; students taking Brown University's PHP 2455a; and participants at Boston University's Biostatistics Seminar, the Harvard TH Chan School of Public Health's Center for Communicable Disease Dynamics seminar (particularly Marc Lipsitch and Kirstin Oliveira Roster), and Indiana University's AI and Public Affairs Workshop (particularly Kevin Bryan, Kosali Simon, and Coady Wing). This work was funded in part by the National Institute of General Medical Sciences (1R35GM155224, AB, JG). The content is solely the responsibility of the authors and should not be construed as views of the funders or others acknowledged here. Since the start of this project, we have had to revise our working example from "Every week, the little girl gives treats to a furry, friendly cat" to "Every week, the little girl and boy give treats to a furry, friendly cat." To that end, we also give due acknowledgment to the natural generative intelligence in our lives, whose growth has paralleled that of the artificial variety and brought much more joy.

## 1. INTRODUCTION

Following the release of ChatGPT in 2022, researchers have increasingly interacted with large language models (LLMs) (OpenAI, 2022). In medicine and health policy, LLMs both support research tasks like coding and classification and are themselves objects of study, with the objective of understanding performance on tasks like exam performance and diagnostic support (Kung et al., 2023; Abbas et al., 2024; Eriksen et al., 2024; Katz et al., 2024). The preeminent LLM architecture is the generative pre-trained transformer (GPT). Popular GPTs include models from OpenAI (e.g., GPT-4o, GPT-o1, GPT-5.2), Meta (Llama), Google (Gemini), and Anthropic (Claude) (OpenAI, 2025a; Meta, 2025; Google, 2025; Anthropic, 2025). The name “GPT” reflects that these models *generate* text after being *pre-trained* on a large corpus (e.g., from the internet) in a self-supervised manner absent explicit labels, employing a neural network architecture called a *transformer* (Radford et al., 2018).

Most researchers interact with GPTs primarily through chat interfaces, and few work directly with their underlying statistical architecture. Furthermore, because GPTs are recent innovations originating in computer science and industry (Vaswani et al., 2017; Radford et al., 2018), which have different jargon, notation, and dissemination practices than health and medicine, few researchers encountered them in formal training. We wrote this paper because, in our attempt to learn about LLMs and teach them to our students, it was difficult to find resources with familiar language that built on our mathematical foundations and intuitions. We hope this translation effort may be useful for others in understanding GPTs.

In this work, we provide an overview of how GPTs work for researchers with a background in biostatistics, epidemiology, health economics or other fields outside computer science for which ordinary least squares remains a mainstay. We start by explaining the objective of GPTs: completing text sequences. We then review OLS and generalized linear models (GLMs) as building blocks for neural networks. Last, we describe features that allowed GPTs, a type of neural network, to achieve text generation at scale. Throughout, we highlight that the mathematical operations underlying GPTs should be familiar to those with an understanding of OLS, even as GPT performance, interpretability, and optimal architecture remain active areas of investigation.

## 2. TEXT GENERATION PROBLEMS

Most users interact with GPTs through text-based chat interfaces, providing a *prompt* and receiving output text, a *completion* (HuggingFace, 2025; OpenAI, 2025b; Community, 2025). Though completions appear as a single block, they in fact represent multiple iterative model runs (Vaswani et al., 2017).

At a high level, GPTs are “text generation” or “completion generation” models: given a sequence of words, they predict what follows (HuggingFace, 2025; Gadesha, 2024; Li et al., 2024) (Figure 1). Readers are likely familiar with the convention of notating a prediction  $\hat{Y}$  for a true outcome  $Y$  based on a vector of inputs  $X$ . Here, functions of the user prompt are used to create the  $X$  vector, and GPTs output a predicted probability distribution over possible next words (or sub-words, called tokens) (Radford et al., 2018). As an illustration, consider predicting the blank in: “Every week, the little girl give treats to a furry, friendly \_\_\_\_\_.” The  $X$  vector should encode attributes of the preceding words (e.g., the adjectives). A high-quality model would assign probability to various friendly animals, perhaps 40% to dog, 30% to cat, 10% to guinea pig. The precise distribution would depend on the training data: a model trained on sitcoms would produce different predictions than one trained on cartoons (e.g., higher probability of “squirrel”) or horror films (e.g., higher probability of “Cujo”).

Using this distribution, we generate a prediction for the output sequence  $\hat{y}$  by generating predictions for sequential tokens in an autoregressive process (Sanderson, 2024). At each step, the model predicts a distribution over the next token given the input  $X$  and all previously generated tokens. The selected token is appended to the sequence, and this process is repeated until a special end-of-sequence (EOS) token is generated, at which point it terminates (Vaswani et al., 2017). For example, given the sequence  $X$  described above, the model may first predict the token *guinea*, then condition on  $X$  augmented with *guinea* to predict *pig*, and finally predict EOS. The resulting decoded sequence  $\hat{Y} = (\{guinea\}, \{pig\})$  constitutes the model’s full response.

Although text completion models have long been subject of research, prior to transformers, they often performed poorly and were difficult to scale (Vaswani et al., 2017). High-quality text generation requires representing word meaning, position, and context while remaining computationally tractable. The following sections describe how modern GPTs achieve this.

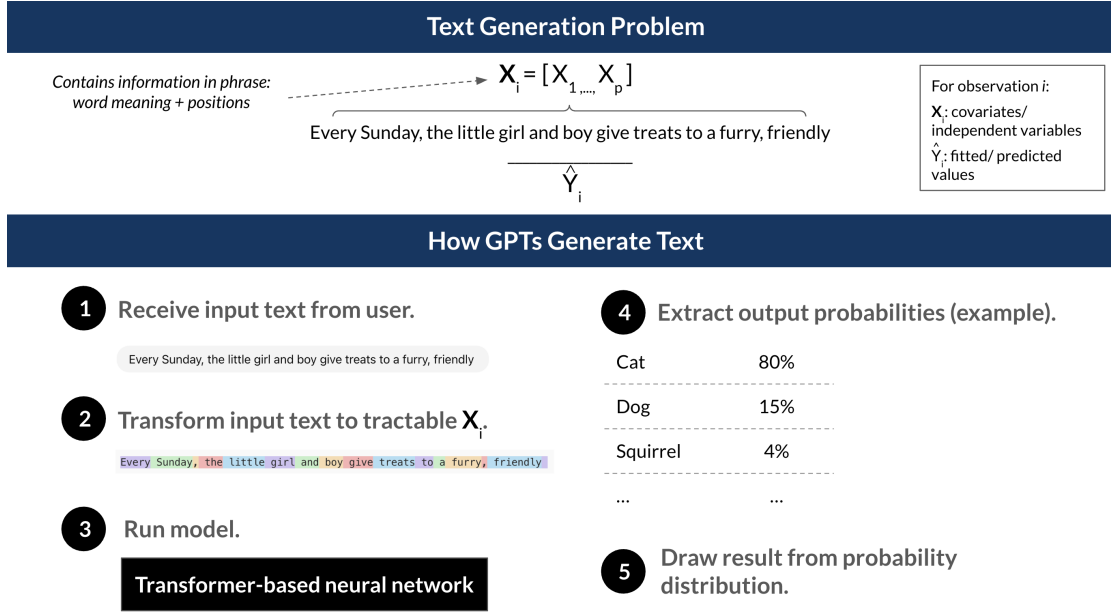


Fig. 1. The text generation problem and an overview of generative pre-trained transformer (GPT) structure.

### 3. BASIC PREDICTION MODELS: OLS, GLMs, AND NEURAL NETWORKS

We begin with a brief review of OLS and GLMs and then use these to explain neural networks, the class of machine learning models that includes GPTs (Figure 2). In the next section, we will detail specific features of GPTs.

#### 3.1. Ordinary least squares (OLS)

As the workhorse for many scientific analyses, OLS estimates parameters (commonly called  $\beta$ s) that minimize mean-squared error of predictions of an outcome  $Y_i$  for observation  $i$  from a linear combination of inputs  $\mathbf{X}_i = [1 \ x_{1i} \ \dots \ x_{pi}]$  (a row vector for observation  $i$ ). We start here not because OLS was ever a realistic option for natural language processing, but as a toy example to ground our discussion in the most intuitive and widely-used method in many fields. Given  $n$  observations ( $i = 1, \dots, n$ ) and  $p$  predictors, we find parameters that minimize:

$$\arg \min_{\beta_0, \beta_1, \dots, \beta_p} \sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi}))^2, \text{ or equivalently,}$$

$$\arg \min_{\beta} \sum_{i=1}^n (Y_i - \mathbf{X}_i \beta)^2,$$

Statistical Methods				
	Ordinary least squares (OLS)	Generalized linear models (GLMs)	Neural networks (NNs)	Generative Pre-trained Transformers: type of neural network with special features
Intuition	Predict $Y_i$ as a linear combination of elements of $X_i$	Predict $Y_i$ as a transformation of a linear combination of elements of $X_i$	Predict $Y_i$ by chaining transformed linear combinations of $X_i$	
Linear predictor	$X_i\beta$	$X_i\beta$	$W^{(0)}a_i^{(0-1)} + b^{(0)}$	
Transformation	-	$g^{-1}(\cdot)$	$\phi(\cdot)$	
Repeat?	No	No	L times	
Prediction of $Y_i$	$X_i\beta$	$g^{-1}(X_i\beta)$	$a_i^{(L)}$ , where $a_i^{(0)} = \phi(W^{(0)}a_i^{(0-1)} + b^{(0)})$	
Fit by minimizing	Squared error	Negative log likelihood	Loss function	
<div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> <p>More functional form flexibility</p> </div> <div style="text-align: center;"> <p>Even more functional form flexibility Works well with large number of parameters</p> </div> </div>				

Fig. 2. An overview of statistical methods for prediction problems. Neural networks are shown abstractly; the final layer depends on the task: regression (identity link), binary classification (sigmoid), or multiclass classification (softmax).

where  $\beta = [\beta_0 \ \beta_1 \ \dots \ \beta_p]^T$ . We obtain a vector of parameter estimates  $\hat{\beta}$  in closed form by taking the derivative of the objective function, setting it equal to 0, and solving for the optimal  $\hat{\beta}$ . With these estimates, our best prediction for the mean of the outcome, conditional on covariates, is  $\hat{Y}_i = X_i\hat{\beta}$ . When there is only one covariate ( $p = 1$ ), we can think of this as generating the best-fit line, in terms of minimizing mean-squared error.

105

To better understand these limitations and motivate GPTs, we consider a stylized example of how OLS might be used for text generation. Letting  $V$  be the number of possible next words, suppose  $X_i$  denotes some simple representation of the preceding text.  $X_i$  could be a basis (one-hot) vector indicating the most recent word, a concatenation of vectors indicating the last several, or counts of word frequency (a “bag of words”). With this representation, we could then fit  $V$  separate OLS regressions, each predicting a binary outcome for whether the next word takes a particular value. This would yield next-word predictions  $\hat{Y}_i = \{\hat{Y}_{i1}, \dots, \hat{Y}_{iV}\}$  for any  $X_i$ . These cannot be interpreted as probabilities, so we would need to normalize them, potentially inspired by linear probability models:

110

$$\hat{p}_{iw} = \frac{\hat{Y}_{iw}}{\sum_{j=1}^V \hat{Y}_{ij}}.$$

115 To generate new text, we could then draw from the resulting distribution. For example, we might deterministically select the word corresponding to the largest  $\hat{p}_{iw}$ . Alternatively, we could randomly draw words with the probability of choosing each word represented by  $\hat{p}_{iw}$ .

Unsurprisingly, this is not a viable approach. While OLS has many desirable properties, it is restricted to linear functions of inputs. When the true relationship between  $\mathbf{X}$  and  $\mathbf{Y}$  is not linear, 120 OLS will yield weak predictions, and language has complex structure that does not lend itself to linear representation. OLS estimates are also unbounded, meaning probability predictions could be negative. Directly modeling probabilities themselves would conceivably yield better predictions. Last, OLS performs poorly with many predictors, with reduced efficiency (i.e., greater variance) as the number of parameters grows relative to  $n$ . When  $p$  exceeds  $n$ , parameter 125 estimates are undefined. To model complex features and interactions of words, we will need a different approach.

### 3.2. Generalized linear models (GLMs)

Generalized linear models (GLMs) allow a partial relaxation of the functional form assumptions imposed by OLS (McCullagh, 2018). This class of models, which includes OLS as well as logistic 130 and Poisson regression, allow a partial relaxation of the functional form assumptions imposed by OLS (McCullagh, 2018). Rather than modeling the mean as a linear function of parameters, GLMs model a monotonic transformation of the mean as a linear function of parameters. We call this transformation the link function  $g(\cdot)$ . We solve:

$$\arg \min_{\boldsymbol{\beta}} \sum_i \left[ -\log p(Y_i | \mu_i(\boldsymbol{\beta})) \right]$$

where  $\mu_i(\boldsymbol{\beta}) = g^{-1}(\mathbf{X}_i \boldsymbol{\beta})$ , and  $p(Y_i | \mu_i(\boldsymbol{\beta}))$  is the likelihood for observation  $i$ . OLS is a special 135 case: with an identity link ( $g(\mu) = \mu$ ) and Gaussian likelihood, minimizing negative log-likelihood is equivalent to minimizing squared error.

Although most GLMs cannot be solved in closed form, parameter estimates can be obtained through gradient descent, which iteratively adjusts parameter estimates along the gradient of the loss (or log-likelihood) until it reaches stable values (McCullagh, 2018). We then predict the mean of  $Y_i$  as  $\hat{Y}_i = g^{-1}(\mathbf{X}_i \hat{\boldsymbol{\beta}})$ . 140

For text completion with  $V$  possible next words, multinomial logistic regression provides a framework for predicting categorical outcomes. Specifying that  $E[Y_{iw} | \mathbf{X}_i] = p_{iw}$ , it uses a linear function  $\mathbf{X}_i \boldsymbol{\beta}_w$  to produce a score (logit) for each next word  $w$ , which is converted to a probability via the softmax function:

$$\hat{p}_{iw} = \text{softmax}(\mathbf{X}_i \hat{\boldsymbol{\beta}})_w = \frac{\exp(\mathbf{X}_i \hat{\boldsymbol{\beta}}^w)}{\sum_{j=1}^V \exp(\mathbf{X}_i \hat{\boldsymbol{\beta}}^j)}$$

Exponentiation ensures all values are positive, and the denominator normalizes them to sum to one. Softmax also amplifies differences between inputs: larger values receive disproportionately more probability, concentrating mass on the most likely outcomes. We return to this property when discussing temperature in Section 4.5. 145

Overall, GLMs enable modeling a wider range of functional forms than OLS. However, they still rely on linear structure relating  $\mathbf{X}_i$  to the transformed mean  $g(E[Y_i | \mathbf{X}_i])$ , which does not hold in the text prediction setting and face similar limitations in terms of performance with high-dimensional inputs. 150

### 3.3. Neural networks

Neural networks can be thought of as chaining GLM-like transformations to more flexibly model outcomes and work well with large  $p$  relative to  $n$  (Ng & Ma, 2023; Nielsen, 2015). The building block is a *neuron*, which has the same basic structure as a GLM: a linear combination of inputs, transformed by a nonlinear function. Given a vector  $\mathbf{a}$ , a single neuron  $f(\mathbf{a})$  computes  $\phi(\mathbf{a}\mathbf{w} + b)$  (Nielsen, 2015). The transformation  $\phi$  is known as the “activation function,” and is analogous to the inverse link function  $g^{-1}$  in a GLM. The “weights”  $\mathbf{w}$  and “bias”  $b$  are analogous to the GLM parameters  $\boldsymbol{\beta}$ . Note that although the general structure of a neuron is similar to that of a GLM, neural networks are not motivated by modeling the distribution of  $Y$ , and the 155  
160

requirements for a GLM need not be strictly met. For example, a neural-network neuron can use non-canonical link/activation functions (e.g., a Rectified Linear Unit (ReLU) function, Appendix Figure S1), and need not assume exponential family error distributions.

Neural networks organize neurons into “layers” (Figure 3). The first layer of a neural network is the “input layer” ( $\ell = 1$ ) and consists of the sample value,  $\mathbf{a} = \mathbf{x}$  (where  $i$  referring to a specific observation is dropped following computer science notation). That layer has  $p$  neurons, with each neuron corresponding to one of the elements in  $\mathbf{x}$ . The next layer of a neural network is created from several neurons, with the output of each mapping onto an input in the next layer ( $\ell = 2$ ). For deeper networks, these neuron then become the inputs that define the third layer, and so on. Layers between the input and output layers are called “hidden” layers because users do not interact with these values. The final layer is called the “output” layer.

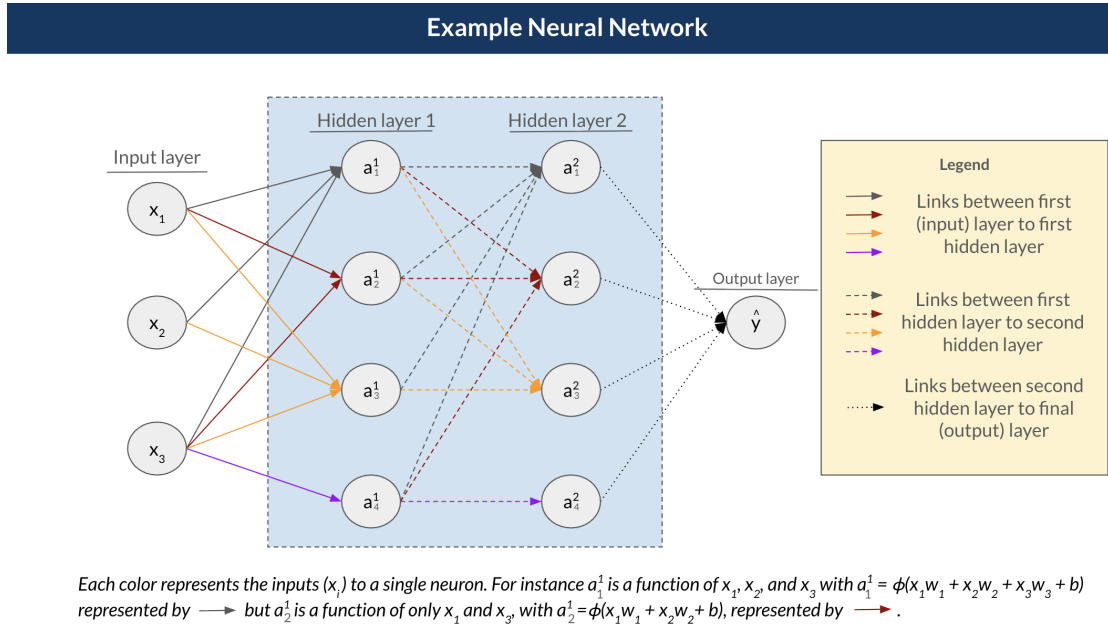


Fig. 3. An example neural network, demonstrating a single input layer, two hidden layers and an output layer. The arrows illustrate how variables are used from one layer to compute the next.

The value of a single node layer  $\ell > 1$  can be written (Nielsen, 2015):

$$a_j^{(\ell)} = \phi \left( \sum_k w_{jk}^{(\ell)} a_k^{(\ell-1)} + b_j^{(\ell)} \right).$$



where  $a_j^{(\ell)}$  is the value of the  $j$ -th neuron in the  $\ell$ -th layer,  $a_k^{(\ell-1)}$  is the value of the  $k$ -th neuron in the  $\ell - 1$ -th layer,  $w_{jk}^{(\ell)}$  is the weight connecting neuron  $k$  in layer  $\ell - 1$  with neuron  $j$  in layer  $\ell$ ,  $b_j^{(\ell)}$  is the bias term for layer  $\ell$ . Alternatively, in vector form, we may write: 175

$$\mathbf{a}^{(\ell)} = \phi(\mathbf{w}^{(\ell)} \mathbf{a}^{(\ell-1)} + \mathbf{b}^{(\ell)})$$

As with GLMs, neural networks are fit using gradient descent. Computing gradients through multiple layers requires the chain rule, working backward from output to input, a procedure called backpropagation (Nielsen, 2015).

In theory, two-layer neural networks with arbitrary hidden dimension can approximate nearly any functional form (Nielsen, 2015; Cybenko, 1989). However, to make learning from data tractable, networks are often provided additional structure and typically stack many hidden layers together (i.e., “deep learning”) (Nielsen, 2015). For example, to process images, researchers use convolutional neural networks (CNNs) that take weighted sums over groups of pixels of different sizes to detect edges and objects (Amidi & Amidi, 2024). However, while CNNs worked well for image recognition, they struggled to encode longer-range dependencies between words. Recurrent neural networks (RNNs) were designed to address this by updating a hidden state when processing each word in each layer, to be passed forward for future computations (Sutskever et al., 2011; Schmidt, 2019). Although RNNs were initially a promising approach to text completion and dominated the natural language space for several years, they struggled to capture relationships across longer context windows and were slow to train because they could not be fully parallelized (Vaswani et al., 2017). In the next section, we discuss the architecture (“the transformer”) and other features of generative pre-trained transformers that dramatically improved next-token prediction at scale. 180 185 190

#### 4. GENERATIVE PRE-TRAINED TRANSFORMERS

195

Generative pre-trained transformers (GPTs) are a class of neural network optimized for predicting the next word in a sequence (Figure 1). In this section, we describe the core components of GPTs in detail.

As a roadmap, GPTs first convert an input sequence of words into smaller units called “tokens.”

200 Tokens are then mapped to vectors called embeddings, which allow researchers to efficiently represent word meaning and position. These components were not unique to transformers or GPTs; tokenization and word embeddings were developed initially with RNNs. Next, embeddings are passed through a neural network with at least one “transformer” block. Transformer blocks include an “attention mechanism”, designed to compute how much each preceding token in a  
205 sequence should influence predictions of what follows. The output from this neural network defines a probability distribution over all possible next tokens in the model’s vocabulary, from which a final output is drawn (with randomness controlled by a temperature parameter). We close the section with discussion model training, as well as factors that determine their performance.

As GPTs have advanced, architectural details have become more proprietary. Therefore, we  
210 describe core structures from published early models here to provide a useful foundation, noting that more recent models may include additional advancements.

#### 4.1. Tokenization

As the first step of processing input text, words are mapped onto “tokens”, either entire words or chunks of words that carry some meaning. Prefixes and suffixes may be tokens; for example,  
215 the word “jumping” may be broken into “jump” and “ing” (OpenAI, 2025c). On average, tokens used by OpenAI as of November 2025 contained 4 characters and represented  $\frac{3}{4}$ th of a word (OpenAI, 2025d).

The set of possible tokens is determined by algorithms that identify frequently occurring sequences of characters and group these into tokens until the total number of unique tokens  
220 reaches the target vocabulary size (Sennrich et al., 2016; Stanford Online, 2024). Tokenization tends to compress the word vocabulary into a smaller set of subwords (Radford et al., 2019; OpenAI, 2026). For instance, GPT-2 used 50,000 tokens and GPT-4o about 200,000, compared to roughly a million words in English (Merriam-Webster, 2025; Radford et al., 2019; OpenAI, 2026). By breaking words into smaller, reusable pieces, tokenization allows related forms to share  
225 information (e.g., *some* in both *something* and *somewhere*).

To ensure that all characters can be mapped to tokens, modern tokenizers operate on bytes (of which there are 256) rather than characters, ensuring that any text, including misspellings, novel words, or emoji, can be represented without requiring an “unknown” token (Hugging Face, 2024).

## 4.2. Embeddings

230

Tokens are then mapped to high-dimensional vectors called “embeddings”, which serve as model inputs. Embeddings are designed so that mathematical operations on them like cosine similarity (a normalized dot product) and vector addition reflect the semantic meaning of tokens and, separately, their position within a sequence. By capturing these relationships in a continuous vector space, models can represent complex interactions using far fewer dimensions than would be required by simpler representations. There are two types of embeddings: semantic and positional embeddings.

### Semantic embeddings

Semantic embeddings represent the meaning of tokens. In a naive setup, we could imagine that we would represent each token as a unique basis vector – e.g.,  $\mathbf{a} = \begin{bmatrix} 1 & 0 & 0 & \dots \end{bmatrix}^T$ , the  $\mathbf{b} = \begin{bmatrix} 0 & 1 & 0 & \dots \end{bmatrix}^T$ . In this setup, the dimensionality of the vector space would have to be at least as large as the vocabulary, and mathematical operations between vectors would not yield linguistically meaningful results. (For example, all dot products would be zero.) Capturing relationships between tokens, such as similarity or compositional meaning, would therefore require explicitly modeling two-way or higher-order interactions between all vectors, requiring many parameters and scaling poorly with sequence length.

In 2013, researchers at Google released a method to address this in the package *word2vec* (Mikolov et al., 2013a). Their approach learned low-dimensional vector representations such that mathematical operations on them reflected semantic relationships. For example, tokens that were similar to each other (e.g., “emperor” and “king”) or found in similar contexts (e.g., “Berlin” and “Germany”) had a higher cosine-similarity score (Mikolov et al., 2013a; Sanderson, 2024). Simple algebraic operations on the vectors could also yield intuitive results: e.g. vector(“King”)

- vector “Man”) + vector(“Woman”) resulted in a vector that was close, as calculated by cosine similarity, to vector(“Queen”) (Mikolov et al., 2013b).

Most GPTs today learn their own embeddings as part of an end-to-end training process. OpenAI’s GPT-3 model, the last model for which they published architectural details, used embeddings of length 12,288 (Brown et al., 2020).

### Positional embeddings

Each token in the sequence is also mapped to a positional embedding of the same length as the semantic embedding, which encode where the token appears in the sequence relative to other tokens. In early models, positional embeddings were constructed using sinusoidal functional so that relative position was recoverable from cosine similarity (Vaswani et al., 2017) (Figure S2). For instance, in the sequence “Every week, the little girl and boy give treats to a furry, friendly”, the dot product of the first position (corresponding to Every) and the second position (corresponding to week) was higher than the dot product of the first position and the fifth position (corresponding to girl).

As with semantic embeddings, more contemporary models (e.g., the OpenAI GPTs) learned positional embeddings, with each position’s embedding a trained parameter (Radford et al., 2018, 2019; Brown et al., 2020). Other models (e.g., LLAMA-2) incorporated position differently, integrating it later into the attention mechanism rather than at the input layer (Su et al., 2024; Touvron et al., 2023). The key intuition nevertheless remains that models seek to represent both meaning and position.

### Model inputs

Transformers generally combine semantic and positional embeddings for each token. While the obvious approach would be to concatenate the two vectors, doing so would double the dimensionality and slow training. Instead, both signals can largely be preserved by adding the embeddings together, as random vectors in high-dimensional space are nearly orthogonal. Per this logic, the original Vaswani et al. (2017) paper proposed adding positional embeddings to the semantic embedding at the first layer of the network. The resulting  $n \times d_{\text{model}}$  matrix (in which  $d_{\text{model}}$  represents the chosen length of the final vector embeddings) would be the model input. This has remained a standard choice, with both absolute and learned positional embeddings (Radford

et al., 2019; Dugas, 2023). Alternatively, some models add positional embeddings into various stages of the attention blocks (Shaw et al., 2018; Raffel et al., 2020; Su et al., 2021; Touvron et al., 2023).

#### 4.3. The attention mechanism

285

Once a sequence is mapped to embeddings, it is then passed into a neural network with a transformer architecture, displayed in Figure 4. The key breakthrough in this design was using an “attention mechanism” without recurrence (i.e., the hidden memory states as had been used in RNNs but hindered parallelization), hence the title of the seminal paper: “Attention is All You Need” (Vaswani et al., 2017). In this section, we describe the standard attention mechanism design before moving to other aspects of the transformer.

290

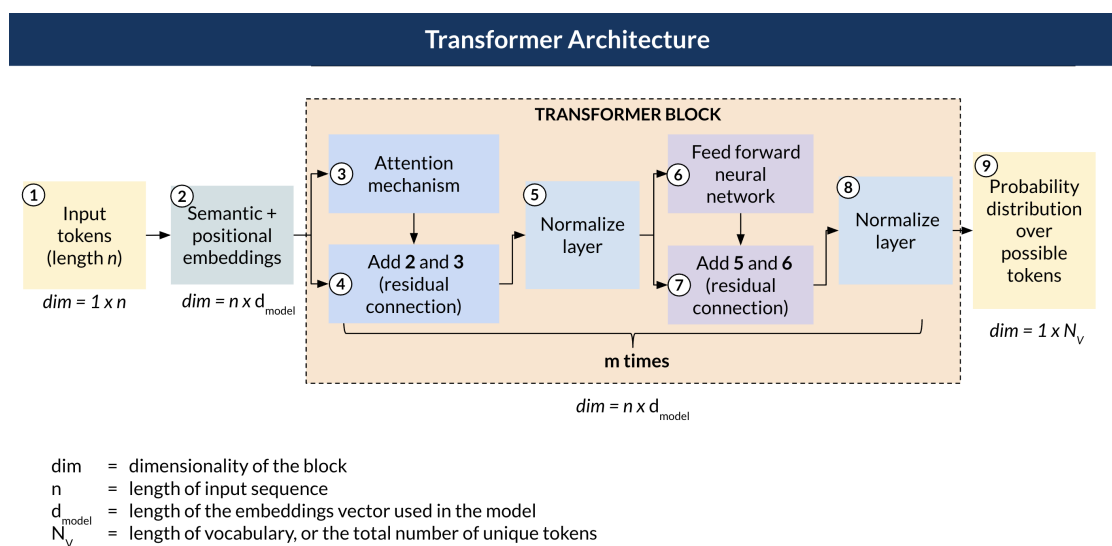


Fig. 4. A simplified transformer architecture that takes in a sequence of tokens as input and produces a probability distribution over all possible next tokens as output.

Consider the text completion task used throughout this paper: “Every Sunday, the little girl and boy give treats to a furry, friendly \_\_\_\_\_.” To predict the missing token, the model must use context. Some tokens are clearly more relevant than others for guessing what comes next. For example, “furry” and “friendly” strongly suggest that the missing token is a

295

noun and, specifically, a pet. Conversely, tokens like “Every” or “Sunday” provide relatively little information to guide this prediction.

The attention mechanism (Figure 5) allows a model to determine which tokens in the input sequence are most relevant for predicting each output (Vaswani et al., 2017; Sanderson, 2024).

300 Importantly, these relevance weights can be computed in parallel, enabling the approach to scale efficiently to long sequences. Broadly, it has a similar structure to search and retrieval algorithms that use key-value pairs and queries (Pichka, 2025). To illustrate, consider a search algorithm in a video website. The value for each video could be the URL associated with a video, and the key is a summary of the video contents (e.g., a title, a series of tags, or both). A query may be what a user searches for in the search bar. For instance, a user may type in “cute cat video” as a query. 305 A search and retrieval algorithm will typically calculate a similarity score between the query and all keys stored within the database, and sort the keys by score. Finally, the algorithm will return the value (in this case the video) associated with the key most similar to the query. When fitting an LLM, researchers estimate parameters that allow us to define a key matrix (characterizing what a token contains) and a query matrix (characterizing what information a token seeks, with 310 the end goal of predicting the next token). The value matrix then characterizes what each token contributes to when matched.

Mathematically, the attention mechanism computes:

$$\text{Attention}(\mathbf{X}_n) = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V}$$

315 where  $\mathbf{Q}$  (the query matrix),  $\mathbf{K}$  (the key matrix), and  $\mathbf{V}$  (the value matrix) are linear transformations of  $\mathbf{X}_n$  (defined below). We will discuss each step of this process in the following sections.

Here is how we can understand these components (Vaswani et al., 2017; Sanderson, 2024).

The key matrix ( $\mathbf{K}$ ) encodes what information each token “contains” in its attributes that might be relevant to other tokens’ queries. It is computed as  $\mathbf{K} = \mathbf{X}_n \mathbf{W}_K$ , where  $\mathbf{W}_K$  is a learned 320 parameter matrix of dimension  $d_k \times d_{model}$ , which we can think of as analogous to  $\beta$  in a traditional linear model. We can also say  $\mathbf{K}$  contains the key vector for each element in the sequence.

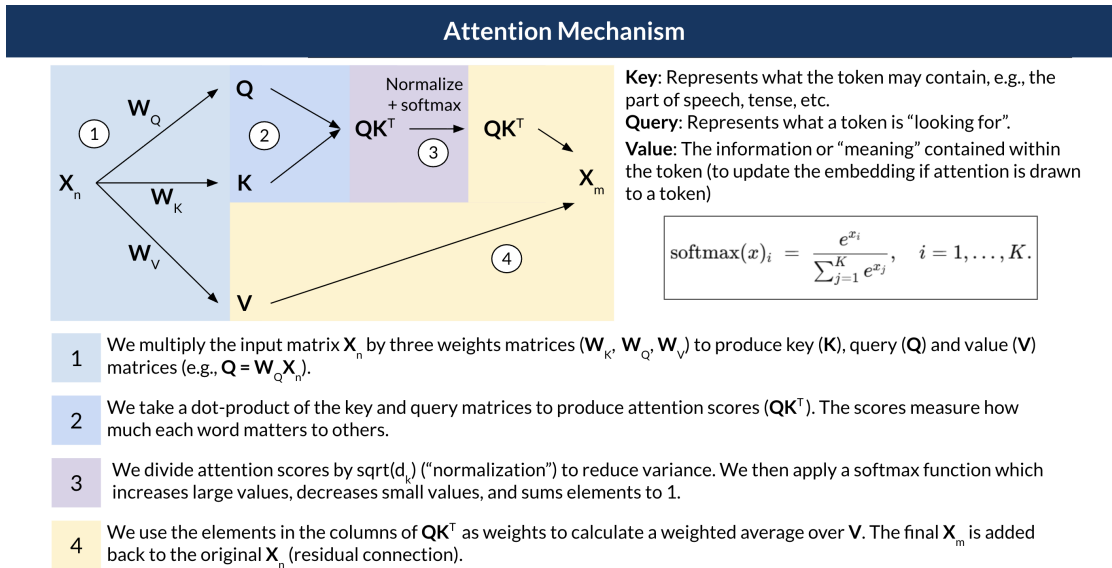


Fig. 5. The attention mechanism described in detail as a sequence of matrix multiplication and other operations.

The query matrix ( $Q$ ) encodes what each token is “looking for” in surrounding context. For instance, the query for “friendly” might encode “looking for a noun to attach to.” It is computed as  $Q = X_n W_Q$ , where  $W_Q$  has the same dimensions as  $W_K$ .

325

*Attention scores* are computed as the product  $QK^T$ , an  $n \times n$  matrix which contains dot product between every query vector and every key vector. Each entry  $(i, j)$  measures how relevant token  $j$  is to token  $i$ —higher values indicate greater relevance. This matrix is normalized by dividing by  $\sqrt{d_k}$ , the square root of the key dimension. This normalization prevents the variance of the dot products from growing with  $d_k$ , which would cause the subsequent softmax to produce extreme values. The softmax function is then applied to each row, transforming the raw scores into a probability distribution  $\left(\phi(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}\right)$ . This exaggerates large values, suppresses small ones, and ensures each row sums to one. The resulting matrix  $\Omega = \text{softmax}(QK^T / \sqrt{d_k})$  contains the *attention weights*.

330

We last define one additional matrix: the value matrix ( $V$ ) contains the information that will be passed forward if a token is deemed relevant. It is computed as  $V = X_n W_V$ , where  $W_V$  is of dimension  $m \times m$ .

335

Multiplying  $\mathbf{\Omega}$  by  $\mathbf{V}$  produces a weighted sum of the value vectors:

$$\mathbf{\Omega V} = \begin{bmatrix} w_{11}\mathbf{V}_{\text{Every}} + w_{12}\mathbf{V}_{\text{Sunday}} + w_{13}\mathbf{V}_{\text{the}} + \dots \\ w_{21}\mathbf{V}_{\text{Every}} + w_{22}\mathbf{V}_{\text{Sunday}} + w_{23}\mathbf{V}_{\text{the}} + \dots \\ \vdots \end{bmatrix} \quad (1)$$

Each row of the output is a weighted combination of all tokens’ value vectors, where the weights reflect how much attention each token pays to the preceding token. Tokens deemed more relevant (via the query-key interaction) contribute more to the output.

Finally, we add this result to update the input in a “residual connection.” This helps prevent vanishing gradients (where gradients shrink toward zero as they pass through many layers during backpropagation, stalling learning) and allows the model to learn incremental refinements to token representations (i.e., updating initial embeddings or results of the prior layer) rather than entirely new representations at each layer.

#### 4.4. *Transformer architecture*

The output of each residual connection is also passed through a layer normalization step (Vaswani et al., 2017; Sanderson, 2024). Layer normalization rescales each token’s representation (i.e., the result of the embedding being added to the output of the attention mechanism and renormalized) to have zero mean and unit variance. Empirically, this improves numerical stability and accelerates convergence during training, particularly in more complex models.

Each block also includes a feedforward network that applies the same nonlinear transformation independently to each token’s representation (Vaswani et al., 2017; Sanderson, 2024). Unlike the attention mechanism, which mixes information across tokens, the feedforward network transforms each representation in isolation, using the standard neural network architecture discussed above. This allows the model to further process each token after contextual information has been incorporated through attention. Finally, models often use another residual connection and layer normalization after the feedforward network.

In practice, the attention mechanism is also employed in “multi-headed” setup. Rather than computing attention weights once, the model computes them multiple times in parallel as part of the same attention mechanism step. For example, the original “Attention Is All You Need” paper



used 8 heads per attention block (Vaswani et al., 2017). Each head has its own query, key, and value matrices, producing a unique set of attention weights  $\text{softmax}(\mathbf{QK}^\top/\sqrt{d_k})$ . This allows the model to capture different types of relationships between tokens simultaneously. In principle, one head might learn to attend to syntactic relationships, another to semantic similarity, another to positional patterns, and so on (although roles may not be as clean in practice to human eyes) (Kissane et al., 2024). The outputs of all heads are concatenated and then projected back to the model dimension  $d_{\text{model}}$  using a learned weight matrix  $\mathbf{W}_O$ . Because each head operates on a lower-dimensional subspace (typically  $d_k = d_{\text{model}}/h$ , where  $h$  is the number of heads), the total computational cost of multi-head attention is similar to that of single-head attention.

In addition to using “multi-headed” attention, models typically stack many such attention blocks in sequence. GPT-3, for instance, used 96 blocks (Brown et al., 2020). Each block refines the token representations further, allowing the model to build increasingly abstract representations of the input.

#### 4.5. Output

As its final output, the model produces an updated matrix of token representations, still of dimension  $n \times m$ . To generate a prediction, the model focuses on the embedding corresponding to the final token in the sequence, which we denote  $\mathbf{y}_n$ . This embedding summarizes the full context observed so far and can be interpreted as the model’s best representation of “what should come next.”

The model converts this representation into a probability distribution over the vocabulary. One approach to do this is simply to compute the dot product of  $\mathbf{y}_n$  with the embedding of every possible token in the language (Dugas, 2023). Recall that embeddings have the property that semantically similar tokens have large dot products. This means that this operation assigns higher value to tokens whose meanings are most compatible with model’s prediction of what comes next given context. For the example sentence

*Every Sunday, the little girl and boy give treats to a furry, friendly \_\_\_\_\_,*

tokens such as *cat* or *dog* will have higher dot products with  $\mathbf{y}_n$  than less related tokens (such as *alligator*). These can then be converted to probabilities using a softmax function. (Alternatively,

models may have another step of an “unembedding” matrix, which allows this process to be a more flexible transformation.)

To make a prediction, the model may select the token with the highest probability (greedy decoding) or sample from this distribution, which introduces randomness and allows for more diverse text generation. A *temperature* parameter controls how concentrated the distribution is: lower temperatures sharpen the distribution toward the highest-probability tokens, producing more deterministic outputs, while higher temperatures flatten it, introducing more randomness (Peeperkorn et al., 2024).

#### 4.6. Model size

So, how many parameters does an LLM have? Consider a single head of attention operating on a sequence of  $n$  tokens with embedding dimension  $m$ . The three weight matrices have dimensions (Vaswani et al., 2017):

$$\dim(\mathbf{W}_Q) = \dim(\mathbf{W}_K) = m \times d_k \quad (2)$$

$$\dim(\mathbf{W}_V) = m \times d_v \quad (3)$$

where  $d_k$  and  $d_v$  are the dimensions of the query/key and value vectors respectively. A single head thus contributes  $2md_k + md_v$  parameters.

In many LLMs like GPT-3,  $d_k = d_v = m/h$ , where  $h$  is the number of heads (Radford et al., 2019; Brown et al., 2020). In this case, the total parameters across all  $h$  heads for the  $m \times m/h$ -dimensional Q, K, and V matrices is  $3m^2$ . An additional matrix  $\mathbf{W}_o$  of dimension  $m \times m$  combines the conjoined head outputs. The attention mechanism therefore contains  $4m^2$  parameters per block. The feedforward network added further parameters. In GPT-3, two layers with an inner dimension of  $4m$  contributed  $8m^2$  parameters (Brown et al., 2020). Thus, each of the  $b$  blocks has  $12m^2$  parameters, yielding a total count of roughly  $12m^2b$ .

The resulting models are enormous, due to large embeddings and a deep stacks of blocks. GPT-3 used  $m = 12,288$  and  $b = 96$ , totaling about  $12m^2b = 174$  billion parameters. Token and position embeddings contributed to the final tally of 175 billion parameters (Brown et al., 2020). Llama 3.1, released by Meta in 2024, is over twice as large. With 16, 384-dimensional embeddings

and 128 blocks, it totals 405 billion parameters (Grattafiori et al., 2024). DeepSeek’s V3 model contains over 600 billion parameters; GPT-4 and GPT-5 are rumored to be significantly larger.

#### 4.7. Training LLMs

420

Training LLMs is feasible for several reasons (Vaswani et al., 2017). First, the core operations in a transformer, matrix multiplications and element-wise nonlinearities, can be parallelized. Second, training examples can also be processed in parallel batches, allowing gradients to be computed across many examples simultaneously. Third, modern GPUs (particularly those from NVIDIA) are optimized for exactly these operations, with training distributed across thousands of GPUs simultaneously.

425

The primary computational bottleneck is the attention mechanism. Computing  $\mathbf{QK}^\top$  produces an  $n \times n$  matrix, meaning that computational and memory costs scale quadratically with sequence length. This is why models specify a maximum *context window*, the total number of tokens the model can process at once. A model’s context window determines how much text it can “see” when generating a response and, in a chat conversation, includes both user input and prior responses. When generating each token, the model considers everything within this window; information beyond it is effectively invisible, as users may encounter when a model appears to forget information from the start of a long conversation.<sup>1</sup>

430

Training typically consists of 3 training phases: pre-training, supervised fine-tuning, and preference learning. This pipeline was standardized by the InstructGPT paper, an important step preceding the release of ChatGPT (Ouyang et al., 2022; Stanford Online, 2024).

435

#### Pre-Training

The first phase of training LLMs is pre-training on vast amounts of data, typically massive corpus of text, comprising books, websites, code repositories, and other webpages. Because pre-training is self-supervised, it does not require labeled data, so any reasonable piece of text may be included in the training corpus; the model simply predicts each next word given the prior sequence. A major innovation of GPT-3 over its predecessor was scaling pre-training: It trained on 300 billion tokens, an order of magnitude bigger than GPT-2 (Radford et al., 2019; Brown

440

<sup>1</sup> As an alternative to GPTs, Mamba is a state-space model whose complexity is linear, rather than quadratic in the sequence length (Gu & Dao, 2023). Although promising, it generally underperforms GPTs; a detailed description is beyond the scope of this work.

et al., 2020). This scaling has only since continued, with contemporary LLMs pre-training on trillions of tokens (Chowdhery et al., 2022; Grattafiori et al., 2024; Deepseek, 2025). Llama 4 claims to use a staggeringly large dataset of 40 trillion tokens (Meta, 2025).

Data curation is an important part of pre-training, and not all pre-training data is equally valuable. For example, academic papers may be higher-quality text than commercial websites (though some may reasonably disagree), and might be upsampled accordingly. In addition, not all text that can be scraped should be used to pretrain a model. Internet datasets are rife with offensive or dangerous content, which researchers attempt to filter out. Models trained on such data are known to learn such behaviors. GPT-3 was found to produce racist, sexist, and violent text based on patterns learned in its pre-training (Brown et al., 2020).

### **Supervised Learning**

Supervised learning adapts a pre-trained model for a range of different tasks. This phase often uses labeled data, with a prompt and a human- or machine-labeled output. Because labeled data are scarce, the total dataset used for supervised fine-tuning is small compared to the pre-training corpus. Pre-training data accounted for 98% of the text used to train InstructGPT (Ouyang et al., 2022).

Pre-trained models are fine-tuned on several tasks at once, each with their own labeled datasets. Tasks might include question answering, document summarization, machine translation, open-ended generation, and rewriting. In the pre-LLM era, models would be trained directly on these datasets. However, pre-training improved performance substantially as it instilled a deep sense of linguistic fluency, background knowledge, and perhaps some degree of reasoning ability. Hallucinations are thought to arise from fine-tuning, as models learn to produce assertive responses even when relevant information was absent from pre-training (Huang et al., 2025; Kalai et al., 2025).

### **Preference Learning**

The third stage, preference learning, is another round of fine-tuning. This phase is responsible for much of the “personality” of modern chatbots, including their confident tone, sycophancy, and tendency toward helpfulness (Ouyang et al., 2022). It was a major breakthrough in the transition from GPT-3 to ChatGPT.

The first preference learning procedure introduced for LLMs was reinforcement learning from human feedback, or RLHF. Originally demonstrated in robotics (Christiano et al., 2017), 475  
RLHF was applied to language models at scale in InstructGPT (Ouyang et al., 2022). In RLHF, human raters compare different model outputs and indicate which response is preferred. These preferences are used to train a reward model that predicts human approval. The language model is then fine-tuned to maximize this reward.

The leading alternative to RLHF is Direct Preference Optimization (Rafailov et al., 2024). 480  
DPO obviates the need for a reward model and RL training. It fine-tunes the LLM directly on the preference data, maximizing the likelihood that the winning responses are preferred. This approach is computationally cheaper and easier to execute, and thus has become common in practice.

More recently, researchers have explored reinforcement learning with verifiable rewards, for 485  
instance, training models on mathematics or coding problems where correctness can be checked automatically, reducing reliance on human judgment (Guo et al., 2025).

#### 4.8. *Why Are Models Getting So Much Better?*

The prior sections have characterized the general structure of LLMs, but models have drastically improved over the past few years. Here, we outline the factors contributing to this, with the aim 490  
of helping readers to understand bottlenecks and anticipate how future models may evolve.

First, models have improved as they have become larger, both in parameter count and training data. GPT-3's 175 billion parameters represented a hundred-fold increase over GPT-2 (1.5 billion) and a thousand-fold increase over GPT-1 (117 million) (Radford et al., 2018, 2019). Subsequent models have continued this trend; the largest publicly documented models now exceed 400 billion 495  
parameters (Grattafiori et al., 2024). Researchers have also become better at optimizing the ratio of model parameters to training data (Hoffmann et al., 2022). For example, GPT-3 was discovered to have an insufficient number of training points relative to its parameter size.

Second, context windows have expanded dramatically. The original transformer had a context window of 512 tokens (Vaswani et al., 2017); GPT-2 extended this to 1,024 tokens (Radford et al., 500  
2019), GPT-3 to 2,048 tokens (Brown et al., 2020), Llama-3 to 128,000 tokens (Grattafiori et al.,

2024) and Gemini reporting up to 10 million in experiments (Team et al., 2024). This expansion increases the model’s effective working memory and allows it to perform better in longer conversations. Because context window size remains a key factor affecting both performance and computational costs, strategies have been developed both to speed computation (e.g., developing  
505 optimization methods that make attention faster on modern hardware (Dao et al., 2022) and improve performance with a given limit (e.g., summarizing or ”compacting” of earlier conversation history to reduce token count or using retrieval systems that fetch relevant information on demand rather than keeping everything in context).

510 Third, the introduction of intermediate reasoning steps, often called chain-of-thought prompting or extended thinking, has substantially improved performance on complex tasks (Wei et al., 2022). This involves generating intermediate reasoning steps before producing a final answer and allows models to tackle multi-step problems that would otherwise exceed their capabilities. This technique has proven particularly effective for mathematical reasoning (Guo et al.,  
515 2025). Interestingly, it performs well even when models are not actually using the reasoning they are stating within intermediate steps (Lanham et al., 2023; Chen et al., 2025).

In addition to direct model improvements, users have also benefited from improvements in supporting architecture, including document processing, integrated code execution environments that make it easy to check computations, and structured output formatting. Still, some challenges  
520 persist; for example, when reading PDFs, layout, tables, and multi-column formatting may be lost or misinterpreted during text extraction.

Overall, these improvements have translated into dramatic gains in real-world performance. METR, an AI safety organization, proposed measuring model capabilities by the length of tasks models can complete autonomously with 50% reliability, where task length is defined by  
525 how long a task takes human professionals (Kwa et al., 2025). By this metric, frontier model capabilities have rapidly improved on software engineering tasks, with task length doubling approximately every seven months from 2019 through early 2025. As of late 2025, Claude Opus 4.5 achieved a time horizon of nearly five hours (METR, 2025). Notably, the models and configurations available to the public at any given time likely understate the frontier: major AI  
530 laboratories typically have more capable models in development that have not yet been released,

deploy longer context windows internally than are publicly available, and support customized agent setups that outperform out-of-the-box tools. Scaling laws suggest that model performance improves predictably with increases in data, parameters, and compute (Kaplan et al., 2020), and labs continue to invest heavily in all three. While the future remains uncertain, researchers can anticipate continued substantial growth in model capabilities.

535

## 5. CONCLUSION

This paper demonstrates that increasingly popular and powerful GPTs can be understood as extensions of familiar statistical tools: like OLS and GLMs, they estimate parameters by minimizing a loss function; like other neural networks, they chain nonlinear transformations to model complex relationships. The key innovations enabling modern text generation, tokenization, learned embeddings, and the attention mechanism, address the specific challenges of representing language and capturing dependencies across sequences. Understanding these foundations has practical value: researchers who grasp that models predict probability distributions over tokens, rather than retrieving facts from a database, are better positioned to anticipate failure modes like hallucination; those who understand context windows can structure prompts more effectively; and those who recognize that fine-tuning shapes behavior through human feedback can better interpret why models respond as they do. For other applications, some model APIs provide the predicted probabilities for output tokens, offering researchers a familiar statistical object to quantify uncertainty, assess model confidence, or construct more principled decision rules, and GPTs may also be applied outside text prediction. As LLMs become more integrated into research and society, continued engagement between researchers and the underlying methodology will help ensure these tools are applied appropriately and their limitations understood.

540

545

550

## A. APPENDIX

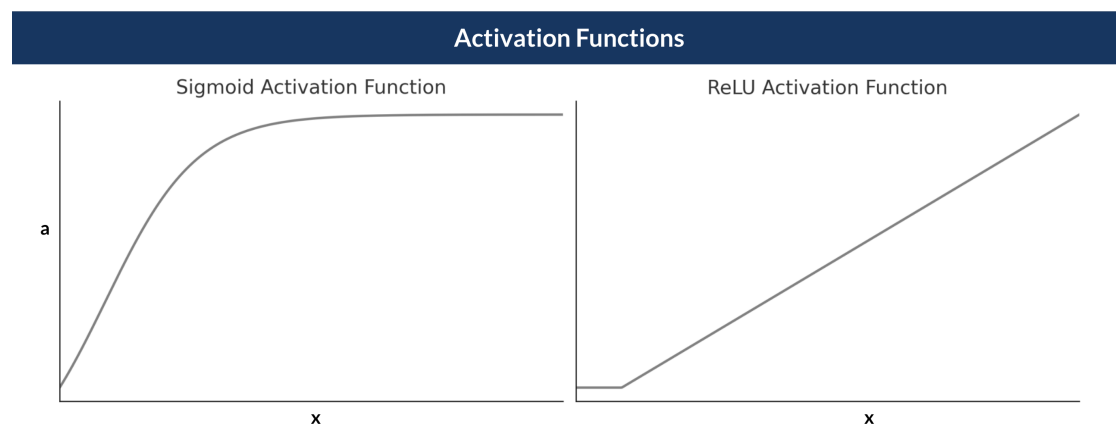


Fig. S1. Sigmoid and ReLU activation functions are both popular for neural networks. While a sigmoid activation function, which is differentiable (left), a ReLU function contains a "kink", below which the value is zero, and the function is linear above the kink.

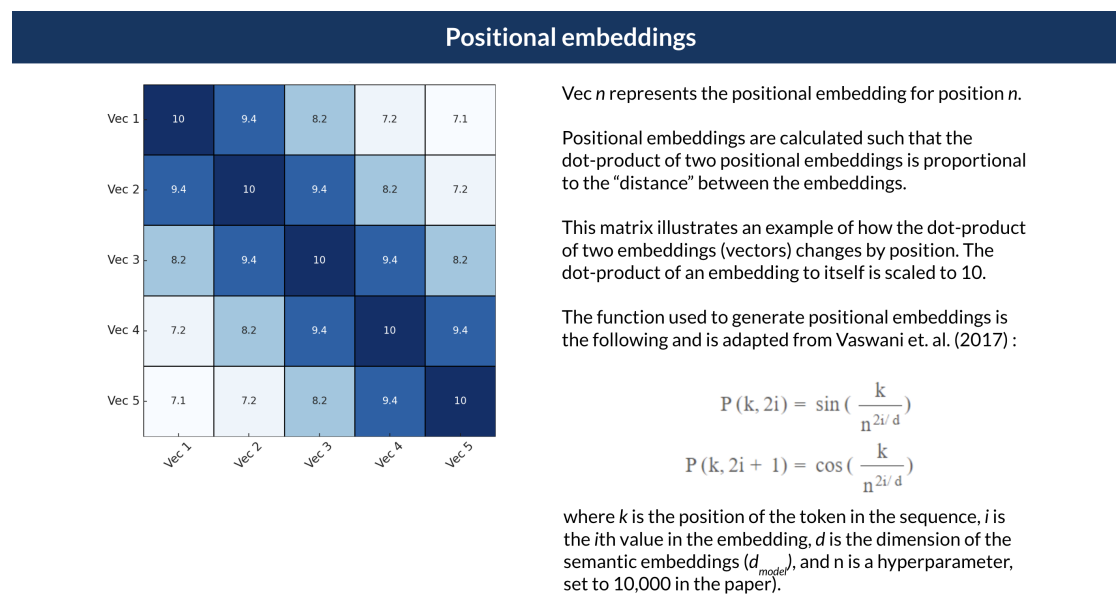


Fig. S2. Positional embeddings, calculated such that the dot product between two embeddings is proportional to the distance between the two positions.



## REFERENCES

- ABBAS, A., REHMAN, M. S. & REHMAN, S. S. (2024). Comparing the Performance of Popular Large Language Models on the National Board of Medical Examiners Sample Questions. *Cureus Publisher: Springer Science and Business Media LLC*. 555
- AMIDI, A. & AMIDI, S. (2024). CS 230 - Convolutional Neural Networks Cheatsheet.
- ANTHROPIC (2025). Models overview.
- BROWN, T., MANN, B., RYDER, N., SUBBIAH, M., KAPLAN, J. D., DHARIWAL, P., NEELAKANTAN, A., SHYAM, P., SASTRY, G., ASKELL, A., AGARWAL, S., HERBERT-VOSS, A., KRUEGER, G., HENIGHAN, T., CHILD, R., RAMESH, A., ZIEGLER, D., WU, J., WINTER, C., HESSE, C., CHEN, M., SIGLER, E., LITWIN, M., GRAY, S., CHESSE, B., CLARK, J., BERNER, C., MCCANDLISH, S., RADFORD, A., SUTSKEVER, I. & AMODEI, D. (2020). Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc. 560
- CHEN, Y., BENTON, J., RADHAKRISHNAN, A., UESATO, J., DENISON, C., SCHULMAN, J., SOMANI, A., HASE, P., WAGNER, M., ROGER, F., MIKULIK, V., BOWMAN, S. R., LEIKE, J., KAPLAN, J. & PEREZ, E. (2025). Reasoning Models Don't Always Say What They Think. *ArXiv:2505.05410 [cs]*. 565
- CHOWDHURY, A., NARANG, S., DEVLIN, J., BOSMA, M., MISHRA, G., ROBERTS, A., BARHAM, P., CHUNG, H. W., SUTTON, C., GEHRMANN, S. et al. (2022). Palm: Scaling language modeling with pathways. *arXiv preprint abs/2204.02311*.
- CHRISTIANO, P. F., LEIKE, J., BROWN, T., MARTIC, M., LEGG, S. & AMODEI, D. (2017). Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan & R. Garnett, eds., vol. 30. Curran Associates, Inc. 570
- COMMUNITY (2025). Terminology evolution? "completion" vs "response". Section: Community.
- CYBENKO, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems* **2**, 303–314. 575
- DAO, T., FU, D., ERMON, S., RUDRA, A. & RÉ, C. (2022). FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness. *Advances in Neural Information Processing Systems* **35**, 16344–16359.
- DEEPSEEK (2025). DeepSeek.
- DUGAS, D. (2023). The GPT-3 Architecture, on a Napkin.
- ERIKSEN, A. V., MÖLLER, S. & RYG, J. (2024). Use of GPT-4 to Diagnose Complex Clinical Cases. *NEJM AI* **1**, A1p2300031. Publisher: Massachusetts Medical Society. 580
- GADESHA, V. (2024). What is text generation?
- GOOGLE (2025). Gemini models | Gemini API.
- GRATTAFIORI, A., DUBEY, A., JAUHRI, A., PANDEY, A., KADIAN, A., AL-DAHLE, A., LETMAN, A., MATHUR, A., SCHELLEN, A., VAUGHAN, A., YANG, A., FAN, A., GOYAL, A., HARTSHORN, A., YANG, A., MITRA, A., SRIVANKUMAR, A., KORENEV, A., HINSVARK, A., RAO, A., ZHANG, A., RODRIGUEZ, A., GREGERSON, A., SPATARU, A., ROZIERE, B., BIRON, B., TANG, B., CHERN, B., CAUCHETEUX, C., NAYAK, C., BI, C., MARRA, C., MCCONNELL, C., KELLER, C., TOURET, C., WU, C., WONG, C., FERRER, C. C., NIKOLAIDIS, C., ALLONSIUS, D., SONG, D., PINTZ, D., LIVSHITS, D., WYATT, D., ESIOBU, D., CHOUDHARY, D., MAHAJAN, D., GARCIA-OLANO, D., PERINO, D., HUPKES, D., LAKOMKIN, E., ALBADAWY, E., LOBANOVA, E., DINAN, E., SMITH, E. M., RADENOVIC, F., GUZMÁN, F., ZHANG, F., SYNNAEVE, G., LEE, G., ANDERSON, G. L., THATTAI, G., NAIL, G., MIALON, G., PANG, G., CUCURELL, G., NGUYEN, H., KOREVAAR, H., XU, H., TOUVRON, H., ZAROV, I., IBARRA, I. A., KLOUMANN, I., MISRA, I., EVTIMOV, I., ZHANG, J., COPET, J., LEE, J., GEFFERT, J., VRANES, J., PARK, J., MAHADEOKAR, J., SHAH, J., LINDE, J. v. d., BILLOCK, J., HONG, J., LEE, J., FU, J., CHI, J., HUANG, J., LIU, J., WANG, J., YU, J., BITTON, J., SPISAK, J., PARK, J., ROCCA, J., JOHNSTUN, J., SAXE, J., JIA, J., ALWALA, K. V., PRASAD, K., UPASANI, K., PLAWIAK, K., LI, K., HEAFIELD, K., STONE, K., EL-ARINI, K., IYER, K., MALIK, K., CHIU, K., BHALLA, K., LAKHOTIA, K., RANTALA-YEARY, L., MAATEN, L. v. d., CHEN, L., TAN, L., JENKINS, L., MARTIN, L., MADAAN, L., MALO, L., BLECHER, L., LANDZAAT, L., OLIVEIRA, L. d., MUZZI, M., PASUPULETI, M., SINGH, M., PALURI, M., KARDAS, M., TSIMPOUKELLI, M., OLDHAM, M., RITA, M., PAVLOVA, M., KAMBADUR, M., LEWIS, M., SI, M., SINGH, M. K., HASSAN, M., GOYAL, N., TORABI, N., BASHLYKOV, N., BOGOYCHEV, N., CHATTERJI, N., ZHANG, N., DUCHENNE, O., ÇELEBI, O., ALRASSY, P., ZHANG, P., LI, P., VASIC, P., WENG, P., BHARGAVA, P., DUBAL, P., KRISHNAN, P., KOURA, P. S., XU, P., HE, Q., DONG, Q., SRINIVASAN, R., GANAPATHY, R., CALDERER, R., CABRAL, R. S., STOJNIC, R., RAILEANU, R., MAHESWARI, R., GIRDHAR, R., PATEL, R., SAUVESTRE, R., POLIDORO, R., SUMBALY, R., TAYLOR, R., SILVA, R., HOU, R., WANG, R., HOSSEINI, S., CHENNABASAPPA, S., SINGH, S., BELL, S., KIM, S. S., EDUNOV, S., NIE, S., NARANG, S., RAPARTHY, S., SHEN, S., WAN, S., BHOSALE, S., ZHANG, S., VANDENHENDE, S., BATRA, S., WHITMAN, S., SOOTLA, S., COLLOT, S., GURURANGAN, S., BORODINSKY, S., HERMAN, T., FOWLER, T., SHEASHA, T., GEORGIU, T., SCIALOM, T., SPECKBACHER, T., MIHAYLOV, T., XIAO, T., KARN, U., GOSWAMI, V., GUPTA, V., RAMANATHAN, V., KERKEZ, V., GONGUET, V., DO, V., VOGETI, V., ALBIERO, V., PETROVIC, V., CHU, W., XIONG, W., FU, W., MEERS, W., MARTINET, X., WANG, X., WANG, X., TAN, X. E., XIA, X., XIE, X., JIA, X., WANG, X., GOLDSCHLAG, Y., GAUR, Y., BABAEI, Y., WEN, Y., SONG, Y., ZHANG, Y., LI, Y., MAO, Y., COUDERT, Z. D., YAN, Z., CHEN, Z., PAKIPIOS, Z., SINGH, A., SRIVASTAVA, A., JAIN, A., KELSEY, A., SHAJNFELD, A., GANGIDI, A., VICTORIA, A., GOLDSTAND, A., MENON, A., SHARMA, A., BOESENBERG, A., BAEVSKI, A., FEINSTEIN, A., KALLET, A., SANGANI, A., TEO, A., YUNUS, A., LUPU, A., ALVARADO, A., CAPLES, A., GU, A., HO, A., POULTON, A., RYAN, A., RAMCHANDANI, A., DONG, A., FRANCO, A., GOYAL, A., SARAF, A., CHOWDHURY, A., GABRIEL, A., BHARAMBE, A., EISENMAN, A., YAZDAN, A., JAMES, B., MAURER, B., LEONHARDI, B., HUANG, B., LOYD, B., PAOLA,

- B. D., PARANJAPPE, B., LIU, B., WU, B., NI, B., HANCOCK, B., WASTI, B., SPENCE, B., STOJKOVIC, B., GAMIDO, B., MONTALVO, B., PARKER, C., BURTON, C., MEJIA, C., LIU, C., WANG, C., KIM, C., ZHOU, C., HU, C., CHU, C.-H., CAI, C., TINDAL, C., FEICHTENHOFER, C., GAO, C., CIVIN, D., BEATY, D., KREYMER, D., LI, D., ADKINS, D., XU, D., TESTUGGINE, D., DAVID, D., PARIKH, D., LISKOVICH, D., FOSS, D., WANG, D., LE, D., HOLLAND, D., DOWLING, E., JAMIL, E., MONTGOMERY, E., PRESANI, E., HAHN, E., WOOD, E., LE, E.-T., BRINKMAN, E., ARCAUTE, E., DUNBAR, E., SMOTHERS, E., SUN, F., KREUK, F., TIAN, F., KOKKINOS, F., OZGENEL, F., CAGGIONI, F., KANAYET, F., SEIDE, F., FLOREZ, G. M., SCHWARZ, G., BADEER, G., SWEE, G., HALPERN, G., HERMAN, G., SIZOV, G., GUANGYI, ZHANG, LAKSHMINARAYANAN, G., INAN, H., SHOJANAZERI, H., ZOU, H., WANG, H., ZHA, H., HABEEB, H., RUDOLPH, H., SUK, H., ASPEGREN, H., GOLDMAN, H., ZHAN, H., DAMLAJ, I., MOLYBOG, I., TUFANOV, I., LEONTIADIS, I., VELICHE, I.-E., GAT, I., WEISSMAN, J., GEBOSKI, J., KOHLI, J., LAM, J., ASHER, J., GAYA, J.-B., MARCUS, J., TANG, J., CHAN, J., ZHEN, J., REIZENSTEIN, J., TEBOUL, J., ZHONG, J., JIN, J., YANG, J., CUMMINGS, J., CARVILL, J., SHEPARD, J., MCPHIE, J., TORRES, J., GINSBURG, J., WANG, J., WU, K., U. K. H., SAXENA, K., KHANDELWAL, K., ZAND, K., MATOSICH, K., VEERARAGHAVAN, K., MICHELENA, K., LI, K., JAGADEESH, K., HUANG, K., CHAWLA, K., HUANG, K., CHEN, L., GARG, L., A. L., SILVA, L., BELL, L., ZHANG, L., GUO, L., YU, L., MOSHKOVICH, L., WEHRSTEDT, L., KHABSA, M., AVALANI, M., BHATT, M., MANKUS, M., HASSON, M., LENNIE, M., RESO, M., GROSHEV, M., NAUMOV, M., LATHI, M., KENEALLY, M., LIU, M., SELTZER, M. L., VALKO, M., RESTREPO, M., PATEL, M., VYATSKOV, M., SAMVELYAN, M., CLARK, M., MACEY, M., WANG, M., HERMOSO, M. J., METANAT, M., RASTEGARI, M., BANSAL, M., SANTHANAM, N., PARKS, N., WHITE, N., BAWA, N., SINGHAL, N., EGEBO, N., USUNIER, N., MEHTA, N., LAPTEV, N. P., DONG, N., CHENG, N., CHERNOGUZ, O., HART, O., SALPEKAR, O., KALINLI, O., KENT, P., PAREKH, P., SAAB, P., BALAJI, P., RITTNER, P., BONTRAGER, P., ROUX, P., DOLLAR, P., ZVYAGINA, P., RATANCHANDANI, P., YUVRAJ, P., LIANG, Q., ALAO, R., RODRIGUEZ, R., AYUB, R., MURTHY, R., NAYANI, R., MITRA, R., PARTHASARATHY, R., LI, R., HOGAN, R., BATTEY, R., WANG, R., HOWES, R., RINOTT, R., MEHTA, S., SIBY, S., BONDU, S. J., DATTA, S., CHUGH, S., HUNT, S., DHILLON, S., SIDOROV, S., PAN, S., MAHAJAN, S., VERMA, S., YAMAMOTO, S., RAMASWAMY, S., LINDSAY, S., LINDSAY, S., FENG, S., LIN, S., ZHA, S. C., PATIL, S., SHANKAR, S., ZHANG, S., ZHANG, S., WANG, S., AGARWAL, S., SAJUYIGBE, S., CHINTALA, S., MAX, S., CHEN, S., KEHOE, S., SATTERFIELD, S., GOVINDAPRASAD, S., GUPTA, S., DENG, S., CHO, S., VIRK, S., SUBRAMANIAN, S., CHOUDHURY, S., GOLDMAN, S., REMEZ, T., GLASER, T., BEST, T., KOEHLER, T., ROBINSON, T., LI, T., ZHANG, T., MATTHEWS, T., CHOU, T., SHAKED, T., VONTIMITTA, V., AJAYI, V., MONTANEZ, V., MOHAN, V., KUMAR, V. S., MANGLA, V., IONESCU, V., POENARU, V., MIHAILESCU, V. T., IVANOV, V., LI, W., WANG, W., JIANG, W., BOUAZIZ, W., CONSTABLE, W., TANG, X., WU, X., WANG, X., WU, X., GAO, X., KLEINMAN, Y., CHEN, Y., HU, Y., JIA, Y., QI, Y., LI, Y., ZHANG, Y., ZHANG, Y., ADI, Y., NAM, Y., YU, WANG, ZHAO, Y., HAO, Y., QIAN, Y., LI, Y., HE, Y., RAIT, Z., DEVITO, Z., ROSNBRICK, Z., WEN, Z., YANG, Z., ZHAO, Z. & MA, Z. (2024). The Llama 3 Herd of Models. ArXiv:2407.21783 [cs].
- GU, A. & DAO, T. (2023). Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*.
- GUO, D., YANG, D., ZHANG, H., SONG, J., WANG, P., ZHU, Q., XU, R., ZHANG, R., MA, S., BI, X., ZHANG, X., YU, X., WU, Y., WU, Z. F., GOU, Z., SHAO, Z., LI, Z., GAO, Z., LIU, A., XUE, B., WANG, B., WU, B., FENG, B., LU, C., ZHAO, C., DENG, C., RUAN, C., DAI, D., CHEN, D., JI, D., LI, E., LIN, F., DAI, F., LUO, F., HAO, G., CHEN, G., LI, G., ZHANG, H., XU, H., DING, H., GAO, H., QU, H., LI, H., GUO, J., LI, J., CHEN, J., YUAN, J., TU, J., QIU, J., LI, J., CAI, J. L., NI, J., LIANG, J., CHEN, J., DONG, K., HU, K., YOU, K., GAO, K., GUAN, K., HUANG, K., YU, K., WANG, L., ZHANG, L., ZHAO, L., WANG, L., ZHANG, L., XU, L., XIA, L., ZHANG, M., ZHANG, M., TANG, M., ZHOU, M., LI, M., WANG, M., LI, M., TIAN, N., HUANG, P., ZHANG, P., WANG, Q., CHEN, Q., DU, Q., GE, R., ZHANG, R., PAN, R., WANG, R., CHEN, R. J., JIN, R. L., CHEN, R., LU, S., ZHOU, S., CHEN, S., YE, S., WANG, S., YU, S., ZHOU, S., PAN, S., LI, S. S., ZHOU, S., WU, S., YUN, T., PEI, T., SUN, T., WANG, T., ZENG, W., LIU, W., LIANG, W., GAO, W., YU, W., ZHANG, W., XIAO, W. L., AN, W., LIU, X., WANG, X., CHEN, X., NIE, X., CHENG, X., LIU, X., XIE, X., LIU, X., YANG, X., LI, X., SU, X., LIN, X., LI, X. Q., JIN, X., SHEN, X., CHEN, X., SUN, X., WANG, X., SONG, X., ZHOU, X., WANG, X., SHAN, X., LI, Y. K., WANG, Y. Q., WEI, Y. X., ZHANG, Y., XU, Y., LI, Y., ZHAO, Y., SUN, Y., WANG, Y., YU, Y., ZHANG, Y., SHI, Y., XIONG, Y., HE, Y., PIAO, Y., WANG, Y., TAN, Y., MA, Y., LIU, Y., GUO, Y., OU, Y., WANG, Y., GONG, Y., ZOU, Y., HE, Y., XIONG, Y., LUO, Y., YOU, Y., LIU, Y., ZHOU, Y., ZHU, Y. X., HUANG, Y., LI, Y., ZHENG, Y., ZHU, Y., MA, Y., TANG, Y., ZHA, Y., YAN, Y., REN, Z. Z., REN, Z., SHA, Z., FU, Z., XU, Z., XIE, Z., ZHANG, Z., HAO, Z., MA, Z., YAN, Z., WU, Z., GU, Z., ZHU, Z., LIU, Z., LI, Z., XIE, Z., SONG, Z., PAN, Z., HUANG, Z., XU, Z., ZHANG, Z. & ZHANG, Z. (2025). DeepSeek-R1 incentivizes reasoning in LLMs through reinforcement learning. *Nature* **645**, 633–638. Publisher: Nature Publishing Group.
- HOFFMANN, J., BORGEAUD, S., MENSCH, A., BUCHATSKAYA, E., CAI, T., RUTHERFORD, E., DE LAS CASAS, D., HENDRICKS, L. A., WELBL, J., CLARK, A., HENNIGAN, T., NOLAND, E., MILLICAN, K., VAN DEN DRIESSCHE, G., DAMOC, B., GUY, A., OSINDERO, S., SIMONYAN, K., ELSSEN, E., RAE, J. W., VINYALS, O. & SIFRE, L. (2022). Training Compute-Optimal Large Language Models. In *Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc.
- HUANG, L., YU, W., MA, W., ZHONG, W., FENG, Z., WANG, H., CHEN, Q., PENG, W., FENG, X., QIN, B. & LIU, T. (2025). A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *ACM Transactions on Information Systems* **43**, 1–55. ArXiv:2311.05232 [cs].
- HUGGING FACE (2024). Byte-Pair Encoding tokenization - Hugging Face LLM Course.
- HUGGINGFACE (2025). What is Text Generation?

- KALAI, A. T., NACHUM, O., VEMPALA, S. S. & ZHANG, E. (2025). Why Language Models Hallucinate. ArXiv:2509.04664 [cs].
- KAPLAN, J., McCANDLISH, S., HENIGHAN, T., BROWN, T. B., CHESSE, B., CHILD, R., GRAY, S., RADFORD, A., WU, J. & AMODEI, D. (2020). Scaling Laws for Neural Language Models. ArXiv:2001.08361 [cs]. 680
- KATZ, U., COHEN, E., SHACHAR, E., SOMER, J., FINK, A., MORSE, E., SHREIBER, B. & WOLF, I. (2024). GPT versus Resident Physicians — A Benchmark Based on Official Board Scores. *NEJM AI* **1**, A1dbp2300192. Publisher: Massachusetts Medical Society.
- KISSANE, C., KRZYZANOWSKI, R., BLOOM, J. I., CONMY, A. & NANDA, N. (2024). Interpreting Attention Layer Outputs with Sparse Autoencoders. ArXiv:2406.17759 [cs]. 685
- KUNG, T. H., CHEATHAM, M., MEDENILLA, A., SILLOS, C., DE LEON, L., ELEPAÑO, C., MADRIAGA, M., AGGABAO, R., DIAZ-CANDIDO, G., MANINGO, J. & TSENG, V. (2023). Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models. *PLOS Digital Health* **2**, e0000198.
- KWA, T., WEST, B., BECKER, J., DENG, A., GARCIA, K., HASIN, M., JAWHAR, S., KINNIMENT, M., RUSH, N., ARX, S. V., BLOOM, R., BROADLEY, T., DU, H., GOODRICH, B., JURKOVIC, N., MILES, L. H., NIX, S., LIN, T., PARIKH, N., REIN, D., SATO, L. J. K., WIJK, H., ZIEGLER, D. M., BARNES, E. & CHAN, L. (2025). Measuring AI Ability to Complete Long Tasks. ArXiv:2503.14499 [cs]. 690
- LANHAM, T., CHEN, A., RADHAKRISHNAN, A., STEINER, B., DENISON, C., HERNANDEZ, D., LI, D., DURMUS, E., HUBINGER, E., KERNION, J., LUKOESIŪTĖ, K., NGUYEN, K., CHENG, N., JOSEPH, N., SCHIEFER, N., RAUSCH, O., LARSON, R., McCANDLISH, S., KUNDU, S., KADAVATH, S., YAO, S. & ZHANG, E. (2025). *of – ThoughtReasoning*. ArXiv : 2307.13702[cs].
- LI, J., TANG, T., ZHAO, W. X., NIE, J.-Y. & WEN, J.-R. (2024). PRE-TRAINED LANGUAGE MODELS FOR TEXT GENERATION: A SURVEY. *ACM COMPUT. SURV.* **56**, 230:1–230:39.
- McCULLAGH, P. (2018). *GENERALIZED LINEAR MODELS*. BOCA RATON: ROUTLEDGE. 695
- MERRIAM-WEBSTER (2025). HOW MANY WORDS ARE THERE IN ENGLISH?
- META (2025). LLAMA BY META.
- METR (2025). MEASURING AI ABILITY TO COMPLETE LONG TASKS. *METR BLOG*.
- MIKOLOV, T., CHEN, K., CORRADO, G. & DEAN, J. (2013A). EFFICIENT ESTIMATION OF WORD REPRESENTATIONS IN VECTOR SPACE. ArXiv:1301.3781 [cs]. 700
- MIKOLOV, T., YIH, W.-T. & ZWEIG, G. (2013B). LINGUISTIC REGULARITIES IN CONTINUOUS SPACE WORD REPRESENTATIONS. IN *PROCEEDINGS OF THE 2013 CONFERENCE OF THE NORTH AMERICAN CHAPTER OF THE ACL: HUMAN LANGUAGE TECHNOLOGIES (NAACL-HLT 2013)*.
- NG, A. & MA, T. (2023). CS229 LECTURE NOTES.
- NIELSEN, M. A. (2015). *NEURAL NETWORKS AND DEEP LEARNING* PUBLISHER: DETERMINATION PRESS. 705
- OPENAI (2022). INTRODUCING CHATGPT.
- OPENAI (2025A). MODELS - OPENAI API.
- OPENAI (2025B). TEXT GENERATION AND PROMPTING.
- OPENAI (2025C). TOKENIZER.
- OPENAI (2025D). WHAT ARE TOKENS AND HOW TO COUNT THEM? 710
- OPENAI (2026). OPENAI/TIKTOKEN. ORIGINAL-DATE: 2022-12-01T23:22:11Z.
- OUYANG, L., WU, J., JIANG, X., ALMEIDA, D., WAINWRIGHT, C. L., MISHKIN, P., ZHANG, C., AGARWAL, S., SLAMA, K., RAY, A., SCHULMAN, J., HILTON, J., KELTON, F., MILLER, L., SIMENS, M., ASKELL, A., WELINDER, P., CHRISTIANO, P., LEIKE, J. & LOWE, R. (2022). TRAINING LANGUAGE MODELS TO FOLLOW INSTRUCTIONS WITH HUMAN FEEDBACK. ArXiv:2203.02155 [cs]. 715
- PEEPPKORN, M., KOUWENHOVEN, T. & BROWN, DAN AND”; JORDANOUS, A. (2024). IS TEMPERATURE THE CREATIVITY PARAMETER OF LARGE LANGUAGE MODELS? IN *PROCEEDINGS OF THE 15TH INTERNATIONAL CONFERENCE ON COMPUTATIONAL CREATIVITY (ICCC 2024)*.
- PICHKA, E. (2025). WHAT IS QUERY, KEY, AND VALUE (QKV) IN THE TRANSFORMER ARCHITECTURE AND WHY ARE THEY USED? 720
- RADFORD, A., NARASIMHAN, K., SALIMANS, T. & SUTSKEVER, I. (2018). IMPROVING LANGUAGE UNDERSTANDING BY GENERATIVE PRE-TRAINING.
- RADFORD, A., WU, J., CHILD, R., LUAN, D., AMODEI, D. & SUTSKEVER, I. (2019). LANGUAGE MODELS ARE UNSUPERVISED MULTITASK LEARNERS.
- RAFAILOV, R., SHARMA, A., MITCHELL, E., ERMON, S., MANNING, C. D. & FINN, C. (2024). DIRECT PREFERENCE OPTIMIZATION: YOUR LANGUAGE MODEL IS SECRETLY A REWARD MODEL. 725
- RAFFEL, C., SHAZEER, N., ROBERTS, A., LEE, K., NARANG, S., MATENA, M., ZHOU, Y., LI, W. & LIU, P. J. (2020). EXPLORING THE LIMITS OF TRANSFER LEARNING WITH A UNIFIED TEXT-TO-TEXT TRANSFORMER. *JOURNAL OF MACHINE LEARNING RESEARCH* **21**, 1–67.
- SANDERSON, G. (2024). TRANSFORMERS, THE TECH BEHIND LLMs | DEEP LEARNING CHAPTER 5. 730
- SCHMIDT, R. M. (2019). RECURRENT NEURAL NETWORKS (RNNs): A GENTLE INTRODUCTION AND OVERVIEW. ArXiv:1912.05911 [cs].
- SENNRICH, R., HADDOW, B. & BIRCH, A. (2016). NEURAL MACHINE TRANSLATION OF RARE WORDS WITH SUBWORD UNITS. IN *PROCEEDINGS OF THE 54TH ANNUAL MEETING OF THE ACL (ACL 2016)*.

- SHAW, P., USZKOREIT, J. & VASWANI, A. (2018). SELF-ATTENTION WITH RELATIVE POSITION REPRESENTATIONS. IN *PROCEEDINGS OF NAACL-HLT*.
- STANFORD ONLINE (2024). STANFORD CS229 I MACHINE LEARNING I BUILDING LARGE LANGUAGE MODELS (LLMs).
- SU, J., AHMED, M., LU, Y., PAN, S., BO, W. & LIU, Y. (2024). RoFormer: ENHANCED TRANSFORMER WITH ROTARY POSITION EMBEDDING. *NEUROCOMPUTING* **568**, 127063.
- SU, J., LU, Y., PAN, S., WEN, B. & LIU, Y. (2021). RoFormer: ENHANCED TRANSFORMER WITH ROTARY POSITION EMBEDDING. *ARXIV PREPRINT ARXIV:2104.09864*.
- SUTSKEVER, I., MARTENS, J. & HINTON, G. E. (2011). GENERATING TEXT WITH RECURRENT NEURAL NETWORKS. IN *PROCEEDINGS OF THE 28TH INTERNATIONAL CONFERENCE ON MACHINE LEARNING (ICML 2011)*.
- TEAM, G., GEORGIEV, P., LEI, V. I., BURNELL, R., BAI, L., GULATI, A., TANZER, G., VINCENT, D., PAN, Z., WANG, S., MARIOORYAD, S., DING, Y., GENG, X., ALCOBER, F., FROSTIG, R., OMERNICK, M., WALKER, L., PADURARU, C., SOROKIN, C., TACCHETTI, A., GAFFNEY, C., DARUKI, S., SERCINOGLU, O., GLEICHER, Z., LOVE, J., VOIGTLAENDER, P., JAIN, R., SURITA, G., MOHAMED, K., BLEVINS, R., AHN, J., ZHU, T., KAWINTIRANON, K., FIRAT, O., GU, Y., ZHANG, Y., RAHTZ, M., FARUQUI, M., CLAY, N., GILMER, J., CO-REYES, J. D., PENCHEV, I., ZHU, R., MORIOKA, N., HUI, K., HARIDASAN, K., CAMPOS, V., MAHDIEH, M., GUO, M., HASSAN, S., KILGOUR, K., VEZER, A., CHENG, H.-T., LIEDEKERKE, R. D., GOYAL, S., BARHAM, P., STROUSE, D. J., NOURY, S., ADLER, J., SUNDARARAJAN, M., VIKRAM, S., LEPIKHIN, D., PAGANINI, M., GARCIA, X., YANG, F., VALTER, D., TREBACZ, M., VODRAHALLI, K., ASAWAROENGCHAI, C., RING, R., KALB, N., SOARES, L. B., BRAHMA, S., STEINER, D., YU, T., MENTZER, F., HE, A., GONZALEZ, L., XU, B., KAUFMAN, R. L., SHAFEEY, L. E., OH, J., HENNIGAN, T., DRIESSCHE, G. v. d., ODOOM, S., LUCIC, M., ROELOFS, B., LALL, S., MARATHE, A., CHAN, B., ONTANON, S., HE, L., TEPLYASHIN, D., LAI, J., CRONE, P., DAMOC, B., HO, L., RIEDEL, S., LENC, K., YEH, C.-K., CHOWDHERY, A., XU, Y., KAZEMI, M., AMID, E., PETRUSHKINA, A., SWERSKY, K., KHODAEI, A., CHEN, G., LARKIN, C., PINTO, M., YAN, G., BADIA, A. P., PATIL, P., HANSEN, S., ORR, D., ARNOLD, S. M. R., GRIMSTAD, J., DAI, A., DOUGLAS, S., SINHA, R., YADAV, V., CHEN, X., GRIBOVSKAYA, E., AUSTIN, J., ZHAO, J., PATEL, K., KOMAREK, P., AUSTIN, S., BORGEAUD, S., FRISO, L., GOYAL, A., CAINE, B., CAO, K., CHUNG, D.-W., LAMM, M., BARTH-MARON, G., KAGOHARA, T., OLSZEWSKA, K., CHEN, M., SHIVAKUMAR, K., AGARWAL, R., GODHIA, H., RAJWAR, R., SNAIDER, J., DOTIWALLA, X., LIU, Y., BARUA, A., UNGUREANU, V., ZHANG, Y., BATSAIKHAN, B.-O., WIRTH, M., QIN, J., DANIHELKA, I., DOSHI, T., CHADWICK, M., CHEN, J., JAIN, S., LE, Q., KAR, A., GURUMURTHY, M., LI, C., SANG, R., LIU, F., LAMPROU, L., MUNOZ, R., LINTZ, N., MEHTA, H., HOWARD, H., REYNOLDS, M., AROYO, L., WANG, Q., BLANCO, L., CASSIRER, A., GRIFFITH, J., DAS, D., LEE, S., SYGNOWSKI, J., FISHER, Z., BESLEY, J., POWELL, R., AHMED, Z., PAULUS, D., REITTER, D., BORSOS, Z., JOSHI, R., POPE, A., HAND, S., SELO, V., JAIN, V., SETHI, N., GOEL, M., MAKINO, T., MAY, R., YANG, Z., SCHALKWYK, J., BUTTERFIELD, C., HAUTH, A., GOLDIN, A., HAWKINS, W., SENTER, E., BRIN, S., WOODMAN, O., RITTER, M., NOLAND, E., GIANG, M., BOLINA, V., LEE, L., BLYTH, T., MACKINNON, I., REID, M., SARVANA, O., SILVER, D., CHEN, A., WANG, L., MAGGIORE, L., CHANG, O., ATTALURI, N., THORNTON, G., CHIU, C.-C., BUNYAN, O., LEVINE, N., CHUNG, T., ELTYSHEV, E., SI, X., LILLICRAP, T., BRADY, D., AGGARWAL, V., WU, B., XU, Y., MCLROY, R., BADOLA, K., SANDHU, P., MOREIRA, E., STOKOWIEC, W., HEMSLEY, R., LI, D., TUDOR, A., SHYAM, P., RAHIMTOROGHI, E., HAYKAL, S., SPRECHMANN, P., ZHOU, X., MINCU, D., LI, Y., ADDANKI, R., KRISHNA, K., WU, X., FRECHETTE, A., EYAL, M., DAFOE, A., LACEY, D., WHANG, J., AVRAHAMI, T., ZHANG, Y., TAROPA, E., LIN, H., TOYAMA, D., RUTHERFORD, E., SANO, M., CHOE, H., TOMALA, A., SAFRANEK-SHRADER, C., KASSNER, N., PAJARSKAS, M., HARVEY, M., SECHRIST, S., FORTUNATO, M., LYU, C., ELSAYED, G., KUANG, C., LOTTES, J., CHU, E., JIA, C., CHEN, C.-W., HUMPHREYS, P., BAUMLI, K., TAO, C., SAMUEL, R., SANTOS, C. N. D., ANDREASSEN, A., RAKIĆEVIĆ, N., GREWE, D., KUMAR, A., WINKLER, S., CATON, J., BROCK, A., DALMIA, S., SHEAHAN, H., BARR, I., MIAO, Y., NATSEV, P., DEVLIN, J., BEHBAHANI, F., PROST, F., SUN, Y., MYASKOVSKY, A., PILLAI, T. S., HURT, D., LAZARIDOU, A., XIONG, X., ZHENG, C., PARDO, F., LI, X., HORGAN, D., STANTON, J., AMBAR, M., XIA, F., LINCIE, A., WANG, M., MUSTAFA, B., WEBSON, A., LEE, H., ANIL, R., WICKE, M., DOZAT, T., SINHA, A., PIQUERAS, E., DABIR, E., UPADHYAY, S., BORAL, A., HENDRICKS, L. A., FRY, C., DJOLONGA, J., SU, Y., WALKER, J., LABANOWSKI, J., HUANG, R., MISRA, V., CHEN, J., SKERRY-RYAN, R. J., SINGH, A., RIJHWANI, S., YU, D., CASTRO-ROS, A., CHANGPINO, B., DATTA, R., BAGRI, S., HRAFNKELSSON, A. M., MAGGIONI, M., ZHENG, D., SULSKY, Y., HOU, S., PAINE, T. L., YANG, A., RIESA, J., ROGOZINSKA, D., MARCUS, D., BADAWY, D. E., ZHANG, Q., WANG, L., MILLER, H., GREER, J., SJOS, L. L., NOVA, A., ZEN, H., CHAABOUNI, R., ROSCA, M., JIANG, J., CHEN, C., LIU, R., SAINATH, T., KRIKUN, M., POLOZOV, A., LESPIAU, J.-B., NEWLAN, J., CANKARA, Z., KWAK, S., XU, Y., CHEN, P., COENEN, A., MEYER, C., TSIHLAS, K., MA, A., GOTTWEIS, J., XING, J., GU, C., MIAO, J., FRANK, C., CANKARA, Z., GANAPATHY, S., DASGUPTA, I., HUGHES-FITT, S., CHEN, H., REID, D., RONG, K., FAN, H., AMERSFOORT, J. v., ZHUANG, V., COHEN, A., GU, S. S., MOHANANEY, A., ILIC, A., TOBIN, T., WIETING, J., BORTSOVA, A., THACKER, P., WANG, E., CAVENESS, E., CHIU, J., SEZENER, E., KASKASOLI, A., BAKER, S., MILLICAN, K., ELHAWATY, M., AISOPPOS, K., LEBSACK, C., BYRD, N., DAI, H., JIA, W., WIETHOFF, M., DAVOODI, E., WESTON, A., YAGATI, L., AHUJA, A., GAO, I., PUNDAK, G., ZHANG, S., AZZAM, M., SIM, K. C., CAELLES, S., KEELING, J., SHARMA, A., SWING, A., LI, Y., LIU, C., BOSTOCK, C. G., BANSAL, Y., NADO, Z., ANAND, A., LIPSCHULTZ, J., KARMARKAR, A., PROLEEV, L., ITTYCHERIAH, A., YEGANEH, S. H., POLOVETS, G., FAUST, A., SUN, J., RRUSTEMI, A., LI, P., SHIVANNA, R., LIU, J., WELTY, C., LEBRON, F., BADDEPUDI, A., KRAUSE, S., PARISOTTO, E., SORICUT, R., XU, Z., BLOXWICH, D., JOHNSON, M., NEYSHABUR, B., MAO-JONES, J., WANG, R., RAMASESH, V., ABBAS, Z., GUEZ, A., SEGAL, C., NGUYEN, D. D., SVENSSON, J., HOU, L., YORK, S., MILAN, K., BRIDGERS, S., GWOREK, W., TAGLIASACCHI, M., LEE-THORP, J., CHANG, M., GUSEYNOV, A., HARTMAN, A. J., KWONG, M.,

ZHAO, R., KASHEM, S., COLE, E., MIECH, A., TANBURN, R., PHUONG, M., PAVETIC, F., CEVEY, S., COMANESCU, R., IVES, R., YANG, S., DU, C., LI, B., ZHANG, Z., IINUMA, M., HU, C. H., ROY, A., BIJWADIA, S., ZHU, Z., MARTINS, D., SAPUTRO, R., GERGELY, A., ZHENG, S., JIA, D., ANTONOGLU, I., SADOVSKY, A., GU, S., BI, Y., ANDREEV, A., SAMANGOEL, S., KHAN, M., KOCISKY, T., FILOS, A., KUMAR, C., BISHOP, C., YU, A., HODKINSON, S., MITTAL, S., SHAH, P., MOUFAREK, A., CHENG, Y., BLONIAZ, A., LEE, J., PEJMAN, P., MICHEL, P., SPENCER, S., FEINBERG, V., XIONG, X., SAVINOV, N., SMITH, C., SHAKERI, S., TRAN, D., CHESUS, M., BOHNET, B., TUCKER, G., GLEHN, T. V., MUIR, C., MAO, Y., KAZAWA, H., SLONE, A., SOPARKAR, K., SHRIVASTAVA, D., COBON-KERR, J., SHARMAN, M., PAVAGADHI, J., ARAYA, C., MISIUNAS, K., GHELANI, N., LASKIN, M., BARKER, D., LI, Q., BRIUKHOV, A., HOULSBY, N., GLAESE, M., LAKSHMINARAYANAN, B., SCHUCHER, N., TANG, Y., COLLINS, E., LIM, H., FENG, F., RECASENS, A., LAI, G., MAGNI, A., CAO, N. D., SIDDHANT, A., ASHWOOD, Z., ORBAY, J., DEHGHANI, M., BRENNAN, J., HE, Y., XU, K., GAO, Y., SAROUFIM, C., MOLLOY, J., WU, X., ARNOLD, S., CHANG, S., SCHRITTWIESER, J., BUCHATSKAYA, E., RADPOUR, S., POLACEK, M., GIORDANO, S., BAPNA, A., TOKUMINE, S., HELLENDORF, V., SOTTIAUX, T., COGAN, S., SEVERYN, A., SALEH, M., THAKOOR, S., SHEFEY, L., QIAO, S., GABA, M., CHANG, S.-Y., SWANSON, C., ZHANG, B., LEE, B., RUBENSTEIN, P. K., SONG, G., KWIATKOWSKI, T., KOOP, A., KANNAN, A., KAO, D., SCHUH, P., STJERNGREN, A., GHIASI, G., GIBSON, G., VILNIS, L., YUAN, Y., FERREIRA, F. T., KAMATH, A., KLIMENKO, T., FRANKO, K., XIAO, K., BHATTACHARYA, I., PATEL, M., WANG, R., MORRIS, A., STRUDEL, R., SHARMA, V., CHOY, P., HASHEMI, S. H., LANDON, J., FINKELSTEIN, M., JHAKRA, P., FRYE, J., BARNES, M., MAUGER, M., DAUN, D., BAATARSUKH, K., TUNG, M., FARHAN, W., MICHALEWSKI, H., VIOLA, F., QUITRY, F. D. C., LAN, C. L., HUDSON, T., WANG, Q., FISCHER, F., ZHENG, I., WHITE, E., DRAGAN, A., ALAYRAC, J.-B., NI, E., PRITZEL, A., IWANICKI, A., ISARD, M., BULANOVA, A., ZILKA, L., DYER, E., SACHAN, D., SRINIVASAN, S., MUCKENHIRN, H., CAI, H., MANDHANE, A., TARIQ, M., RAE, J. W., WANG, G., AYOUB, K., FITZGERALD, N., ZHAO, Y., HAN, W., ALBERTI, C., GARRETTE, D., KRISHNAKUMAR, K., GIMENEZ, M., LEVSKAYA, A., SOHN, D., MATAK, J., ITURRATE, I., CHANG, M. B., XIANG, J., CAO, Y., RANKA, N., BROWN, G., HUTTER, A., MIRROKNI, V., CHEN, N., YAO, K., EGYED, Z., GALILEE, F., LIECHTY, T., KALLAKURI, P., PALMER, E., GHEMAWAT, S., LIU, J., TAO, D., THORNTON, C., GREEN, T., JASAREVIC, M., LIN, S., COTRUTA, V., TAN, Y.-X., FIEDEL, N., YU, H., CHI, E., NEITZ, A., HEITKAEMPER, J., SINHA, A., ZHOU, D., SUN, Y., KAED, C., HULSE, B., MISHRA, S., GEORGAKI, M., KUDUGUNTA, S., FARABET, C., SHAFRAN, I., VLASIC, D., TSITSULIN, A., ANANTHANARAYANAN, R., CARIN, A., SU, G., SUN, P., V. S., CARVAJAL, G., BRODER, J., COMSA, I., REPINA, A., WONG, W., CHEN, W. W., HAWKINS, P., FILONOV, E., LOHER, L., HIRNSCHALL, C., WANG, W., YE, J., BURNS, A., CATE, H., WRIGHT, D. G., PICCININI, F., ZHANG, L., LIN, C.-C., GOG, I., KULIZHSKAYA, Y., SREEVATSA, A., SONG, S., COBO, L. C., IYER, A., TEKUR, C., GARRIDO, G., XIAO, Z., KEMP, R., ZHENG, H. S., LI, H., AGARWAL, A., NGANI, C., GOSHVADI, K., SANTAMARIA-FERNANDEZ, R., FICA, W., CHEN, X., GORGOLEWSKI, C., SUN, S., GARG, R., YE, X., ESLAMI, S. M. A., HUA, N., SIMON, J., JOSHI, P., KIM, Y., TENNEY, I., POTLURI, S., THIET, L. N., YUAN, Q., LUISIER, F., CHRONOPOULOU, A., SCELLATO, S., SRINIVASAN, P., CHEN, M., KOVERKATHU, V., DALIBARD, V., XU, Y., SAETA, B., ANDERSON, K., SELLAM, T., FERNANDO, N., HUOT, F., JUNG, J., VARADARAJAN, M., QUINN, M., RAUL, A., LE, M., HABALOV, R., CLARK, J., JALAN, K., BULLARD, K., SINGHAL, A., LUONG, T., WANG, B., RAJAYOGAM, S., EISENSCHLOS, J., JIA, J., FINCHLSTEIN, D., YAKUBOVICH, A., BALLE, D., FINK, M., AGARWAL, S., LI, J., DVIJOTHAM, D., PAL, S., KANG, K., KONZELMANN, J., BEATTIE, J., DOUSSE, O., WU, D., CROCKER, R., ELKIND, C., JONNALAGADDA, S. R., LEE, J., HOLTMANN-RICE, D., KALLARACKAL, K., LIU, R., VNUKOV, D., VATS, N., INVERNIZZI, L., JAFARI, M., ZHOU, H., TAYLOR, L., PRENDKI, J., WU, M., ECCLES, T., LIU, T., KOPPARAPU, K., BEAUFAYS, F., ANGERMUELLER, C., MARZOCA, A., SARCAR, S., DIB, H., STANWAY, J., PERBET, F., TRDIN, N., STERNECK, R., KHORLIN, A., LI, D., WU, X., GOENKA, S., MADRAS, D., GOLDSHTEIN, S., GIERKE, W., ZHOU, T., LIU, Y., LIANG, Y., WHITE, A., LI, Y., SINGH, S., BAHARGAM, S., EPSTEIN, M., BASU, S., LAO, L., OZTUREL, A., CROUS, C., ZHAI, A., LU, H., TUNG, Z., GAUR, N., WALTON, A., DIXON, L., ZHANG, M., GLOBERSON, A., UY, G., BOLT, A., WILES, O., NASR, M., SHUMAILOV, I., SELVI, M., PICCINNO, F., AGUILAR, R., MCCARTHY, S., KHALMAN, M., SHUKLA, M., GALIC, V., CARPENTER, J., VILLELA, K., ZHANG, H., RICHARDSON, H., MARTENS, J., BOSNJAK, M., BELLE, S. R., SEIBERT, J., ALNAHLAWI, M., MCWILLIAMS, B., SINGH, S., LOUIS, A., DING, W., POPOVICI, D., SIMICICH, L., KNIGHT, L., MEHTA, P., GUPTA, N., SHI, C., FATEHI, S., MITROVIC, J., GRILLS, A., PAGADORA, J., MUNKHDALAI, T., PETROVA, D., EISENBUD, D., ZHANG, Z., YATES, D., MITTAL, B., TRIPURANENI, N., ASSAEL, Y., BROVELLI, T., JAIN, P., VELIMIROVIC, M., AKBULUT, C., MU, J., MACHEREY, W., KUMAR, R., XU, J., QURESHI, H., COMANICI, G., WIESNER, J., GONG, Z., RUDDOCK, A., BAUER, M., FELT, N., GP, A., ARNAB, A., ZELLE, D., ROTHFUSS, J., ROSGEN, B., SHENOY, A., SEYBOLD, B., LI, X., MUDIGONDA, J., ERDOGAN, G., XIA, J., SIMSA, J., MICH, A., YAO, Y., YEW, C., KAN, S., CASWELL, I., RADEBAUGH, C., ELISSEFF, A., VALENZUELA, P., MCKINNEY, K., PATERSON, K., CUI, A., LATORRE-CHIMOTO, E., KIM, S., ZENG, W., DURDEN, K., PONNAPALLI, P., SOSEA, T., CHOQUETTE-CHOO, C. A., MANYIKA, J., ROBENEK, B., VASHISHT, H., PEREIRA, S., LAM, H., VELIC, M., OWUSU-AFRIYIE, D., LEE, K., BOLUKBASI, T., PARRISH, A., LU, S., PARK, J., VENKATRAMAN, B., TALBERT, A., ROSIQUE, L., CHENG, Y., SOZANSCHI, A., PASZKE, A., KUMAR, P., AUSTIN, J., LI, L., SALAMA, K., PERZ, B., KIM, W., DUKKIPATI, N., BARYSHNIKOV, A., KAPLANIS, C., SHENG, X., CHERVONYI, Y., UNLU, C., CASAS, D. D. L., ASKHAM, H., TUNYASUVUNAKOOL, K., GIMENO, F., PODER, S., KWAK, C., MIECNIKOWSKI, M., MIRROKNI, V., DIMITRIEV, A., PARISI, A., LIU, D., TSAI, T., SHEVLANE, T., KOURIDI, C., GARMON, D., GOEDECKEMEYER, A., BROWN, A. R., VIJAYAKUMAR, A., ELQURSH, A., JAZAYERI, S., HUANG, J., CATHY, S. M., HOOVER, J., KIM, L., KUMAR, S., CHEN, W., BILES, C., BINGHAM, G., ROSEN, E., WANG, L., TAN, Q., ENGEL, D., PONGETTI, F., CESARE, D. D., HWANG, D., YU, L., PULLMAN, J., NARAYANAN, S., LEVIN, K., GOPAL, S., LI, M., AHARONI, A., TRINH, T., LO, J., CASAGRANDE,

- N., VIJ, R., MATTHEY, L., RAMADHANA, B., MATTHEWS, A., CAREY, C. J., JOHNSON, M., GORANOVA, K., SHAH, R.,  
 860 ASHRAF, S., DASGUPTA, K., LARSEN, R., WANG, Y., VUYYURU, M. R., JIANG, C., IJAZI, J., OSAWA, K., SMITH, C.,  
 BOPANA, R. S., BILAL, T., KOIZUMI, Y., XU, Y., ALTUN, Y., SHABAT, N., BARIACH, B., KORCHEMNIY, A., CHOO, K.,  
 RONNEBERGER, O., IWUANYANWU, C., ZHAO, S., SOERGEL, D., HSIEH, C.-J., CAI, I., IQBAL, S., SUNDERMEYER, M.,  
 CHEN, Z., BURSSTEIN, E., MALAVIYA, C., BIADSY, F., SHROFF, P., DHILLON, I., LATKAR, T., DYER, C., FORBES, H.,  
 865 NICOSIA, M., NIKOLAEV, V., GREENE, S., GEORGIEV, M., WANG, P., MARTIN, N., SEDGHI, H., ZHANG, J., BANZAL,  
 P., FRITZ, D., RAO, V., WANG, X., ZHANG, J., PATRAUCEAN, V., DU, D., MORDATCH, I., JURIN, I., LIU, L., DUBEY,  
 A., MOHAN, A., NOWAKOWSKI, J., ION, V.-D., WEI, N., TOJO, R., RAAD, M. A., HUDSON, D. A., KESHAVA, V.,  
 AGRAWAL, S., RAMIREZ, K., WU, Z., NGUYEN, H., LIU, J., SEWAK, M., PETRINI, B., CHOI, D., PHILIPS, I., WANG, Z.,  
 BICA, I., GARG, A., WILKIEWICZ, J., AGRAWAL, P., LI, X., GUO, D., XUE, E., SHAIK, N., LEACH, A., KHAN, S. M.,  
 WIESINGER, J., JEROME, S., CHAKLADAR, A., WANG, A. W., ORNDUFF, T., ABU, F., GHAFKARKHAH, A., WAINWRIGHT,  
 870 M., CORTES, M., LIU, F., MAYNEZ, J., TERZIS, A., SAMANGOUEI, P., MANSOUR, R., KEPA, T., AUBET, F.-X., ALGYMR,  
 A., BANICA, D., WEISZ, A., ORBAN, A., SENEGES, A., ANDREJCZUK, E., GELLER, M., SANTO, N. D., ANKLIN, V.,  
 MEREY, M. A., BAEUML, M., STROHMAN, T., BAI, J., PETROV, S., WU, Y., HASSABIS, D., KAVUKCUOGLU, K., DEAN,  
 J. & VINYALS, O. (2024). GEMINI 1.5: UNLOCKING MULTIMODAL UNDERSTANDING ACROSS MILLIONS OF TOKENS OF  
 CONTEXT. ArXiv:2403.05530 [cs].
- 875 TOUVRON, H., MARTIN, L., STONE, K., ALBERT, P., ALMAHAIRI, A., BABAEI, Y., BASHLYKOV, N., BATRA, S., BHARGAVA,  
 P., BHOSALE, S., BIKEL, D., BLECHER, L., FERRER, C. C., CHEN, M., CUCURULL, G., ESIÖBU, D., FERNANDES, J.,  
 FU, J., FU, W., FULLER, B., GAO, C., GOSWAMI, V., GOYAL, N., HARTSHORN, A., HOSSEINI, S., HOU, R., INAN, H.,  
 KARDAS, M., KERKEZ, V., KHABSA, M., KLOUMANN, I., KORENEV, A., KOURA, P. S., LACHAUX, M.-A., LAVRIL, T.,  
 880 LEE, J., LISKOVICH, D., LU, Y., MAO, Y., MARTINET, X., MIHAYLOV, T., MISHRA, P., MOLYBOG, I., NIE, Y., POULTON,  
 A., REIZENSTEIN, J., RUNGTA, R., SALADI, K., SCHELLEN, A., SILVA, R., SMITH, E. M., SUBRAMANIAN, R., TAN, X. E.,  
 TANG, B., TAYLOR, R., WILLIAMS, A., KUANG, J. X., XU, P., YAN, Z., ZAROV, I., ZHANG, Y., FAN, A., KAMBADUR,  
 M., NARANG, S., RODRIGUEZ, A., STOJNIC, R., EDUNOV, S. & SCIALOM, T. (2023). LLAMA 2: OPEN FOUNDATION  
 AND FINE-TUNED CHAT MODELS. ArXiv:2307.09288 [cs].
- VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER, U. & POLOSUKHIN, I.  
 885 (2017). ATTENTION IS ALL YOU NEED. IN *ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS*, vol. 30. CURRAN  
 ASSOCIATES, INC.
- WEI, J., WANG, X., SCHUURMANS, D., BOSMA, M., ICHTER, B., XIA, F., CHI, E., LE, Q. V. & ZHOU, D. (2022). CHAIN-  
 OF-THOUGHT PROMPTING ELICITS REASONING IN LARGE LANGUAGE MODELS. *ADVANCES IN NEURAL INFORMATION  
 PROCESSING SYSTEMS* **35**, 24824–24837.