



Large Language Models (LLMs)

Large language models (LLMs) are very large neural network models trained on massive text datasets. They are designed to understand context and generate human-like language. Thanks to their scale (often billions of parameters) and extensive training data, LLMs can perform a wide range of language tasks – they can **infer meaning from context, produce coherent and relevant responses, translate between languages, summarize text, answer questions, and even assist in writing or coding** ¹. These models have essentially revolutionized natural language processing by achieving a level of fluency and contextual understanding that was previously unattainable with smaller models ¹. Notably, the trend in recent years has been to increase model size dramatically: for example, OpenAI's GPT-1 in 2018 had about 117 million parameters, GPT-2 in 2019 grew to 1.5 billion, and GPT-3 in 2020 jumped to **175 billion parameters** ². This rapid scaling of model size (and training data) has been a key factor in improving LLM capabilities.

The Transformer Architecture

Modern LLMs are almost all based on the **Transformer** architecture introduced by Vaswani et al. in 2017. The Transformer's key innovation is the **self-attention mechanism**, which allows the model to consider relationships between all words (tokens) in a sequence regardless of their positions ³. In each self-attention layer, every token generates a set of vectors – a **Query, Key, and Value** – and the model computes attention scores by matching the query against all keys to decide how much focus to give to each word's value ⁴. In essence, attention produces a *weighted average of the value vectors*, where the weights (attention scores) indicate relevance of other words to the current word ⁵. This mechanism lets the Transformer effectively capture context: each word's representation is updated by looking at *every other word* in the sentence, enabling the model to, for example, understand that in *"I arrived at the bank after crossing the river,"* the word "bank" likely refers to a river bank ³.

Another important aspect is **multi-head attention**. Instead of computing a single set of attention scores, the Transformer uses multiple attention "heads" in parallel. Each head can attend to different patterns or aspects of the input, which means the model learns a richer set of relationships within the sentence ⁶. The outputs of these heads are then concatenated and processed, allowing the model to integrate various contextual cues. The benefit is a more nuanced understanding – one head might focus on syntactic structure while another captures long-range dependencies, for instance ⁶.

The original Transformer architecture has an **Encoder-Decoder** structure. The **encoder** stack reads and encodes the input sequence (for example, a sentence in the source language for translation) into an intermediate representation. The **decoder** stack then takes that representation and generates an output sequence (e.g. the translated sentence in the target language), one token at a time. This encoder-decoder design proved highly effective for tasks like machine translation ³. However, large language models today often use **simplified variants** of this architecture depending on their goals: - Some models use **encoder-only** Transformers (just the stack of self-attention encoder layers). These are good for producing a contextualized understanding of text, which is useful in tasks like classification or question answering. **BERT** is a prime example of an encoder-only model ⁷. - Other models use **decoder-only** Transformers, which

generate text by predicting the next token. These models internalize language patterns and excel at text generation. **GPT (Generative Pre-trained Transformer)** models are decoder-only Transformers ⁸.

In a decoder-only model, the self-attention is typically masked so that the model can only attend to earlier positions (preventing it from “cheating” by looking ahead at future tokens during generation). In an encoder-only model, self-attention is bidirectional (unmasked), allowing the model to look at context on both left and right of a given token to fully understand it ⁹. The architecture choice (encoder vs. decoder) thus aligns with the use-case: encoders for understanding, decoders for generating.

Transfer Learning in NLP

A major reason for the rapid progress in LLMs is the **transfer learning** paradigm. Instead of training a separate model from scratch for every new NLP task, researchers train a large **general-purpose model** on an extremely large corpus of text (this is called **pre-training**). The model learns a broad understanding of language – it picks up grammar, facts, and some reasoning abilities from this pre-training. Then, for a specific task (say, sentiment analysis or named-entity recognition), the pre-trained model is **fine-tuned** on a much smaller, task-specific dataset. This fine-tuning adjusts the model's knowledge to the task at hand.

This approach is powerful for several reasons: - **Data Efficiency:** The pre-trained model has already learned general language features, so it needs far less task-specific data to achieve high performance ¹⁰. - **Faster Development:** Rather than training a huge model from scratch for each task, one can fine-tune an existing model in a shorter time (fewer parameters need updating) ¹⁰. - **Improved Performance:** Models that have “seen” a broad range of language during pre-training often generalize better and achieve **better results** on downstream tasks than models trained only on the smaller task data ¹⁰.

In practice, the fine-tuning process usually involves adding a small task-specific layer (often called a “head”) on top of the pre-trained model and then training the model on the task data for a few epochs. For example, if we have a pre-trained Transformer and we want to do **text classification**, we add a classification layer on the final hidden representation of the [CLS] token (in BERT, a special token for “classification”) and fine-tune on labeled examples. The key is that the **language knowledge learned during pre-training is transferred** to the new task ¹¹. This is why even with limited data, fine-tuned LLMs can perform impressively well on tasks like **sentiment analysis**, **question answering**, and **named entity recognition** ¹¹.

It's worth noting that **BERT** itself was an early success story of transfer learning in NLP. BERT was pre-trained on a huge corpus (more on this below) and then fine-tuned to achieve state-of-the-art results on a variety of benchmarks with relatively little task-specific data ¹¹. This demonstrated the benefits of transfer learning and kicked off the era of large pre-trained language models.

(Transfer learning has additional benefits in practice, such as leveraging models across domains or languages, but at its core, the idea is as described: pre-train on massive general data, then fine-tune on the specific task.)

BERT: Bidirectional Encoder Representations from Transformers

BERT is a landmark LLM released by Google in 2018. The name stands for *Bidirectional Encoder Representations from Transformers*. As the name implies, BERT uses the **encoder** part of the Transformer

architecture and is designed to learn deep bidirectional representations of text (meaning it considers context from both left and right of each word) ⁷ . This bidirectional understanding allows BERT to capture nuances of language that unidirectional models might miss.

Training data and method: BERT was trained on an extremely large text corpus: notably, **the entire English Wikipedia (about 2.5 billion words) plus a huge collection of books (BookCorpus, ~800 million words)** ⁷ . The training was *unsupervised*, using two novel pre-training tasks: - *Masked Language Modeling (MLM)*: BERT randomly masks some percentage of the words in each input sentence and learns to predict those masked words from context ¹² . For example, "The cat sat on the [MASK]." BERT should predict "mat." To do this, it must use both left context ("cat sat on the") and right context (".") – hence the bidirectional nature. - *Next Sentence Prediction (NSP)*: BERT is also trained on pairs of sentences and tasked with predicting if the second sentence is the actual next sentence in the original text or just a random unrelated sentence ¹³ . This teaches BERT to understand relationships between sentences, which is useful for tasks like question answering and text coherence.

Through these tasks, BERT develops a strong grasp of language structure and meaning. Google released BERT in two model sizes: **BERT Base** (110 million parameters) and **BERT Large** (340 million parameters) ⁷ . These were among the largest models of their time (especially BERT Large).

Fine-tuning and impact: Once BERT is pre-trained, it can be fine-tuned and applied to a variety of NLP tasks with excellent results. For example, BERT achieved state-of-the-art accuracy on question answering (SQuAD dataset) and natural language inference tasks by fine-tuning with just a few epochs of training. Importantly, fine-tuning BERT for **sentiment analysis** or other classification tasks is straightforward: you just add a classification layer on top of the encoder's output for the [CLS] token and train on your labeled dataset. BERT's pre-trained knowledge makes it extremely effective even with limited data ¹¹ .

One real-world application of BERT is in **Google Search**. In late 2019, Google announced it had incorporated BERT into its search engine to better understand user queries ¹⁴ ⁹ . This was a big deal – for the first time, Google Search could grasp the context of words in queries, rather than treating queries as a bag of keywords. For instance, in the query "2019 brazil traveler to USA need a visa," the word "to" is critical for understanding the meaning (the query is about a Brazilian traveling to the USA, not the other way around). Previously, the search algorithm might have overlooked such a common word, but with BERT, it understands that "to" indicates direction and completely changes the intent. BERT allowed Google to return more relevant results by interpreting the query in a more human-like way ¹⁵ . In fact, Google reported that BERT improved the understanding of **10% of all search queries** in English (a very significant improvement in such a mature product) ¹⁶ .

In summary, BERT is an **encoder-based LLM** that excels in understanding text. It demonstrated the power of pre-training on vast data and fine-tuning for specific tasks. BERT's success paved the way for an explosion of other pre-trained models (RoBERTa, ALBERT, DistilBERT, etc.), and it remains a strong baseline for many NLP applications ¹⁷ ⁷ .

OpenAI's GPT Models and the Rise of Generative LLMs

Around the same time BERT was introduced for understanding tasks, **OpenAI** was developing the GPT series focused on **text generation**. **GPT** stands for *Generative Pre-Trained Transformer*. Unlike BERT, which

uses an encoder, GPT models use the **Transformer's decoder** mechanism and are trained to predict the next word in a sequence (making them inherently good at generating fluent text) ⁸ .

The progression of GPT models is a story of scaling up: - **GPT-1 (2018)**: 117 million parameters. This model demonstrated the viability of the approach: even a 117M parameter Transformer, pre-trained on a large chunk of the internet (BooksCorpus and Wikipedia, among other text), could generate paragraphs of text that were somewhat coherent. - **GPT-2 (2019)**: 1.5 billion parameters. With an order of magnitude more parameters and more training data, GPT-2 could produce even more convincing text. It could generate news-like articles given a prompt, often requiring an expert to tell they were machine-written. OpenAI initially hesitated to release the full GPT-2 model out of concern for misuse (e.g. generating fake news), highlighting how powerful text generation was becoming. - **GPT-3 (2020)**: 175 billion parameters. This was a massive leap in scale and is the model that truly vaulted LLMs into public awareness. GPT-3 was trained on hundreds of billions of words (essentially most of the internet, plus books, Wikipedia, etc.), and the result was a model that can generate extremely fluent and often impressively relevant text on almost any topic. With GPT-3, users found that the model could not only write essays and answer questions, but also perform tasks like writing basic code, composing poems, or completing knowledge tasks – all without any task-specific training (this is called *zero-shot* or *few-shot learning*, where the model is prompted with some examples and completes the task).

Why GPT-3 was remarkable: Beyond the sheer size (175 billion parameters ¹⁸), GPT-3 demonstrated emergent abilities. It could handle tasks it wasn't explicitly trained for, just by understanding the task described in the prompt. For instance, if prompted with "Translate the following English sentence to French: 'Where is the library?'", GPT-3 would produce a correct translation, despite not being a dedicated translation model. The quality of its English text generation was a big step up from GPT-2 – it uses more natural vocabulary, makes fewer grammatical errors, and stays on topic more reliably. As one researcher noted when comparing GPT-3 to its predecessor, *scaling up the model made it "a little better at English grammar, a little better at trivia questions," and eventually these improvements start to look like general reasoning abilities* ¹⁹ . In other words, by the time we scale to GPT-3, the model's outputs often feel startlingly coherent and even knowledgeable. Many users interacting with GPT-3 for the first time described the experience as borderline uncanny, because it could produce paragraphs of text that read as if a human wrote them on a wide variety of subjects.

It's important to note that GPT-3 (like GPT-2 and GPT-1) is a **decoder-only Transformer** model ⁸ . During training, it learns to predict the next token given all previous tokens in the sequence. There is no bidirectional context in the decoder (it can't peek at future words), which is appropriate for generation tasks. However, because it's so large and trained on so much data, GPT-3 acquires a broad base of world knowledge and linguistic context. When you prompt GPT-3 with some text, it continues that text in whatever style or context you set. This makes it extremely flexible – essentially a *general-purpose text generation engine*. People have used GPT-3 (and its successors like GPT-3.5 and GPT-4) for an astonishing array of applications: writing creative fiction, generating dialog for chatbots, summarizing documents, creating website copy, drafting emails, answering questions, tutoring in various subjects, and more.

Quality of English: The user's note specifically mentioned GPT-3 "improved the quality of English." Indeed, observers have commented that GPT-3's outputs are more fluent and grammatically correct than those of smaller models. The improvement is not just in surface-level grammar; GPT-3 is better at maintaining context over a long passage and producing relevant, sensible continuations. To give an example highlighted by Vox media: GPT-2 was "*pretty good*" at generating a few paragraphs of text, but it might lose track of the

narrative or make obvious mistakes ²⁰. GPT-3, with its greater capacity, is “a lot smarter” – it can carry on for longer, stick to a given style or story, and even make reasonable arguments or analogies ²¹. This doesn't mean GPT-3 is perfect (it certainly can produce errors or nonsensical answers, and it has no true understanding or conscience behind its words), but the leap in quality from GPT-2 to GPT-3 was widely recognized as a **milestone** in AI. It became much harder to distinguish GPT-3's writing from a human's in many cases ¹⁹.

Beyond GPT-3: Following GPT-3, the field saw even larger models and specialized variants: - **Microsoft/Nvidia's Megatron-Turing NLG (530B parameters, 2021)**, - **DeepMind's Gopher (280B, 2021)**, - **Google's PaLM (540B, 2022)**, etc., and eventually OpenAI's **GPT-4 (2023)** which is not just larger but also more fine-tuned and aligned with user intent. Each of these built on the idea that *scaling up* and *more training data* can unlock new levels of performance. There are also **open-source LLMs** like Meta's *LLaMA* series (with variants ranging from 7B to 65B parameters) showing that even smaller LLMs, if trained on the right data, can be very effective.

In summary, the GPT series exemplifies the generative side of LLMs. GPT models use the Transformer decoder to **generate text**, and by pre-training on large swaths of the internet, they acquire the ability to produce human-like text for a variety of tasks. GPT-3, in particular, highlighted how increasing model size dramatically improved the **fluency and correctness of generated English** ¹⁹, which in turn has had huge implications for applications of AI in writing and communication.

Concluding Remarks

To reorganize the key points from these notes: - **Transformers** are the backbone of modern large language models, introducing self-attention and a flexible architecture (encoder/decoder) that can capture context effectively. - **Large Language Models (LLMs)** leverage this architecture at an unprecedented scale – they are pre-trained on enormous datasets, which gives them a broad competence in language. - **Transfer Learning** is essential: LLMs are first pre-trained (which is computationally expensive and data-intensive), then fine-tuned to specific tasks, enabling faster development and strong performance even with limited task data ¹⁰. - **BERT** is an example of an encoder-based LLM aimed at understanding language. It's bidirectional, trained on masked word prediction and next-sentence prediction, and has been used in everything from research benchmarks to improving Google's search results ¹⁵. - **GPT** exemplifies decoder-based LLMs for text generation. The jump to GPT-3's 175B parameters showed that with enough data and size, a model can generate remarkably human-like text and handle diverse tasks without explicit training for each ¹⁹ ¹⁸.

As you continue your AI learning, remember that this field is evolving quickly. New models (like those incorporating **multimodal inputs** or improved training techniques) are emerging, but the foundation you've captured in these notes – Transformers, large-scale pre-training, and fine-tuning – will remain central to understanding how AI is able to read and write language so effectively. Good luck with your studies!

Sources:

- Vaswani et al., “*Attention Is All You Need*” (NeurIPS 2017) – introduced the Transformer architecture.
- Google AI Blog – “*Transformer: A Novel Neural Network Architecture for Language Understanding*” ³.
- Google Blog – “*Understanding searches better than ever before*” (Oct 2019), on applying BERT to search ¹⁵.

- Devlin et al., “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding” (2018) ^{7 12} .
- OpenAI – “GPT-3” announcement and various analyses (2020); Vox coverage of GPT-3’s capabilities ¹⁹ .
- IBM Developer – “What are large language models?” (2023), overview of LLM capabilities ¹ .
- Medium (Anil George) – “Visualizing the size of Large Language Models” (2023), data on model sizes ² .

¹ What Are Large Language Models (LLMs)? | IBM

<https://www.ibm.com/think/topics/large-language-models>

² Visualizing the size of Large Language Models | by Anil George | Medium

<https://medium.com/@georgeanil/visualizing-size-of-large-language-models-ec576caa5557>

³ Transformer: A Novel Neural Network Architecture for Language Understanding

<https://research.google/blog/transformer-a-novel-neural-network-architecture-for-language-understanding/>

^{4 5} neural networks - What exactly are keys, queries, and values in attention mechanisms? - Cross Validated

<https://stats.stackexchange.com/questions/421935/what-exactly-are-keys-queries-and-values-in-attention-mechanisms>

⁶ Exploring Multi-Head Attention: Why More Heads Are Better Than One | by Hassaan Idrees | Medium

<https://medium.com/@hassaanidrees7/exploring-multi-head-attention-why-more-heads-are-better-than-one-006a5823372b>

^{7 11 12 13 17} BERT (language model) - Wikipedia

[https://en.wikipedia.org/wiki/BERT_\(language_model\)](https://en.wikipedia.org/wiki/BERT_(language_model))

^{8 18} GPT-3 - Wikipedia

<https://en.wikipedia.org/wiki/GPT-3>

^{9 15 16} Understanding searches better than ever before

<https://blog.google/products/search/search-language-understanding-bert/>

¹⁰ Transfer Learning vs. Fine-Tuning: Unlocking the Power of Pretrained Models in Modern Machine Learning

<https://www.linkedin.com/pulse/transfer-learning-vs-fine-tuning-unlocking-power-models-mohanty-vp1hf>

¹⁴ Timeline of AI and language models – Dr Alan D. Thompson – LifeArchitect.ai

<https://lifearchitect.ai/timeline/>

^{19 20 21} GPT-3, explained: OpenAI’s new language AI is uncanny, funny- and a big deal | Vox

<https://www.vox.com/future-perfect/21355768/gpt-3-ai-openai-turing-test-language>