

The Emergence of a New Computational Paradigm in Natural Language Processing: A Review of Architectures, Adaptation, and Applications of Large Language Models

Surajo Nuhu Umar^{1*}; Aliyu Ishaq Abdullahi²;
Muhammad Abdulrazak Rabi³; Abdulrahman Rabi⁴ Umar⁴

^{1,2,3}Department of Computer Science and Engineering, Shobhit Institute of Engineering, India.

⁴Department of Computer Science and Engineering, Kalinga University Raipur, India.

Corresponding Author: Surajo Nuhu Umar^{1*}

Publication Date: 2026/02/14

Abstract: Large Language Models have become a transformational element in Natural Language Processing because they introduce new approach for understanding and generating languages. This paper is a formal review of the development of Large Language Models from a variety of perspectives, including the architectural advances, pre training strategies, and adaptation techniques. The paper focuses on the process of moving from early contextual word representations to large scale transformer based systems trained using very large collections of written language, describing significant advancements in model architecture, pretraining methods, and techniques to adapt the models for future tasks. Furthermore, the major applications, including text summarization, translation, dialogue systems, information extraction, and question answering are discussed. The paper further analyzes critical challenges such as computational scalability, data requirements, model alignment, inference efficiency, ethical concerns, and deployment limitations.

Keywords: Large Language Models; Natural Language Processing, Transformer Based Systems, Model Architecture.

How to Cite: Surajo Nuhu Umar; Aliyu Ishaq Abdullahi; Muhammad Abdulrazak Rabi; Abdulrahman Rabi Umar (2026) The Emergence of a New Computational Paradigm in Natural Language Processing: A Review of Architectures, Adaptation, and Applications of Large Language Models. *International Journal of Innovative Science and Research Technology*, 11(2), 509-517. <https://doi.org/10.38124/ijisrt/26feb382>

I. INTRODUCTION

Natural language processing (NLP) has seen a revolution thanks to large language models (LLMs), which have demonstrated previously unheard-of performance on a variety of tasks. These models can understand and produce language that is similar to that of humans because they have billions of parameters that have been trained on large text corpora [1].

The paradigm shift started in the late 2010s when deep contextualized word representations such as ELMo (Peters et al., 2018) showed that rich linguistic features that greatly enhance downstream NLP tasks could be produced through pre-training on large text data [2]. Shortly after, a new era was ushered in by the introduction of pre-trained Transformers: BERT (Devlin et al., 2018) demonstrated how a bidirectional Transformer trained on unsupervised text could be optimized to attain cutting-edge performance on language

comprehension tasks [3]. Simultaneously, the GPT series investigated unidirectional (autoregressive) Transformers for language production. OpenAI's GPT-2 (2019) and GPT-3 (2020) significantly increased the size of their models (to 1.5 billion and then 175 billion parameters), and GPT-3 showed impressive few-shot learning, which is the capacity to complete tasks using only prompt examples and no fine-tuning [4]. GPT-3's 175B model, which was trained on nearly the entire Internet, performed well on tasks like question answering and translation just by being prompted [4]. These results demonstrated that LLMs have an emergent capability, meaning that they can "learn" from context and adapt to new tasks when given only natural language prompts. This led to the idea of "foundation models," which are big models that can be tailored to a wide range of tasks after being trained on extensive data [5].

The Stanford CRFM report from 2021 formally defined these as models such as GPT-3, BERT, and others, pointing out their significant but unfinished nature: at scale, they demonstrate new capabilities (like reasoning and coherent generation), but their behavior can be unpredictable and inherit any defects in the training data [5]. Since then, there has been an arms race in the field's scaling, with models growing from tens of billions to hundreds of billions of parameters. For example, Google's PaLM (540B) in 2022 achieved groundbreaking few-shot results on numerous benchmarks and even outperformed average human performance on some reasoning tasks [6]. The community sought openness in addition to closed-source initiatives, with

initiatives like LLaMA and BLOOM (176B) making large models available to researchers. Both enthusiasm about LLMs' potential applications in NLP and worries about their drawbacks and moral ramifications have accompanied their quick development. In this review, we thoroughly look at the development of large language models, the technical frameworks that make them possible, the diverse range of applications they support, the advantages they offer, and the difficulties they present, such as concerns about safety, bias, and accuracy. In order to present an academic viewpoint on the state-of-the-art for LLMs, we base our analysis on recent literature (2018-2024) and reference peer-reviewed research and reliable reports.

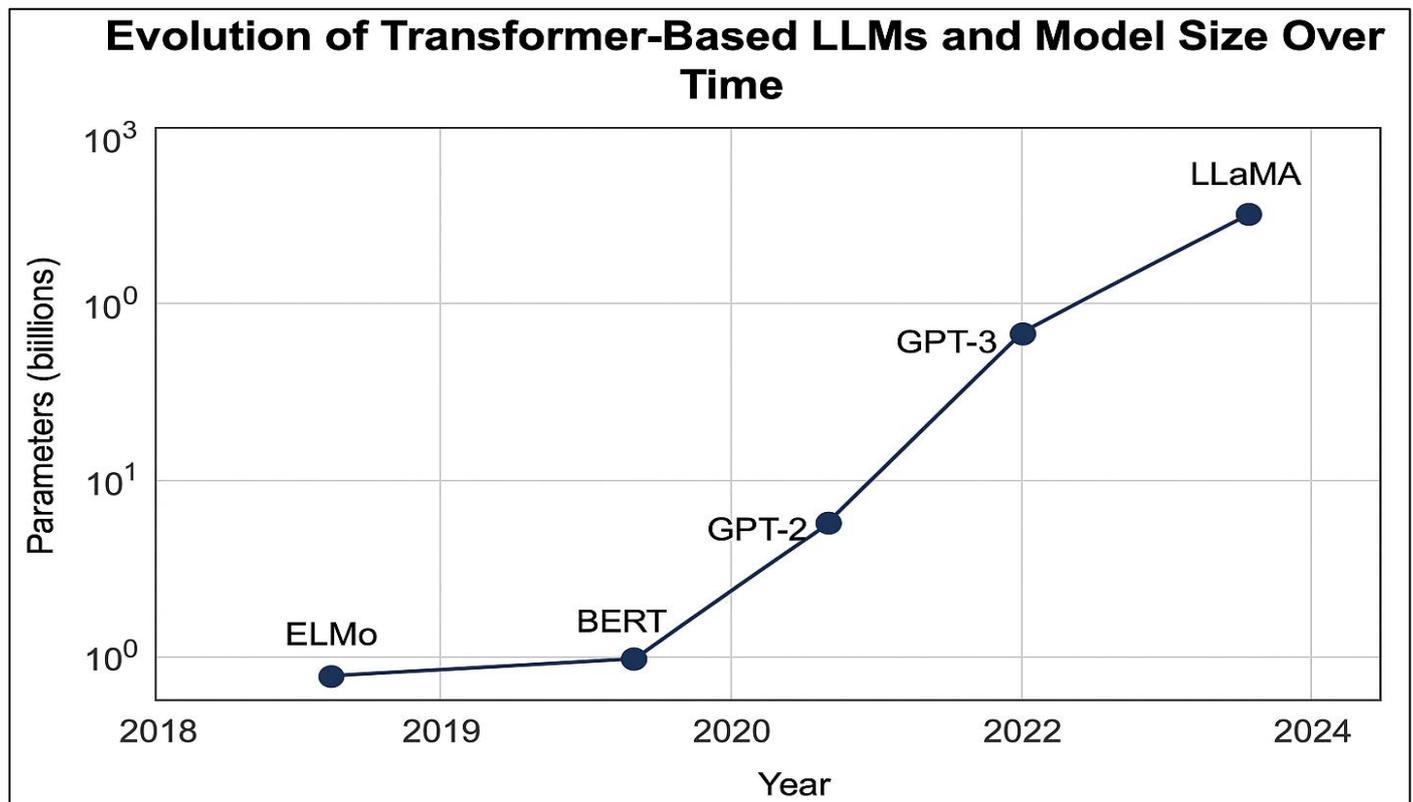


Fig 1 The Evolution of LLMs

II. FOUNDATIONAL ADVANCES IN LARGE LANGUAGE MODELS

➤ *Rise of Pre-trained LLMs and Foundation Models*

Modern LLMs are built on a foundation of prior to training on large text corpora and job adaption. Models that used unlabeled data to develop deep latent representations of language were the first to achieve success. BERT was one of the first to implement bidirectional Transformer encoding of text. BERT can collect rich context from both the left and right of a word because it is taught with a masked language modeling target to predict hidden words [3]. With very minor architecture modifications, BERT could be optimized upon release to attain cutting-edge performance on a variety of applications, including natural language inference and question answering [3]. A surge of research into ever-larger and better pre-trained models was sparked by BERT's success, which validated the paradigm of large-scale pre-training + fine-tuning.

The generative capabilities of autoregressive Transformers were showcased by OpenAI's GPT models. The scale of GPT-2 (2019) (1.5B parameters, the largest of its kind at the time) and its unexpected capacity to produce coherent, fluid text passages made it noteworthy. GPT-2 was dubbed a "unsupervised multitask learner" after it was demonstrated to complete tasks like reading comprehension and translation in a zero-shot fashion (without gradient updates) after being trained on 40 GB of Internet text. The capabilities of GPT-3 were significantly increased. After scaling the model to 175 billion parameters, Brown et al. (2020) found that GPT-3 performs well regardless of the task: without any fine-tuning and with just a task description or a few examples provided in the prompt, GPT-3 matched or outperformed previous state-of-the-art results on numerous benchmarks [7]. By using prompting, for example, GPT-3 showed proficient language translation, open-domain question responding, and even basic math and reasoning [7]. Known as "in-context learning," this few-shot prompt-based

learning was a startling emergent trait of scale. The distinction between pre-training and "learning" at inference time becomes hazy, as it was proposed that very large models

can encode both language information and the capacity to modify behavior based on context.

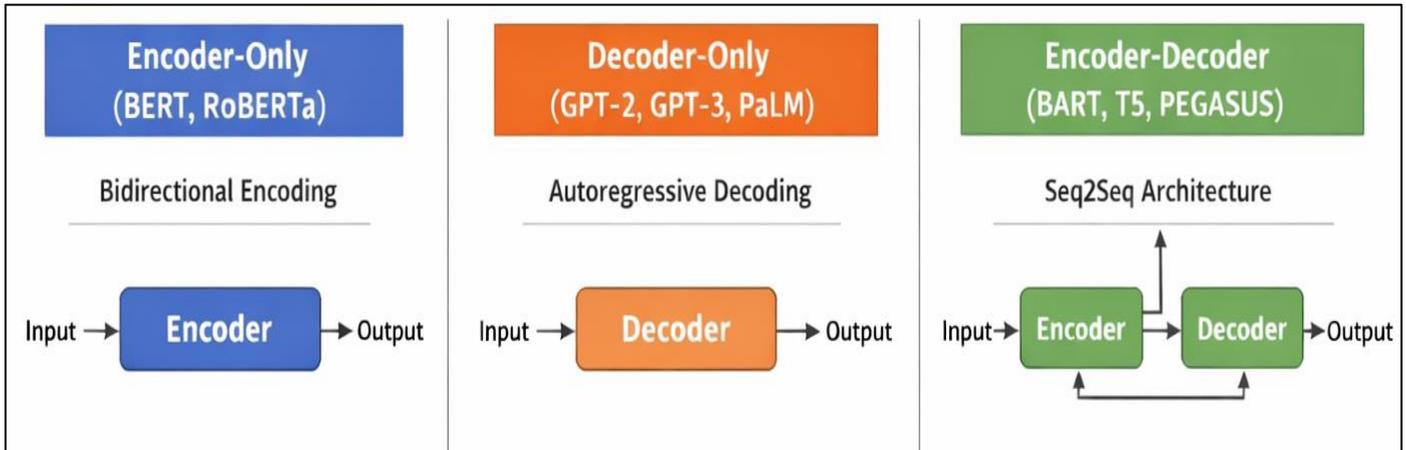


Fig 2 Comparison of Transformer Architecture Types Used in LLMs

The race for scale went on. Researchers looked into scaling laws for LMs. Kaplan et al. (2020) shown that performance increases predictably with an increase in model parameters, training data, and compute [8]. Importantly, they concluded that decreasing returns only occurred gradually and that the greatest outcomes might be obtained by training the largest models within a compute budget [8]. This knowledge was furthered in 2022 by Hoffmann et al. (DeepMind) with the Chinchilla study, where they found that there is an ideal trade-off between model size and data for a given compute budget and that many large models were under-trained in comparison to their size [4].

As a result of scaling up Transformer networks and utilizing massive unlabeled datasets, LLMs have become increasingly popular, producing models that function as general-purpose language learners. With the right adaption, these foundation models may be used for almost every language task, bringing NLP together. Previously task-specific architectures were needed for certain tasks, but today the same pre-trained model can frequently tackle them. How these models are specifically modified and matched to certain applications is covered in the following paragraph.

Table 1 Summary of Key Large Language Models Discussed in the Review

Model	Architecture Type	Parameters	Key Contribution / Feature	Source
ELMo	BiLSTM Contextual Encoder		Introduced deep contextualized word representations	[2]
BERT	Encoder-only Transformer	340M	Bidirectional masked language modeling	[3]
GPT-2	Decoder-only Transformer	1.5B	Unsupervised multitask generation	[4]
GPT-3	Decoder-only Transformer	175B	Few-shot & in-context learning	[4]
PaLM	Decoder-only Transformer	540B	Advanced reasoning, large-scale performance	[6]
Chinchilla	Optimally trained LM	70B	Compute-optimal training using 1.4T tokens	[15]
NLLB-54B	Mixture-of-Experts LLM	54B	Multilingual translation for 200 languages	[10]

➤ *Adaptation and Alignment of LLMs*

Pre-trained LLMs collect a lot of factual and linguistic information, but an equally important aspect of the LLM paradigm is aligning and modifying these models to meet the demands (and values) of users. The simplest adaptation method is fine-tuning, which involves using supervised learning to further train an already-trained LLM on a task-specific dataset. By adding a task-specific output layer and training on labeled data, BERT was successful through fine-tuning [3], and when there is a sufficient amount of supervised data available, it continues to work well for LLMs. For instance, huge models optimized on translation corpora or summary datasets continue to perform well in those tasks [9] [10]. However, because of the computational expense and the possibility of overfitting or catastrophic forgetting, it is frequently unfeasible to fine-tune a massive model for every task.

Thus, researchers have looked into several adaption strategies. One family of techniques learns tiny extra parameters, like adaptor layers or prompt embeddings, that guide the model on a new job while mostly or completely freezing the overall model's parameters. This comprises LoRA (Low-Rank Adaptation), which introduces trainable low-rank matrices into the Transformer layers, and prompt tuning/prefix tuning, which involves optimizing a set of virtual tokens or a prefix supplied into the model to condition it for the job. Even for 100B+ models on consumer hardware, fine-tuning is possible because to these techniques, which enable effective adaption of LLMs with only a small portion of the original model's parameters requiring changes. Prompt-based learning can perform well, although it occasionally falls short of complete fine-tuning. The ability to condition LLMs through prompts is a distinct strength, though, as

demonstrated by GPT-3, which shows that an LLM can frequently complete a task without any parameter updates when prompts are cleverly designed (or a few examples in context) [4]. Thus, prompt engineering has evolved into an art

and science that involves figuring out the best ways to induce the desired behavior from a fixed model through examples or instructions.

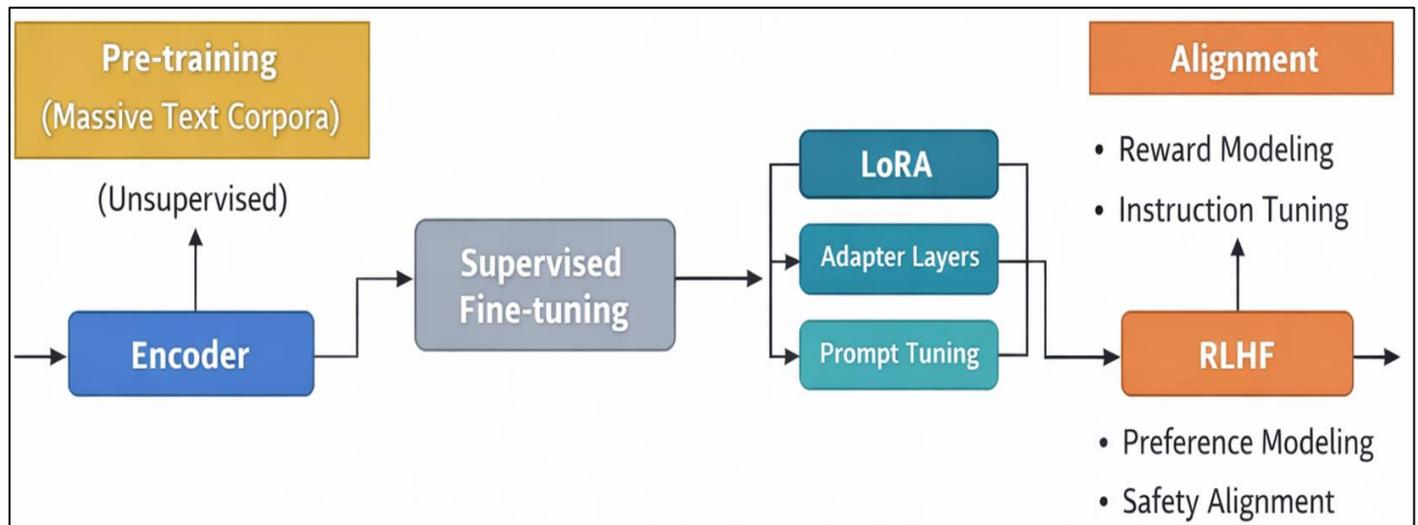


Fig 3 Adaptation and Alignment Pipeline for LLMs

III. LLM-BASED APPLICATIONS IN NLP

➤ Abstractive Text Summarization

The goal of summarization is to reduce a document's length while maintaining its essential content. The generative power of LLMs and their capacity to capture distant context are advantageous for this task. News articles were used to train sequence-to-sequence networks for the first neural summarization models, but these networks frequently had problems with coherence and factual accuracy. Abstractive summarization has significantly improved with large pre-trained models. Pre-training a model designed for summarization is one method, while fine-tuning a general LLM (such as GPT-3 or T5) on summarization datasets is another. The latter was best illustrated by PEGASUS (Zhang et al., 2020), which added a pre-training objective in which entire sentences are masked (as "gap sentences") so that the model can produce those as a summary [9]. Following this pre-training, PEGASUS improved on 12 summarization tasks and achieved state-of-the-art results on all of them, even in low-resource environments. Its summaries were frequently rated similarly to those written by humans [9]. Another well-known model is BART [11], a massive Transformer with 400 million parameters that has been trained to function as a denoising autoencoder by corrupting text and then learning to reconstruct it. Utilizing a bidirectional encoder and autoregressive decoder, BART expands on the concepts of BERT and GPT and has demonstrated remarkable efficacy when optimized for text generation tasks, such as summarization [11].

➤ Machine Translation

In the long-standing NLP problem of machine translation (MT), specialized models (such as seq2seq with attention) have demonstrated high accuracy for numerous language pairs. How can general LLMs help with this? It's interesting to note that even in the absence of explicit training,

large language models exhibit a startling level of translation ability. GPT-3, for instance, was tested on translating between English, French, Spanish, and other languages using only a few sample translations in its prompt. It generated translations that were reasonably accurate, occasionally coming close to the caliber of supervised machine translation systems [4].

This implies that cross-lingual mappings are acquired by LLMs as a result of training on online multilingual or translated content. Nevertheless, until recently, specialized MT models continued to have an advantage, particularly for high-resource languages or complex text. As researchers construct massively multilingual LLMs, the gap is narrowing. One notable endeavor is the No Language Left Behind (NLLB) project by Meta AI, which produced a single Transformer model for 200 languages [10]. To handle numerous low-resource languages by utilizing transfer learning from high-resource ones, NLLB used a 54B-parameter model with a Mixture-of-Experts architecture [10]. A breakthrough in multilingual machine translation resulted from this: the NLLB model significantly improved translation quality for many under-served languages, achieving an average BLEU score 44% higher than the previous state-of-the-art on a benchmark of 40,000 translation directions [10]. By training data mining techniques and making the model open-source, NLLB established a basis for research and useful translation tools that span languages from Yoruba to Asturian.

➤ Dialogue Systems and Chatbots

The use of LLMs in dialogue systems AI chatbots that can have natural, open-ended conversations with users may be the most well-known application of LLMs. AI has long sought to create a human-like open-domain chatbot, but previous efforts (Cleverbot, Microsoft XiaoIce, etc.) were constrained by small models and pre-written responses. By

offering the capacity to produce contextually relevant, cohesive, and varied responses on almost any topic, LLMs have radically altered this landscape. Google's Meena (Adiwardana et al., 2020) was a 2.6 billion-parameter seq2seq model that was trained on public social media conversations and served as a forerunner to modern chatbots. Meena was one of the first chatbots to use multi-turn dialogue data to train the model end-to-end to minimize perplexity (predict next reply) [12]. Meena made significant progress in conversational quality; measured by a new human metric called SSA (Sensibleness and Specificity Average), it achieved 79% SSA, which was 23% better than previous chatbots and not far from human-level (86% SSA) [12].

➤ *Information Extraction and Understanding*

An LLM can be trained to carry out extraction in a few-shot prompting scenario. Take the following example: "Extract all the person names from the following text:" and include a paragraph. The LLM may be persuaded to list entities or fill in slots by a well-crafted prompt, perhaps accompanied by a few examples. This makes use of the model's strong pattern recognition, which allows it to generalize to the extraction task because it saw a lot of lists of names, sentences with similar structures, etc. during pre-training. With only a few of the prompt's demonstration examples, research has shown that GPT-3 can perform close to state-of-the-art on tasks like named entity recognition (NER) and even relation extraction [13].

LLMs excel at tasks like question answering and reading comprehension. They can take a paragraph and answer questions about it in natural language, often reasoning through the answer. In multi-choice QA, for instance, LLMs can use their vast knowledge to eliminate wrong options and choose the correct answer [14]. They also demonstrate a form of common-sense understanding answering questions that require implicit world knowledge or logical reasoning [14].

Large models are especially helpful for extraction tasks that call for common sense or prior knowledge. Even though a sentence may not state "X causes Y" directly, an LLM may deduce this from context and prior knowledge. This is an example of extracting cause-effect relations from text. Such inferences can be made by LLMs that have been pre-trained on a variety of data to instill common sense in them. Another example is knowledge base completion, which involves predicting a missing fact given some entity-related facts. By taking advantage of the fact that LLMs function as enormous parametric knowledge repositories, they can be encouraged to do this [13]. On some factual recall questions, Petroni et al. demonstrated that BERT was as accurate as conventional knowledge base techniques [13]. Because of their sequence reasoning skills, more recent models, such as GPT-4, can even solve problems like reading a short story and extracting a complex set of information (characters, their relationships, timeline of events), effectively performing reading comprehension and IE together.

IV. SCALABILITY AND DEPLOYMENT CONSIDERATIONS

➤ *Model Scaling: Compute, Data, and Performance*

Generally, scaling up an LLM entails using more computing power to train a larger model on more data. As previously mentioned, empirical scaling laws imply that performance (as determined by downstream accuracy or cross-entropy loss) increases smoothly with increases in model size, data size, and compute budget [8]. It is not easy to implement such scaling, though. It takes distributed computing across numerous GPUs or TPUs to train a model with 100B+ parameters. These massive models are fitted in memory and computed using a combination of techniques such as data parallelism (replicating the model on multiple devices, each processing different data) and model parallelism (splitting the model's layers or parameters across devices). In order to train Transformer models with up to 530 billion parameters, for example, NVIDIA's Megatron-LM library introduced a combination of model and pipeline parallelism, distributing the computation and weights among dozens of GPUs in parallel. Similarly, PaLM's 540B training was able to run on 6144 TPU chips in parallel thanks to Google's Pathways system [6]. Pushing the envelope has been made possible by these advancements in distributed training software and algorithms (such as XLA, DeepSpeed, and Mesh-TensorFlow).

Scaling on the data side has meant guaranteeing diversity and quality in addition to using more tokens. An extensive but largely uncurated web scrape was used by early LLMs such as GPT-2/GPT-3 (WebText, Common Crawl). Subsequent models avoided excessively noisy or redundant data by selecting high-quality subsets (such as the Pile dataset or MassiveText). A variety of data sources, including books, Wikipedia, news, code, scholarly articles, and more, are essential because they provide a range of styles and expertise. The Chinchilla study revealed an intriguing conclusion: a lot of big models were trained on a disproportionately small number of tokens [15]. The fact that GPT-3 (175B) saw roughly 300 billion tokens while Chinchilla (70B) was trained on approximately 1.4 trillion suggests that GPT-3 was undertrained and could have used a lot more data [15]. Newer models then try to use the compute-optimal data size, which means that the number of training tokens should roughly scale in proportion to the model parameters [15]. You eventually run out of high-quality text on the internet, so of course, getting trillions of tokens has its limits. Creating synthetic training data (possibly using smaller models or human-in-the-loop) and adding multilingual text to increase volume are two methods for extending data.

To sum up, scaling LLMs involves striking a balance between compute optimization, data quantity and quality, and model size. The largest models continue to push the limits of what is technically feasible, despite the community's advanced techniques for training multi-billion-parameter Transformers at a reasonable cost. We are probably at a point where increasing the number of parameters is no longer as effective as it once was; instead, future improvements might come from creative training paradigms (such as multimodal

training or training on better objectives) and architectural changes. In fact, according to a 2023 result, by cleverly reusing some training compute, one can train U-PaLM (an upgraded PaLM) to achieve the same performance as a standard PaLM that used twice as much compute with just 0.1% more compute [16]. This suggests that intelligent training can outperform blind scaling.

➤ *Efficient Inference and Serving of LLMs*

The difficulty lies in effectively providing an LLM to users after it has been trained. These models are very large; for example, a 175B parameter model in FP16 uses about 350 GB of memory for weights alone. Because of the size and depth of the model, running inference for a single input can be computationally demanding, particularly if the outputs are lengthy. It would be impractical to use LLMs in real-time applications (chatbots, interactive assistants, etc.) without optimization. As a result, various methods are used to speed up inference and compress models.

Model compression through knowledge distillation or smaller, specialized models is one important strategy. Using a given dataset, distillation entails teaching a smaller model (student) to mimic the outputs of the larger model (teacher). This worked well with BERT-Distiller. About 97% of BERT's language understanding performance was retained when it was compressed by 40% (Sanh et al., 2019). During training, it matched the teacher's intermediate representations and logits to accomplish this. The end result was a model that is nearly as good on downstream tasks but 60% faster and much smaller [17].

Quantization is an additional method that lowers the model weights and activations' numerical precision. In order to cut memory in half, 16-bit floats (FP16/BF16) are typically used for training rather than 32-bit. It is frequently possible to use 8-bit or even 4-bit integers for weights in inference. Recent advances in 8-bit approximation (Dettmers et al., 2022) demonstrate that, with careful handling of outlier features, a model such as GPT-3 can be compressed to 8-bit weights with minimal loss in accuracy [18]. Their LLM.int8 () approach achieved 0% performance degradation for 175B model inference at 8-bit by introducing a vector-wise quantization scheme and a fallback to higher precision for a small fraction of outlier values [14].

In effect, this reduces the amount of memory and bandwidth needed by half. In fact, using 8-bit weight loading, they showed how to run the 175B model on a single machine with consumer GPUs [18]. More aggressively, 4-bit quantization (and even 3-bit in research) has been investigated. Although some accuracy is lost, the community has released 4-bit quantized versions of LLaMA-65B, etc., which significantly reduce the hardware barrier for use. Quantization is often used post-hoc with some calibration; it doesn't require retraining the model. The drawback is that very low precision may affect the quality of the model; however, this can be lessened by methods like mixed precision (keeping some layers at a higher precision) or quantization-aware training.

Table 2 Reported Performance Metrics of Major LLMs Across NLP Tasks

Model	Task Evaluated	Performance Summary	Source
GPT-3	Few-shot NLI, QA, Translation	Matched/exceeded previous SOTA with prompts only	[4]
PEGASUS	Abstractive Summarization	Achieved SOTA, human-level summary quality	[9]
BART	Generation & Denoising Tasks	Strong performance on abstractive summarization	[11]
Meena	Open-domain Chatbot (SSA metric)	79% SSA close to human baseline	[12]
BERT	Knowledge Probing (LAMA)	Comparable to symbolic knowledge retrieval	[13]
NLLB-54B	Multilingual MT (40k directions)	44% BLEU improvement over previous SOTA	[10]
PaLM	Complex Reasoning Tasks	Outperformed human average on reasoning benchmarks	[6]

V. REAL-WORLD USE CASES OF LLMs

After talking about their potential and difficulties, we move on to well-known real-world uses for large language models. Across all industries, LLMs are being used, frequently as the main ingredient in brand-new goods and services. We focus on three main areas: how LLMs are enabling search engines and information retrieval tools; how they improve productivity software and content creation; and how domain-specific LLMs are applied in specialized fields such as law and medicine.

➤ *Search Engines and Information Access*

For many years, using a search engine involved entering keywords and perusing a list of results. By allowing search engines to converse and provide direct answers to queries, LLMs are revolutionizing this paradigm and making it harder to distinguish between chatbots and search engines. For

instance, GPT-4 has been incorporated into Microsoft's Bing search engine. Users can ask sophisticated questions in natural language and receive a synthesized response that includes citations to online sources. This is a significant change because the search engine now understands and presents information in a conversational way in addition to retrieving it. Aggressive summarization and contextual integration of data from various documents are skills that LLMs excel at.

By initiatives like Google Bard (which is based on LaMDA) and experimental search features that display an AI snapshot of the response above links, Google is also integrating LLMs. Given that LLMs have a propensity for hallucinations, these engines take care to guarantee factual accuracy and steer clear of incorrect responses. With LLMs permitting follow-up inquiries, clarifications, and even multi-turn search sessions, conversational search is evidently

becoming a reality. This makes accessing information more natural because, like in a conversation with an expert, one can ask more specific questions to hone their query.

One novel method for enhancing LLMs with search is Retrieval-Augmented Generation (RAG). As an example, startups are developing RAG-based assistants that can answer questions from users, conduct real-time web searches or database queries, and more. The Bing/Google integration essentially does this, but it's also used in customer service (LLM searches a knowledge base) and even personal assistants (LLM searches your calendar, emails, etc., to answer questions about your schedule, for example). Answers are more current and factual when the LLM is grounded in retrieved data [19]. For instance, since its training data has a cutoff, an LLM cannot be aware of current events without retrieval; however, it can be informed about today's news and then have a conversation about it by using a search query. The LLM and search engine work well together because the search engine offers new, targeted information and the LLM offers reasoning and language fluency [19].

➤ *Content Creation and Productivity Tools*

The way people create text, code, and other media has changed dramatically as a result of the quick adoption of LLMs in content creation tools. LLMs are sophisticated writing assistants because of their ability to produce text that is human-like. When writing an email, for instance, a user can include bullet points, and an LLM will expand them into formal, courteous language. In actuality, programs like Microsoft Word's AI integration or GPT-3-powered plugins can generate whole paragraphs on a particular theme, rephrase awkward sentences, and suggest sentence completions. This significantly speeds up writing assignments and helps people who have trouble writing by offering a first draft or a creative spark. According to one survey, foundation models demonstrate exceptional capabilities for producing and modifying content in a variety of fields, ranging from technical documentation to imaginative storytelling [20].

In software development, LLMs are changing the game. With a natural language description of the problem, models such as OpenAI Codex (and its successor in GitHub Copilot) can automatically generate code because they have been trained on vast source code corpora [20]. Along with swiftly producing boilerplate code, they also help with code explanation, bug fixes, and even test case writing. In this way, LLMs function as an AI pair-programmer.

In reality, Copilot works with code editors and provides developers with contextually relevant code recommendations while they type. Research on the effects of Copilot revealed that developers could finish tasks much more quickly. Participants who used Copilot finished a coding task 55.8% faster than those who did not in a controlled experiment [21]. Email and document drafting are two more common use cases. More sophisticated LLM integration is now emerging, thanks to features like Gmail's Smart Compose, which suggests next words using a smaller language model. In order to enable tasks like "Draft a job offer letter for a product

manager position," Google announced integrating LLMs into Docs and Gmail. The model will generate a complete draft, which the user can then edit.

➤ *Domain-Specific Applications (Medicine, Law, etc.)*

Large language models are also being customized for specialized fields, where their extensive knowledge and cognitive powers open up new avenues. Medicine and law are two well-known examples, but other industries that demand expert-level language comprehension also exhibit comparable patterns, including academia and finance.

LLMs are used in medicine to assist patients and medical professionals. GPT-4's performance on the United States Medical Licensing Examination (USMLE), a difficult set of tests for medical professionals, was a notable accomplishment. GPT-4 achieved scores in the top quartile (approximately the 90th percentile) of human test-takers in addition to passing all three exam steps [22]. This demonstrates the capacity to retain and apply sophisticated medical knowledge (anatomy, physiology, pathology, etc.) at an expert level. On some sections of the USMLE, GPT-3.5 had previously achieved near-passing scores (~60% accuracy, when passing was ~60%), but GPT-4's improvement was notable [22], [23]. To function as clinical assistants, LLMs are being trained on medical texts (developing Med-PaLM, BioGPT, etc.) in addition to exams. As a diagnostic decision support, for instance, an LLM can assist physicians by combining patient data and making recommendations for potential diagnoses or treatment regimens.

LLMs have made waves in the legal field by passing sections of the Multistate Bar Exam. On the bar exam, GPT-4 significantly outperformed GPT-3.5, which was in the bottom 10% of test-takers, scoring in the top 10% [24]. This implies the model has learned a great deal about the law and is able to apply it to hypothetical fact patterns that resemble exams. Legal research tasks, such as posing queries regarding case law and receiving responses with citations, are being investigated by law firms and startups. An LLM could find and summarize pertinent cases or statutes in seconds, citing the original texts, as opposed to a junior associate spending hours doing so. It's even possible to speed up the process of creating legal documents (such as contracts, wills, and NDAs) by using templates. The LLM creates a customized document after the user specifies the important terms. Products already exist that use AI to help draft or review contracts and identify potentially dangerous language.

Other fields, besides law and medicine, are doing the same:

- *Finance:*

Applications for LLM can also be found in the finance sector, including automated customer support in banking, financial report analysis, and market data inquiries. An LLM with a finance concentration that was further trained on financial texts (such as SEC filings and analyst reports) was presented in a recent case study [14].

- *Scientific Research:*

Elicit and other tools use LLMs to assist researchers in locating pertinent papers and even compiling evidence from them. By using the literature, an LLM trained on scientific papers can respond to factual queries, such as "What do studies say about the effect of X on Y?" This can expedite literature reviews, but it cannot take the place of real experiments or human interpretation.

- *Education:*

Personal tutors in all subjects are provided by LLMs such as Khanmigo (by Khan Academy). When learning a foreign language or programming language, for example, LLMs can create exercises, adjust to the student's level, and offer feedback.

- *Customer Support and Domain-Specific Assistance:*

To manage complicated inquiries, numerous businesses are incorporating LLMs into their chatbots for customer service. A tax preparation service might, for instance, have an AI that guides users through filing by providing detailed answers to tax code questions (provided the AI is properly trained or connected to the tax knowledge base).

VI. CONCLUSION

Large Language Models have revolutionized the natural language processing field tremendously. These models have opened doors towards unified approaches for various natural language understanding and production tasks. This paper focuses on discussing the architectural basics of large language models. Here, we highlighted the contributions of large-scale Transformer models, effective pre-training approaches, and the adaptations for these models towards realizing significant achievements. In addition, this survey highlighted various effective application areas of large language models. These areas include text summarization, translation, dialogue systems, and information extraction. According to the survey, large language models have greatly reduced the necessity of specialized architectures for task implementation. However, various challenges still need to be addressed. These challenges include effective computation, large-scale application, and various ethical considerations. They are critical to ensuring the sustainable progress of large language models. Future large language models are likely to focus more on effective computation, inference, and alignment. In conclusion, large language models are a great milestone in natural language processing. They are likely to influence this field significantly in the future.

REFERENCE

- [1]. R. Qureshi et al., Large Language Models: A Comprehensive Survey of its Applications, Challenges, Limitations, and Future Prospects. 2024. doi: 10.36227/techrxiv.23589741.v7.
- [2]. M. E. Peters et al., "Deep contextualized word representations," Mar. 22, 2018, arXiv: arXiv:1802.05365. doi: 10.48550/arXiv.1802.05365.
- [3]. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional

Transformers for Language Understanding," May 24, 2019, arXiv: arXiv:1810.04805. doi: 10.48550/arXiv.1810.04805.

- [4]. T. B. Brown et al., "Language Models are Few-Shot Learners," July 22, 2020, arXiv: arXiv:2005.14165. doi: 10.48550/arXiv.2005.14165.
- [5]. R. Bommasani et al., "On the Opportunities and Risks of Foundation Models," July 12, 2022, arXiv: arXiv:2108.07258. doi: 10.48550/arXiv.2108.07258.
- [6]. A. Chowdhery et al., "PaLM: Scaling Language Modeling with Pathways," Oct. 05, 2022, arXiv: arXiv:2204.02311. doi: 10.48550/arXiv.2204.02311.
- [7]. M. Chen et al., "Evaluating Large Language Models Trained on Code," July 14, 2021, arXiv: arXiv:2107.03374. doi: 10.48550/arXiv.2107.03374.
- [8]. J. Kaplan et al., "Scaling Laws for Neural Language Models," Jan. 23, 2020, arXiv: arXiv:2001.08361. doi: 10.48550/arXiv.2001.08361.
- [9]. J. Zhang, Y. Zhao, M. Saleh, and P. J. Liu, "PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization," July 10, 2020, arXiv: arXiv:1912.08777. doi: 10.48550/arXiv.1912.08777.
- [10]. N. Team et al., "No Language Left Behind: Scaling Human-Centered Machine Translation".
- [11]. M. Lewis et al., "BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension," Oct. 29, 2019, arXiv: arXiv:1910.13461. doi: 10.48550/arXiv.1910.13461.
- [12]. D. Adiwardana et al., "Towards a Human-like Open-Domain Chatbot," Feb. 27, 2020, arXiv: arXiv:2001.09977. doi: 10.48550/arXiv.2001.09977.
- [13]. F. Petroni et al., "Language Models as Knowledge Bases?," Sept. 04, 2019, arXiv: arXiv:1909.01066. doi: 10.48550/arXiv.1909.01066.
- [14]. P. Kumar, "Large language models (LLMs): survey, technical frameworks, and future challenges," *Artif. Intell. Rev.*, vol. 57, no. 10, p. 260, Aug. 2024, doi: 10.1007/s10462-024-10888-y.
- [15]. J. Hoffmann et al., "Training Compute-Optimal Large Language Models," Mar. 29, 2022, arXiv: arXiv:2203.15556. doi: 10.48550/arXiv.2203.15556.
- [16]. Y. Tay et al., "Transcending Scaling Laws with 0.1% Extra Compute," Nov. 16, 2022, arXiv: arXiv:2210.11399. doi: 10.48550/arXiv.2210.11399.
- [17]. V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter," Mar. 01, 2020, arXiv: arXiv:1910.01108. doi: 10.48550/arXiv.1910.01108.
- [18]. T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer, "LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale," Nov. 10, 2022, arXiv: arXiv:2208.07339. doi: 10.48550/arXiv.2208.07339.
- [19]. H. Xiong et al., "When Search Engine Services meet Large Language Models: Visions and Challenges," June 28, 2024, arXiv: arXiv:2407.00128. doi: 10.48550/arXiv.2407.00128.

- [20]. J. Schneider, C. Meske, and P. Kuss, "Foundation Models: A New Paradigm for Artificial Intelligence," *Bus. Inf. Syst. Eng.*, vol. 66, no. 2, pp. 221-231, Apr. 2024, doi: 10.1007/s12599-024-00851-0.
- [21]. S. Peng, E. Kalliamvakou, P. Cihon, and M. Demirer, "The Impact of AI on Developer Productivity: Evidence from GitHub Copilot," Feb. 13, 2023, arXiv: arXiv:2302.06590. doi: 10.48550/arXiv.2302.06590.
- [22]. A. Meyer, J. Riese, and T. Streichert, "Comparison of the Performance of GPT-3.5 and GPT-4 With That of Medical Students on the Written German Medical Licensing Examination: Observational Study," *JMIR Med. Educ.*, vol. 10, p. e50965, Feb. 2024, doi: 10.2196/50965.
- [23]. Y. Chen et al., "Performance of ChatGPT and Bard on the medical licensing examinations varies across different cultures: a comparison study," *BMC Med. Educ.*, vol. 24, no. 1, p. 1372, Nov. 2024, doi: 10.1186/s12909-024-06309-x.
- [24]. D. M. Katz, M. J. Bommarito, S. Gao, and P. Arredondo, "GPT-4 Passes the Bar Exam," *SSRN Electron. J.*, 2023, doi: 10.2139/ssrn.4389233.