



Foundation Models and Their Transformative Impact on Machine Learning

Srashti Farkya

Department of Computer Science and Engineering
Mandsaur University
Mandsaur, India
srashtifarkya2005@gmail.com

Priyanka Khabiya

Assistant Professor
Department of CSE
Mandsaur University
Mandsaur, India
khabiya198727@gmail.com

Abstract—Artificial Intelligence (AI) has made remarkable progress in recent years, mostly enabled by the growth of large-scale ML models based on abundant and various datasets. Although conventional models are meant for just a single job, foundation models are adaptable and may be applied in fields such as natural language processing, computer vision and even the creation of programming code. All these models, for example, BERT, GPT-4, Claude and Gemini, work with transformers and have been trained in an unsupervised manner. Because of this approach, models are able to figure out what's in the data and how the components relate which requires very little labeled information. Models trained on a large dataset can be customized or guided to do different tasks with only a little more data. The paper looks into the concept, architecture and how foundation models are applied. It describes the upsides of ML such as its adaptable setup, speed and versatility and it similarly notes the main issues such as potential bias in data, moral aspects, huge computing requirements and privacy concerns. The objective is to explain how foundation models are changing ML and what issues should be kept in mind when they are used or introduced.

Keywords—Foundation Models, Artificial Intelligence, Transformer Architecture, Self-supervised Learning, versatile in nature Models, GPT-4, BERT, Claude, Gemini, Ethical AI, Machine learning Applications

I. INTRODUCTION

A. Paradigm Shift in Artificial Intelligence

The rise of foundation models is a major shift in how AI is developing today [1]. Foundation models can do more than traditional AI systems because they are trained to work in many different areas and tasks [2]. The models are mainly transformer-based and are pre-trained using huge and varied data set through self-supervised learning [3].

Three main advances make possible this major shift in how it thinks about computing:

- **Scalable Architectures:** Parallelized sequence modeling is made possible in Transformers by using self-attention, helping them get rid of the sequential drawbacks seen in RNNs and LSTMs [4]. Because of their scalability, they can deal with deep and high-dimensioned training, identifying long-range ones more efficiently.
- **Massive Data and Compute Resources:** The availability of large-scale corpora, combined with powerful computing infrastructure (GPUs, TPUs) [5], makes it feasible to train models on trillions of tokens and optimize parameters in the billions or trillions [6].

- **Unified Self-Supervised Objectives:** Basic language modelling (e.g., BERT) and next-token prediction (e.g., GPT) are examples of pretraining goals. Minimize the requirement for task-specific datasets using annotations. This allows foundation models to be applied broadly with minimal downstream fine-tuning.

Together, these advances signify a movement away from siloed models toward highly generalizable systems that learn foundational representations transferable across multiple domains.

B. Historical Evolution of Foundation Models

One way to think about the evolution of foundation models is as a multi-phase process, marked by key innovations and paradigm shifts [7]. Figure 1 shows the Bar graph of the impacts of model training on environment are as follows:

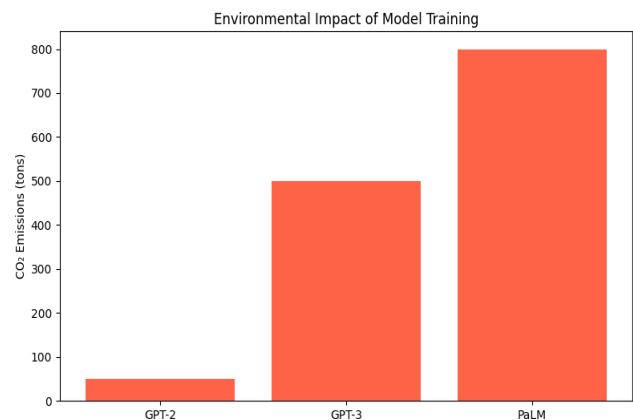


Fig. 1. Bar graph of the impacts of model training on environment

1) Precursor Phase (2017–2018):

- **Transformer Introduction:** Vaswani et al. introduced the transformer architecture, replacing recurrence with attention mechanisms.
- **BERT:** Devlin et al. developed BERT, introducing deep bidirectional context via masked language modeling.
- **GPT-1:** OpenAI's Generative Pre-trained Transformer showcased autoregressive pretraining, laying the groundwork for generative AI.

2) First Wave (2019–2020):

- **GPT-2:** Demonstrated coherent text generation and few-shot learning capabilities.
- **RoBERTa and ALBERT:** Refined BERT through improved training strategies and parameter efficiency.

3) *Emergence of Capabilities (2021–2022):*

- **GPT-3:** With 175 billion parameters, GPT-3 revealed emergent properties such as reasoning and translation without task-specific training.
- **CLIP and DALL-E:** Enabled multimodal learning, linking text to images and generating creative content.

4) *The Modern Era (2023–Present):*

- **Advanced Models:** GPT-4, Claude, Gemini, and LLaMA-3 pushed boundaries in safety, alignment, and multimodality.
- **Widespread Integration:** Foundation models now power tools in productivity, education, healthcare, and beyond, becoming embedded in everyday life [8][9].

This historical journey highlights the exponential progression in scale, capability, and impact of foundation models over just a few years.

C. *Societal Impact and Industrial Relevance*

Foundation models are reshaping industries by enhancing productivity, transforming job roles, and enabling new forms of creativity and automation shown in Table I.

TABLE I. SECTOR-WISE IMPACT OF FOUNDATION MODELS

Sector	Productivity Gain	Job Impact	Innovation Potential
Healthcare	+35%	Augments diagnostics	Drug discovery medical image analysis
Education	+40%	Personalized tutoring	Adaptive learning automated content creation
Legal	+45%	Contract generation	Case law summarization, legal research
Creative Industries	+60%	Co-creation with AI	AI art, storytelling, music composition
Customer Service	+50%	Virtual assistance	Conversational agent, AI-based support
Finance	+30%	Decision support	Fraud detection, algorithmic trading

These models enable automation of repetitive tasks, augmentation of complex decision-making, and creation of entirely new workflows. However, their widespread deployment also raises serious concerns regarding [10]:

- Misinformation amplification
- Bias and discrimination
- Surveillance and privacy invasion
- Labor market disruption

This dual-edged impact demands responsible development, deployment, and regulation of foundation models.

D. *Motivation and Scope of This Paper*

Given the growing centrality of foundation models in modern AI, this paper aims to explore their [11]:

- **Architectural Foundations:** Transformer innovations, scalability techniques, and attention optimizations.
- **Training Paradigms:** Self-supervised learning, scaling laws, and emergent capabilities.
- **Adaptation Strategies:** Fine-tuning, prompting, adapters, and instruction alignment.
- **Deployment Challenges:** Efficiency, interpretability, and environmental cost. safety,

- **Ethical Dimensions:** Fairness, transparency, societal risk, and governance mechanisms.

By providing an in-depth analysis, it aims to foster a deeper learning of both the technical underpinnings and the real world implications of foundation models.

II. RELATED WORK

The evolution of foundation models builds upon decades of ML research, the categorized prior work into three key areas[12]:

A. *Early Language Models*

Traditional approaches relied on statistical methods like n-gram models, later surpassed by neural architectures such as Word2Vec and ELMo [13][14]. These models demonstrated the value of distributed representations but lacked contextual awareness.

B. *Transformer Revolution*

A paradigm change occurred with the release of the Transformer, enabling parallel processing of sequential data. Follow-up work adapted this architecture for[15]:

- **Bidirectional Contexts:** BERT introduced masked language modeling
- **Generative Tasks:** GPT series pioneered autoregressive pretraining
- **Efficiency:** Sparse attention mechanisms addressed quadratic scaling

C. *Scaling and Generalization*

Recent breakthroughs emerged from understanding neural scaling laws. Key findings include[16]:

- Emergent abilities appear at sufficient scale
- Multitask performance improves predictably with compute

Model alignment techniques (e.g., RLHF) become critical at scale.

D. *Open Challenges*

Prior surveys identify unresolved issues:

- Energy efficiency vs. capability trade-offs
- Evaluation frameworks for general-purpose AI
- Societal impacts of model centralization

The work extends these analyses with updated architectural comparisons and ethical considerations [17].

III. FOUNDATION MODELS: OVERVIEW

A. *Defining Foundation Models*

The term "foundation models" describes extensive ML algorithms that have been trained on a variety of data and are able to adjust to a broad range of downstream tasks [18]. They are typically based on transformer-based structures and are trained using self-supervised learning on massive datasets, allowing them to adapt to other domains with little task-specific modification [19]. These models "serve as the basis" for a number of applications, including text production, image captioning, code synthesis, and speech translation. Examples include BERT, GPT, CLIP, and DALL-E shown in Table II [20].

B. Key Characteristics

- **Scale:** Trained on billions of tokens/images/code/data points using billions of parameters.
- **Transferability:** Excellent zero-shot, few-shot, and fine-tuned performance.
- **Multimodality:** Integration across text, image, audio, video, and code modalities.
- **Generalization:** Ability to perform well across unseen tasks and languages.

C. Historical Milestones

TABLE II. KEY MILESTONES IN FOUNDATION MODELS

Year	Milestone
2018	Contextual embeddings in NLP
2020	(BERT) Few-shot reasoning (GPT-3)
2021	Vision-language fusion (CLIP, DALL·E)
2022	Open-ended instruction tuning (InstructGPT, FLAN)
2023	Multimodal interactions and agents (Gemini, Kosmos, Claude)

D. Categorization of Foundation Models

1) Paradigm Shift in AI

In Table III shows the foundation models represent a change from conventional task-specific models to models that are more general in nature, where the same model architecture can be adapted to numerous applications. This reduces engineering overhead, promotes reusability, and democratizes access to powerful AI systems.

TABLE III. CATEGORIES OF FOUNDATION MODELS

Type	Examples	Focus
Language Models	GPT, BERT, T5	Text understanding and generation
Vision Models	ViT, SAM	Image classification, segmentation
Multimodal Models	CLIP, Flamingo	Text-image, image-video tasks
Code Models	Codex, CodeGen	Code generation, synthesis
Speech Models	Whisper, wav2vec	Speech recognition and synthesis

IV. ARCHITECTURES AND TECHNIQUES

Foundation models derive their strength not only from scale but also from the architectural innovations that enable them to handle vast, diverse datasets and generalize across tasks [21]. Over the years, these architectures have evolved to support better learning, efficiency, and modality alignment.

A. Transformer: The Foundational Backbone

The Transformer architecture, first presented in the groundbreaking work Attention Is All You Need, has grown to be the cornerstone of contemporary foundation models. The self-attention mechanism, which enables the model to focus on various segments of the input sequence and capture intricate relationships regardless of distance, is its special strength [22]. Transformers provide significant efficiency benefits by processing input sequences in parallel, as opposed to recurrent models [23].

Transformers include feedforward layers, multi-head self-attention layers, and residual connections with layer normalization. Each of these components contributes to stabilizing deep training and enhancing expressiveness.

Transformer structure variants serve as the basis for the majority of foundation models:

- **BERT** The encoder-only architecture of (Bidirectional Encoder Representations from Transformers) is optimized for job comprehension.
- **GPT** (Generative Pretrained Transformer) adopts a decoder-only structure, fine-tuned for generative tasks like language modeling and code generation.
- **T5** and **BART** implement encoder-decoder frameworks, making them ideal for jobs involving sequences, such as summarization or translation.

B. Emerging Architectural Innovations

As models scale, researchers have introduced novel modifications to improve efficiency, performance, and specialization [24]:

- **Sparse Attention and Long Sequence Handling:** Standard Transformers scale quadratically with input length. Models like Long former, Big Bird, and Performer introduce sparse or linear attention mechanisms to process longer sequences without compromising on context.
- **Mixture-of-Experts (MoE):** This approach, seen in models like G-Shard and Switch Transformer, selectively activates a subset of model parameters (experts) per input, thereby drastically increasing capacity without proportional compute cost.
- **Retrieval-Augmented Models:** Instead of relying solely on internal parameters, models like RETRO and REALM access external databases during inference, blending memory and generation for more factually accurate outputs.

C. Cross-modal and Multimodal Models

With the growing need to model information across multiple modalities, Architectures have developed to handle pictures, music, and video in addition to text. These multimodal models use modality-specific encoders (e.g., CNNs for vision or spectrogram-based encoders for audio) alongside shared Transformer backbones [25].

- **CLIP** (Contrastive Language-Image Pretraining) jointly trains on image-caption pairs to align visual and textual representations.
- **DALL·E**, **Flamingo**, and **PaLM-E** further expand the capability to generate images, captions, or interpret inputs across modalities.
- Models like **Gemini** (from Google DeepMind) unify vision, text, audio, and even video under one model interface, pushing the boundaries of general intelligence.

D. Training Techniques and Adaptations

To improve task generalization and adaptability, foundation models incorporate several advanced training methods [26]:

- **Prompting and Instruction Tuning:** These approaches guide model behavior using carefully designed input phrases or examples. Instruction tuning, used in models like FLAN and Instruct GPT, fine-tunes the model with human-readable instructions to make outputs more reliable and controllable.
- **Reinforcement Learning from Human Feedback (RLHF):** These techniques, which powers Chat-GPT and similar models, aligns model outputs with human

preferences[27], enhancing factual correctness and social alignment.

- **Parameter-Efficient Fine-Tuning (PEFT):** Fine-tuning with a small number of trainable parameters is made possible by methods including BitFit, LoRA, as well as adapters. These methods are particularly useful when computational resources are limited or when models need to be customized across multiple tasks or domains.

E. Beyond Traditional Architectures

The witnessing a move toward hybrid and modular architectures [28]:

- **Perceiver and Perceiver IO models** generalize attention mechanisms to handle any data type and dimensionality.
- **Composable Models** allow for flexible assembly of components, each fine-tuned for a specific function (e.g., reasoning, summarizing, translating).

This ongoing innovation in architectures reflects a core theme of foundation models: not just bigger, but smarter and more adaptable. Figure 2 shows a flowchart of the foundation model and challenges and mitigation are given :

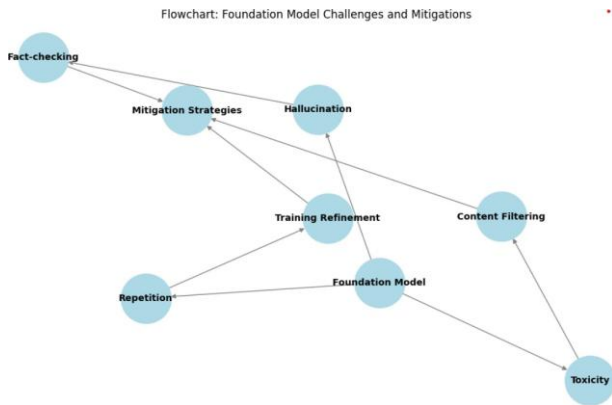


Fig. 2. Foundation model and challenges and mitigation

V. TRAINING STRATEGIES AND DATA CONSIDERATIONS

The success of foundation models is rooted not only in architectural sophistication but also in the strategies used during training, as well as the calibers, volume, and variety of data they encounter [29]. To Designing effective training methodologies involves a complex interplay of data collection, preprocessing, optimization techniques, and resource management.

A. Data Scale and Diversity

Foundation models are typically trained on massive corpora that span multiple domains, styles, and modalities. The idea is to introduce models to a broad range of ideas, language structures, and understanding representations, allowing them to generalize effectively across tasks [30].

- **Text Data Sources:** Common text sources include web crawls (e.g., Common Crawl), books, academic papers, news articles, social media content, and conversational data. Models like GPT-3, PaLM, and LLaMA utilize hundreds of billions of tokens from curated datasets [31].
- **Multimodal Datasets:** For models that go beyond text, datasets such as ImageNet, LAION, Audio Set,

and YouTube-8M provide aligned visual and auditory data.

- **Data Curation:** Cleaning, deduplication, and filtering are essential steps. Poor-quality or redundant data can lead to overfitting, hallucination, or biased outputs. Human-in-the-loop data filtering, like in the Pile or Open Web Text, ensures higher relevance and appropriateness.

B. Optimization and Scaling Techniques

Training large-scale models is computationally intensive. Several techniques help manage the complexity and cost [32]:

- **Gradient Accumulation:** Allows training on large batch sizes even with memory constraints [33].
- **Mixed Precision Training:** Lower-precision (FP16, for example) arithmetic speeds up training as well as uses less memory without compromising accuracy.
- **Distributed Training:** Data parallelism and model parallelism strategies, like Megatron-LM and DeepSpeed [34], allow model components and data batches to be spread across multiple GPUs or nodes.
- **Curriculum Learning:** Some approaches prioritize simpler examples early in training, gradually introducing more complex data. This can stabilize learning and accelerate convergence.

C. Data Augmentation and Synthetic Data

To enhance generalization, especially in low-resource domains, data augmentation strategies are employed [35]:

- Backtranslation, paraphrasing, and noising augment textual data.
- Synthetic datasets, generated through rules or existing models, are often used to expand data availability for rare tasks or low-resource languages.

In multimodal models, techniques like image-text alignment through caption generation or text-to-image pair filtering improve quality and reduce noise.

D. Human Feedback and Alignment

The increasing significance of matching model results to human values as well as preferences is shown by recent developments:

- Reinforcement Learning from Human Feedback (RLHF) plays a vital role in fine-tuning models post-pretraining. Human evaluators rank model outputs, and in order to develop a reward model that directs reinforcement learning, these rankings are utilized.
- Instruction tuning, where models are trained on datasets with clear task instructions, has led to better usability and alignment with user intent.

E. Resource and Environmental Considerations

Developing foundation models is very expensive in terms of resources and impact on the environment [36]:

- Training one large model uses a huge amount of GPU which can pollute tons of CO depending on the infrastructure.
- Techniques like model distillation, sharing parameters and using efficient algorithms work to address these problems by making training more productive while still keeping performance high.

Activities in the training phase of foundation models depend on having sufficient resources, following an effective strategy and focusing on responsibility, fairness and efficiency.

VI. EVALUATION METRICS AND BENCHMARKS

Foundation models must be judged using a variety of methods since their powers cover various subjects, not just the regular narrow-task ones. This requires their evaluation to cover their use of language, their ability to reason, as well as their reliability, impartiality and agreement with what society finds important[37].

A. Traditional NLP Benchmarks

A wide range of standardized datasets are commonly employed to study how well foundational models use language and reason. The benchmark known as GLUE comprises tasks including phrase similarity, sentiment analysis, as well as natural language inference. SuperGLUE is an enhanced and more challenging version of GLUE that tests deeper reasoning and understanding. SQuAD (Stanford Question Answering Dataset) evaluates reading comprehension

by requiring models to extract precise answers from textual passages. MNLI (Multi-Genre Natural Language Inference) assesses the model's capability to understand relationships between sentence pairs [38]. Datasets like CoQA and QuAC focus on conversational question answering and help gauge a model's performance in multi-turn dialogue. Although these benchmarks offer standardized evaluation, they are limited by the static nature and size of their tasks.

B. Multimodal Benchmarks

Multimodal benchmarks have become essential with the development of algorithms that can interpret several input formats, including text, pictures, audio, as well as video[39]. A model's capacity to decipher and respond to queries based on visuals is assessed using VQA (Visual Question Answering). Datasets like Flickr30k and MSCOCO are used for image captioning and evaluating text-to-image alignment [40]. In jobs requiring text-to-image creation, the CLIP Score can be utilized to evaluate how well the created pictures match the written descriptions. These benchmarks test the model's capacity for accurate cross-modal understanding.

C. Instruction Following and Alignment Benchmarks

As foundation models are increasingly applied in interactive systems, their ability to follow instructions and align with human values becomes essential. BIG-Bench (Beyond the Imitation Game) is a broad benchmark covering diverse areas including reasoning, mathematics, and commonsense understanding. HELMa (Helpful, Honest, and Harmless) focuses on ethical alignment by assessing safety and truthfulness in model outputs. TruthfulQA is designed to measure how often a model gives factually correct responses. Particularly after reinforcement learning with human feedback (RLHF), these standards often depend on human preferences as well as controlled reminders.

D. Emerging Trends in Evaluation

New trends in model evaluation are pushing beyond static benchmarks. Dynamic evaluation involves testing models on evolving datasets, such as real-time question answering or trending topics[41]. Human-in-the-loop evaluation includes real-world human feedback to measure the helpfulness,

coherence, and safety of responses. For creative tasks like story generation or open-ended dialogue, traditional accuracy metrics fall short. Metrics like BLEU, ROUGE, and METEOR provide some structure, but ultimately, human judgment plays a pivotal role.

VII. APPLICATIONS AND USE CASES

Foundation models, due to their immense scale and generalization capabilities, have significantly impacted a variety of sectors, transforming not only how tasks are executed but also redefining what is possible through AI. NLP is among the most well-known application fields[42]. These models are capable of generating text, summarizing, translating, analyzing sentiment, responding to queries, and having chat-based discussions. A key example is automated customer support, where platforms like Intercom and Duolingo utilize OpenAI's models to provide real-time [43], multilingual support with contextual understanding and sentiment detection. This has resulted in reduced operational costs and highly personalized, 24/7 customer experiences.

Text-to-image generation, object identification, captioning, picture categorization, as well as visual reasoning are all handled by foundation models in the area of computer vision [44]. In healthcare, vision-language models are deployed to interpret radiology reports, detect anomalies in X-rays and MRI scans, and assist in early diagnostics, especially in under-resourced areas leading to increased diagnostic accuracy and faster medical assessments [45]. Multimodal applications, combining text, images, and some- times even audio or video, are expanding rapidly. Tools like DALL·E and MidJourney are revolutionizing the creative industry by generating images from text, making art and design more accessible and significantly speeding up creative processes. Similarly, code generation and software development have seen a transformation with AI pair programming tools like GitHub Copilot, which supports real-time code completion, debugging, and test generation, enhancing developer productivity and supporting new learners.

In Education and E-learning, foundation models power personalized tutors, auto-grading systems, and dynamic content generation. Adaptive tutoring in language and coding apps now tailors' explanations based on individual learner progress, promoting personalized learning paths and improved engagement [46]. The Healthcare and Biomedical Research domain benefits from models like AlphaFold2 though not purely language-based, it's built on foundational architectures accurately predicting protein structures, which has accelerated drug discovery and opened possibilities for rare disease treatments.

For Business Intelligence and Operations, these models streamline document automation, trend analysis, forecasting, and report generation [47]. Financial institutions now utilize AI to process contracts and generate compliance documentation, resulting in considerable time savings and enhanced decision-making. In Law and Governance, legal AI assistants help draft and analyze documents, summarize case law, and offer contextual suggestions, which reduces workloads and increases accessibility to legal aid.

In Scientific Discovery and Research, AI models assist researchers by summarizing papers, generating hypotheses, and interpreting data, fostering faster innovation and interdisciplinary collaboration. Lastly, foundation models contribute meaningfully to Accessibility and Inclusion. The

power as- strive technologies such as real-time captions for the hearing impaired, visual aids for those with vision impairments, and personalized assistants that adapt to individual needs making digital environments more inclusive and empowering for all users.

VIII. CHALLENGES AND LIMITATIONS

While foundation models have revolutionized the field of artificial intelligence, they come with a host of challenges and limitations that must be acknowledged and addressed. One of the most prominent concerns is the high computational and environmental cost associated with training these models[48],

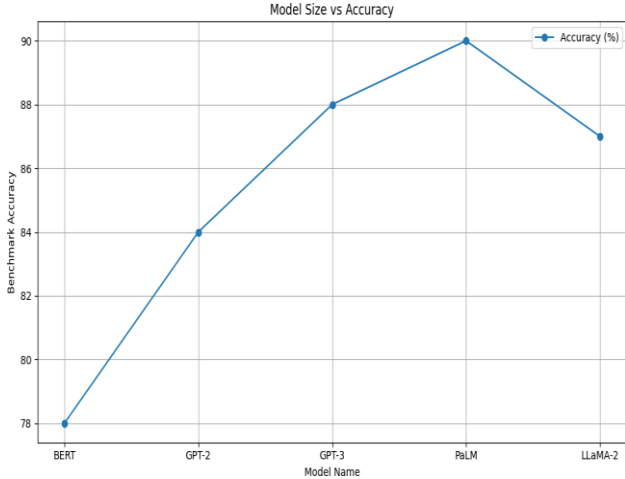


Fig. 3. Performance of the accuracy curve vs model size

For example, training GPT-3 consumed several hundred megawatt-hours of electricity, highlighting the significant energy demands and raising sustainability concerns around scaling such systems [49], (see in Figure 3) shows the accuracy curve vs model size. This has sparked an urgent need to balance innovation with environmental responsibility.

Equally pressing are issues related to data bias and representation gaps. Since foundation models learn from vast datasets that often mirror real-world stereotypes and imbalances, they can produce outputs that are discriminatory, exclusionary, or culturally insensitive. Subtle and systemic biases such as associating certain professions with specific genders or ethnicities can propagate unchecked, especially without rigorous auditing of training data. This underscores the importance of fairness and inclusivity in dataset design and model evaluation. Another core limitation is the opacity and lack of interpretability inherent in many foundation models. Functioning largely as black boxes, these models make it difficult to trace or explain their decision-making processes. This limits their trustworthiness, particularly in critical sectors like healthcare, finance, and law, where accountability and transparency are non-negotiable. The field of explainable AI (XAI) is actively exploring ways to bridge this gap, but significant progress is still needed. The economic barrier to entry poses another challenge, as the cost of training or fine-tuning these models is prohibitively high for most academic institutions, startups, and smaller organizations. This contributes to the centralization of AI capabilities within a handful of large tech corporations, raising questions about equity, innovation access, and model democratization. Efforts to create open-source alternatives are crucial to leveling the playing field and fostering inclusive

participation in AI development. Here are Table IV of the challenge mapping of foundation models as given below:

TABLE IV. CHALLENGE MAPPING OF FOUNDATION MODELS

Challenge	Domain Affected	Severity	Suggested Mitigation
Energy	Environmental	High	Efficient model design
Consumption Bias	Social, Ethical	High	Diverse dataset curation
Hallucination	Technical, Safety	High	Fact-checking mechanisms
Security Risks	Operational	Medium	Adversarial training
Lack of Interpretability	Trust & Regulation	Medium	XAI integration

In essence, while foundation models have remarkable capabilities, their unchecked use may do more harm than good. As a result, having a careful, well-structured and diverse approach to regulation using this approach is vital for dealing with these various challenges.

IX. FUTURE DIRECTIONS AND RESEARCH OPPORTUNITIES

As foundation models keep transforming AI, the next years are expected to offer many opportunities and sets of problems. Getting rid of bulk and instead relying on compact and efficient models is becoming a main direction. As people worry more about accessibility and environmental consequences, many are interested in developing lightweight models that function the same as bigger ones. Experts in areas like model compression, quantization and distillation work to allow models to be used on edge devices and in settings where resources are scarce.

Multimodal foundation models are bringing even more exciting advances. Besides text, these models now focus on images, videos and audio. CLIP, DALL·E and Flamingo represent the trend towards this kind of technology. The final aim is to design models that can think across different types of information and better understand their situation. In addition, the use of neuromyotonic integration is increasing, letting deep learning and symbolic reasoning work together. Thanks to this hybrid way, models may now carry out logical thinking and structured decisions and this would allow them to be used in planning, verification and discovery in science. Focus on sustainability is becoming more urgent. Large models need a lot of processing power to train which can hurt the environment [50]. Records of growth in Green AI, data-efficient learning and examining a model's total impacts mean that new models are both effective and also, paying attention to protecting the environment.

AI models are still challenged by inequity when it comes to language and culture. Most of today's foundation models are trained on English data and this usually introduces Western bias [51]. For this reason, FLORES, Masakhane and IndicNLP are creating models that address various languages, expressions and local aspects, placing AI within reach of more people all over the globe.

Robustness is still a major problem because models are often fooled by minor changes in the input that leads to unintended outcomes. Improving foundation models so they can tolerate these kinds of changes is important to guarantee their usefulness in practical situations.

X. CONCLUSION

Foundation models show a major change in AI, taking AI from systems designed for single tasks to engines built for general use and flexibility. In studying this technological advancement carefully, it finds several main lessons about it.

Architects see transformer-based towers as being superior at storing and transporting power compared to other types. cross-modal interaction of different elements in the brain. The self-attention mechanism, above all, has allowed models to work on and generate information with a comparable level of understanding across images, text, code and sensory data. Through this advance, along with access to a lot of data and computing, machines can now display capabilities that were not directly programmed.

Second, using self-supervised learning has completely altered the economics of developing AI. Because they lower the need for huge labels, foundation models are now allowing more people to access AI technology, yet it is now also challenging to create new ones a trend that will be important for years to come. Advances in few-shot learning and prompt engineering show that they are still only learning how to make the most of these systems.

Still, there are many challenges that it discovered and must guide future work in this area:

- The environmental impact of training ever-larger models is unsustainable without breakthroughs in energy-efficient architectures and training methods. Recent work suggests diminishing returns from pure scale, pointing toward more sophisticated approaches like mixture-of-experts and modular architectures.
- The black-box nature of these systems creates trust barriers for critical applications in healthcare, law, and finance. While techniques like attention visualization and concept activation vectors provide partial explanations, the lack comprehensive frameworks for model interpretability at scale.
- The concentration of development resources in a handful of organizations raises concerns about equitable access and the potential for monocultures in AI development. Open-source alternatives and public-private partnerships may help mitigate these risks.
- The alignment problem remains fundamentally unsolved, with models still prone to hallucination, bias amplification, and unpredictable behavior. The success of RLHF is promising but incomplete, requiring more robust approaches to value alignment.

Looking ahead, three trajectories appear particularly significant:

- **Multimodal integration** will likely dominate the next generation of foundation models, moving beyond separate text, image, and audio systems toward truly unified cognitive architectures. Early examples like GPT-4V and Gemini demonstrate the potential of this direction.
- **Neuro-symbolic hybridization** may address current limitations in reasoning and factual grounding, Combining structured knowledge representation, logical reasoning, as well as the pattern detection capabilities of foundation models
- **Efficiency breakthroughs** in areas like sparse attention, dynamic computation, and model distillation

could dramatically reduce the resource requirements while maintaining capability, potentially enabling localized deployment and edge applications.

The societal implications of foundation models are profound and multifaceted. These systems are not merely tools but collaborators that will reshape education, creative work, scientific discovery, and knowledge work across all sectors. Their development must therefore be guided by interdisciplinary collaboration involving not just computer scientists but also cognitive scientists, ethicists, policymakers, and domain experts across all affected fields.

As stand at this inflection point in AI development, the choices it makes about how to develop, deploy, and govern foundation models will have lasting consequences. The technical brilliance demonstrated in these systems must be matched by equal innovation in safety, ethics, and human-centered design. Only through such balanced advancement can it realizes the full potential of foundation models as amplifiers of human potential rather than as sources of disruption or harm.

This study has offered a thorough assessment of foundation models' present status as well as a path forward for their responsible development. The coming years will test whether it can harness this remarkable technology while mitigating its risks - a challenge that will require sustained effort from the entire AI community.

REFERENCES

- [1] R. Bommasani *et al.*, "On the Opportunities and Risks of Foundation Models," *Cent. Res. Found. Model.*, 2021.
- [2] G. Maddali, "Enhancing Database Architectures with Artificial Intelligence (AI)," *Int. J. Sci. Res. Sci. Technol.*, vol. 12, no. 3, pp. 296–308, 2025.
- [3] R. Q. Majumder, "Machine Learning for Predictive Analytics: Trends and Future Directions," *Int. J. Innov. Sci. Res. Technol.*, vol. 10, no. 4, pp. 3557–3564, 2025.
- [4] S. Singamsetty, "Dynamic Stock Price Prediction Leveraging LSTM, ARIMA, and Sparrow Search Algorithm," *Int. J. Comput. Math. IDEAS*, vol. 16, no. 3, pp. 3031–3051, 2024.
- [5] P. Choudhary, R. Choudhary, and S. Garaga, "Enhancing Training by Incorporating ChatGPT in Learning Modules: An Exploration of Benefits, Challenges, and Best Practices," *Int. J. Innov. Sci. Res. Technol.*, vol. 9, no. 11, 2024.
- [6] T. B. Brown *et al.*, "Language Models are Few-Shot Learners," *Adv. Neural Inf. Process. Syst.*, 2020.
- [7] C. Raffel *et al.*, "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer," *J. Mach. Learn. Res.*, vol. 21, pp. 1–67, 2020.
- [8] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," *NAACL HLT 2019 - 2019 Conf. North Am. Chapter Assoc. Comput. Linguist. Hum. Lang. Technol. - Proc. Conf.*, vol. 1, no. Mlm, pp. 4171–4186, 2019.
- [9] R. Tarafdar and Y. Han, "Finding Majority for Integer Elements," *J. Comput. Sci. Coll.*, vol. 33, no. 5, pp. 187–191, 2018.
- [10] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, and I. Akkaya, "GPT-4 Technical Report," vol. 4, pp. 1–100, 2024.
- [11] J. Kaplan *et al.*, "Scaling Laws for Neural Language Models," 2020.
- [12] A. Radford and T. Salimans, "Improving Language Understanding by Generative Pre-Training," pp. 1–12, 2018.
- [13] L. Floridi and J. Cowls, "A Uni ed Framework of Five Principles for AI in Society," no. June 2019, 2021, doi: 10.1162/99608f92.8cd550d1.
- [14] S. Pahune and M. Chandrasekharan, "Several Categories of Large Language Models (LLMs): A Short Survey," *Int. J. Res. Appl. Sci.*

- Eng. Technol.*, vol. 11, no. 7, pp. 615–633, 2023, doi: 10.22214/ijraset.2023.54677.
- [15] Y. Chen and B. Dolan, “GPT: Large-Scale Generative Pre-training for Conversational Response Generation,” 2018.
- [16] I. Raji and J. Buolamwini, “Actionable Auditing: Investigating the Impact of Publicly Naming Biased Performance Results of Commercial AI Products,” 2019, pp. 429–435. doi: 10.1145/3306618.3314244.
- [17] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language Models are Unsupervised Multitask Learners,” 2018.
- [18] A. Vaswani *et al.*, “Attention is all you need,” *Adv. Neural Inf. Process. Syst.*, 2017.
- [19] N. Prajapati, “The Role of Machine Learning in Big Data Analytics: Tools, Techniques, and Applications,” *ESP J. Eng. Technol. Adv.*, vol. 5, no. 2, pp. 16–22, 2025, doi: 10.56472/25832646/JETA-V5I2P103.
- [20] K. Kishorkumar and R. K. Nimavat, “Machine learning ranking system,” 2025.
- [21] A. H. Anju, “Extreme Gradient Boosting using Squared Logistics Loss function,” *Int. J. Sci. Dev. Res.*, vol. 2, no. 8, pp. 54–61, 2017.
- [22] R. Kumar, “Evaluating and Enhancing Spatial Reasoning in Large Language Models,” in *ICIDA*, 2025.
- [23] Y. Liu *et al.*, “RoBERTa: A Robustly Optimized BERT Pretraining Approach,” *ArXiv*, vol. abs/1907.1, 2019, [Online]. Available: <https://api.semanticscholar.org/CorpusID:198953378>
- [24] A. Ramesh *et al.*, “Zero-Shot Text-to-Image Generation,” vol. 200, 2020.
- [25] J. Wei *et al.*, “Emergent Abilities of Large Language Models,” pp. 1–30, 2022.
- [26] A. Chowdhery *et al.*, “PaLM: Scaling Language Modeling with Pathways,” pp. 1–87.
- [27] H. Kali, “Diversity, Equity, and Inclusion Analytics in HR: How Workday Enables Data-Driven Decision-Making,” vol. 3, no. 2, pp. 162–170, 2023, doi: 10.56472/25832646/JETA-V3I2P113.
- [28] H. Touvron *et al.*, “LLaMA: Open and Efficient Foundation Language Models,” 2022.
- [29] S. S. U., M. S. Nalluri, V. M. Bindu, and M. S. Babu, “Deep Guard: Fortifying Digital Authenticity with Deep Q-Learning and Gorilla Troop Optimization,” *African J. Biol. Sci.*, vol. 6, no. 7, pp. 3868–3887, 2024.
- [30] S. Bubeck *et al.*, “Sparks of Artificial General Intelligence: Early experiments with GPT-4,” 2023. doi: 10.48550/arXiv.2303.12712.
- [31] S. Nokhwal, P. Chilakalapudi, P. Donekal, S. Nokhwal, S. Pahune, and A. Chaudhary, “Accelerating Neural Network Training: A Brief Review,” in *2024 8th International Conference on Intelligent Systems, Metaheuristics & Swarm Intelligence (ISMSI)*, New York, NY, USA: ACM, Apr. 2024, pp. 31–35. doi: 10.1145/3665065.3665071.
- [32] J. Hoffmann *et al.*, “Training Compute-Optimal Large Language Models,” vol. 3, no. 2020, pp. 1–36, 2023.
- [33] A. Balasubramanian, “Improving Air Quality Prediction Using Gradient Boosting,” *Int. J. Sci. Technol.*, vol. 13, no. 2, pp. 1–9, 2022.
- [34] M. Shah and A. Gogineni, “Distributed Query Optimization for Petabyte-Scale Databases,” *Int. J. Recent Innov. Trends Comput. Commun.*, vol. 10, no. 10, pp. 223–231, 2022.
- [35] L. Ouyang *et al.*, “Training language models to follow instructions with human feedback,” 2022.
- [36] R. Thoppilan *et al.*, “LaMDA: Language Models for Dialog Applications,” 2022.
- [37] R. Kumar, “Assessing Maryland Watermark Robustness Against Paraphrasing Attacks,” in *ICIDA 2024*, 2025.
- [38] D. Ganguli *et al.*, “Red Teaming Language Models to Reduce Harms: Methods, Scaling Behaviors, and Lessons Learned,” 2022. doi: 10.48550/arXiv.2209.07858.
- [39] A. Zeng *et al.*, “GLM-130B: An Open Bilingual Pre-trained Model,” *ICLR*, 2023.
- [40] S. R. P. Madugula and N. Malali, “AI-Powered Life Insurance Claims Adjudication Using LLMs and RAG Architectures,” *Int. J. Sci. Res. Arch.*, vol. 15, no. 1, pp. 460–470, Apr. 2025, doi: 10.30574/ijrsra.2025.15.1.0867.
- [41] P. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei, “Deep reinforcement learning from human preferences,” 2017, doi: 10.48550/arXiv.1706.03741.
- [42] E. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, “On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?,” 2021, pp. 610–623. doi: 10.1145/3442188.3445922.
- [43] S. R. Thota, S. Arora, and S. Gupta, “Quantum-Inspired Data Processing for Big Data Analytics,” in *2024 4th International Conference on Advancement in Electronics & Communication Engineering (AECE)*, 2024, pp. 502–508. doi: 10.1109/AECE62803.2024.10911758.
- [44] R. Bommasani and P. Liang, “Holistic Evaluation of Language Models,” *Ann. N. Y. Acad. Sci.*, vol. 1525, 2023, doi: 10.1111/nyas.15007.
- [45] S. Arora and S. R. Thota, “Ethical Considerations and Privacy in AI-Driven Big Data Analytics,” *Int. Res. J. Eng. Technol.*, vol. 11, no. 05, 2024.
- [46] S. Murri, M. Bhoyar, G. P. Selvarajan, and M. Malaga, “Transforming Decision-Making with Big Data Analytics: Advanced Approaches to Real-Time Insights, Predictive Modeling, and Scalable Data Integration,” *Int. J. Commun. Networks Inf. Secur.*, vol. 16, no. 5, pp. 506–519, 2024.
- [47] J. Wei *et al.*, “Finetuned Language Models Are Zero-Shot Learners,” 2021. doi: 10.48550/arXiv.2109.01652.
- [48] V. Sanh *et al.*, “Multitask Prompted Training Enables Zero-Shot Task Generalization,” 2021. doi: 10.48550/arXiv.2110.08207.
- [49] T. Eloundou, S. Manning, P. Mishkin, and D. Rock, “GPTs are GPTs: An Early Look at the Labor Market Impact Potential of Large Language Models,” 2023. doi: 10.48550/arXiv.2303.10130.
- [50] N. Malali, “Adversarial Robustness of AI-Driven Claims Management Systems,” *Int. J. Adv. Res. Sci. Commun. Technol.*, 2025.
- [51] I. Solaiman *et al.*, “Release Strategies and the Social Impacts of Language Models,” 2019. doi: 10.48550/arXiv.1908.09203.