# Foundations and Frontiers of Multimodal Transformer Based Operator Learning: Toward Unified Foundation Models for Language, Vision, Physics, and Robotics

[1] Prof. Rafael D. Costa
[1] Department of Computer Science, Technical University of Munich, Germany

## Abstract

*The rapid evolution of transformer architectures has fundamentally reshaped machine learning across language, vision, multimodal perception, and scientific computing. Initially developed for natural language processing, transformers demonstrated unprecedented few shot learning capabilities, scaling behavior, and contextual reasoning, thereby establishing the paradigm of foundation models. Subsequent adaptations extended the transformer framework to images, multimodal tasks, robotics, and increasingly to scientific domains involving partial differential equations and operator learning. This article synthesizes theoretical and methodological developments spanning large language models, vision transformers, multimodal systems, neural operators, and physics informed transformer architectures. Drawing exclusively on recent foundational works, it develops a unified conceptual and methodological framework that interprets operator learning, in context reasoning, and multimodal integration as manifestations of a broader representational paradigm grounded in tokenization, attention, and scale.*

*The study begins by examining the emergence of few shot language modeling as a foundation paradigm, highlighting its implications for data efficiency, transfer learning, and contextual generalization. It then traces architectural adaptations into visual domains and embodied multimodal systems, emphasizing cross modal alignment and representation sharing. Building upon universal approximation theorems for nonlinear operators and the introduction of Fourier neural operators, the discussion transitions into operator learning for parametric partial differential equations. Recent advances in multimodal PDE foundation models, physics informed token transformers, and efficient PDE specific foundation architectures are analyzed in depth. Special attention is given to unsupervised pretraining, in context operator learning, continuous number encoding, and knowledge distillation in scientific forecasting contexts.*

*Through extensive theoretical elaboration, the article argues that foundation models for physics and robotics represent not merely domain specific adaptations but structural generalizations of transformer based representation learning. Results are presented through descriptive comparative analysis of architectural paradigms, training strategies, and generalization behaviors. The discussion critically examines scalability, interpretability, computational cost, domain shift, and epistemic uncertainty. Finally, the work outlines future directions toward unified multimodal operator foundation models capable of integrating language, perception, physical reasoning, and control in coherent computational systems.*

Keywords: Transformer models, foundation models, neural operators, multimodal learning, partial differential equations, in context learning, robotics.

## 1. Introduction

The development of transformer architectures has fundamentally altered the theoretical and practical landscape of artificial intelligence. Originally introduced for sequence modeling in natural language processing, transformers achieved a breakthrough in representation learning by replacing recurrence with attention mechanisms capable of modeling long range dependencies. The decisive turning point occurred when large scale autoregressive language models demonstrated few shot learning capabilities without task specific fine tuning, thereby inaugurating the era of foundation models (Brown et al., 2020). These models exhibited the ability to generalize across tasks using in context examples, suggesting that sufficiently large networks trained on diverse corpora acquire abstract reasoning capacities emergent from scale.

The implications of this paradigm extend beyond language. The introduction of bidirectional pretraining strategies further advanced contextual representation learning, enabling models to capture deep semantic structure (Devlin et al., 2019). Subsequent analyses framed pretraining as a universal transfer mechanism, positioning large scale self supervised learning as the cornerstone of modern AI systems (Han et al., 2021). The unifying characteristic of these approaches lies in their reliance on tokenization, self attention, and large scale pretraining across heterogeneous datasets.

Parallel developments emerged in the visual domain. The insight that images could be partitioned into patches and processed analogously to word tokens led to the vision transformer architecture, demonstrating competitive performance at scale (Dosovitskiy, 2020). The unification of vision and language representations followed naturally, as models such as VisualBERT introduced cross modal attention mechanisms capable of integrating textual and visual information within a shared representational space (Li et al., 2019). These multimodal extensions broadened the scope of foundation models from purely linguistic tasks to grounded perception.

The extension of transformer architectures into embodied systems further underscored their generality. Multimodal language models capable of integrating perception and action, such as embodied models that couple textual reasoning with sensorimotor data, revealed the feasibility of general purpose robotic reasoning (Driess et al., 2023). Comprehensive surveys of foundation models in robotics highlight both the promise and challenges of scaling transformer architectures to control and physical interaction tasks (Firoozi et al., 2025). These developments suggest that the transformer paradigm may provide a unifying computational substrate for perception, language, and action.

Concurrently, scientific computing witnessed a profound transformation through neural operator learning. Traditional machine learning models approximate functions, mapping finite dimensional inputs to outputs. However, many physical systems are governed by operators mapping functions to functions, as in partial differential equations. The theoretical foundation for approximating nonlinear operators using neural networks was established decades ago (Chen and Chen, 1995). Yet practical architectures capable of efficiently learning such operators emerged only recently. The Fourier neural operator introduced a scalable framework for learning mappings between function spaces using spectral representations, achieving strong generalization across discretizations (Li et al., 2020). This innovation marked a conceptual shift from pointwise function approximation to operator learning.

Recent years have witnessed the convergence of transformer architectures and operator learning. Transformer based operator networks have been proposed for solving partial differential equations through tokenization of spatial and temporal domains (Lorsung et al., 2024). Efficient foundation models tailored for PDEs demonstrate that scale and pretraining strategies analogous to language modeling can be leveraged for scientific domains (Herde et al., 2024). Multimodal PDE foundation models integrate symbolic representations, time series data, and field observations within unified transformer architectures (Liu et al., 2024a; Liu et al., 2024b). These models extend the concept of in context learning from language tasks to operator forecasting and multi physics systems (Cao et al., 2024; Chen et al., 2024).

Despite rapid progress, significant theoretical and methodological gaps remain. First, the conceptual relationship between few shot language models and operator learning architectures has not been systematically articulated. Second, multimodal integration across language, vision, and physics remains fragmented, lacking a unified representational theory. Third, the role of continuous numerical encoding, search strategies, and knowledge distillation in scientific foundation models requires deeper examination (Golkar et al., 2023; Kulikov et al., 2018). Fourth, applications in weather and climate modeling highlight domain specific challenges such as stability, uncertainty quantification, and long horizon forecasting (de Burgh Day and Leeuwenburg, 2023).

This article addresses these gaps by constructing a comprehensive theoretical synthesis grounded exclusively in the referenced works. It proposes that transformer based foundation models and neural operator architectures share a common representational principle: tokenized contextual approximation of operators across heterogeneous modalities. By elaborating the architectural, theoretical, and training paradigms in detail, the article seeks to provide a coherent framework for understanding the next generation of multimodal foundation models that integrate language, perception, physical reasoning, and robotic control.

## 2. Methodology

The methodological approach of this research article is conceptual and synthetic rather than experimental. It systematically analyzes and integrates theoretical contributions from transformer based foundation models, multimodal architectures, and neural operator learning frameworks. The objective is to construct a unified descriptive framework that explains how these paradigms converge and how their underlying principles can be interpreted within a common representational theory.

The first methodological component involves architectural analysis. Transformer architectures are examined through the lens of scaling behavior, attention mechanisms, tokenization strategies, and pretraining objectives. The analysis begins with autoregressive language modeling, emphasizing the role of self attention in enabling contextual inference and few shot generalization (Brown et al., 2020). Bidirectional pretraining is then considered as an alternative representation learning paradigm (Devlin et al., 2019). Comparative examination reveals how attention based architectures facilitate the learning of contextual embeddings that generalize across tasks.

The second methodological component addresses modality extension. Vision transformers are analyzed as a structural generalization of language transformers, replacing word tokens with image patches (Dosovitskiy, 2020). Multimodal extensions are examined through cross modal attention and joint embedding strategies (Li et al., 2019). Embodied multimodal language models are studied to understand how textual reasoning integrates with sensory and motor inputs (Driess et al., 2023). This component identifies the mechanisms enabling representational alignment across heterogeneous modalities.

The third methodological component focuses on operator learning theory. The universal approximation of nonlinear operators by neural networks provides the theoretical foundation (Chen and Chen, 1995). Fourier neural operators are examined as practical architectures for learning mappings between function spaces using spectral parameterization (Li et al., 2020). Bayesian extensions and stochastic gradient Langevin diffusion methods are analyzed for uncertainty quantification in noisy parametric PDEs (Lin et al., 2021). Physics informed token transformers are studied as attention based operator learners that incorporate physical constraints directly into token representations (Lorsung et al., 2024).

The fourth methodological component examines foundation models for PDEs and multi physics systems. Efficient PDE specific foundation architectures demonstrate the feasibility of scaling operator models through pretraining (Herde et al., 2024). Multimodal operator models that integrate symbolic expressions, time series data, and visual representations are analyzed in depth (Liu et al., 2024a; Liu et al., 2024b). Unsupervised pretraining and in context learning strategies are evaluated for data efficient operator generalization (Chen et al., 2024). Knowledge distillation and refinement mechanisms are explored in forecasting tasks (Jollie et al., 2024).

The fifth methodological component addresses numerical encoding and search strategies. Continuous number encoding mechanisms for large language models are analyzed to understand how transformers handle numerical data (Golkar et al., 2023). Search and evaluation strategies in neural modeling are examined as mechanisms that influence generation quality and inference reliability (Kulikov et al., 2018). Regression based loss functions for time series forecasting are reviewed to contextualize training objectives in operator forecasting tasks (Jadon et al., 2024).

Finally, domain specific applications such as numerical weather and climate modeling are incorporated to ground the theoretical synthesis in real world scientific challenges (de Burgh Day and Leeuwenburg, 2023). Multimodal motion prediction in autonomous systems and object detection tasks are analyzed to illustrate how stacked transformers manage temporal and spatial prediction (Liu et al., 2021; Feng et al., 2020).

The methodology thus proceeds through layered conceptual integration. It identifies shared architectural primitives, compares training paradigms, evaluates theoretical guarantees, and synthesizes application specific adaptations. Rather than conducting empirical experiments, it constructs a coherent theoretical narrative supported by rigorous citation of foundational works.

## 3. Results

The results of this conceptual synthesis reveal several unifying principles that cut across language models, multimodal transformers, and neural operator architectures.

First, contextual tokenization emerges as a universal representational mechanism. In language models, tokens correspond to discrete lexical units whose embeddings capture semantic relationships (Brown et al., 2020). In vision transformers, tokens correspond to image patches encoding spatial information (Dosovitskiy, 2020). In operator learning, tokens may represent spatial grid points, spectral modes, or discretized function segments (Li et al., 2020; Lorsung et al., 2024). Across these domains, attention mechanisms compute pairwise interactions between tokens, effectively approximating integral operators through weighted aggregation. This structural similarity suggests that transformer attention can be interpreted as a learnable kernel operator acting on tokenized domains.

Second, scale and pretraining confer emergent generalization properties. Few shot learning in large language models demonstrates that exposure to diverse data distributions enables in context adaptation without gradient updates (Brown et al., 2020). Analogously, unsupervised pretraining for operator learning improves data efficiency and generalization across PDE families (Chen et al., 2024). Efficient PDE foundation models trained on broad distributions of parametric equations exhibit transfer capabilities reminiscent of linguistic foundation models (Herde et al., 2024). These observations indicate that scaling laws may extend beyond language to functional and physical domains.

Third, multimodal alignment enhances operator reasoning. Models that integrate symbolic expressions, time series data, and field observations within a unified transformer demonstrate improved multi operator forecasting (Liu et al., 2024a; Liu et al., 2024b). Embodied multimodal language models show that shared representations can bridge textual instructions and physical actions (Driess et al., 2023). In robotics, foundation models leverage cross modal embeddings to facilitate perception action loops (Firoozi et al., 2025). These results suggest that operator learning benefits from multimodal grounding, particularly when symbolic and numerical information are jointly encoded.

Fourth, spectral and token based operator architectures exhibit complementary strengths. Fourier neural operators efficiently capture global interactions through spectral convolution (Li et al., 2020). Physics informed token transformers incorporate local attention and physical priors directly within token interactions (Lorsung et al., 2024). Bayesian DeepONet variants enhance robustness under noisy conditions (Lin et al., 2021). The coexistence of these approaches indicates that operator learning may require hybrid architectures combining spectral efficiency with attention based adaptability.

Fifth, domain specific applications reveal both potential and limitations. In weather and climate modeling, machine learning models must address stability, long horizon forecasting, and physical consistency (de Burgh Day and Leeuwenburg, 2023). Multimodal motion prediction tasks in autonomous driving highlight challenges in integrating spatial and temporal signals (Liu et al., 2021). Deep multimodal object detection systems underscore the complexity of aligning heterogeneous sensor modalities (Feng et al., 2020). These findings emphasize that while transformer based operator models are powerful, domain constraints impose stringent requirements on reliability and interpretability.

## 4. Discussion

The synthesis presented in this article supports the thesis that transformer based foundation models and neural operator architectures represent different instantiations of a shared representational paradigm. Attention mechanisms approximate context dependent operators over tokenized domains. Pretraining across diverse datasets induces emergent generalization. Multimodal integration aligns heterogeneous representations within shared embedding spaces.

However, several limitations and open challenges persist. Computational cost remains a central concern. Large scale transformer models demand significant memory and energy resources, raising sustainability issues. Efficient PDE specific foundation models attempt to mitigate this challenge through architectural specialization (Herde et al., 2024), yet scalability remains constrained by quadratic attention complexity.

Interpretability constitutes another challenge. While spectral operators provide some analytical transparency, attention weights in large transformers are often difficult to interpret meaningfully. Physics informed token transformers partially address this by embedding physical constraints directly into the architecture (Lorsung et al., 2024). Nevertheless, rigorous theoretical understanding of generalization in operator transformers remains incomplete.

Uncertainty quantification and robustness are particularly

critical in scientific applications. Bayesian approaches and stochastic gradient diffusion methods offer promising directions (Lin et al., 2021). Knowledge distillation strategies may enhance efficiency and stability in forecasting tasks (Jollie et al., 2024). Yet comprehensive frameworks for uncertainty propagation in multimodal operator models are still emerging.

Another unresolved issue concerns numerical representation. Continuous number encoding mechanisms aim to improve the handling of quantitative data in language models (Golkar et al., 2023). Their integration into operator foundation models could enhance numerical precision, particularly in physics simulations. The interplay between discrete tokenization and continuous numerical domains warrants further theoretical investigation.

Finally, the integration of language, vision, physics, and robotics into unified foundation models raises philosophical questions about abstraction and embodiment. Embodied multimodal models demonstrate the feasibility of linking linguistic reasoning with physical action (Driess et al., 2023). Robotics surveys indicate that foundation models may serve as general policy learners (Firoozi et al., 2025). However, ensuring safety, reliability, and domain adaptability remains a substantial challenge.

Future research should explore hybrid architectures combining spectral operators with attention mechanisms, develop scalable pretraining strategies for multi physics datasets, integrate continuous numerical encodings, and establish rigorous theoretical guarantees for generalization across function spaces. Additionally, domain specific benchmarks in weather, climate, and robotics should be expanded to evaluate multimodal operator foundation models under realistic constraints.

## 5.   Conclusion

Transformer architectures have evolved from language specific models into general purpose computational frameworks capable of representing text, images, physical fields, and robotic actions. Neural operator learning extends this paradigm into function space approximation, enabling data driven solutions to partial differential equations. Multimodal integration further broadens applicability, linking symbolic reasoning with perception and control.

This article has provided an extensive theoretical synthesis of these developments, demonstrating that contextual tokenization, attention based operator approximation, and large scale pretraining constitute a unifying foundation. By bridging language modeling, vision transformers, multimodal robotics systems, and PDE foundation models, it articulates a coherent vision of next generation AI systems capable of integrated reasoning across modalities and domains.

The convergence of foundation models and operator learning suggests that the future of artificial intelligence lies in unified multimodal architectures that learn operators over structured domains while leveraging contextual reasoning and scale. Realizing this vision will require continued advances in efficiency, interpretability, robustness, and theoretical understanding. The works synthesized herein provide the conceptual and methodological building blocks for this transformative trajectory.

## References

1.   Brown, T. et al. Language models are few shot learners. Advances in Neural Information Processing Systems. 33:1877 to 1901, 2020.

2.   Cao, S. Choose a transformer: Fourier or Galerkin. Advances in Neural Information Processing Systems. 34:24924 to 24940, 2021.

3.   Cao, Y., Liu, Y., Yang, L., Yu, R., Schaeffer, H., and Osher, S. VICON: Vision in context operator networks for multi physics fluid dynamics prediction. arXiv:2411.16063, 2024.

4.   Chen, T. and Chen, H. Universal approximation to nonlinear operators by neural networks with arbitrary activation functions and its application to dynamical systems. IEEE Transactions on Neural Networks. 6(4):911 to 917, 1995.

5.   Chen, W., Song, J., Ren, P., Subramanian, S., Morozov, D., and Mahoney, M. W. Data efficient operator learning via unsupervised pretraining and in context learning. arXiv:2402.15734, 2024.

6.   de Burgh Day, C. O. and Leeuwenburg, T. Machine learning for numerical weather and climate modelling: A review. Geoscientific Model Development. 16(22):6433 to 6477, 2023.

7.   Devlin, J., Chang, M. W., Lee, K., and Toutanova, K. BERT: Pre training of deep bidirectional transformers for language understanding. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. 4171 to 4186, 2019.

8.   Dosovitskiy, A. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv:2010.11929, 2020.

9.   Driess, D. et al. PaLM E: An embodied multimodal language model. arXiv:2303.03378, 2023.

10.   Feng, D., Haase Schutz, C., Rosenbaum, L., Hertlein,

H., Glaeser, C., Timm, F., Wiesbeck, W., and Dietmayer, K. Deep multi modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges. IEEE Transactions on Intelligent Transportation Systems. 22(3):1341 to 1360, 2020.

11. Firoozi, R. et al. Foundation models in robotics: Applications, challenges, and the future. International Journal of Robotics Research. 44(5):701 to 739, 2025.

12. Golkar, S. et al. xVal: A continuous number encoding for large language models. arXiv:2310.02989, 2023.

13. Han, X. et al. Pre trained models: Past, present and future. AI Open. 2:225 to 250, 2021.

14. Herde, M., Raonic, B., Rohner, T., Kappeli, R., Molinaro, R., Bezenac, E., and Mishra, S. Poseidon: Efficient foundation models for PDEs. arXiv:2405.19101, 2024.

15. Jadon, A., Patil, A., and Jadon, S. A comprehensive survey of regression based loss functions for time series forecasting. International Conference on Data Management, Analytics and Innovation. 117 to 147, 2024.

16. Jollie, D., Sun, J., Zhang, Z., and Schaeffer, H. Time series forecasting, knowledge distillation, and refinement within a multimodal PDE foundation model. arXiv:2409.11609, 2024.

17. Kulikov, I., Miller, A. H., Cho, K., and Weston, J. Importance of search and evaluation strategies in neural dialogue modeling. arXiv:1811.00907, 2018.

18. Li, L. H., Yatskar, M., Yin, D., Hsieh, C. J., and Chang, K. W. VisualBERT: A simple and performant baseline for vision and language. arXiv:1908.03557, 2019.

19. Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., and Anandkumar, A. Fourier neural operator for parametric partial differential equations. arXiv:2010.08895, 2020.

20. Lin, G., Moya, C., and Zhang, Z. Accelerated replica exchange stochastic gradient Langevin diffusion enhanced bayesian DeepONet for solving noisy parametric PDEs. arXiv:2111.02484, 2021.

21. Liu, Y., Zhang, J., Fang, L., Jiang, Q., and Zhou, B. Multimodal motion prediction with stacked transformers. Proceedings of the IEEE CVF Conference on Computer Vision and Pattern Recognition. 7573 to 7582, 2021.

22. Liu, Y., Zhang, Z., and Schaeffer, H. PROSE: Predicting multiple operators and symbolic expressions using multimodal transformers. Neural Networks. 180:106707, 2024.

23. Liu, Y., Sun, J., He, X., Pinney, G., Zhang, Z., and Schaeffer, H. PROSE FD: A multimodal PDE foundation model for learning multiple operators for forecasting fluid dynamics. arXiv:2409.09811, 2024.

24. Lorsung, C., Li, Z., and Farimani, A. B. Physics informed token transformer for solving partial differential equations. Machine Learning: Science and Technology. 5(1):015032, 2024.