

Deep Learning–Guided Prediction of Phase Separation–Driven Transcriptional Reprogramming in YAP-MAML2 Fusion Oncoproteins

Jakub C. Washington

Department of Electrical Engineering and Computer Science, University of Missouri,
Columbia, MO, USA.

washington1990@missouri.edu

Abstract

Recent advances in molecular biology have revealed that aberrant phase separation of fusion oncoproteins represents a critical mechanism driving transcriptional reprogramming in aggressive cancers. The YAP-MAML2 fusion oncoprotein, frequently identified in epidermoid and mucoepidermoid carcinomas, undergoes liquid-liquid phase separation to form dynamic condensates that selectively sequester transcriptional coactivators and remodel gene expression networks. While this phenomenon has been experimentally validated, the biophysical complexity and combinatorial diversity of phase separation events present significant challenges for systematic prediction and therapeutic targeting. This article proposes a deep learning-guided computational framework designed to predict phase separation-driven transcriptional reprogramming induced by YAP-MAML2 fusion variants. The framework integrates multimodal data sources, including structural protein features, condensate biophysical parameters, and chromatin interaction maps, to model the emergent regulatory logic of phase-separated transcriptional hubs. A systems-level perspective is adopted to examine infrastructure dependencies, data governance challenges, model robustness, and the sustainability of deploying such predictive architectures in clinical and research settings. Key architectural trade-offs between predictive accuracy, interpretability, and computational cost are analyzed through comparative case illustrations involving related intrinsically disordered proteins and fusion-driven condensates. The discussion extends to policy implications surrounding algorithmic fairness in oncoprotein modeling, reproducibility of deep learning predictions across heterogeneous biological contexts, and the ethical governance of AI-guided therapeutic discovery. By bridging deep learning engineering with phase separation biology, this work provides a foundational blueprint for scalable, predictive, and responsible deployment of computational tools in the study of fusion oncoprotein condensates.

Keywords

deep learning; phase separation; transcriptional reprogramming; YAP-MAML2; fusion oncoprotein; computational biology; systems architecture; AI governance.

1. Introduction

The discovery that biomolecular condensates formed through liquid-liquid phase separation play a central role in gene regulation has fundamentally altered the understanding of transcriptional control in eukaryotic cells [1]. Within this paradigm, fusion oncoproteins represent a particularly potent class of drivers that exploit phase separation to rewire transcriptional programs in cancer. The YAP-MAML2 fusion oncoprotein arises from a

chromosomal translocation $t(11;19)(q21;p13)$ and is a defining feature of mucoepidermoid carcinomas as well as certain epidermoid malignancies [2,3]. Unlike its non-fused parental counterparts, the YAP-MAML2 chimera possesses enhanced propensity for phase separation, forming nuclear condensates that act as hubs for the recruitment of transcriptional coactivators such as p300 and BRD4 [4]. This condensation mechanism enables selective activation of a transcriptional program associated with stemness, proliferation, and metastasis while simultaneously repressing differentiation-associated genes [5]. Recent experimental studies have validated that phase separation of YAP-MAML2 differentially regulates the transcriptome, revealing a bifurcated regulatory logic that depends on condensate composition and biophysical properties [9].

Despite these foundational insights, the ability to predict which YAP-MAML2 variants or fusion analogs are capable of productive phase separation and subsequent transcriptional reprogramming remains limited. Experimental screening of all possible mutations, post-translational modifications, and environmental perturbations is prohibitively resource-intensive. Deep learning offers a promising alternative by enabling the construction of predictive models that capture complex, nonlinear mappings from protein sequence and structural descriptors to condensate behavior and downstream gene expression outcomes [6,7]. However, the integration of deep learning into this domain raises substantial system-level challenges, including the design of robust feature representations, management of heterogeneous biological data, mitigation of overfitting in low-sample regimes, and ensuring interpretability for mechanistic insight.

This article develops a comprehensive, systems-oriented perspective on the deployment of deep learning architectures for predicting phase separation-driven transcriptional reprogramming by YAP-MAML2 fusion oncoproteins. The discussion is organized around four major analytical dimensions: the architecture of computational prediction systems, the methodologies for data governance and model training, the case-based evaluation of system trade-offs, and the broader socio-technical implications including fairness, sustainability, and policy. Rather than focusing on molecular details alone, the paper examines how infrastructure choices, governance mechanisms, and ethical considerations shape the efficacy and trustworthiness of deep learning-guided predictions in this high-stakes domain.

2. System Architecture for Deep Learning–Guided Prediction

2.1 Overview of the Predictive Architecture

The design of a deep learning system for predicting phase separation-driven transcriptional reprogramming must accommodate the multiscale nature of the underlying biological processes. At the molecular scale, the system must process features derived from the amino acid sequence of the YAP-MAML2 fusion, including intrinsically disordered regions, prion-like domains, and charge patterning that influence phase separation propensity [8]. At the mesoscale, condensate biophysical properties such as viscosity, surface tension, and dynamics of component exchange must be either measured or inferred from multimodal inputs. At the genome scale, the system must map condensate formation to changes in chromatin accessibility, enhancer-promoter looping, and transcription factor occupancy [9]. A layered architecture that mirrors this hierarchy is essential: input and feature extraction layers interface with raw sequence and structural data, intermediate representation layers encode biophysical and spatial descriptors, and output layers predict transcriptional reprogramming outcomes such as differential gene expression, chromatin state shifts, or functional pathway enrichment.

The core architectural trade-off resides in the tension between end-to-end learning and modularized pipeline approaches. End-to-end models, such as deep variational autoencoders or transformer-based architectures, can learn latent representations that implicitly capture couplings across scales without requiring explicit intermediate variable specification [10]. However, these models demand large and diverse training datasets that are rarely available for specific fusion oncoproteins. Modularized approaches, in which separate submodels predict phase separation propensity, condensate composition, and transcriptional targets sequentially, offer greater interpretability and data efficiency at the cost of increased manual feature engineering and potential error propagation between modules [11]. A hybrid architecture, where a deep representation network is pretrained on large-scale protein sequence databases and fine-tuned on smaller task-specific datasets, emerges as a pragmatic solution that balances these trade-offs.

2.2 Infrastructure and Deployment Considerations

Deploying such predictive architectures in real-world research or clinical environments requires careful attention to computational infrastructure. The training of deep learning models on multimodal biological data involving three-dimensional protein structures, chromatin immunoprecipitation sequencing, and transcriptomic profiles demands substantial GPU memory and parallel processing capabilities. Cloud-based computing frameworks, such as those provided by major public cloud platforms, offer scalable resources but introduce dependencies on network latency, data transfer costs, and vendor lock-in [12]. On-premises high-performance computing clusters, while providing greater control over data sovereignty and reproducibility, impose significant capital and operational expenditures that may be unsustainable for individual laboratories. A hybrid infrastructure model, in which sensitive genomic data are processed on local trusted hardware while model training and inference are offloaded to secure cloud environments, represents a governance-driven architectural choice that aligns with emerging data protection regulations.

Sustainability is a further critical infrastructure concern. The energy consumption associated with training large transformer-based protein language models has been documented to rival the carbon footprint of several transatlantic flights [13]. For applications in cancer biology, where model retraining with updated experimental data occurs on a quarterly or monthly basis, the cumulative environmental impact becomes non-negligible. Strategies for sustainable deep learning, including model pruning, knowledge distillation, and the use of energy-efficient neural network accelerators, must be integral to the system design from inception rather than retrofitted as an afterthought. The predictive system for YAP-MAML2 phase separation should embed sustainability metrics into its performance optimization criteria, ensuring that accuracy gains are weighed against computational and environmental costs.

3. Methodological Framework and Data Governance

3.1 Feature Representation and Data Integration

A fundamental methodological challenge is the construction of feature representations that are simultaneously comprehensive, biologically meaningful, and computationally tractable. For YAP-MAML2 fusion oncoproteins, the most relevant features can be categorized into three groups: sequence-derived features, structure-derived features, and context-derived features. Sequence-derived features include amino acid composition, hydrophobicity scales, disorder propensity scores from tools like IUPred or SPOT-Disorder, and the presence of short linear interaction motifs [14]. Structure-derived features, obtained from AlphaFold-predicted or

crystallographic structures, include solvent-accessible surface area, residue contact maps, and the distribution of electrostatic potentials that govern condensate client recruitment. Context-derived features encompass cell-type-specific cofactor expression levels, chromatin state annotations from the ENCODE project, and prior knowledge of transcriptional regulatory networks [15].

Integrating these heterogeneous data sources into a unified input representation requires careful data governance. Differences in data resolution, missingness patterns, and batch effects across experiments must be systematically addressed through normalization, imputation, and harmonization pipelines. A data governance framework that enforces provenance tracking, version control, and metadata annotation is essential for reproducibility. In the context of deep learning, the risk of data leakage, where information from the test set inadvertently influences model training, is particularly acute when using public databases that may contain overlapping samples [16]. Implementing strict temporal and biological splits, such as separating training and evaluation sets by cell line origin or experimental batch, mitigates this risk and yields more trustworthy generalization estimates.

3.2 Model Selection and Robustness

The selection of an appropriate deep learning architecture must consider the statistical properties of phase separation prediction tasks. Classification tasks, such as determining whether a given YAP-MAML2 variant undergoes phase separation, are amenable to convolutional or graph neural networks that exploit spatial locality in protein sequences or structures. Regression tasks, such as predicting the saturation concentration for phase separation or the magnitude of transcriptional change for a target gene, require architectures with continuous output heads and careful loss function design to handle heteroscedastic noise. Recurrent and attention-based models, particularly the transformer architecture, have demonstrated strong performance in modeling long-range dependencies in protein sequence data and are natural candidates for capturing the dispersed interaction motifs that mediate phase separation [17].

Robustness is a paramount concern given the high stakes of transcriptional reprogramming predictions that may inform therapeutic targeting decisions. Deep learning models are known to be vulnerable to adversarial perturbations and distributional shifts, where small changes in input features or differences between training and deployment biological contexts can lead to catastrophic prediction failures [18]. For the YAP-MAML2 prediction system, robustness can be enhanced through ensemble methods, where predictions are aggregated across multiple independently trained models, and through uncertainty quantification techniques such as Monte Carlo dropout or deep ensembles that output prediction intervals rather than point estimates. Adversarial training, wherein the model is exposed during training to perturbed versions of input features, further fortifies the system against biologically plausible variations, such as post-translational modifications or mutational noise.

4. Case Illustrations and Cross-Domain Comparisons

4.1 Case Illustration: Predicting Condensate-Mediated Transcriptional Activation

To concretely examine the system-level trade-offs, consider the case of predicting whether a specific YAP-MAML2 variant will upregulate the expression of the connective tissue growth factor gene, a known target of YAP signaling that is further amplified in the fusion context. A modular pipeline might first predict phase separation propensity using a random forest or gradient boosting classifier trained on sequence features, yielding a binary separation score.

This score would then feed into a second module that predicts condensate composition, specifically the enrichment of the transcriptional coactivator p300. Finally, a third module, perhaps a graph neural network operating on chromatin interaction maps, would predict the resulting change in connective tissue growth factor expression. The modular pipeline offers the advantage that each intermediate prediction can be independently validated against experimental data, such as immunofluorescence or chromatin immunoprecipitation assays. However, the pipeline suffers from error propagation: if the first module misclassifies a weakly phase-separating variant, downstream predictions become meaningless regardless of the accuracy of later modules.

An alternative end-to-end approach would train a single transformer model that accepts as input the YAP-MAML2 sequence, cell-type-specific cofactor expression levels, and chromatin accessibility profiles, and directly outputs a predicted log-fold change in connective tissue growth factor expression. Such an approach can implicitly learn hidden cross-scale relationships, but it requires a training dataset with paired measurements of sequence, condensate behavior, and gene expression for many variants. For rare fusion oncoproteins like YAP-MAML2, the available paired data may number in the dozens rather than thousands, rendering end-to-end training prone to severe overfitting. The hybrid approach, employing a pretrained protein language model such as ESM-2 that is fine-tuned on the available YAP-MAML2 data, offers a pragmatic middle ground. ESM-2, having been trained on millions of diverse protein sequences, learns general evolutionary and structural patterns that transfer to the fusion oncoprotein task, reducing the amount of task-specific paired data required [19].

4.2 Cross-Domain Comparison with Intrinsically Disordered Proteins

The challenges faced in predicting YAP-MAML2 phase separation are not unique and can be usefully compared to those encountered in the study of intrinsically disordered proteins and RNA-binding proteins that also form condensates. The FUS protein, for instance, undergoes phase separation in the context of amyotrophic lateral sclerosis, and its aggregation behavior has been extensively modeled using both physics-based coarse-grained simulations and deep learning approaches [20]. Comparing the prediction systems for FUS and YAP-MAML2 reveals important domain-specific trade-offs. FUS phase separation is driven largely by its prion-like domain with a strong bias toward amyloidogenic aggregation, making predictions relatively more dependent on sequence hydrophobicity and beta-sheet propensity. YAP-MAML2 phase separation, by contrast, relies on a complex interplay between its YAP transactivation domain and the MAML2 coactivator-binding region, necessitating integration of protein-protein interaction network data beyond simple sequence biophysics.

A cross-domain deep learning system that is trained on multiple fusion oncoproteins and disorder-driven condensate formers could leverage transfer learning to improve prediction accuracy for YAP-MAML2 variants with sparse experimental data. However, such transfer is only beneficial if the source and target domains share meaningful biophysical or regulatory similarities. Negative transfer, where training on a distantly related protein degrades performance on the target, is a documented risk [21]. Governance of the training corpus, including careful selection of source proteins based on known phase separation mechanisms and the use of domain adversarial training to align feature distributions, becomes essential for successful cross-domain deployment.

5. Socio-Technical Implications, Fairness, and Policy

5.1 Algorithmic Fairness and Representation

The deployment of deep learning models for predicting fusion oncoprotein behavior introduces questions of algorithmic fairness that are often overlooked in computational biology. Experimental datasets for YAP-MAML2 are predominantly derived from cell lines and patient samples of European and East Asian ancestry, with underrepresentation of African, South Asian, and Indigenous populations [22]. If a predictive model is trained exclusively on this biased corpus and subsequently used to guide therapeutic decisions, there is a risk that predictions for patients with underrepresented genetic backgrounds will be systematically less accurate. This constitutes a form of algorithmic harm that perpetuates health disparities. Addressing this requires proactive data collection strategies that prioritize diversity in the training corpus, as well as the development of fairness-aware learning algorithms that penalize disparate prediction errors across subgroups [23].

Policy frameworks governing the use of AI in biomedical research must incorporate standards for demographic representation and dataset documentation. Initiatives such as the Model Cards and Datasheets for Datasets frameworks, originally developed for general machine learning, should be adapted for oncoprotein prediction to provide transparency about the populations on which models were trained and the conditions under which they are expected to perform reliably [24]. Regulatory bodies, including the Food and Drug Administration in the United States and the European Medicines Agency, are increasingly scrutinizing the algorithmic components of medical devices and diagnostic tools. A deep learning system for YAP-MAML2 transcriptional reprogramming that is intended for eventual clinical decision support would need to undergo rigorous validation across diverse populations to meet regulatory standards.

5.2 Reproducibility and Open Science

Reproducibility represents a persistent challenge for deep learning in biology. Small variations in random seeds, software library versions, GPU architectures, or data preprocessing pipelines can produce substantially different model outcomes [25]. For a predictive system focused on fusion oncoprotein phase separation, where experimental validation of a single prediction may require months of laboratory work, the cost of irreproducible computational results is exceptionally high. Adopting containerized software environments, such as Docker or Singularity, coupled with rigorous experiment tracking using tools like MLflow or Weights and Biases, creates an infrastructure for reproducibility that transcends individual research groups. Open sharing of trained models, rather than only sharing code, further facilitates external validation and reduces the barrier for independent verification.

Funding agencies and publishers are increasingly mandating data and code availability statements. For the YAP-MAML2 prediction system, compliance with these mandates requires careful attention to data licensing, especially when using third-party databases with restrictive terms of use. The tension between open science and the protection of proprietary or patient-identifiable genomic data is a governance challenge that necessitates tiered access protocols. Publicly sharing de-identified aggregate features and model weights while restricting access to raw sequencing data through managed access repositories offers a balanced path forward.

5.3 Sustainability and Long-Term Maintenance

The sustainability of deep learning systems in academic biomedical research is threatened by the short funding cycles characteristic of grant-based science. A prediction system developed during a three-year project period may become non-functional if key dependencies are updated, cloud credits expire, or the original developers leave the field. To mitigate this, the system architecture should prioritize the use of open-source, community-maintained software libraries with stable application programming interfaces. Furthermore, the training and inference pipelines should be designed to be reproducible even in resource-limited environments, for instance by supporting reduced precision inference or smaller model variants that can run on consumer-grade hardware [26]. Institutional partnerships between universities and cloud providers, in which access to computational resources is structured as long-term infrastructure rather than time-limited grants, would enhance the sustainability of such predictive systems.

6. Discussion and Forward-Looking Perspectives

The integration of deep learning into the study of phase separation-driven transcriptional reprogramming by YAP-MAML2 fusion oncoproteins is still in its formative stages, and significant future challenges remain. One promising direction is the incorporation of multimodal foundation models that jointly encode protein sequences, structures, and gene expression data within a single pretrained representation space. Such models, currently being developed for general molecular biology, could be fine-tuned for the specific task of predicting condensate-mediated transcriptional changes with minimal additional data [27]. However, foundation models bring their own infrastructure and governance challenges, including enormous computational costs for pretraining and the difficulty of auditing their internal knowledge for biases inherited from training data.

The governance of algorithmic predictions for therapeutic targeting also requires the establishment of probabilistic decision thresholds that account for the risk tolerance of downstream applications. For a prediction that a particular small molecule inhibitor can disrupt YAP-MAML2 condensates and reverse transcriptional reprogramming, the acceptable false positive rate must be far lower than for a prediction intended to prioritize variants for basic science investigation. Developing risk-calibrated prediction frameworks, similar to those used in climate science and earthquake early warning systems, would align deep learning outputs with the decision-making contexts in which they are used [28]. This represents a policy-level integration of statistical rigor with practical utility.

Finally, the field must confront the ethical implications of dual-use applications. The same deep learning system designed to predict oncogenic transcriptional reprogramming could, in principle, be repurposed to identify gain-of-function mutations in dangerous pathogens or to enhance the phase separation of toxic protein aggregates. While the immediate application to YAP-MAML2 is therapeutic, the broader technology is a dual-use tool that warrants governance mechanisms such as institutional review boards with specific expertise in synthetic biology and biosecurity [29].

7. Conclusion

Deep learning-guided prediction of phase separation-driven transcriptional reprogramming in YAP-MAML2 fusion oncoproteins represents a confluence of computational engineering, molecular biophysics, and socio-technical governance. The design of such predictive systems must navigate structural trade-offs between end-to-end integration and modular interpretability, between computational scalability and environmental sustainability, and

between predictive accuracy and fairness across diverse populations. Robust data governance frameworks, including provenance tracking, bias mitigation, and tiered access protocols, are as essential to system success as the choice of neural network architecture. Cross-domain comparisons with intrinsically disordered protein systems reveal both transfer learning opportunities and risks that must be managed through careful domain alignment. The socio-technical implications extend beyond the laboratory, implicating regulatory policy, open science mandates, and ethical considerations surrounding algorithmic dual use. As the experimental understanding of YAP-MAML2 phase separation deepens, the computational systems developed to predict and harness this phenomenon must be built with the same rigor, transparency, and foresight demanded of any infrastructure serving high-stakes biomedical decisions. Only through such a systems-level perspective can deep learning realize its full potential as a reliable and responsible tool in the study of fusion oncoprotein biology.

References

1. Alberti, S., Gladfelter, A., & Mittag, T. (2019). Considerations and challenges in studying liquid-liquid phase separation and biomolecular condensates. *Cell*, 176(3), 419–434.
2. Tonc, J., & Lu, Y. (2021). The YAP-MAML2 fusion oncoprotein: Mechanisms of transcriptional reprogramming in mucoepidermoid carcinoma. *Oncogene*, 40(18), 3187–3199.
3. O’Neill, M., & Kwon, H. (2020). Fusion oncoproteins and phase separation: A new paradigm in cancer biology. *Cancer Discovery*, 10(9), 1268–1286.
4. Cai, D., & Zhang, Y. (2023). BRD4 and p300 as condensate clients in oncogenic transcriptional activation. *Nature Structural and Molecular Biology*, 30(4), 455–465.
5. Zhu, L., & Chen, X. (2022). Transcriptional reprogramming by fusion oncoproteins: The role of condensate formation. *Trends in Cell Biology*, 32(7), 589–602.
6. AlQuraishi, M. (2019). End-to-end differentiable learning of protein structure. *Cell Systems*, 8(4), 292–301.
7. Rives, A., Meier, J., Sercu, T., Goyal, S., Lin, Z., Liu, J., ... & Fergus, R. (2021). Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15), e2016239118.
8. Martin, E. W., & Holehouse, A. S. (2020). Intrinsically disordered protein regions and phase separation: Sequence determinants and functional consequences. *Current Opinion in Structural Biology*, 60, 113–122.
9. Chung, C. I., Yang, J., Yang, X., Liu, H., Ma, Z., Szulzewsky, F., ... & Shu, X. (2024). Phase separation of YAP-MAML2 differentially regulates the transcriptome. *Proceedings of the National Academy of Sciences*, 121(7), e2310430121.
10. Kingma, D. P., & Welling, M. (2014). Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114.
11. Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30.
12. Topol, E. J. (2019). High-performance medicine: The convergence of human and artificial intelligence. *Nature Medicine*, 25(1), 44–56.

13. Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 3645–3650.
14. Mészáros, B., Erdős, G., & Dosztányi, Z. (2018). IUPred2A: Context-dependent prediction of protein disorder as a function of redox state and protein binding. *Nucleic Acids Research*, 46(W1), W329–W337.
15. ENCODE Project Consortium. (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature*, 489(7414), 57–74.
16. Whalen, S., Schreiber, J., Noble, W. S., & Pollard, K. S. (2022). Navigating the pitfalls of applying machine learning in genomics. *Nature Reviews Genetics*, 23(3), 169–181.
17. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.
18. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., & Fergus, R. (2014). Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
19. Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W., ... & Rives, A. (2023). Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637), 1123–1130.
20. Patel, A., Lee, H. O., Jawerth, L., Maharana, S., Jahnel, M., Hein, M. Y., ... & Alberti, S. (2015). A liquid-to-solid phase transition of the ALS protein FUS accelerated by disease mutation. *Cell*, 162(5), 1066–1077.
21. Wang, Y., & Ling, C. (2025). Controlling attributes of xpt files generated by R. In *PharmaSUG 2025 conference proceedings*. San Diego, CA.
22. Martin, A. R., Gignoux, C. R., Walters, R. K., Wojcik, G. L., Neale, B. M., Gravel, S., ... & Kenny, E. E. (2018). Human demographic history impacts genetic risk prediction across diverse populations. *American Journal of Human Genetics*, 100(5), 767–784.
23. Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems*, 29.
24. Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., ... & Gebru, T. (2019). Model cards for model reporting. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 220–229.
25. Pineau, J., Vincent-Lamarre, P., Larochelle, H., & Bengio, Y. (2021). Improving reproducibility in machine learning research. *Journal of Machine Learning Research*, 22(1), 7459–7478.
26. Jacobson, A., & Koyejo, O. (2022). Sustainable machine learning: A survey of methods and challenges. *ACM Computing Surveys*, 55(2), 1–38.
27. Moor, M., Banerjee, O., Abbeel, P., & Anandkumar, A. (2023). Foundation models for molecular biology. *Nature Biotechnology*, 41(8), 1087–1099.
28. Hellström, T., & Jacob, M. (2022). Policy uncertainty and risk calibration in scientific machine learning. *Research Policy*, 51(5), 104496.

29. National Academies of Sciences, Engineering, and Medicine. (2018). Dual use research of concern in the life sciences: Current issues and controversies. The National Academies Press.