# **Generative Biology**

This manuscript (<u>permalink</u>) was automatically generated from <u>In-Vivo-Group/generative-biology@9e31e96</u> on November 7, 2023.

#### **Authors**

- Alexander J. Titus <sup>™</sup>
  - **(D** 0000-0002-0145-9564 **· (7** alexandertitus

In Vivo Group, Washington, DC, USA; International Computer Science Institute, Berkeley, CA, USA · Funded by Grant TBD

- Matthew E. Walsh

Center for Health Security, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA; Department of Environmental Health and Engineering, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

#### **Abstract**

The rapid pace of progress in generative artificial intelligence (AI) techniques like deep learning, reinforcement learning, and transformer neural networks is transforming the life sciences and biomedicine. This living review paper provides an updatable, comprehensive overview and analysis of the latest literature on generative biology – the application of cutting-edge generative AI methods to accelerate insights and innovation across the life sciences and healthcare. All authors are welcome to contribute to this review via pull requests.

The review synthesizes key developments in using generative models for de novo biomedical discovery, design, and decision support. It examines techniques and applications including deep learning on omics data for personalized medicine, generative chemistry for drug development, protein structure prediction for molecular engineering, image synthesis for pathology, language models for clinical decision support, robotic simulation for prosthetics, and generative networks for cell programming.

The review highlights representative studies and benchmarks in each area while contextualizing progress, limitations, emerging best practices, and directions for future work. It also discusses social and ethical challenges raised by generative biology applications, such as compounding bias, system opacity, and dual-use risks, alongside proposed solutions.

As a living review, this paper will be continually updated as the field rapidly advances to provide researchers and practitioners with an up-to-date reference on the state of the art in employing generative AI to accelerate biomedicine for the collective good.

# **Executive Summary**

The goal with be a 2 page TL;DR of the review after v1 is complete. Need more content.

## Introduction

This is the start of the Generative Biology living review!

## **Recent Advances in Generative Al**

#### Introduction

Brief background on rise of generative models like GPT-3, DALL-E, AlphaFold, etc. Summary of scope and goals of chapter focused on key advances in last 1-2 years.

### **Transformer-Based Language Models**

#### **GPT-3 and Foundation Models**

- Overview of GPT-3 architecture and self-supervised training on massive text corpus
- Discussion of GPT-3 capabilities and limits, including few-shot learning
- Concept of foundation models as basis for many downstream applications

#### Other Notable Models

- Summary of other major transformer language models like Google's PaLM, DeepMind's Gopher, Meta's OPT, Anthropic's Claude etc.
- · Comparison of model sizes, architectures, training approaches
- Benchmarking of models on various NLP tasks

#### **Multimodal Generative Models**

### **DALL-E 2 and Text-to-Image Generation**

- Explain DALL-E 2 architecture and training methodology
- Discuss capabilities in text-to-image generation
- Issues around bias, appropriate use cases

#### Other Multimodal Models

- Overview of models like Imagen, Parti, Flamingo for text-to-image
- Discussion of video generation models like Googles Imagen Video
- Models for text to 3D shapes, text to music etc.

#### **Outlook**

• Key challenges and limitations of current generative models

- Likely future advances building on these models
- Broader societal impact of widely available generative models

## **Advances in Generative AI for Proteins**

#### Introduction

- Background on proteins as key molecular machines in biology
- Promise of generative AI to accelerate protein discovery and engineering
- Overview of scope covering recent advances in last 1-2 years

#### **Structure Prediction**

- AlphaFold2 as a breakthrough method for structure prediction
- Novel model architecture and training methodology
- Examples of new biological insights from predicted structures

#### **Function Prediction**

- Using predicted structures to infer protein functions
- Structure-based identification of catalytic and binding sites
- Case studies of novel enzyme functions discovered

### **Designing Novel Proteins**

- Generative models for designing functional protein sequences
- Leveraging structural constraints for optimized protein engineering
- · Applications in industrial enzymes, therapeutics, biomaterials

#### **Interaction Prediction**

- Modeling protein-protein interactions with graph networks
- Structure-based prediction of protein-drug bindings
- Applications in drug discovery and toxicity screening

#### Outlook

- Challenges and next steps in improving accuracy
- Hybrid physics- and data-driven approaches
- Ethical considerations in synthetic protein design

# **Advances in Generative AI for RNA**

#### Introduction

- Background on key roles of RNA in biology
- · Promise of generative models to advance RNA research
- Scope focused on latest advances in RNA prediction and design

#### **Structure Prediction**

- Transformer-based prediction of RNA folding
- Novel architectures for incorporating chemical constraints
- Improved accuracy on complex structures like ribosomes

#### **Function Prediction**

- Inferring RNA functions from sequence and structure
- Identifying motifs, domains, and atomic binding sites
- Applications in understanding long noncoding RNAs

#### Interaction Prediction

- Graph neural networks for RNA-protein interactions
- Structure-augmented modeling of splice sites
- Predicting RNA base editing targets

### **Design of RNA Therapeutics**

- Generative models for optimizing siRNA, antisense design
- Reinforcement learning for chemical modification patterns
- Progress in computational RNA-targeted drug design

### **Outlook**

- Key challenges in prediction of long RNA structures
- Design of RNA for self-assembly and scaffolding
- Ethical use of synthetic RNA technologies

### Advances in Generative AI for DNA

### Introduction

- Background on the role of DNA as a carrier of genetic information
- Promise of generative models to advance DNA research
- Scope focused on the latest advances in DNA prediction and design

### **Sequence Modeling**

### Transformer architectures for modeling DNA sequences

The HyenaDNA paper introduces a new genomic foundation model called HyenaDNA that can process DNA sequences at single nucleotide resolution with ultralong context lengths up to 1 million base pairs - a 500x increase over previous transformer models. HyenaDNA uses a parameter-efficient convolutional architecture that allows it to scale subquadratically with sequence length, enabling the use of full genome-scale context. On a range of regulatory genomics prediction tasks, HyenaDNA matches or exceeds the performance of previous state-of-the-art models while using 1500x fewer parameters and 3200x less pretraining data. HyenaDNA also demonstrates the ability to perform

challenging species classification using the full mutational profile visible at 1 million base pairs of context. The authors explore new training techniques to enable ultralong sequence modeling as well as prompt-based tuning methods for rapidly adapting to new tasks without updating pre-trained weights [1].

### Pretraining on large genomic datasets

### Applications in variant calling and annotation

### **Regulation Prediction**

- Graph neural networks for modeling 3D genome architecture
- Predicting enhancer-promoter interactions and expression
- Design of synthetic promoters and enhancers

### **Genome Editing**

- · Generative models for CRISPR guide design
- Contextual prediction of on-target editing efficacy
- Modeling of off-target effects during optimization

#### **DNA Data Generation**

- Variational autoencoders for realistic DNA sequences
- Generating paired genomic-transcriptomic data
- Applications in training genome interpretation models

#### **Outlook**

### Challenges in predicting long-range chromatin interactions

The development of transformer architectures like HyenaDNA that can leverage genome-scale context creates new opportunities for understanding long-range chromatin interactions, gene regulation, and intercellular networks. However, significant challenges remain in improving model accuracy and uncertainty quantification. Hybrid physics- and data-driven approaches may help address these gaps. Safety considerations around synthetic genome design also warrant further research to ensure responsible innovation as these generative capabilities advance. Overall, the future looks bright for generative models that can capture both local mutations and global patterns critical to biology [1].

### Responsible design of synthetic genomes

Ethical considerations for human genome editing

# Generative Al for Autonomous Experimentation

### Introduction

- Promise of generative models to accelerate scientific discovery
- Rise of autonomous labs and robotics for automated experimentation
- Overview of generative Al's role in self-driving research

### **Closed-Loop Systems**

- Integrating computational hypothesis generation with robotic wet lab testing
- Reinforcement learning pipelines for autonomous optimization
- Case studies in materials science, drug discovery

### **Automated Experiment Design**

- Using generative models to design novel compounds, genes
- Leveraging simulations to predict experimental outcomes
- Robotic execution of designed experiments

### **Adaptive Sampling**

- Active learning to iteratively select most informative experiments
- · Bayesian optimization powered by neural networks
- · Applications in probing molecular design spaces

### **Real-Time Learning**

- Deploying models on lab edge devices
- Online learning from experimental data streams
- Improving models and experiment plans on-the-fly

#### Outlook

- Key challenges around model accuracy and integration
- The future of data-driven, self-driving laboratories
- Risks of full automation and need for human oversight

### **Conclusion**

- Generative Al as a powerful tool for autonomous experimentation
- Accelerating discovery alongside human researchers
- Responsible implementation will maximize benefit

# Generative AI and the Biosecurity Landscape

#### Introduction

- Background on biosecurity threats from natural, accidental, and intentional pathogens
- Rise of generative AI as a dual use technology for biodefense and misuse

### **Enhanced Risks**

- Automated bioweapon design with generative models
- Relatively low computing needs to generate dangerous agents
- Challenges detecting artificially generated sequences/organisms

### **Enhanced Response Capabilities**

- Generative models for vaccine and therapeutic design
- High-throughput testing of countermeasures with synthetic data
- Al for early detection of emergent pathogens and outbreaks

#### Recommendations

- Increased oversight for generative model development/release
- Expanding biosecurity legislation and regulations
- Fostering open research and global cooperation

#### **Outlook**

- Trajectory toward increasingly powerful generative biological capabilities
- Need for preventative ethics research and guidance
- · Maintaining public trust and avoiding overreaction

#### Conclusion

- · Balancing generative Al's benefits and risks in biology
- Importance of thoughtful governance and responsible innovation
- Staying ahead of the curve on biosecurity

# Policy Responses to Generative AI in Biology

#### Introduction

- Background on rise of powerful generative models for biology
- Overview of risks like bioweapons, environmental damage
- Need for governance to ensure responsible development

### **Self-Governance by Developers**

- Voluntary guidelines on ethical AI by corporations
- Limiting access to certain capabilities like human editing
- Issues with self-regulation and transparency

### **Governmental Regulations**

- New biosecurity regulations on certain AI technologies
- Restrictions on use of synthetic biology IP
- Challenges with fast pace of technology change

### **International Governance**

- Proposals for global observatory to monitor risks
- Treaties restricting development of bioweapons
- Difficulty achieving consensus and compliance

#### **Public Deliberation and Ethics**

- Activities to involve broader stakeholders
- Gathering public attitudes on acceptable uses of generative bio Al
- Informing policy with deliberative democracy

#### **Outlook**

- Likely increase in debate and policy activity in this space
- Balancing innovation and security will be a key challenge
- Importance of thoughtful multidisciplinary discourse

### Conclusion

- No easy policy solutions, but inaction also carries risks
- Policy should enable innovation but promote responsible use
- This will require sustained public deliberation and coordination

## References

1. **HyenaDNA: Long-Range Genomic Sequence Modeling at Single Nucleotide Resolution** Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Callum Birch-Sykes, Michael Wornow, Aman Patel, Clayton Rabideau, Stefano Massaroli, Yoshua Bengio, ... Chris Ré *arXiv* (2023) <a href="https://doi.org/gs3v5j">https://doi.org/gs3v5j</a>

DOI: 10.48550/arxiv.2306.15794