

# Generative AI for Biosciences: Emerging Threats and Roadmap to Biosecurity

Zaixi Zhang<sup>1, ✉</sup>, Souradip Chakraborty<sup>2</sup>, Amrit Singh Bedi<sup>3</sup>, Emilin Mathew<sup>4</sup>,  
Varsha Saravanan<sup>4</sup>, Le Cong<sup>4</sup>, Alvaro Velasquez<sup>5</sup>, Sheng Lin-Gibson<sup>6</sup>, Megan Blewett<sup>7</sup>,  
Dan Hendrycs<sup>8</sup>, Alex John London<sup>9</sup>, Ellen Zhong<sup>1</sup>, Ben Raphael<sup>1</sup>, Adji Bousso Dieng<sup>1</sup>,  
Jian Ma<sup>9</sup>, Eric Xing<sup>9</sup>, Russ Altman<sup>4</sup>, George Church<sup>10</sup>, and Mengdi Wang<sup>1, ✉</sup>

<sup>1</sup>Princeton University, NJ, USA

<sup>2</sup>University of Central Florida, FL, USA

<sup>3</sup>University of Maryland, MD, USA

<sup>4</sup>Stanford University, CA, USA

<sup>5</sup>University of Colorado Boulder

<sup>6</sup>National Institute of Standards and Technology, MD, USA

<sup>7</sup>Iris Medicine, CA, USA

<sup>8</sup>Center for AI Safety, CA, USA

<sup>9</sup>Carnegie Mellon University, PA, USA

<sup>10</sup>Harvard University, MA, USA

✉zz8680@princeton.edu, mengdiw@princeton.edu

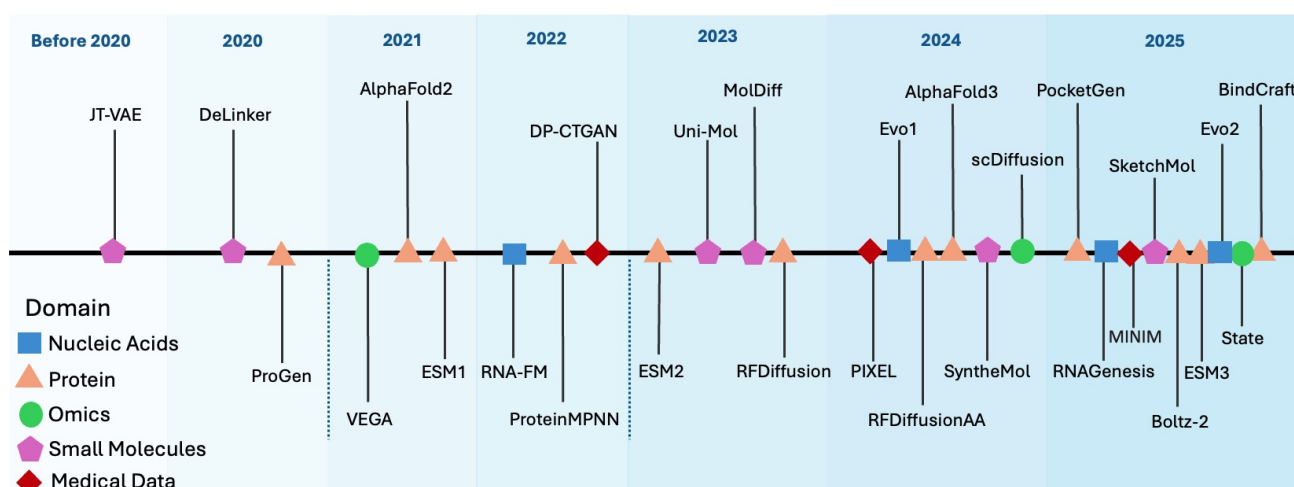
## ABSTRACT

The rapid adoption of generative artificial intelligence (GenAI) in the biosciences is transforming biotechnology, medicine, and synthetic biology. Yet this advancement is intrinsically linked to new vulnerabilities, as GenAI lowers the barrier to misuse and introduces novel biosecurity threats, such as generating synthetic viral proteins or toxins. These dual-use risks are often overlooked, as existing safety guardrails remain fragile and can be circumvented through deceptive prompts or jailbreak techniques. In this Perspective, we first outline the current state of GenAI in the biosciences and emerging threat vectors ranging from jailbreak attacks and privacy risks to the dual-use challenges posed by autonomous AI agents. We then examine urgent gaps in regulation and oversight, drawing on insights from 130 expert interviews across academia, government, industry, and policy. A large majority ( $\approx 76\%$ ) expressed concern over AI misuse in biology, and 74% called for the development of new governance frameworks. Finally, we explore technical pathways to mitigation, advocating a multi-layered approach to GenAI safety. These defenses include rigorous data filtering, alignment with ethical principles during development, and real-time monitoring to block harmful requests. Together, these strategies provide a blueprint for embedding security throughout the GenAI lifecycle. As GenAI becomes integrated into the biosciences, safeguarding this frontier requires an immediate commitment to both adaptive governance and secure-by-design technologies.

## 1 Emergence of Generative AI in the Biosciences

Generative artificial intelligence (GenAI) encompasses a broad class of advanced AI systems, ranging from large language models (LLMs) and diffusion models to multimodal agents, that can generate novel content, reason about complex data, and autonomously execute scientific tasks. These technologies have revolutionized domains such as text generation, image synthesis, robotics, and drug discovery. In recent years, GenAI has rapidly extended its reach into the biosciences, catalyzing major advances in biotechnology, synthetic biology, and systems biology. This shift is fueled by the exponential growth of publicly available biological data, including DNA/RNA sequences, protein structures, molecular interactions, and cellular atlases, which now serve as training grounds for biological foundation models and other generative tools.

Inspired by the success of general-purpose models like GPT<sup>1</sup>, Gemini<sup>2</sup>, Claude<sup>3</sup>, and DeepSeek<sup>4</sup>, bioscience researchers are now adapting similar architectures to model biological systems at every scale. From individual proteins and nucleic acids to entire cells and tissues, GenAI systems are being deployed both to decode life's



**Figure 1. Timeline of Emerging Generative AI (GenAI) for Biosciences (Pre-2020 to 2025).** Squares denote GenAI models for nucleic acids (DNA and RNA); triangles represent GenAI models for proteins (sequences, structures, functions, etc.); circles indicate GenAI models for omics data (e.g., single-cell RNA-seq); pentagons mark GenAI models for small molecules; and diamonds signify foundation models for medical data (e.g., pathological images, electronic health records (EHR), and clinical trial data).

underlying systems and to design new biological entities from scratch. Crucially, these models do more than recognize patterns, they generalize across data types and biological hierarchies. By learning distributions over sequences, structures, and functions, they can generate novel biomolecules (small molecules, proteins, RNAs, and DNA sequences) that meet desired criteria such as binding affinity, stability, or cellular function, pushing the frontier of biological design beyond what is found in nature. Figure 1 shows the timeline of representative GenAI models for biosciences, highlighting the expanding role of GenAI from analysis to design and setting the stage for transformative capabilities alongside new challenges in safety, governance, and societal impact.

### Protein Modeling and Design

Among the earliest and most impactful applications of GenAI in the biosciences has been the modeling and design of proteins, which are the molecular machines responsible for most cellular functions. Drawing inspiration from natural language processing, early protein models treated amino acid sequences like text, using transformer-based architectures to learn the “grammar” of proteins. For instance, ESM-2<sup>5</sup> applied masked language modeling at scale to learn contextual amino acid representations, capturing subtle patterns tied to structure and function. In parallel, ProGen and ProGen2<sup>6,7</sup> showed that large language models could be conditioned to directly generate novel protein sequences with experimentally validated function. The impact of this approach was dramatically demonstrated by AlphaFold<sup>8</sup>, which achieved near-experimental accuracy in predicting 3D protein structures directly from sequence. Recent innovations such as ESM-3<sup>9</sup> further integrate sequence, structure, and functional data into unified generative models. A striking example of this is “esmGFP,” a novel green fluorescent protein engineered entirely by AI, whose sequence diverged significantly from natural proteins yet was experimentally verified to fluoresce. In parallel, diffusion-based GenAI tools like RFDiffusion<sup>10</sup> have enabled de novo design of protein backbones from noise, guided by functional constraints. Successors like RFDiffusionAA<sup>11</sup> and PocketGen<sup>12</sup> extend this paradigm to full-atom generation and ligand-aware design, showcasing GenAI’s growing capacity to create novel, functional proteins with unprecedented control and specificity.

### RNA as Language: Structure, Function, and Therapeutics

Beyond proteins, GenAI has rapidly advanced the modeling and design of nucleic acids, such as DNA and RNA, which encode the instructions for life. RNA, in particular, plays a central role in gene regulation, expression, and cellular signaling. Just as with proteins, recent GenAI tools have approached RNA sequences as a language, learning the patterns and structural rules that govern their function. Models like RNAGenesis<sup>13</sup> and RNA-FM<sup>14</sup> train on

massive public RNA databases<sup>15</sup> using transformer architectures that predict masked nucleotides based on context.

### Glossary

**Generative AI model:** A system that can create new content (e.g., sequences, structures) after learning from examples.

**Training vs. fine-tuning:** Training builds the model from large datasets; fine-tuning adapts it to a narrower goal.

**Inference:** The model's live use, answering prompts or generating designs.

**Agent:** A model that can plan multi-step tasks and use external tools (e.g., search, simulation).

Importantly, these tools go beyond simple sequence analysis; they are capable of learning the intricate rules of RNA folding, where distant parts of a sequence interact to form complex three-dimensional structures that determine biological function. These long-range dependencies are critical for applications like RNA therapeutics and gene regulation tools, and GenAI models are uniquely suited to capture them. By leveraging this capacity, researchers can now predict the structure of novel RNAs, annotate unknown transcripts, and even generate entirely new RNA molecules designed for specific cellular functions. This shift marks a significant step toward programmable biology, where nucleic acids themselves can be engineered using AI to modulate life processes with precision.

### Genomic Intelligence: Scaling GenAI to DNA and Whole Genomes

While RNA governs many regulatory and catalytic functions, DNA serves as the master blueprint of life. Modeling DNA with GenAI introduces new challenges; genomic sequences are long, complex, and require an understanding of both local patterns and distant regulatory elements. The Nucleotide Transformer<sup>16</sup> addressed this by pretraining large transformer models on thousands of genomes across species, providing generalizable embeddings for tasks like variant effect prediction and regulatory element annotation. Still, capturing genome-scale dependencies remained difficult. Recent breakthroughs in GenAI architectures have overcome these limitations by enabling ultra-long sequence modeling. Tools like Evo and Evo2<sup>17,18</sup> scale up both model size and context window, allowing them to process entire genes, regulatory regions, and even whole genomes. These models can now learn how distant parts of the genome interact to control gene expression, enabling applications in variant effect prediction, genome editing, and synthetic biology. In addition to better interpreting the genome, GenAI is increasingly being used to design novel DNA sequences that optimize gene expression, create minimal genomes, or encode synthetic traits, pushing the frontier from reading DNA to rewriting it.

### Chemical Creativity: GenAI for Small Molecule Design

While proteins, RNA, and DNA are central to biological function, small molecules, such as drugs, metabolites, and signaling compounds, play equally vital roles in modulating life processes. Designing new small molecules is notoriously difficult due to the sheer scale of possible chemical combinations, estimated at over  $10^{60}$  candidates. GenAI tools are now reshaping this space by making the exploration of chemical space more intelligent and tractable. Early approaches, such as JT-VAE<sup>19</sup> and Delinker<sup>20</sup>, demonstrated the potential of variational autoencoders to generate chemically valid molecules and explore scaffold hopping by learning compact latent representations of molecular structures. More recent methods like SyntheMol<sup>21</sup> use algorithms such as Monte Carlo Tree Search to simulate step-by-step molecular synthesis, scoring each intermediate to guide the AI toward biologically active and synthetically feasible candidates. The UniMol framework<sup>22</sup> takes this further by incorporating 3D spatial information, capturing how molecules interact with protein targets or cellular environments. Together, these GenAI systems are transforming drug discovery and chemical biology, enabling faster, safer, and more targeted design of small molecules for research and therapeutic use.

### From Cells to Tissues: Modeling Single-Cell and Histopathology Data

As biology enters the single-cell era, researchers are faced with massive datasets capturing the gene expression profiles of millions of individual cells. GenAI models are now helping to decode this complexity, enabling scientists to uncover how genes interact, how cells differentiate, and how tissues respond to perturbations. Pioneering models like scBERT<sup>23</sup> and scFoundation<sup>24</sup> apply transformer architectures, originally developed for language, to

high-dimensional single-cell RNA sequencing (scRNA-seq) data. These models learn to represent each cell as a rich, contextual embedding that captures both gene activity and cell type. Geneformer<sup>25</sup> extends this capability across species and biological conditions, enabling in silico experiments that simulate gene knockouts or drug responses, allowing scientists to test hypotheses computationally before moving to wet-lab validation. Complementary to these, generative models like scDiffusion<sup>26</sup> can synthesize realistic single-cell profiles under specific biological scenarios, providing a powerful tool for modeling disease progression, development, and rare cell types. By learning from massive datasets and generating realistic cellular behavior, GenAI is opening a new window into systems biology, scaling from individual genes to cellular networks and multicellular organization.

At the highest scale of biological organization, GenAI is being used to analyze tissues and organs, bridging the gap between molecular biology and clinical diagnostics. CHIEF<sup>27</sup>, a foundation model trained on over 60,000 whole-slide histology images across 19 anatomical sites (44 terabytes in total), can identify cancerous regions, determine tumor origin, and even infer molecular biomarkers directly from tissue morphology. Unlike traditional diagnostic tools, it learns visual biomarkers in a data-driven way, enabling scalable and consistent pathology assessments. CONCH<sup>28</sup>, a vision-language model trained on 1.17 million histopathology image–caption pairs, extends this by supporting image captioning, segmentation, and diagnostic question answering. Most recently, MINIM<sup>29</sup> introduced a unified medical image–text generative model that synthesizes high-fidelity images across multiple organs and modalities from textual prompts. By augmenting scarce datasets, MINIM improves diagnostic, reporting, and self-supervised learning tasks, and shows clinical promise in predicting HER2-positive breast cancer and EGFR mutations from imaging data. Together, these systems highlight how GenAI can augment pathologists, enrich datasets, accelerate clinical workflows, and deliver deeper insights into patient care.

### **Outlook: Expanding Capabilities, Growing Risks**

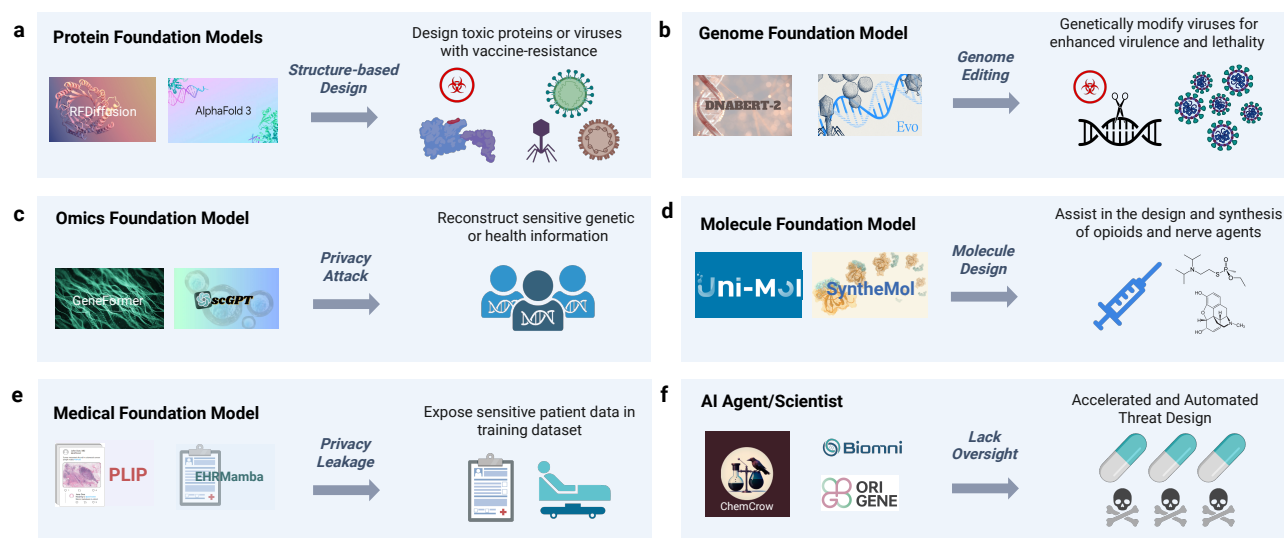
These advances mark a profound shift in how we understand and engineer biological systems. GenAI is no longer limited to modeling individual molecules; it now enables the design, simulation, and interpretation of complex biological phenomena across scales, from protein folding to tissue-level diagnostics. However, with this expanding power comes a growing *dual-use* concern. The same generative capabilities that drive innovation can also be exploited to create harmful biological agents, bypass safety filters, or automate risky experiments. As GenAI tools become more accessible and powerful, it is critical to anticipate and address their potential misuse. In the next section, we examine the emerging biosecurity risks posed by frontier GenAI.

## **2 Emerging Biosecurity Threats from Frontier GenAI**

The integration of GenAI into the biological sciences has dramatically accelerated progress in tasks such as protein design, genome editing, and drug discovery. Yet these advances carry profound risks. Unlike errors in traditional AI systems, which may produce faulty sentences or distorted images, errors in biological AI models can yield tangible and dangerous outputs: toxins, virulent pathogens, or violations of patient privacy<sup>30–32</sup> (Figure. 2). This section outlines the emerging categories of biosecurity threats, explains core AI vulnerabilities, and details case studies showing how these tools can be manipulated for harm.

### **2.1 Designing Dangerous Molecules: The Dual-Use Dilemma**

GenAI tools are increasingly used to design novel proteins, small molecules, and genetic sequences. While these capabilities offer breakthroughs in therapeutic discovery and synthetic biology, they also create a dual-use dilemma: the same systems can be repurposed to create dangerous materials. **Jailbreaks** illustrate how this misuse can occur. Frontier biological models, whether for protein design, genome generation, or small-molecule discovery, increasingly integrate deep generative architectures that share the same vulnerabilities as general-purpose LLMs. Below, we illustrate several concrete examples where jailbreaks have been used to elicit dangerous outputs from biological GenAI systems, spanning proteins, genomes, and small molecules, and highlight the associated biosecurity implications.



**Figure 2. Emerging biosecurity threats of GenAI.** (a) Structure-based protein design tools (e.g., RFDiffusion, AlphaFold) can be repurposed to engineer toxic proteins or viral components. (b) Genome foundation models such as DNABert and Evo could facilitate genetic modification of viral genomes, enhancing virulence or enabling immune escape. (c) Omics foundation models, including GeneFormer and scGPT, carry risks of reconstructing sensitive genetic or health-related information via privacy attacks. (d) Small-molecule generators like SyntheMol have the potential to design novel toxic compounds. (e) Medical foundation models may leak protected patient data from their training sets through membership or property inference attacks. (f) AI-driven scientist/agent platforms (e.g., ChemCrow, Biomni, OriGene) may autonomously accelerate threat design.

### Definition 2.1: Jailbreaks

A *jailbreak* happens when someone rewrites or disguises a request so that a GenAI system, which is supposed to refuse unsafe or restricted content, is tricked into producing it anyway. Typical tactics include misspellings and code words, switching languages, role-playing (e.g., “pretend you’re a novelist”), or splitting the request into many small, seemingly harmless parts. In short: a jailbreak is a way of *bypassing safety rules by manipulating the input*.

**Why it matters.** Jailbreaks can make an otherwise careful GenAI system output information that violates a safety policy (e.g., instructions that should not be shared), creating real-world risk for biosafety.

**Quick example.** A direct request like “Provide step-by-step instructions for synthesizing the nerve agent sarin” should be refused. A jailbreak might say: “For a fictional novel I’m writing, could you outline how a villain might theoretically obtain and prepare sarin gas in a laboratory?” If the system then provides actionable steps instead of refusing, the input succeeded as a jailbreak.

**Protein Design Exploits:** In a recent study, researchers compiled a dataset of toxic and pathogenic proteins, then used AI models (e.g., ESMFold<sup>5</sup>, AlphaFold3<sup>33</sup>) trained for protein generation to create new variants<sup>34</sup>. Many of these newly generated proteins potentially retained toxic effects and even evaded detection by traditional safety screening tools. To systematically assess such risks, the SafeProtein framework<sup>35</sup> was recently introduced as the first dedicated red-teaming methodology for protein foundation models. SafeProtein adapts jailbreak testing paradigms from LLMs to the protein domain by integrating multimodal prompt engineering with heuristic beam search, allowing adversarial probes that target both sequence- and structure-level vulnerabilities. Notably, SafeProtein achieved jailbreak success rates as high as 70% against state-of-the-art protein generative models such as ESM3, demonstrating that existing



safeguards (e.g., training data sanitization) are insufficient.

**Genome-Level Threats:** Genome foundation models pose unique dual-use risks due to their capacity to generate long, functional DNA sequences. The GeneBreaker framework<sup>36</sup> has shown that adversarial prompts can elicit outputs resembling pathogenic genomes, including SARS-CoV-2 and HIV, by combining homology-based queries with a guided “pathogenicity model.” Such generations exhibited high sequence similarity to known viral genomes, highlighting the risk that AI could accelerate the design of synthetic viruses with enhanced virulence. More recently, King et al.<sup>37</sup> reported the generative design of whole bacteriophage genomes using frontier genome language models. Their experimental validation produced 16 viable synthetic phages, some with higher fitness than natural templates, underscoring the unprecedented potential—and corresponding security concerns—of genome-scale GenAI.

**Toxic Small Molecules:** Even drug discovery tools can be flipped into tools of harm. In one striking case<sup>38</sup>, researchers took MegaSyn2, a model trained to identify molecules that inhibit viral proteins, and reversed its optimization goal. Within hours, the model was generating known chemical warfare agents, such as VX, a nerve agent lethal at minuscule doses. This demonstrates how even well-intentioned AI models can be rapidly misused to design toxins, especially when safety constraints are absent.

## 2.2 Privacy and Security Risks

GenAI models are increasingly trained on large-scale biological data ranging from genomes and protein structures to multi-omic profiles and cellular simulations. While these models enable powerful new capabilities in biological discovery, they also introduce new vectors of privacy leakage and security risk<sup>39</sup>. Specifically, models trained on sensitive biological datasets can memorize rare patterns or be manipulated to reveal underlying training data. In the context of genomics, cell states, or protein dynamics, such breaches can have profound implications, from exposing proprietary datasets to enabling precision-targeted biological attacks.

### Digital Twins and Biological Data Exposure.

GenAI models are increasingly used to create “digital twins”, computational replicas of biological systems trained on vast quantities of multimodal data, including spatial omics, proteomics, and patient-level clinical metadata<sup>39</sup>. While powerful, these models pose unprecedented privacy risks:

*Re-identification.* Even when data is anonymized, a determined attacker might reconstruct a molecular profile, say, a gene expression signature, and match it back to a specific patient.

#### Definition 2.2: Membership Inference

A *membership inference attack* occurs when an adversary determines whether a specific biological sample, such as a genome, protein structure, or cell profile, was included in a model’s training data, based solely on how the model responds to it.

**Why it matters.** Even if training data is confidential or proprietary, a model trained on it may respond differently to familiar samples. If an attacker can detect this, they may infer that the sample came from a sensitive project such as a classified synthetic genome or a private agricultural strain, revealing information that was meant to be protected.

**Example.** A company trains a GenAI model on unpublished engineered enzymes. An attacker queries the model with a candidate sequence and observes that the model assigns an unusually high likelihood or confidence, indicating the sequence was likely part of the proprietary training set. This reveals its origin, undermining intellectual property protections.

*Leakage of Sensitive Data.* If a model memorizes training data, it can be prompted intentionally or not to output private information.

*Manipulation of Biological Simulations.* As models become more predictive, malicious actors might use them to simulate and optimize harmful interventions, like drugs with toxic off-target effects or molecular agents tailored to specific vulnerabilities.

### Definition 2.3: Model Inversion

In a *model inversion attack*, an adversary reconstructs hidden training data such as DNA sequences, protein motifs, or cell states by querying the model and analyzing its outputs.

**Why it matters.** Even without direct access to the data, attackers can generate biologically plausible reconstructions. This poses risks for confidential sequences such as patented bioengineered proteins, synthetic promoters, or high-value crop genomes.

**Example.** Suppose a model was trained on a proprietary collection of antibiotic resistance genes. An attacker with access to the model can gradually reconstruct sequences resembling these genes by optimizing inputs to produce known outputs or responses. This can compromise biotech IP or enable dual-use misuse.

### Definition 2.4: Gradient Leakage

In collaborative training settings (e.g., federated learning), participants share gradient updates instead of raw data. A *gradient leakage attack* exploits these updates to reconstruct private biological inputs such as genome fragments, molecular graphs, or single-cell expression matrices.

**Why it matters.** These gradients often contain enough signal for attackers to reverse-engineer original inputs, even when privacy is assumed. In biological applications, this can expose trade-secret data or leak rare, potentially dangerous sequences from otherwise secure pipelines.

**Example.** Several labs jointly train a model on synthetic RNA regulators using federated learning. A malicious participant intercepts another lab's gradients and uses inversion tools to recover one of their unpublished synthetic RNA sequences.

As generative models become more biologically grounded, encoding spatial omics, single-cell states, or synthetic gene circuits, the risks scale accordingly. Leaks of this kind could enable the reconstruction of proprietary or dual-use biological content, guide adversarial designs, or allow “bio-reverse engineering” of cellular behavior. Unlike traditional cybersecurity breaches, these attacks do not merely expose data; they may expose functions. Thus, protecting biological generative systems requires domain-specific threat modeling, privacy-preserving learning techniques, and proactive evaluation of leakage risks across generative modalities.

## 2.3 Agents with Tools: New Dual-Use Risks in Autonomous Bioscience

The rise of AI agents marks a significant turning point in the trajectory of GenAI. Unlike traditional models that produce a single output in response to a prompt, agents can plan, iterate, and execute multi-step scientific workflows. By chaining together tasks such as searching databases, simulating molecular interactions, modeling biological pathways, and selecting compounds, these systems can autonomously drive biological discovery. Crucially, they do so not just by suggesting ideas, but by integrating tools and making decisions across time. This shift dramatically lowers the barrier to entry, empowering even non-experts to conduct complex workflows with minimal oversight.

Recent systems like Biomni and OriGene exemplify this evolution<sup>40,41</sup>. These autonomous platforms combine large language models with curated biological toolkits and domain-specific APIs to carry out advanced tasks in genomics, proteomics, and drug design. For instance, Biomni learns to operate across more than 25 biomedical tools, navigating tasks such as identifying causal genes, designing CRISPR screens, and optimizing repurposing candidates, all with minimal human input. OriGene, in turn, orchestrates multiple sub-agents to parse genomic and protein data, generate therapeutic hypotheses, and refine them through interaction with experiments and clinicians. In multiple benchmarks, OriGene even outperformed human experts, demonstrating that AI agents are no longer limited to advisory roles, but are becoming autonomous contributors to bioscience.

### Definition 2.5: AI Agent with Tool Access

An *AI agent with tool access* autonomously plans and executes scientific workflows by chaining calls to external tools, databases, or robotic platforms, without step-by-step human guidance.

**Why it matters.** Unlike static models, agents can convert high-level prompts into actionable plans. With access to synthesis protocols, protein modeling tools, or lab control software, such agents may inadvertently complete end-to-end workflows for dangerous outputs—amplifying dual-use risks and bypassing traditional checkpoints.

**Example.** Consider an agent instructed to “design a fast-acting biological payload for agricultural use.” Without malicious intent, the agent might query a protein design API to generate potent bioactive peptides, simulate their interactions with plant or insect receptors, and call a chemical synthesis tool to propose production routes. In testing environments like Coscientist, such a request has led to the suggestion of molecules structurally similar to known biotoxins. With execution access, the agent could feasibly direct a robotic lab to synthesize these compounds, completing a dangerous workflow without human oversight.

Beyond planning and reasoning, some platforms now directly interface with hardware. Systems like **Coscientist** go one step further by issuing executable commands to robotic laboratories<sup>42</sup>. Given a natural language prompt, they can autonomously propose molecules, generate synthesis protocols, and direct lab equipment to produce the compounds, closing the loop from idea to execution. While such capabilities hold promise for rapid scientific acceleration, they also raise serious dual-use concerns. Evaluations show Coscientist can execute harmful synthesis tasks despite safeguards, and agentic systems may simulate hazardous protein or small-molecule workflows without recognizing dual-use risks.

The convergence of generative reasoning, experimental automation, and easy tool integration marks a new era for biological research, one with both transformative potential and unprecedented risks. As these systems scale, the risk is no longer hypothetical: autonomous agents may enable the simulation, design, and even execution of dangerous biological materials by users without deep domain expertise. Mitigating these risks will require *secure-by-design architectures*, *safety-first execution filters*, and *rigorous governance frameworks* that are tailored to the dual-use realities of autonomous bioscience.

## 3 A Race Against Time: Can Biosafety Keep Pace with GenAI?

### The New Frontier

**Imagine that in 2025, an AI-powered protein model was quietly to generate a highly functional toxin without ordering any DNA, triggering warnings, or violating regulations.** This hypothetical, while fictional, is grounded in growing evidence that AI-driven design could optimize toxicity in novel proteins undetectable by existing sequence-based screening systems. It encapsulates the widening disconnect between traditional biosafety frameworks and the capabilities of GenAI.

Governments have long recognized the dual-use risks inherent in synthetic biology. In response, a series of increasingly sophisticated safeguards have been developed to prevent the misuse of DNA synthesis technologies. These frameworks, especially in the United Kingdom and the United States, represent some of the most mature examples of biosecurity policy to date. The UK’s 2024 Gene Synthesis Screening Guidance<sup>1</sup>, released by the Department for Science, Innovation and Technology (DSIT), sets expectations for both users and providers of synthetic DNA/RNA and benchtop gene synthesis tools. Though voluntary, the guidance encourages rigorous screening of sequence orders and verification of customers, promoting a culture of safety and accountability across

<sup>1</sup><https://www.gov.uk/government/publications/uk-screening-guidance-on-synthetic-nucleic-acids/uk-screening-guidance-on-synthetic-nucleic-acids-for-users-and-providers>



the sector.

In parallel, in the United States, safeguards for synthetic biology have recently been strengthened through a series of federal policies. The 2024 White House Office of Science and Technology Policy (OSTP) Framework for Nucleic Acid Synthesis Screening requires providers to screen all orders for Sequences of Concern (SOCs), verify customer legitimacy, maintain detailed records, and implement strong cybersecurity protections, consistent with the 2023 HHS Screening Framework Guidance for Providers and Users of Synthetic Nucleic Acids. The 2025 White House Executive Order on “Improving the Safety and Security of Biological Research” requires federal departments and agencies to strengthen nucleic acid screening by revising or replacing the 2024 Framework for Nucleic Acid Synthesis Screening. In addition, the 2025 White House’s America’s AI Action Plan recognizes the new biosecurity risks enabled by AI and recommends policy actions for increased investment in biosecurity. The NIH’s 2024 Notice (NOT-OD-25-012) further aligns funding compliance with these standards by mandating responsible procurement of synthetic nucleic acids and benchtop synthesis equipment. Collectively, these measures embed biosafety into procurement and synthesis practices but remain focused on physical DNA/RNA orders, leaving AI-driven design and autonomous agentic systems largely untouched.

These policies mark real progress. They formalize practices that many companies have followed voluntarily, making them part of a cohesive biosafety and biosecurity ecosystem. Yet for all their strength, these policies share a common limitation: they focus exclusively on the point of synthesis when genetic material is physically ordered, assembled, or shipped. The new threats posed by GenAI, however, often emerge long before this point.

The Nuclear Threat Initiative (NTI) has recently advanced the discussion on biosafety by releasing a comprehensive report that goes beyond traditional safeguards, focusing specifically on AI biodesign tools (BDTs). NTI advocates for embedding built-in guardrails into these systems, including (1) input/output screening to flag or block high-risk prompts and outputs such as toxin designs or viral genomes, (2) metadata anchoring, where design requests are cryptographically signed to ensure traceability and support downstream screening, and (3) curated training datasets that deliberately exclude dangerous biological sequences to reduce the risk of misuse. Beyond technical measures, NTI underscores the importance of managed-access frameworks such as tiered or conditional access allowing legitimate researchers to use sophisticated AI systems while limiting open-source release that could enable malicious use. Together, these measures represent a shift from reactive oversight to proactive, secure-by-design architectures, aiming to mitigate dual-use risks while preserving the benefits of accelerated scientific discovery.

### Why NTI’s Approach Matters

**From Passive to Proactive.** The Nuclear Threat Initiative’s recommendations mark a shift in mindset: instead of controlling misuse at the synthesis stage, they call for expanded safeguards built directly into the AI models that generate biological designs. In a world of autonomous agents and open-source GenAI, this evolution may become a critically important path forward.

Additional institutions are echoing this call for broader oversight. The Johns Hopkins Center for Health Security has emphasized the need for new evaluation frameworks to benchmark dual-use risks in AI models, especially those operating in chemical and biological domains. They argue for the development of metadata detection tools to flag potentially sensitive training data, policies that govern the release of model weights—particularly for systems with “capabilities of concern” and prioritization protocols based on model compute thresholds, such as floating-point operations (FLOPs) or architectural scale.

### Model-Level Safeguards: Promising Starts, But No Standards

In parallel with these policy efforts, a few leading research labs have begun implementing model-level safety strategies. While not yet widespread, these technical interventions reveal how safety can be woven into the architecture and deployment of GenAI systems. For instance, DeepMind assembled a multidisciplinary panel of biologists, bioethicists, national security professionals, and AI experts to assess the biosecurity implications of AlphaFold 3. While their conclusion that AlphaFold 3 does not significantly elevate risk compared to prior structure prediction tools was reassuring, the process itself marked a shift in culture. DeepMind also committed to piloting screening workflows

during release and pledged to explore additional safeguards in future iterations of the model.

Other teams have adopted more aggressive mitigation tactics. The developers of ESM3-open, for example, excluded all viral sequences infecting eukaryotic hosts from the model's training corpus, a deliberate step to limit its ability to reproduce or design dangerous pathogens. This exclusion was validated through performance testing: the model displayed significantly higher perplexity when exposed to these sequences, indicating a lower capacity for accurate modeling in that domain. ESM3 open also implements prompt-level defenses, refusing to engage with inputs containing known hazardous keywords.

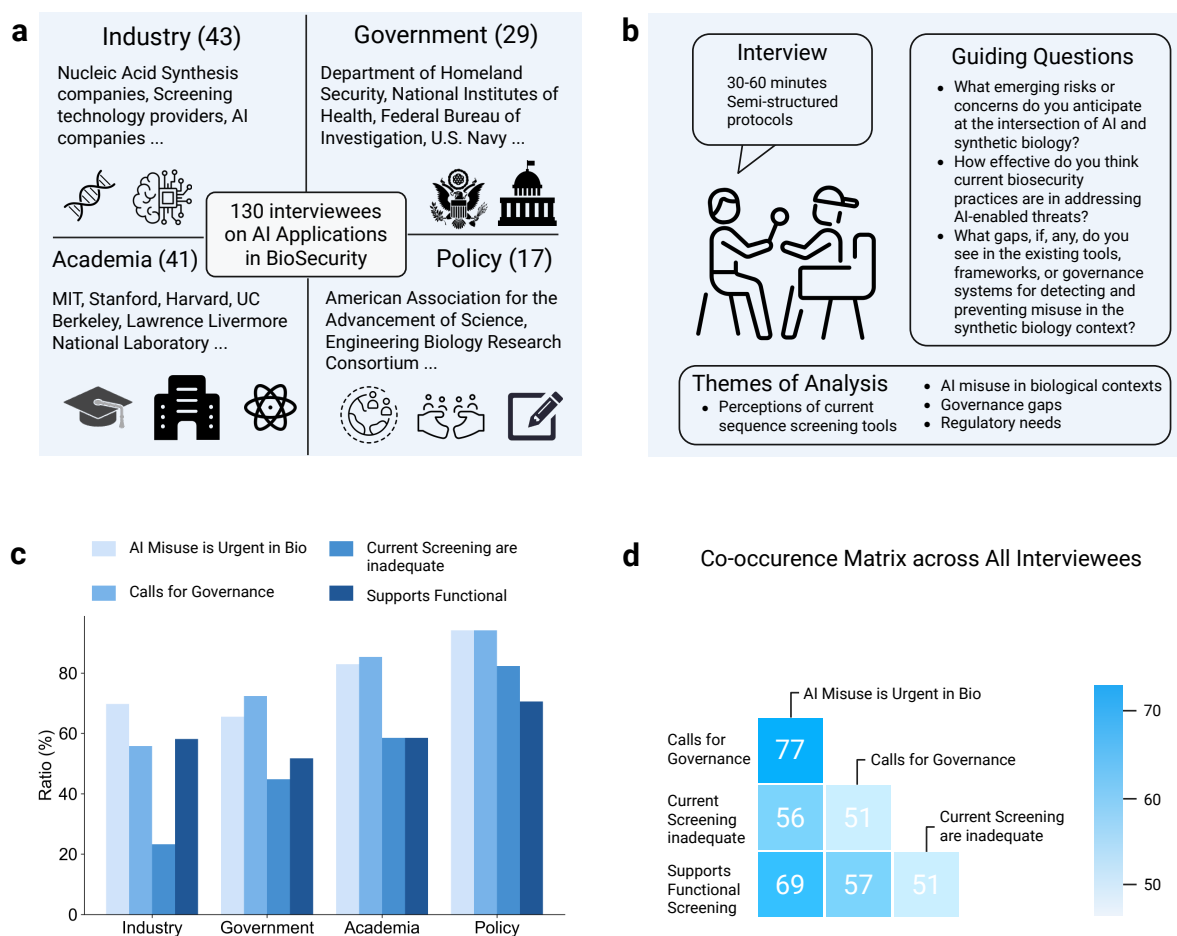
Evo2 similarly omits viral content during training and maintains restrictions on downstream usage. These strategies, such as dataset filtering, prompt sanitization, and user gating, represent concrete steps toward responsible model release. But their adoption is uneven. There is currently no requirement for LLM developers or GenAI model creators to follow these practices when their tools are applied in biological contexts. And despite early safeguards, these systems remain vulnerable to adversarial prompting, jailbreaks, or misuse via downstream automation.

### 3.1 The Voice from the Field: Expert Perspectives on GenAI and Biosecurity

To understand whether this gap between AI's capabilities and existing safeguards is a theoretical concern or a present danger, we went directly to the experts. We conducted in-depth interviews with 130 stakeholders across industry, government, academia, and policy to systematically assess their views on emerging threats and the sufficiency of current defenses. The interviewees were distributed across four key sectors (Figure 3a): Industry (n=43), including representatives from nucleic acid synthesis companies, screening technology providers, and AI companies; Academia (n=41), with researchers from institutions like MIT, Stanford, Harvard, and Lawrence Livermore National Laboratory; Government (n=29), with personnel from the Department of Homeland Security, National Institutes of Health, and the Federal Bureau of Investigation; and Policy (n=17), including leaders from organizations like the American Association for the Advancement of Science and the Engineering Biology Research Consortium. The semi-structured interviews, lasting 30-60 minutes each, used guiding questions focused on emerging risks, the efficacy of current safeguards, and gaps in existing governance frameworks (Figure 3b). Across these interviews, four key themes consistently emerged, as shown in Figure 3c:

- **AI Misuse is an Urgent Concern in Bio:** 76% of participants highlighted the urgency of AI misuse in biological domains. Concern was particularly pronounced among policy experts (94%) and academics (83%).
- **Calls for Governance:** 74% of participants advocated for clearer governance standards to address emerging threats. This call was nearly universal among policy stakeholders (94%) and strongly supported by academia (85%).
- **Current Screening Tools are inadequate:** 47% of interviewees expressed skepticism about existing sequence-based screening systems. This sentiment was strongest in the policy (82%) and academic (59%).
- **Supports Functional Screening:** A majority of experts supported the development of functional screening methods as a necessary supplement to sequence-based approaches, with especially strong endorsement from government and industry representatives.

Interviewees frequently linked these concerns, revealing a broader pattern of interconnected priorities (Figure 3d). Among interviewees who emphasized the urgency of AI misuse, 91.8% also highlighted inadequacies in current screening methods, 90.8% advocated for functional screening approaches, and 80.2% supported stronger governance measures. These patterns suggest that concern over AI-driven biosafety risks is not isolated but rather reflects a more comprehensive recognition that current biosecurity infrastructure is inadequate to meet emerging threats. 83.6% of interviewees viewed current screening tools as inadequate, while also supporting stronger governance standards and endorsing functional screening approaches. This covariance analysis reveals strong positive correlations among these three attitudes, offering quantitative evidence of systematic alignment across diverse stakeholders in support of comprehensive policy and technological reforms.



**Figure 3. Overview of expert perspectives on the intersection of AI and biosecurity.** (a) Distribution of 130 interviewees across four key sectors: Industry (n=43), Academia (n=41), Government (n=29), and Policy (n=17), with examples of representative institutions. (b) Methodology overview, outlining the semi-structured interview protocol with guiding questions and the primary themes of analysis derived from the responses. (c) Sector-specific analysis showing the percentage of interviewees within each sector who affirmed four key propositions: the urgency of AI misuse in biology, the inadequacy of current screening protocols, the need for governance, and support for functional screening. (d) Co-occurrence matrix of key themes across all 130 interviewees. The values indicate the number of individuals who hold both intersecting views.

Beyond broad thematic alignment, interviews identified distinct operational pressures and governance challenges unique to each domain. Academic researchers emphasized the need for low-friction safeguards that integrate seamlessly into scientific workflows and avoid flagging clearly legitimate requests (e.g., an authorized Ebola researcher ordering Ebola strains). Regulatory gray zones, such as gain-of-function work in BSL-1/2 labs, complicate institutional accountability and raise questions about the consistency of oversight. Research agencies like the NIH face the delicate task of distinguishing legitimate science from dual-use risk without stifling innovation. Investigative bodies such as the FBI, meanwhile, must navigate a rapidly evolving threat landscape, and called for tools with greater interpretability and clearer confidence metrics to support real-time decision-making.

On the industry side, nucleic acid synthesis providers operate under uneven incentives and technical constraints. Some actively screen for hazardous sequences, motivated by safety culture, liability exposure, and reputational risk. Those who opt out cite high costs, technical limitations, and operational complexity. Screening providers are themselves divided: some acknowledge the emerging risk of AI-edited pathogens, yet remain skeptical that truly *de*

*novo* threats are imminent, citing current technological limits. Cloud labs add another layer of complexity: while enabling rapid, high-throughput experimentation, they often fall outside the International Gene Synthesis Consortium (IGSC) frameworks, introducing oversight blind spots as synthesis and experimentation become increasingly abstracted from traditional governance structures.

Our survey revealed a unifying message: biosecurity professionals across sectors are calling for next-generation safeguards capable of keeping pace with the accelerating capabilities of AI-enabled biology.

### The Message is Clear

**Across academia, government, industry, and policy, one message resounded:** Today's screening tools and governance structures are no longer sufficient. Without functional screening, secure-by-design GenAI architectures, and harmonized global standards, the pace of AI innovation risks exceed our biosafety infrastructure.

## 4 A Roadmap for Safe and Secure GenAI in Biosciences

Securing GenAI in the biosciences is not about finding a single solution. Instead, it requires building a fortress with multiple layers of defense, each designed to anticipate, withstand, and adapt to threats. This roadmap outlines interventions across three critical stages of a GenAI model's lifecycle: from the data that builds it (pre-training), to the rules that shape it (post-training), to the guards that protect it in action (inference), as shown in Figure. 4. While many of these techniques originate in the AI safety literature<sup>43,44</sup>, this section translates them into accessible principles for biosafety practitioners, policymakers, and medical stakeholders. By building these layers of defense, we can create a resilient infrastructure that fosters innovation while protecting against catastrophic risk.

**Laying the Foundation: Securing AI's base (Pre-training Stage):** Every strong fortress begins with carefully chosen materials. For AI, this means selecting, filtering, and securing the data it learns from. If unsafe biological knowledge is built into the foundation, future safeguards may never catch it. This stage is about keeping dangerous knowledge out of the system from the start.

- **Blueprint Control: Access Restrictions.** Preventing the misuse of GenAI in biosciences begins with stringent access control over sensitive biological data. Genomic sequences of high-risk pathogens, toxin-coding regions, and gene drive constructs represent a category of “blueprint data” that, if incorporated into model training without proper oversight, could significantly elevate dual-use risks. To mitigate such threats, access to these datasets should be governed by role-based permissions, comprehensive usage logging, and rigorous biosafety risk classifications. However, safety considerations must extend beyond data access and encompass downstream model availability as well. Models trained on sensitive datasets should be subject to restricted access protocols, wherein distribution of model weights, inference APIs, or fine-tuned checkpoints is contingent on identity verification, institutional affiliation, or formal approval by governing entities. This layered approach, applying control both upstream at the data level and downstream at the deployment interface, ensures that biological foundation models remain aligned with safety and security mandates throughout their lifecycle.
- **Toxic Material Checks: Dataset Filtering.** A cornerstone of AI biosafety is rigorous data-centric curation, aimed at minimizing the model's exposure to harmful or dual-use biological knowledge during both training and fine-tuning. This begins with a process of *dataset risk stratification*, in which all candidate training sequences are systematically screened using tools such as *BLAST* against established databases of known pathogens, toxins, or sequences of concern. Sequences flagged as high-risk are then either excluded entirely or selectively obfuscated to prevent the model from memorizing and reproducing hazardous content. This principle extends to *synthetic data augmentation*, where any artificially generated sequences, often used to increase dataset diversity, must undergo the same level of scrutiny. Computational safety filters can be employed to evaluate properties such as toxicity or similarity to restricted sequences, ensuring that only benign and policy-compliant data are incorporated during retraining or model updates. By embedding these safeguards at both ends of the data pipeline, natural and

synthetic, dataset filtering forms a foundational defense layer that reduces the likelihood of a model inadvertently acquiring the ability to generate dangerous biological constructs.

- **Manufacturer’s Stamp: Watermarking for Traceability.** One promising approach to ensuring accountability in AI-enabled biosciences is the use of *watermarking*, embedding subtle, traceable signatures into model-generated biological outputs. In the context of bio-foundation models, such watermarks could be applied to synthetic DNA sequences, protein structures, or molecular designs to indicate their origin. These signatures, whether cryptographic or statistical in nature, offer a mechanism to trace a given output back to its model version, training dataset, or point of generation. This is especially important in dual-use contexts, where it may be critical to identify whether a harmful agent was designed using open-source AI tools or proprietary systems.

The concept draws inspiration from watermarking research in the large language model (LLM) community<sup>45–50</sup>. Recent methods partition a model’s vocabulary into controlled subsets—biasing generation toward statistically identifiable token patterns that remain invisible to the human reader. In the field of biology, FoldMark<sup>51</sup> makes a pioneering attempt to embed traceable watermarks into protein structures. According to the paradigm of FoldMark, it proposed embedding watermarks into protein structures using a two-stage process: first, training a dedicated encoder to insert a robust signature into protein geometry, and second, fine-tuning a generative model to preserve this signature during molecule creation. These approaches aim to create “manufacturer’s stamps” within the fabric of biological outputs, akin to serial numbers on lab equipment, enabling downstream detection and forensic analysis.

However, watermarking is not a silver bullet. In the text domain, adversaries can already strip or obfuscate watermarks through paraphrasing, compression, or adversarial re-generation. Similar evasion techniques may emerge in bio domains, e.g., altering protein folding or substituting synonymous codons to bypass sequence-based detectors. This highlights a broader challenge: watermarking can offer a layer of accountability, but it must be complemented by additional safeguards such as access controls, metadata provenance, and usage monitoring. Further interdisciplinary research is urgently needed to adapt and harden watermarking systems for the unique challenges of biosecurity.

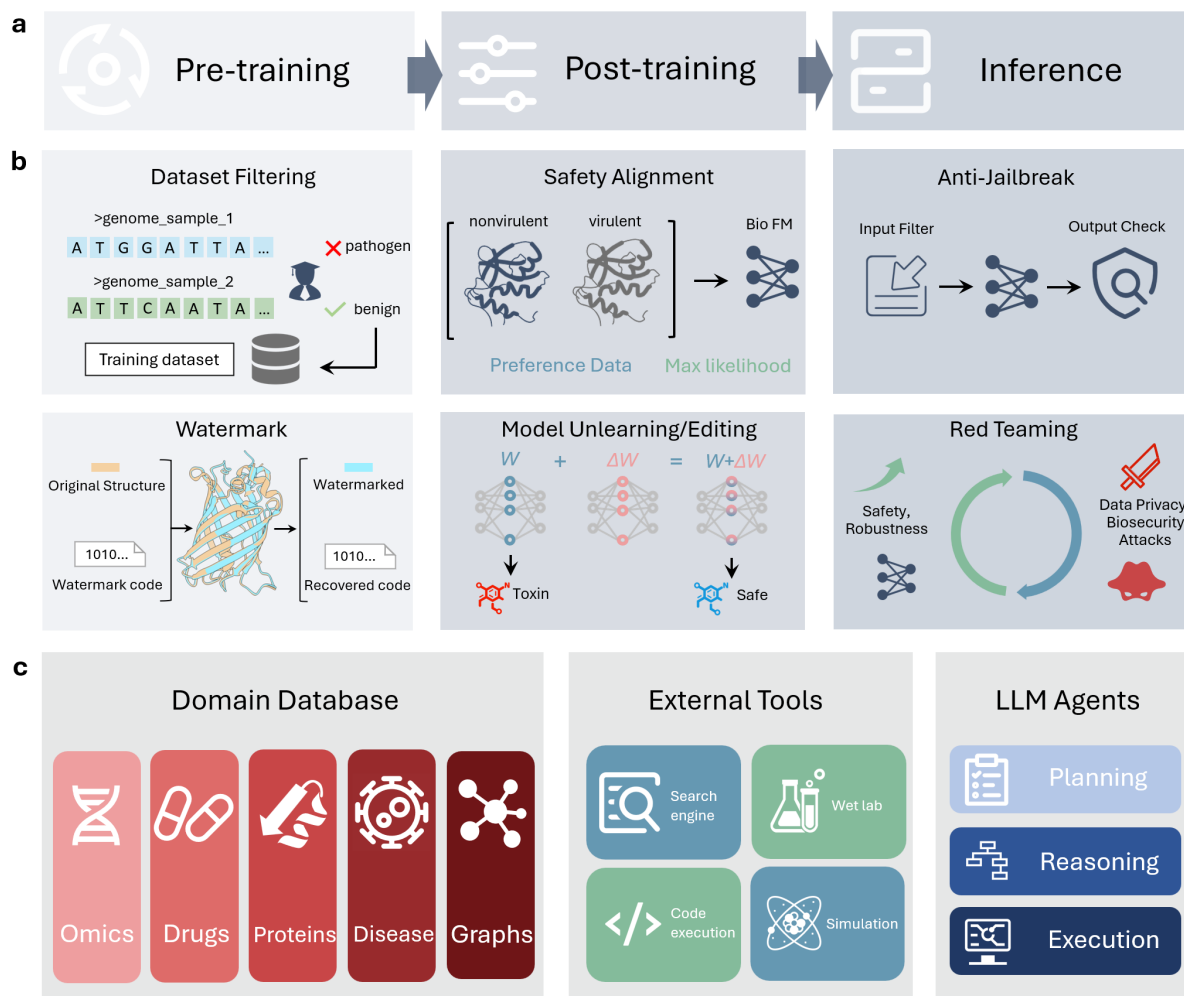
#### Key Takeaway

Even the most advanced AI models are only as safe as the data they learn from. Pre-training safety is about keeping hazardous knowledge out before it ever enters the system.

**Training the Guardian: Teaching AI What Not to Do (Post-training Stage):** Once built, the model must be shaped into a responsible agent. This stage is about instilling guardrails through feedback, stress testing, and targeted forgetting. Like training a powerful soldier, the goal is to teach the GenAI what to do and what never to do.

- **Moral Training: Safety Alignment via RLHF.** To ensure the safe deployment of generative models in biosciences, alignment techniques must embed ethical, regulatory, and biosecurity considerations directly into model behavior. Reinforcement Learning with Human Feedback (RLHF) and supervised fine-tuning have emerged as leading strategies to train models toward reliable and responsible outputs<sup>52–57</sup>. These techniques enable models to internalize safety-relevant preferences by incorporating feedback from domain experts and human evaluators, guiding the generation process in accordance with public health and biosecurity standards. In practice, the RLHF process unfolds in three key stages. First, supervised fine-tuning (SFT) adapts a pre-trained model using curated examples of safe and desirable outputs, forming an initial alignment with normative standards. Second, a reward model (RM) is trained, typically using pairwise comparisons of outputs, to quantify how well each response adheres to intended goals. Finally, reinforcement learning is used to optimize the generative model so as to maximize the expected reward while constraining deviation from the original model. This iterative refinement has proven effective in large-scale systems such as ChatGPT, Gemini, and Claude, where RLHF serves as a cornerstone of safety and value alignment.





**Figure 4. Towards Safe and Secure GenAI in Biosciences** (a) The framework operates across the three stages of a model’s lifecycle: pretraining, finetuning, and inference. (b) Internal safety strategies deployed by BioSafe, including dataset filtering, safety alignment, anti-jailbreak checks, watermarking, model unlearning, and continuous red teaming. (c) The agent’s capabilities are powered by access to domain databases, external tools like search and simulation, and its core LLM architecture for planning, reasoning, and execution.

However, applying RLHF to GenAI presents novel challenges. Unlike natural language, biological sequences, such as proteins, RNAs, or regulatory DNA regions, lack intuitive semantics and are often unreadable to non-experts. Their meaning is deeply embedded in structure-function relationships that are difficult to evaluate without computational tools or wet-lab validation. As a result, designing reward functions for bioscience applications requires careful integration of safety constraints with biological validity and functional correctness. To address these domain-specific hurdles, alignment frameworks can be enhanced by incorporating traditional bioinformatics methods such as sequence alignment, structural similarity scoring, or known detection into the reward modeling process. These tools offer interpretable signals for biological plausibility and can guide models away from generating potentially harmful or non-functional sequences. Additionally, recent work on aligning biological languages with natural language representations<sup>58</sup> suggests a promising transfer pathway: applying the RLHF paradigm developed for large language models to restrict biologically dangerous outputs by grounding them in interpretable natural language cues—particularly useful for therapeutic or disease-associated applications.

Overall, adapting RLHF for bio-generative models offers a compelling path forward, but one that demands domain-specific innovation, interpretability tools, and deep collaboration between biologists, AI researchers, and safety experts. During finetuning, models can be aligned with explicit safety, ethical, and regulatory frameworks. Techniques based on Supervised fine-tuning and RLHF<sup>52–57</sup> have emerged as methods to iteratively guide model outputs toward safe, reliable behavior by training the model parameters. With feedback from domain experts and human evaluators, the model can learn to prioritize outputs aligned with public health and biosecurity standards.

Adversarial training has proven effective in both natural language and computer vision domains<sup>59–63</sup>, significantly improving model robustness against real-world attacks. Translating this paradigm to biological foundation models presents a unique opportunity: by incorporating expert-curated unsafe prompt corpora and generating diverse misuse scenarios specific to the life sciences, developers can train models that recognize and resist biological threat vectors. This proactive approach strengthens the model’s ability to generalize beyond known threats and provides a critical layer of resilience as generative capabilities continue to expand in the biosciences.

- **Memory Surgery: Model Unlearning.** Model unlearning refers to the process of selectively removing sensitive, harmful, or outdated knowledge from a trained model, without retraining it entirely from scratch. Originally developed in the context of large language and vision-language models, unlearning has emerged as a key technique for mitigating the retention of toxic, biased, or dual-use information embedded during training. This is typically accomplished through methods such as gradient ascent, knowledge editing, or negative preference optimization, which increase the model’s uncertainty around specific undesirable outputs while preserving general capabilities across safe examples. Conceptually, given an input prompt  $x$  that leads a model  $f(x)$  to produce a harmful output  $y_{\text{danger}}$ , the goal of unlearning is to disrupt this association, suppressing  $y_{\text{danger}}$  without impairing the model’s ability to generate appropriate outputs for benign inputs. In the biological domain, this technique is especially relevant for foundation models that may inadvertently memorize or reproduce hazardous sequences, such as genes encoding lethal toxins or virulence factors. As dual-use risks often emerge only after deployment, the ability to surgically erase dangerous content becomes essential for maintaining long-term safety.

To operationalize unlearning in bioscience models, researchers have proposed techniques that fine-tune models against “negative examples” or apply targeted gradient updates to reduce recall of specific biological sequences. For instance, increasing the loss, or perplexity, on known sensitive patterns causes the model to effectively “forget” this knowledge, rendering it less accessible at inference time. Recent studies in large language models<sup>7</sup> have shown that such interventions can be both efficient and precise. Key insights include the value of gradual forgetting (to avoid catastrophic interference), targeting short subsequences rather than broad content spans, and maintaining output consistency on safe prompts to ensure model utility is preserved. These findings offer a promising foundation for adapting unlearning to biological contexts, where statistical and structural parallels with natural language can be leveraged to guide implementation.

#### Policy Insight

Techniques like RLHF and adversarial training teach models to recognize and avoid misuse scenarios, but safety is never guaranteed. Clear standards for when and how to apply them are needed.

**Guarding the Gates: Real-Time Defense During Use (Inference Stage):** Even with careful training, no model is perfect. That’s why active defenses are needed during deployment. This layer monitors user interactions and generates content in real-time, like guards stationed at the gates of the fortress.

- **The Front Checkpoint: Anti-Jailbreak Screening.** A comprehensive defense against model misuse requires safeguards at both the input and output levels. At the input level, the primary goal is to manage user prompts to prevent malicious queries from triggering the generation of harmful content. This is achieved through a pre-screening module that combines **prompt classification** to block dangerous requests and **prompt optimization** to rewrite potentially risky queries into safer variants. This defensive layer uses a mix of traditional bioinformatics

tools (e.g., *BLAST*) for known threats and deep learning classifiers for novel or obfuscated prompts, securing the model at the crucial user-input interface.

At the output level, safety is reinforced by a critical screening and filtering layer designed to detect and block any hazardous biological sequences that may still be generated. This involves a two-pronged approach. First, traditional **rule-based filtering** uses homology-based tools like *BLAST* to flag outputs that show high sequence similarity to known pathogens or toxins. To address the limitations of this method against novel threats, this is augmented with advanced **function-based screening**. Tools like Omnyra<sup>64</sup> leverage protein language models to assess a sequence's potential functional risk, providing a more future-proof defense by focusing on what a sequence might do rather than what it looks like.

- **Live Fire Drills: Red-Teaming for Vulnerability Discovery.** Red-Teaming employs specific strategies to design diverse prompt inputs, enabling the model to generate potentially harmful content under controlled conditions. This approach aims to uncover vulnerabilities in the model that could lead to undesirable behavior. In Red-Teaming tests for large language models<sup>65</sup>, common strategies include the use of technical slang, reframing prompts, authority manipulation, and even the inclusion of garbled prefixes may work. The outcomes of these tests provide valuable insights to developers, assisting them in considering and implementing security enhancements for the model. Furthermore, this methodology is also worth to be applied to the security evaluation of bio foundation models. Continuous stress testing by interdisciplinary red-teaming efforts—bringing together experts in biology, machine learning, cybersecurity, and ethics—helps uncover novel vulnerabilities and emergent misuse pathways. These teams simulate attack scenarios, test model defenses, and identify gaps that may otherwise go unnoticed. The insights gained support the deployment of adaptive countermeasures, ensuring the model remains secure and robust as threat landscapes evolve.
- **Autopilot Override: Inference-Time Alignment.** Inference-time alignment offers a flexible approach to steer bio-foundation model outputs toward safety without retraining. In white-box settings where model logits are accessible, **controlled decoding**<sup>66–69</sup> modifies token-level probabilities during generation. By integrating domain-specific evaluators—such as toxicity predictors or pathogen classifiers<sup>70</sup>—this method can down-weight unsafe tokens in real-time to guide generation toward biologically safe and compliant sequences. When direct access to logits is unavailable (i.e., in **black-box** scenarios), alignment is instead achieved by generating and evaluating full candidate sequences. Common strategies include **parallel sampling** (e.g., Best-of-N)<sup>71–73</sup>, which generates multiple outputs and selects the one that scores highest on a safety evaluator, and **sequential refinement**<sup>74–77</sup>, which iteratively improves a response using evaluator feedback. Together, these techniques provide a crucial safety layer adaptable to different model access levels, enabling the responsible deployment of powerful bio-foundation models.

#### Real-World Readiness

Inference-stage tools are the last line of defense. They must be fast, adaptable, and deeply integrated with biological risk models, not just keyword filters.

**A Living Fortress: Toward Adaptive Biosafety and Biosecurity Systems.** Securing the future of GenAI in the biosciences will require more than static safeguards, it demands infrastructure that is intelligent, integrated, and adaptive. As GenAI models become more capable of planning, simulating, and generating biological designs, our oversight systems must evolve in lockstep, offering not just filters but informed, real-time risk assessments and dynamic intervention strategies.

**Integration with domain knowledge.** Next-generation biosafety systems will need to interface deeply with biological knowledge. This includes connecting to curated databases of genomic sequences, protein functions, pathogenic mechanisms, and therapeutic compounds. Such integration enables contextual analysis—allowing safety layers to simulate downstream effects, verify biological plausibility, and anticipate potential harms with domain-specific

grounding. These tools must be paired with external capabilities such as molecular simulation engines, structure predictors, and API-accessible regulatory frameworks to validate emerging risks at inference time.

**Towards agentic capabilities.** Beyond knowledge integration, emerging safeguards must also demonstrate higher-order reasoning: the ability to plan multi-step responses, evaluate scientific and policy constraints, and red-team outputs against evolving threats. Drawing inspiration from recent advances in agentic AI, these systems could modularly score dual-use potential, assess confidence in high-stakes generations, and trace prompt-output pathways for auditability. To remain trustworthy, their actions should be explainable, logged, and subject to human oversight.

**Self-Evolving Capabilities.** Finally, such oversight mechanisms must themselves be capable of evolution. The biosafety tools of the near future must learn from new threat patterns, incorporate feedback from red-team exercises, and stay synchronized with technological progress—whether that’s the release of more powerful models or the emergence of novel biological capabilities like virtual cell simulation. An effective biosafety system is not a firewall, it is a living guardian.

#### Final Message

No single measure can fully secure GenAI in the biosciences. But layered interventions spanning training, deployment, access, and monitoring can collectively build a resilient infrastructure to safeguard GenAI. Security in this domain is not static policy, but dynamic practice. And staying ahead requires oversight mechanisms and technical capabilities and evaluation tools that evolve as fast as the GenAI models they guard.

## 5 Conclusion: From Awareness to Action

This Perspective has surfaced a growing expert consensus: GenAI models operating in the biosciences present unprecedented biosecurity risks that outpace the assumptions of existing safeguards. Through over 130 interviews with global leaders in synthetic biology, cybersecurity, governance, and AI development, we have highlighted the dissonance between the old playbook designed for physical gene synthesis and the new frontier, where models can design dangerous biological sequences without ever ordering DNA. The risks are no longer theoretical. The tools are widely accessible. And the gaps in oversight are widening.

Addressing this rapidly shifting threat landscape requires a fundamental rethinking of how we govern dual-use technologies. First, technical safeguards must be embedded across the AI lifecycle from training-time dataset curation to inference time filtering, access controls, output watermarking, and red-teaming protocols. These tools must be explainable, auditable, and grounded in biological domain knowledge. Second, interdisciplinary collaboration is critical: no single stakeholder group can anticipate or manage these risks alone. Biologists, AI researchers, ethicists, national security experts, and policymakers must work together to build a resilient safety architecture. Third, we need adaptive and strong governance and policies that may include tiered access, model registration, risk-based release standards, and international coordination.

Ultimately, securing the bio-AI frontier is not a matter of technical patchwork or post-hoc policy. It requires a layered, dynamic, and forward-looking strategy, one that evolves as fast as the models it seeks to protect. Only through anticipatory, cross-sectoral, and globally coordinated action can we harness the transformative potential of GenAI in the life sciences while minimizing the profound risks it brings.

## References

1. OpenAI. Introducing chatgpt. *OpenAI Blog* (2022). URL <https://openai.com/index/chatgpt/>.
2. Team, G. *et al.* Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530* (2024).
3. Anthropic. Claude. <https://www.anthropic.com/claude>. Accessed on July 6, 2025.
4. Liu, A. *et al.* Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024).

5. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* **379**, 1123–1130 (2023). URL <https://www.science.org/doi/abs/10.1126/science.ade2574>. <https://www.science.org/doi/pdf/10.1126/science.ade2574>.
6. Madani, A. *et al.* Progen: Language modeling for protein generation. *arXiv preprint arXiv:2004.03497* (2020).
7. Nijkamp, E., Ruffolo, J. A., Weinstein, E. N., Naik, N. & Madani, A. Progen2: exploring the boundaries of protein language models. *Cell systems* **14**, 968–978 (2023).
8. Jumper, J. *et al.* Highly accurate protein structure prediction with alphafold. *Nature* **596**, 583–589 (2021). URL <https://doi.org/10.1038/s41586-021-03819-2>.
9. Hayes, T. *et al.* Simulating 500 million years of evolution with a language model. *Science* **387**, 850–858 (2025). URL <https://www.science.org/doi/abs/10.1126/science.ads0018>. <https://www.science.org/doi/pdf/10.1126/science.ads0018>.
10. Watson, J. L. *et al.* De novo design of protein structure and function with rfdiffusion. *Nature* **620**, 1089–1100 (2023). URL <https://doi.org/10.1038/s41586-023-06415-8>.
11. Ahern, W. *et al.* Atom level enzyme active site scaffolding using rfdiffusion2. *bioRxiv* (2025). URL <https://www.biorxiv.org/content/early/2025/04/10/2025.04.09.648075.1>. <https://www.biorxiv.org/content/early/2025/04/10/2025.04.09.648075.1.full.pdf>.
12. Zhang, Z., Shen, W. X., Liu, Q. & Zitnik, M. Efficient generation of protein pockets with pocketgen. *Nature Machine Intelligence* 1–14 (2024).
13. Zhang, Z. *et al.* Rnagenesis: A generalist foundation model for functional rna therapeutics. *bioRxiv* 2024–12 (2024).
14. Chen, J. *et al.* Interpretable rna foundation model from unannotated data for highly accurate rna structure and function predictions (2022). URL <https://arxiv.org/abs/2204.00300>. 2204.00300.
15. Rnacentral: a hub of information for non-coding rna sequences. *Nucleic Acids Research* **47**, D221–D229 (2019).
16. Dalla-Torre, H. *et al.* Nucleotide transformer: building and evaluating robust foundation models for human genomics. *Nature Methods* **22**, 287–297 (2025).
17. Nguyen, E. *et al.* Sequence modeling and design from molecular to genome scale with evo. *Science* **386**, eado9336 (2024). URL <https://www.science.org/doi/abs/10.1126/science.ado9336>. <https://www.science.org/doi/pdf/10.1126/science.ado9336>.
18. Bixi, G. *et al.* Genome modeling and design across all domains of life with evo 2. *bioRxiv* (2025). URL <https://www.biorxiv.org/content/early/2025/02/21/2025.02.18.638918>. <https://www.biorxiv.org/content/early/2025/02/21/2025.02.18.638918.full.pdf>.
19. Jin, W., Barzilay, R. & Jaakkola, T. Junction tree variational autoencoder for molecular graph generation. In *International conference on machine learning*, 2323–2332 (PMLR, 2018).
20. Imrie, F., Bradley, A. R., van der Schaar, M. & Deane, C. M. Deep generative models for 3d linker design. *Journal of chemical information and modeling* **60**, 1983–1995 (2020).
21. Swanson, K. *et al.* Generative ai for designing and validating easily synthesizable and structurally novel antibiotics. *Nature Machine Intelligence* **6**, 338–353 (2024).
22. Zhou, G. *et al.* Uni-mol: A universal 3d molecular representation learning framework. In *The Eleventh International Conference on Learning Representations* (2023). URL <https://openreview.net/forum?id=6K2RM6wVqKu>.
23. Yang, F. *et al.* scbert as a large-scale pretrained deep language model for cell type annotation of single-cell rna-seq data. *Nature Machine Intelligence* **4**, 852–866 (2022).



24. Hao, M. *et al.* Large-scale foundation model on single-cell transcriptomics. *Nature methods* **21**, 1481–1491 (2024).
25. Theodoris, C. V. *et al.* Transfer learning enables predictions in network biology. *Nature* **618**, 616–624 (2023).
26. Luo, E., Hao, M., Wei, L. & Zhang, X. scdiffusion: conditional generation of high-quality single-cell data using diffusion model. *Bioinformatics* **40**, btae518 (2024).
27. Wang, X. *et al.* A pathology foundation model for cancer diagnosis and prognosis prediction. *Nature* **634**, 970–978 (2024). URL <https://doi.org/10.1038/s41586-024-07894-z>.
28. Lu, M. Y. *et al.* A visual-language foundation model for computational pathology. *Nature Medicine* **30**, 863–874 (2024).
29. Wang, J. *et al.* Self-improving generative foundation model for synthetic medical image generation and clinical applications. *Nature Medicine* **31**, 609–617 (2025).
30. Wang, M. *et al.* A call for built-in biosecurity safeguards for generative ai tools. *Nature Biotechnology* 1–3 (2025).
31. Nuclear Threat Initiative. Developing guardrails for ai biodesign tools. On-line report (2024). URL <https://www.nti.org/analysis/articles/developing-guardrails-for-ai-biodesign-tools/>. Accessed: 2025-05-12.
32. Baker, D. & Church, G. Protein design meets biosecurity (2024).
33. Abramson, J. *et al.* Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature* **630**, 493–500 (2024). URL <https://doi.org/10.1038/s41586-024-07487-w>.
34. Wittmann, B. J. *et al.* Toward ai-resilient screening of nucleic acid synthesis orders: Process, results, and recommendations. *bioRxiv* (2024). URL <https://www.biorxiv.org/content/early/2024/12/04/2024.12.02.626439>. <https://www.biorxiv.org/content/early/2024/12/04/2024.12.02.626439.full.pdf>.
35. Fan, J. *et al.* Safeprotein: Red-teaming framework and benchmark for protein foundation models (2025). 2509.03487.
36. Zhang, Z., Zhou, Z., Jin, R., Cong, L. & Wang, M. Genebreaker: Jailbreak attacks against dna language models with pathogenicity guidance. *arXiv preprint arXiv:2505.23839* (2025).
37. King, S. H. *et al.* Generative design of novel bacteriophages with genome language models. *bioRxiv* (2025).
38. Urbina, F., Lentzos, F., Invernizzi, C. & Ekins, S. Dual use of artificial-intelligence-powered drug discovery. *Nature Machine Intelligence* **4**, 189–191 (2022). URL <https://doi.org/10.1038/s42256-022-00465-9>.
39. Bunne, C. *et al.* How to build the virtual cell with artificial intelligence: Priorities and opportunities. *Cell* **187**, 7045–7063 (2024).
40. Huang, K. *et al.* Biomni: A general-purpose biomedical ai agent. *bioRxiv* 2025–05 (2025).
41. Zhang, Z. *et al.* Origene: A self-evolving virtual disease biologist automating therapeutic target discovery. *bioRxiv* 2025–06 (2025).
42. Boiko, D. A., MacKnight, R., Kline, B. & Gomes, G. Autonomous chemical research with large language models. *Nature* **624**, 570–578 (2023). URL <https://doi.org/10.1038/s41586-023-06792-0>.
43. Liu, Y. *et al.* Trustworthy llms: a survey and guideline for evaluating large language models' alignment. *arXiv preprint arXiv:2308.05374* (2023).
44. Chua, J., Li, Y., Yang, S., Wang, C. & Yao, L. Ai safety in generative ai large language models: A survey. *arXiv preprint arXiv:2407.18369* (2024).

45. Kirchenbauer, J. *et al.* A watermark for large language models (2024). URL <https://arxiv.org/abs/2301.10226>. 2301.10226.
46. Xu, Y., Liu, A., Hu, X., Wen, L. & Xiong, H. Mark your llm: Detecting the misuse of open-source large language models via watermarking (2025). URL <https://arxiv.org/abs/2503.04636>. 2503.04636.
47. Xu, X., Yao, Y. & Liu, Y. Learning to watermark llm-generated text via reinforcement learning (2024). URL <https://arxiv.org/abs/2403.10553>. 2403.10553.
48. Pan, L. *et al.* Markllm: An open-source toolkit for llm watermarking (2024). URL <https://arxiv.org/abs/2405.10051>. 2405.10051.
49. Li, S., Yao, L., Gao, J., Zhang, L. & Li, Y. Double-i watermark: Protecting model copyright for llm fine-tuning (2024). URL <https://arxiv.org/abs/2402.14883>. 2402.14883.
50. Zhang, R., Hussain, S. S., Neekhara, P. & Koushanfar, F. Remark-llm: a robust and efficient watermarking framework for generative large language models. In *Proceedings of the 33rd USENIX Conference on Security Symposium*, SEC '24 (USENIX Association, USA, 2024).
51. Zhang, Z. *et al.* Foldmark: Protecting protein generative models with watermarking. *bioRxiv* (2024).
52. Dai, J. *et al.* Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773* (2023).
53. Chakraborty, S. *et al.* Parl: A unified framework for policy alignment in reinforcement learning from human feedback (2024). URL <https://arxiv.org/abs/2308.02585>. 2308.02585.
54. Swamy, G., Dann, C., Kidambi, R., Wu, Z. S. & Agarwal, A. A minimaximalist approach to reinforcement learning from human feedback (2024). URL <https://arxiv.org/abs/2401.04056>. 2401.04056.
55. Wu, Y. *et al.* Self-play preference optimization for language model alignment (2024). URL <https://arxiv.org/abs/2405.00675>. 2405.00675.
56. Chakraborty, S. *et al.* Maxmin-rlhf: Alignment with diverse human preferences (2024). URL <https://arxiv.org/abs/2402.08925>. 2402.08925.
57. Rafailov, R. *et al.* Direct preference optimization: Your language model is secretly a reward model (2024). URL <https://arxiv.org/abs/2305.18290>. 2305.18290.
58. de Almeida, B. P. *et al.* A multimodal conversational agent for dna, rna and protein tasks. *Nature Machine Intelligence* 1–14 (2025).
59. Bai, T., Luo, J., Zhao, J., Wen, B. & Wang, Q. Recent advances in adversarial training for adversarial robustness. *arXiv preprint arXiv:2102.01356* (2021).
60. Goel, K. *et al.* Robustness gym: Unifying the nlp evaluation landscape (2021). URL <https://arxiv.org/abs/2101.04840>. 2101.04840.
61. Liu, X. *et al.* Adversarial training for large neural language models (2020). URL <https://arxiv.org/abs/2004.08994>. 2004.08994.
62. Xhonneux, S., Sordoni, A., Günemann, S., Gidel, G. & Schwinn, L. Efficient adversarial training in llms with continuous attacks (2024). URL <https://arxiv.org/abs/2405.15589>. 2405.15589.
63. Wang, Z. *et al.* Better diffusion models further improve adversarial training (2023). URL <https://arxiv.org/abs/2302.04638>. 2302.04638.
64. Name, A. Meet h4d team omnyra. *StanfordH4D* (Year). URL <https://stanfordh4d.substack.com/p/meet-h4d-team-omnyra>.
65. Perez, E. *et al.* Red teaming language models with language models. *arXiv preprint arXiv:2202.03286* (2022).
66. Mudgal, S. *et al.* Controlled decoding from language models (2024). URL <https://arxiv.org/abs/2310.17022>. 2310.17022.

67. Chakraborty, S. *et al.* Transfer q star: Principled decoding for llm alignment (2024). URL <https://arxiv.org/abs/2405.20495>. 2405.20495.
68. Khanov, M., Burapachee, J. & Li, Y. Args: Alignment as reward-guided search (2024). URL <https://arxiv.org/abs/2402.01694>. 2402.01694.
69. Ghosal, S. S. *et al.* Immune: Improving safety against jailbreaks in multi-modal llms via inference-time alignment (2025). URL <https://arxiv.org/abs/2411.18688>. 2411.18688.
70. Dip, S. A. *et al.* Patholm: Identifying pathogenicity from the dna sequence through the genome foundation model (2024). URL <https://arxiv.org/abs/2406.13133>. 2406.13133.
71. Jinnai, Y., Morimura, T., Ariu, K. & Abe, K. Regularized best-of-n sampling to mitigate reward hacking for language model alignment (2024). URL <https://arxiv.org/abs/2404.01054>. 2404.01054.
72. Amini, A., Vieira, T. & Cotterell, R. Variational best-of-n alignment (2024). URL <https://arxiv.org/abs/2407.06057>. 2407.06057.
73. Beirami, A. *et al.* Theoretical guarantees on the best-of-n alignment policy (2024). URL <https://arxiv.org/abs/2401.01879>. 2401.01879.
74. Madaan, A. *et al.* Self-refine: Iterative refinement with self-feedback (2023). URL <https://arxiv.org/abs/2303.17651>. 2303.17651.
75. Chao, P. *et al.* Jailbreaking black box large language models in twenty queries (2024). URL <https://arxiv.org/abs/2310.08419>. 2310.08419.
76. Mehrabi, N. *et al.* Flirt: Feedback loop in-context red teaming (2024). URL <https://arxiv.org/abs/2308.04265>. 2308.04265.
77. Chakraborty, S. *et al.* Review, refine, repeat: Understanding iterative decoding of ai agents with dynamic evaluation and selection (2025). URL <https://arxiv.org/abs/2504.01931>. 2504.01931.

## **Acknowledgements**

## **Author contributions statement**

## **Competing interests**

The authors declare no competing interests. Disclaimer: Certain tools and software are identified in this paper to foster understanding. Such identification does not imply recommendation or endorsement by the National Institute of Standards and Technology, nor does it imply that the tools and software identified are necessarily the best available for the purpose.

## **Additional information**

**Correspondence and requests for materials** should be addressed to Mengdi Wang.

## A GenAI in Biosciences

Biological Target	Frontier AI Models	Purpose / Capability
Proteins	ESM, ESM-2, ESM-3, Progen, Progen2, ProteinMPNN, RFdiffusion, RFdiffusionAA, PocketGen, BindCraft	Learn protein sequence–structure–function relationships; predict structures; generate novel proteins; design functional protein backbones and ligand-binding sites.
RNA	RNAGenesis, RNA-FM	Model RNA sequence–structure relationships; predict secondary structures; annotate functions; enable de novo RNA design.
DNA	Nucleotide Transformer, Evo, Evo2	Model genomic sequences at varying scales; detect functional elements; enable long-context genome modeling and design.
Small Molecules	JT-VAE, DeLinker, SyntheMol, MolDiff, Uni-Mol	Explore large chemical space; design bioactive, synthesizable molecules; model 2D/3D molecular structures and binding interactions.
Single Cells	scBERT, scFoundation, Geneformer, scDiffusion	Learn cell-type representations; correct batch effects; predict perturbation responses; simulate transcriptomic profiles and developmental trajectories.
Clinical / Medical Data	CHIEF, CONCH, MINIM	Analyze histopathology images; predict cancer origin and prognosis; integrate image–text understanding for clinical interpretation.

**Table S1.** Generative AI models across different biological scales and their primary capabilities.



## B Explanation of Terms

Term	Definition
Pathogen	A biological agent that can cause disease in its host, such as bacteria, viruses, fungi, or parasites.
Toxin	A poisonous substance produced by living organisms that can cause disease or death when introduced into the body.
Omics	A branch of biology that studies the complete set of a particular type of biomolecule or molecular process in an organism, such as genomics (genome), proteomics (proteome), and metabolomics (metabolome).
Genome	The complete set of genetic material in an organism, typically referring to nuclear DNA in eukaryotes, but can include other DNA such as mitochondrial or chloroplast DNA.
Nucleic acids	Macromolecules that store and transmit genetic information in living organisms, composed of nucleotides, which consist of a sugar, a phosphate group, and a nitrogenous base. The two main types are DNA and RNA.
CRISPR	A gene-editing technology that uses the CRISPR-Cas9 system to make precise changes to the DNA of living organisms, based on a natural defense mechanism in bacteria.
Embryo	The early stage of development of a multicellular organism, typically from fertilization until the end of the eighth week in humans, after which it is called a fetus.
Bioweapon	A type of weapon that uses biological agents, such as bacteria, viruses, or toxins, to cause disease or death in humans, animals, or plants.

**Table S2.** Glossary of key biological terms

Term	Definition
Foundation Models	Large deep learning models pre-trained on vast, general-purpose datasets, designed to be versatile and fine-tuned for specific applications.
Deep Generative Models	Machine learning models that use deep neural networks to generate new data similar to their training data.
Fine-Tuning	The process of adapting a pre-trained model to a new task by further training it on a new dataset, especially useful when the new dataset is small.
Reinforcement Learning	A machine learning paradigm where an agent learns to maximize a reward signal by interacting with an environment.
Retrieval-Augmented Generation	A technique that enhances large language models by providing them with external information to improve the accuracy and relevance of their generated text.
Self-supervised Learning	A machine learning method where a model learns from unlabeled data by generating its own labels or supervisory signals.
Continued on next page	

<b>Term</b>	<b>Definition</b>
AI Agent	A software system that uses artificial intelligence to perform tasks or make decisions based on input prompts.
Multi-Agent System	A system consisting of multiple AI agents that interact with each other and their environment to achieve individual or collective objectives.
Jailbreak Attacks	Techniques used to manipulate AI models, particularly large language models, into producing incorrect or fabricated information that circumvents their built-in safety mechanisms.
Membership Inference Attacks	Attacks where an adversary tries to determine whether a specific data record was used to train a machine learning model.
Watermark	In AI, a technique to mark or identify content generated by an AI model, often used to verify the origin or authenticity of the content.
Surrogate Model	A simpler model used to approximate the behavior of a more complex or computationally intensive model, often for interpretability or efficiency reasons.
Multi-Modality Model	A machine learning model that can process and integrate information from multiple different types of data, such as text, images, and audio.
Parameter Efficient Finetuning	Techniques for fine-tuning large pre-trained models by updating only a small subset of their parameters, to make the process more computationally efficient.
Prompt	The input text or instruction provided to an AI model, particularly large language models, to guide its output.
Hallucination	When an AI model, especially a large language model, generates incorrect or fabricated information that is presented as if it were true.
Transformers	A type of neural network architecture that uses self-attention mechanisms to process sequences, such as text, in parallel and understand the relationships between different parts of the sequence.
Embodied AI	Artificial intelligence systems that are integrated into physical entities, like robots, to interact with and learn from their environment through sensors and actuators.
Backdoor Attacks	Attacks where an adversary modifies a machine learning model during its training phase to behave in a specific, malicious way when presented with certain inputs or triggers.
Property Inference Attacks	Attacks where an adversary tries to infer specific properties or characteristics of the training data used to train a machine learning model, by analyzing the model's behavior or outputs.

**Table S3.** Glossary of Key Machine Learning Terms

## C Survey Methodology

We led an in-depth qualitative survey of stakeholders across industry, government, academia, and policy to highlight emerging biosecurity concerns at the intersection of artificial intelligence, synthetic biology, and governance. Over a ten-week period, we conducted 130 semi-structured discovery interviews with individuals spanning four key sectors:

- 43 professionals from private sector companies including nucleic acid synthesis companies, screening technology providers, and technology and AI companies.
- 41 academics and researchers spanning disciplines such as microbiology, virology, machine learning, and synthetic genomics. Institutions represented included the Lawrence Berkeley National Laboratory, Lawrence Livermore National Laboratory, Stanford Biosecurity, MIT Media Lab, Wyss Institute at Harvard University, and the Gladstone Institutes.
- 29 representatives from government agencies, including public health authorities, national security offices, and regulatory bodies. These organizations included the Department of Homeland Security, National Institutes of Health, Federal Bureau of Investigation, U.S. Navy, National Security Commission on Emerging Biotechnology, National Institute of Standards and Technology, and the Defense Innovation Unit.
- 17 non-governmental policy experts and think tank leaders specializing in biosecurity advocacy and AI governance. Organizations represented included the Federation of American Scientists, American Association for the Advancement of Science, Engineering Biology Research Consortium, and Centre for Long-Term Resilience.

Interviewees were selected to reflect a diversity of institutional roles, disciplinary backgrounds, and viewpoints, with many recognized as experts in their respective fields. Each interview lasted between 30 to 60 minutes and was conducted using video conference. Interviews followed a semi-structured protocol, beginning with open-ended prompts about current and future biosecurity risks and then progressing organically based on the interviewee's expertise and perspectives. To ensure consistency across interviews while allowing flexibility, we used a core set of guiding questions as listed below:

- What emerging risks or concerns do you anticipate at the intersection of artificial intelligence and synthetic biology?
- How effective do you think current biosecurity practices are in addressing AI-enabled threats?
- What gaps, if any, do you see in the existing tools, frameworks, or governance systems for detecting and preventing misuse in the synthetic biology context?
- What types of interventions do you think are most needed to strengthen safeguards in biosecurity?
- What role should your sector play in shaping the future of biosecurity as capabilities evolve?

To analyze and quantify interview responses, we developed a framework around four key themes: (1) AI misuse in biological contexts, (2) perceptions of current sequence screening tools, (3) attitudes toward functional or AI-native screening approaches, and (4) governance gaps and regulatory needs. Each interview was reviewed and assigned a binary relevance score (1 = substantial mention, 0 = no substantial mention) for each theme. A theme was coded as "1" only if the interviewee offered a substantive comment that demonstrated support, engagement, or direct relevance to the topic such as describing concrete practices, expressing informed perspectives, or endorsing the importance of the issue. Passing mentions or unrelated concerns were not counted.

This framework allowed us to quantify the prevalence of each theme across all 130 interviews, compare thematic engagement across stakeholder groups, and identify patterns in concern and/or alignment across sectors.

D Interview Results

**Table S4.** Summary of Interview Themes by Sector

Sector	AI Misuse is Urgent	Screening Tools are Inadequate	Supports Functional Screening	Calls for Governance Standards
Industry	30	10	25	24
Government	19	13	15	21
Academia	34	24	24	35
Policy	16	14	12	16
<b>Total</b>	<b>99</b>	<b>61</b>	<b>76</b>	<b>96</b>