

Active Inference AI Systems for Scientific Discovery

Karthik Duraisamy

*Michigan Institute for Computational Discovery & Engineering,
University of Michigan, Ann Arbor.*

Abstract

The rapid evolution of artificial intelligence has led to expectations of transformative scientific discovery, yet current systems remain fundamentally limited by their operational architectures, brittle reasoning mechanisms, and their separation from experimental reality. Building on earlier work—including foundational contributions in automated scientific discovery [58], closed-loop laboratory systems [54, 14, 92], and causal machine learning [75, 85]—this perspective contends that progress in AI-driven science now depends on closing three fundamental gaps—the abstraction gap, the reasoning gap, and the reality gap—rather than on model size/data/test time compute. While prior systems have made substantial progress on individual gaps, we argue that genuinely novel scientific discovery—defined here as the identification of previously unknown causal mechanisms, physical laws, or theoretical frameworks that generalize beyond training distributions—requires their integrated resolution. Scientific reasoning demands internal representations that support simulation of actions and response, causal structures that distinguish correlation from mechanism, and continuous calibration. Active inference AI systems for scientific discovery are defined as those that (i) maintain long-lived research memories grounded in causal self-supervised foundation models, (ii) employ symbolic or neuro-symbolic planners equipped with Bayesian guardrails, (iii) grow persistent knowledge graphs where thinking generates novel conceptual nodes, reasoning establishes causal edges, and real-world interaction prunes false connections while strengthening verified pathways, and (iv) refine their internal representations through closed-loop interaction with both high-fidelity simulators and automated laboratories—an operational loop where mental simulation guides action and empirical surprise reshapes understanding. In essence, this work outlines design principles for an architecture in which discovery arises from the interplay between internal models that enable counterfactual reasoning and external validation that grounds hypotheses in reality. It is also argued that the inherent ambiguity in feedback from simulations and experiments, and underlying uncertainties makes human judgment indispensable, not as a temporary scaffold but as a permanent architectural component.

1 Present day AI Systems and Scientific Discovery

Over the past decade, the evolution of AI foundation model research has followed a clear sequence of discrete jumps in capability. The advent of the Transformer[96] marked a phase dominated by architectural innovations, which was rapidly succeeded by scaling demonstrations such as GPT-2[81]. The maturation of large-language-model pre-training then gave way to the “usability turn”: the shift from models optimized purely for benchmark performance to chat-oriented systems fine-tuned for alignment and safety that enabled direct human interaction[72]. The current frontier is characterised by reasoning-emulation systems that incorporate tool use, scratch-pad planning, or program-synthesis objectives[70]. A fifth, still-incipient phase points toward autonomous agents which are software systems capable of perceiving their environment, making decisions, and taking actions to achieve specified goals without continuous human supervision [103]that can decompose tasks, invoke external software or laboratories, and learn from the resulting feedback. Scientific applications of AI have

echoed each of these transitions at a compressed cadence. As examples, SchNet translated architectural advances to quantum chemistry[87]; AlphaFold combined evolutionary multiple-sequence-alignment features with SE(3)-equivariant neural architectures [48] leveraging decades of accumulated protein structure data to achieve high accuracy, though still requiring petascale compute and vast unlabeled evolutionary sequences rather than dramatically fewer examples; ChemBERTa [23] and FourCastNet [74] adapted language and vision innovations to molecular and climate domains; and AlphaGeometry applied reasoning-centric objectives to symbolic mathematics[94]. Collectively, recent works [38, 9, 11] chart a shift from single, specialized pre-trained model to workflow orchestration, suggesting that future breakthroughs may hinge on integrating heterogeneous, domain-aware agents capable of planning experiments, steering simulations, and iteratively refining hypotheses across scales.

This highlights a deeper challenge for scientific discovery, which must reason across stacked layers of abstraction: the emergence of unexpected phenomena at higher scales, just as local atmospheric equations do not directly predict large-scale El Niño patterns [20, 19]. This challenge of multi-scale abstraction has been recognized and discussed extensively in the automated scientific discovery literature [58]. To address this challenge, it may be required to deliberately architect systems with built-in mechanisms for hierarchical inference, equipping them with specialized components that can navigate between reductionist details and emergent phenomena. A compelling counter-argument posits that such abstract reasoning is not a feature to be explicitly engineered, but an emergent property that will arise from sufficient scale and diverse data. Proponents of this view might point to tools such as AlphaGeometry [94], where complex, formal reasoning appears to emerge from a foundation model trained on vast synthetic data. However, we contend that while scaling can master any pattern present in the training distribution—even highly complex ones—it is fundamentally limited to learning correlational structures. Scientific discovery, in contrast, hinges on understanding interventional and counterfactual logic [75, 43]: *what happens when the system is deliberately perturbed?* This knowledge cannot be passively observed in static data; it must be actively acquired through interaction with the world or a reliable causal model thereof though we note that causal assumptions can sometimes be justified from observational data when appropriate structural constraints hold [18, 6]. The ‘reality gap’ thus remains a significant barrier that pure scaling may not cross.

It is also pertinent to examine the nature of present-day scientific discovery before speculating the role of AI. Modern science has moved beyond the romanticized vision of solitary geniuses grappling with nature’s mysteries. It may be difficult to generalize or even define the nature of discovery, but it is safe to assume that many of today’s discoveries emerge from vast collaborations parsing petabytes of data from instruments such as the Large Hadron Collider or from distributed sensor networks or large-scale computations and most importantly, refining hypothesis in concurrence with experiments and simulations. In fields such as high-energy physics, the bottleneck has shifted toward complexity management, whereas in data-constrained arenas such as fusion-plasma diagnostics, insight scarcity remains dominant; any general framework must therefore account for both regimes. Even if one possesses the raw data to answer profound questions, we often lack the cognitive architecture to navigate the combinatorial explosion of hypotheses, interactions, and emergent phenomena. The connection between cognitive architecture and the capacity for rich abstraction is direct: effective reasoning over scientific concepts requires working memory structures that can hold, manipulate, and compose abstract representations [57], capabilities that current transformer architectures approximate but do not fully realize. This creates an opportunity for AI systems—to excel precisely where human cognition fails, in maintaining consistency across very high-dimensional parameter spaces, identifying and reasoning about subtle patterns in noisy data. At this juncture, it has to be emphasized that generating novel hypotheses might be the easy part [39]: the challenge

is in rapidly assessing the impact of a hypothesis or action in an imaginary space. Thus AI systems have to be equipped with rich world models that can rapidly explore vast hypothesis spaces, and integrated with efficient computations and experiments to provide valuable feedback.

It is essential to distinguish between *engineering discovery* -the optimization of known systems toward well-defined objectives- and *scientific discovery* -the identification of novel causal mechanisms, physical laws, or theoretical frameworks. Engineering problems in materials science or drug design can often be addressed by closed-loop optimization systems operating over predefined search spaces [14, 92]. Scientific discovery, by contrast, requires the capacity for causal reasoning and theoretical insight that may restructure the problem space itself. Prior closed-loop systems have demonstrated genuine discoveries: Adam, an automated scientist for yeast metabolism, combined ontological reasoning with active learning to design and execute falsifiable experiments, leading to novel findings in genomics [54]; its successor Eve identified that triclosan is effective against malaria [100]. These systems embody the active inference principles we advocate: maintaining scientific memory while engaging in closed-loop interaction with physical laboratories. The contribution of this perspective is not to claim that such systems are impossible, but to argue that scaling beyond well-defined domains requires addressing the abstraction, reasoning, and reality gaps in an integrated manner, and to propose architectural principles for doing so.

Modern science also operates within a myriad of constraints that exist in economic, legal and social dimensions rather than physical laws. These constraints favor certain types of AI-driven discovery while effectively prohibiting others. Additionally, the structure of the modern scientific enterprise—with its emphasis on incremental, publishable units and citation metrics—may be fundamentally misaligned with the kind of patient, integrative thinking required for paradigm shifts. AI systems could theoretically ignore these social pressures, pursuing research programs too risky or long-term for humans. But this same freedom from social constraint raises the possibility of AI systems optimizing for discovery without broader goals accounted for. Perhaps more concerning, AI systems trained on existing scientific literature risk amplifying current biases and narrowing the space of explored ideas. In other words, AI systems might converge on well-studied paths [17, 65], reducing the rich variety of research directions to a handful of statistically probable avenues. Scientific progress demands not convergence but divergence—an explosion of hypotheses, methodologies, and frameworks that challenge orthodoxy. The challenge is in designing AI systems that expand rather than constrain the landscape of scientific imagination.

Against this backdrop, the remainder of this perspective piece is organized around three interlocking hurdles that scientific discovery architectures must clear: (i) the abstraction gap, which separates low-level statistical regularities from the mechanistic concepts on which scientists actually reason; (ii) the reasoning gap, which limits today’s models to correlation-driven pattern completion rather than causal, counterfactual inference; and (iii) the reality gap, which isolates computation from the empirical feedback loops that ultimately arbitrate truth. Each gap both constrains and amplifies the others: without rich abstractions there is little substrate for reasoning, and without tight coupling to reality even the most elegant abstractions may drift toward irrelevance. While we treat these gaps in separate sections for expository clarity, they are deeply interconnected: the abstraction gap concerns *what* representations the system can manipulate, the reasoning gap concerns *how* those representations are transformed, and the reality gap concerns *whether* those transformations correspond to the external world. Closing one gap in isolation yields limited progress; for instance, rich abstractions without causal reasoning remain correlational, while causal reasoning without empirical grounding may drift into unfalsifiable speculation.

Philosophical Foundations

The quest for scientific discovery via computation confronts a fundamental paradox. Gödel [37] proved that formal systems are incomplete, and cannot self-consistently prove all truths, while Wolfram [101] demonstrates that computational irreducibility pervades nature. Additionally, Penrose [76, 77] contends that human insight transcends algorithms. Yet scientific theories and computations are found to be highly effective in many cases. Insight can be gained from Wolfram’s recent comment [102]: *The very presence of computational irreducibility necessarily implies that there must be pockets of computational reducibility, where at least certain things are regular and predictable. It is within these pockets of reducibility that science fundamentally lives.*

In most cases, these pockets cannot be deduced a priori—they require empirical discovery. This connects to Popper’s [78] falsificationism: we cannot prove we have found true reducibility, but we can discover boundaries through experiments that challenge assumptions. Empirical feedback escapes Gödel’s constraints while delineating where nature permits shortcuts. Kuhn’s analysis [56] adds temporal dynamics: Science alternates between normal science within established pockets and paradigm shifts that restructure understanding. AI systems must balance exploiting known regularities with flexibility to reconceptualize when evidence demands.

This synthesis directly informs our architecture. Thinking explores for new pockets and tests boundaries; reasoning exploits discovered regularities. World models encode provisional maps of known pockets, subject to Popper’s falsification and Kuhn’s paradigm shifts. Human steering proves essential. Humans provide non-computational insight for recognizing genuine understanding, value judgments for directing exploration, and navigation through paradigm shifts where evaluation criteria themselves transform. Humans can shape the search process by encoding domain knowledge, identifying significant anomalies, and recognizing connections that form larger frameworks. When Faraday discovered electromagnetic induction, he did not deduce it from Maxwell’s equations (which did not yet exist)—he found it through experiment. Thus, productive collaborations can implement the complete scientific method: AI generates and tests hypotheses at scale; humans provide insight and judgment and empirical feedback provides critical steering.

2 The Abstraction Gap

While early models largely manipulated tokens and pixels, recent advances in concept-bottleneck networks[55], symmetry-equivariant graph models[93], and neuro-symbolic hybrids[63] show preliminary evidence that contemporary AI can already represent and reason over higher-order scientific concepts and principles. Yet a physicist reasons in conservation laws and symmetry breaking, whereas language models still operate on surface statistics. Closing this abstraction gap requires addressing several intertwined weaknesses.

Modern transformer variants assemble chain-of-thought proofs[99] by replaying patterns observed during pre-training; they do not build explicit causal graphs or exploit formal logic engines except in narrow plug-in pipelines. As a result they fail at problems that demand deep compositionality. Several other shortcomings have also been pointed out [66]: (i) AI systems can create “illusions of understanding” where researchers mistake fluent outputs for genuine insight; (ii) they may narrow the diversity of research directions by channeling attention toward well-represented topics in training data; and (iii) they can erode researchers’ own reasoning skills through over-reliance on automated suggestions.

Thinking and reasoning

A critical distinction emerges between thinking and reasoning: Thinking can be operationalized as an iterative, exploratory process—searching for partial solutions in the form of patterns without guaranteed convergence. It is the slow, generative phase where new connections form and novel patterns emerge from a number of possibilities. Reasoning, by contrast, represents the fast, deterministic traversal of established knowledge structures—building the most expressive path through a graph of already-discovered patterns. This dichotomy [47] may explain why current AI systems excel at certain tasks while failing at others. Large language models can reason impressively when the requisite patterns exist in their training data—they rapidly traverse their learned knowledge graphs to construct seemingly intelligent responses. Yet they struggle with genuine thinking: the patient, iterative discovery of patterns that do not yet exist in their representational space. Scientific discovery demands both capabilities in careful balance.

Thinking generates hypotheses by discovering new patterns through mental simulation and exploration; reasoning then rapidly tests these patterns against existing knowledge and empirical constraints. The purpose of thinking therefore is not to solve problems directly but to expand the pattern vocabulary available for subsequent reasoning. Each thinking cycle potentially adds new nodes and edges to the knowledge graph, creating shortcuts and abstractions that make previously intractable reasoning paths suddenly accessible. This is perhaps why breakthrough discoveries often seem obvious in retrospect—the thinking phase has restructured the problem space so thoroughly that the reasoning path becomes trivial. We note that this terminology differs from Kahneman’s [49] System 1/System 2 framework, where System 2 denotes slow, deliberate reasoning. Our usage follows Johnson-Laird’s [47] mental models tradition, where “thinking” refers to the creative construction of new mental models (slow, exploratory), while “reasoning” refers to drawing inferences within established models (fast once patterns exist). The key insight is that scientific discovery requires both: the slow construction of novel conceptual frameworks *and* the rapid verification of their logical consequences. Both frameworks recognize the importance of dual cognitive modes; the difference lies in which dimension (effort vs. novelty) defines the slow/fast distinction. Recent neuro-symbolic RL agents hint at this synergy: the survey of Acharya et al.[1] chronicles agents that fuse neural perception with first-order symbolic planners while Mao et al.[63] demonstrate compositional question-answering by training a neural concept learner that hands off logic programs to a symbolic executor.

The gap between correlation and causation represents perhaps the most fundamental challenge in automated scientific discovery. While current models excel at finding statistical regularities, scientific understanding requires the ability to reason about interventions—to ask not just “what correlates with what?” but “what happens when we change this?”

Pearl’s causal hierarchy [75] distinguishes three levels of cognitive ability: association (seeing), intervention (doing), and counterfactuals (imagining). Current AI systems operate primarily at the associative level, occasionally reaching intervention through experimental design. True scientific reasoning requires all three, particularly the counterfactual ability to imagine alternative scenarios that violate observed correlations. This claim finds support in the history of science: the link between smoking and cancer was established through epidemiological reasoning about counterfactuals without randomized controlled trials [27]. Similarly, Darwin’s theory of evolution was fundamentally a counterfactual framework: imagining what would happen to populations under different selection

pressures. This connects directly to ethologist Konrad Lorenz’s insight [61]—first tied to learning systems by Scholkopf [85], building on Pearl’s foundational work [75]—that thinking is fundamentally about acting in imaginary spaces where we can violate the constraints of observed data. This *mental experimentation*—impossible in physical reality but accessible in the imagination—forms the basis of scientific law formation.

3 Unhobbling Intelligence

A certain level of consensus appears to be forming in the community that incremental scaling of present architectures may not deliver the qualitative leap that scientific discovery demands. Progress hinges on unhobbling—removing the design constraints that keep today’s models predictable, yet fundamentally limited—through concurrent advances in algorithms, speculation control, hardware co-design, and access models.

Algorithmic Gains Future systems must balance the complementary modes of thinking and reasoning as first-class architectural principles. Thinking—or slow, iterative discovery of new patterns—demands (i) world-model agents that can explore counterfactual spaces through mental simulation [41]; (ii) curiosity-driven mechanisms that reward pattern novelty over immediate task performance [84, 73]; and (iii) mechanisms for preventing premature convergence, such as entropy regularization [42] or explicit exploration bonuses that maintain hypothesis diversity during search [7]. Reasoning—the fast, deterministic traversal of pattern graphs—demands (i) efficient knowledge graph architectures with learned traversal policies [29]; (ii) neuro-symbolic stacks that maintain both continuous representations and discrete logical structures [63]; and (iii) caching mechanisms that transform expensive thinking outcomes into rapid reasoning primitives. The interplay between these modes mirrors how scientists alternate between exploratory experimentation (thinking) and theoretical derivation (reasoning).

The notion that “*thinking is acting in an imaginary space*”—as Konrad Lorenz observed [61]—provides a foundational principle for understanding how world models enable scientific discovery. Just as biological organisms evolved the capacity to simulate actions internally before committing physical resources, AI systems with rich world models can explore vast hypothesis spaces through mental simulation. This capability enables a qualitatively different mode of inference compared to pure pattern matching: it supports counterfactual reasoning, experimental design optimization, and the anticipation of empirical surprises before they manifest in costly real-world experiments. World models can serve as the substrate for this imaginary action space, encoding not just correlations but causal structures that permit intervention and manipulation. The fidelity of these mental simulations—their alignment with physical reality—determines whether the system’s thoughts translate into valid discoveries.

Scientific progress thrives on disciplined risk: venturing beyond received wisdom while remaining falsifiable. Current alignment protocols deliberately dampen exploratory behaviour, biasing models toward safe completion of well-trodden trajectories. Controlled speculation frameworks—for example, curiosity-driven reinforcement learning [71] combined with Bayesian epistemic guards—could allow systems to seek novel hypotheses, flag them with calibrated uncertainty, and propose targeted experiments for arbitration. Mechanisms such as self-consistency voting [98], adversarial peer review, and tool-augmented chain-of-thought audits offer additional scaffolding to keep high-variance reasoning tethered to empirical reality.

Recent empirical work by Buehler [13, 12] demonstrates that graph-based knowledge representations can bridge the abstraction gap. Specifically, recursive graph expansion experiments show that autonomous systems naturally develop hierarchical, scale-free networks mirroring human scientific knowledge structures. Without predefined ontologies, these systems spontaneously form

conceptual hubs and persistent bridge nodes, maintaining both local coherence and global integration—addressing precisely the limitations that prevent current AI from connecting low-level patterns to high-level scientific concepts. Indeed, success in one class of problems does not guarantee translation to other problems, domains and disciplines, but these works show that with appropriate graph-based representations, AI systems can discover novel conceptual relationships.

Computational Inefficiency Scaling laws show that models get predictably better with more data, parameter count and test time compute, yet every small gain might come at a great expense in time and/or energy. Such brute-force optimization contrasts sharply with biological economies in which sparse, event-based spikes[30] and structural plasticity[50] deliver continual learning at milliwatt scales. Bridging the gap will demand both algorithmic frugality—latent-variable models, active-learning curricula, reversible training—and hardware co-design. State-of-the-art foundation models require months of GPU time and $> 10^{25}$ FLOPs to reach acceptable performance on long-horizon benchmarks. Memory-reversible Transformers [62, 107] and curriculum training [97] have recently reduced end-to-end *training* costs by 30–45 %, without loss of final accuracy. Similar level of cost reductions have been reported [25] leveraging energy and power draw scheduling.

The von Neumann bottleneck—shuttling tensors between distant memory and compute—now dominates energy budgets [64]. Processing-in-memory fabrics [53], spiking neuromorphic cores that exploit event sparsity, analog photonic accelerators for low-latency matrix products, quantum samplers for combinatorial sub-routines [2] could open presently unreachable algorithmic spaces. Realising their potential outside of niche applications, however, will require co-design of hardware, software and algorithms and extensive community effort.

Evaluations Current leaderboards—e.g. MathBench[60], ARC[24], GSM8K[26]—scarcely probe the generative and self-corrective behaviours central to science. A rigorous suite should test whether a model can (i) identify when empirical data violate its latent assumptions, (ii) propose falsifiable hypotheses with quantified uncertainty, and (iii) adapt its internal representation after a failed prediction. Concretely, this may involve closed-loop benchmarks[51] in which the system picks experiments from a simulated materials lab, updates a dynamical model, and is scored on discovery efficiency; or theorem-proving arenas where credit is given only for proofs accompanied by interpretable lemmas. Without such stress-tests, superficial gains risk being mistaken for conceptual breakthroughs. Future evaluations can also assess the human-AI-reality-discovery feedback loop itself. Early exemplars such as DiscoveryWorld [46], PARTNR [21] and SciHorizon [80] represent steps towards this direction.

4 Architecture for the Era of Experience

Empirical feedback complements formal reasoning by supplying information inaccessible to purely deductive systems, thereby expanding—rather than mechanically escaping—the set of testable scientific propositions. The interplay between formal systems and empirical validation creates a bootstrap mechanism that circumvents incompleteness and irreducibility constraints. This suggests that AI systems for discovery must be fundamentally open—not just to new data, but to surprise from reality itself. Scientific history abounds with internally coherent theories that later failed empirical tests, underscoring the indispensability of continuous validation against data. Current AI systems excel at interpolation within their training distributions but struggle with the extrapolation that defines discovery. This is exacerbated by the fact that many scientific domains are characterized by sparse, expensive data and imperfect simulators. Unlike language modeling where data is abundant, a single protein crystallography experiment might take months and cost thousands of dollars. Simulations help but introduce their own biases—the “sim-to-real gap” that plagues robotics extends

to all of computational science. Our architecture must therefore implement a hybrid loop: physics priors guide ML surrogates, which direct active experiments, which update our understanding in continuous iteration.

Scientific Intelligence and the ‘Bitter Lesson’

The interplay between imagination and experimental validation creates a synergistic process at the heart of scientific discovery. Thinking provides the hypothesis generation engine, while empirical feedback provides the selection pressure that refines these mental models toward truth. These ideas have also been popularized by Yann LeCun [59] in the context of general machine intelligence. True understanding requires recognizing which causal relationships remain invariant as we move between the imaginary space of simplified models and the reality of full-scale experiments. The divergence between mental simulation and empirical outcome provides the richest learning signal.

Rich Sutton’s ‘bitter lesson’ [90] observes that scalable, general methods often outperform handcrafted heuristics, yet recent domain-aware models like AlphaFold illustrate that judiciously chosen scientific priors can accelerate learning within compute-intensive regimes. AlphaFold’s use of SE(3)-equivariant networks and multiple-sequence-alignment priors enabled high accuracy with ≈ 170000 labeled structures, but success still relied on petascale compute and vast unlabeled evolutionary data rather than a dramatic reduction from trillions of examples.

We propose a ‘complete lesson’: intelligent systems arise from computation plus constraints with feedback from reality. The ‘bitter lesson’ correctly observes that general learning methods eventually outperform human-engineered heuristics within a fixed computational domain. However, it implicitly assumes that all relevant information is already encoded in the data. We contend that the most critical constraint for scientific intelligence is not a clever, human-designed inductive bias, but the information-rich feedback from the physical world, a qualitatively different signal that contains latent information about causality, invariance, and the validity of the model’s core assumptions that may be absent from any tractably finite, static dataset.

The key insight is knowing when to leverage existing knowledge versus when to remain agnostic: When exploring many engineering domains (drug design, materials science), encode known physics. When questioning the nature of reality (quantum gravity, consciousness), avoid premature constraints. It is worth noting that much of science and engineering relies on the former scenario, and thus the art lies in choosing the right inductive biases for the problem at hand. This does not represent a retreat from the bitter lesson but its fulfillment: learning what to search for and how to search, guided by accumulated scientific knowledge and empirical feedback.

Causal Models The current paradigm of domain-specific foundation models—from protein language models to molecular transformers—represents significant progress in encoding domain knowledge. However, these models fundamentally learn correlational patterns rather than causal mechanisms. ChemBERTa can predict molecular properties through pattern matching but cannot simulate how modifying a functional group alters reaction pathways. AlphaFold predicts protein structures through evolutionary patterns but does not model the physical folding process.

Scientific discovery demands models that transcend pattern recognition to capture causal dynamics. A causal molecular model would not just recognize that certain molecular structures corre-

late with properties—it would explain how electron density distributions cause reactivity, and how thermodynamic gradients drive reactions. This causal understanding enables the counterfactual reasoning essential to science: predicting outcomes of novel interventions never seen in training data. This architectural choice has profound implications: foundation models scale with data and compute, but causal models scale with understanding. As we accumulate more structural data, foundation models improve at interpolation. As we refine causal mechanisms, foundation models improve at extrapolation—the essence of scientific discovery.

Physics priors While generative models like Sora create visually compelling outputs, they lack physical consistency—objects appear and disappear, gravity works intermittently, and causality is merely suggested rather than enforced. Mitchell [67] states that without biases to prefer some generalizations over others, a learning system cannot make the inductive leap necessary to classify instances beyond those it has already seen. Such inductive biases or physics priors—can be built-in to ensure generated realizations obey conservation laws, maintain object permanence, and support counterfactual reasoning about physical interactions.

Recent implementations demonstrate that world models can also discover physical laws through interaction. The joint embedding predictive architecture [3, 4] learns to predict object movements without labeled data, suggesting that the feedback loop between mental simulation and empirical observation can be implemented through self-supervised learning objectives that reward accurate forward prediction. Current world models and coceptualizations thereof, however, remain limited to relatively simple physical scenarios. While they excel at rigid body dynamics and basic occlusion reasoning, they are generally insufficient to describe complex phenomena like fluid dynamics or emergent collective behaviors. This gap between toy demonstrations and the full complexity of scientific phenomena represents the next frontier.

Active Inference AI Systems to Navigate Complex Scientific Questions Many scientific phenomena exhibit chaotic dynamics, multiscale interactions, and emergent properties that defy reductionist analysis. Climate systems, biological networks, and turbulent flows operate across scales from molecular to planetary. Traditional ML approaches that assume smooth, well-behaved functions fail catastrophically in these domains. We need architectures that can reason across scales, identify emergent patterns, and know when deterministic prediction becomes impossible. No single formal or informal computational system can accomplish these tasks, and hence we propose an AI stack. An exemplar architecture is shown in Figure 1. Some of the components of the architecture include:

1. *Base reasoning model suite with inference-tunable capabilities:* This top-layer component comprises large reasoning models that can dynamically adjust their inference strategies based on the problem context. In contrast to being optimized for next-token prediction, these models support extended thinking times, systematic exploration of solution paths, and explicit reasoning chains. The suite has the ability to recognize which mode of reasoning is appropriate. Value specifications from humans guide the reasoning process, ensuring that resources are allocated to scientifically meaningful directions rather than arbitrary pattern completion.
2. *Multi-modal domain foundation models with shared representations:* These are effectively world models that maintain causal representations of scientific domains. These models allow the system to mentally simulate interventions, test counterfactuals, and explore hypothesis spaces before committing to physical experiments. These function as oracles or world models, serving as the substrate for both pattern discovery (thinking) and rapid inference (reasoning). These domain-specific models must share embeddings that enable cross-pollination of insights.

3. *Dynamic knowledge graphs as evolving scientific memory:* Unlike static knowledge bases, these graphs function as cognitive architectures that grow through the interplay of thinking, reasoning, and experimentation. Nodes represent concepts ranging from raw observations to abstract principles, while weighted edges encode causal relationships with associated uncertainty. The graphs expand as thinking discovers new patterns (adding nodes), reasoning establishes logical connections (adding edges), and experiments validate or falsify relationships (adjusting weights). Version-controlled evolution allows the system to maintain competing hypotheses, track conceptual development, and recognize when anomalies demand fundamental restructuring rather than incremental updates. This persistent, growing memory enables genuine scientific progress rather than mere information retrieval.
4. *Reality tethering through verification layers:* The verification layer partitions scientific claims into formally provable statements and empirically testable hypotheses. Mathematical derivations, algorithmic properties, and logical arguments can be decomposed into proof obligations for interactive theorem provers (Lean [68], Coq [8]), creating a growing corpus of machine-verified knowledge that future reasoning can build upon. For claims beyond formal correctness—predictions about physical phenomena, chemical reactions, or biological behaviors—the system generates targeted computational simulations and experimental protocols. This dual approach acknowledges that scientific knowledge spans from mathematical certainty to empirical contingency. Crucially, failed verifications become learning opportunities, updating the system’s confidence bounds and identifying gaps between its world model and reality. When formal and empirical verification diverge for instance, when a mathematically valid prediction fails experimentally, the system must flag this inconsistency for human review, as the divergence may indicate modeling assumptions that require revision rather than simple parameter updates.
5. *Human-steerable orchestration:* Humans excel at recognizing meaningful patterns and making creative leaps; AI can perform exhaustive search and maintaining consistency across vast knowledge spaces; Well-understood computational science tools (e.g. optimal experimental design) can execute efficient agentic actions in a reliable manner. This symbiotic relationship ensures that the system’s powerful reasoning capabilities remain tethered to meaningful scientific questions, and existing algorithms are efficiently leveraged. We note that optimal experimental design methods can be inadequate when conditioned on misspecified models [89]; human oversight is therefore essential not just for value alignment but for detecting when the system’s foundational assumptions require revision.
6. *Proactive exploration engines:* Rather than passively responding to queries (the primary mode in which language models are used currently), these systems work persistently in the background to generate hypotheses, identify gaps in knowledge, and propose experiments. Driven by uncertainty quantification and novelty detection algorithms, these engines can maintain a priority queue of open questions ranked by their potential to achieve specified goals versus resource requirements. This layer enables the system to operate across multiple time horizons—pursuing rapid experiments vs long-term research campaigns that systematically map uncharted territories in the knowledge space. The ranking of scientific goals reflects value judgments that must ultimately derive from human priorities; the system provides estimates of expected information gain and resource costs, but the weighting between scientific impact, feasibility, and broader societal considerations remains a human responsibility. The tension between uncertainty-driven exploration and novelty-driven exploration remains an open challenge [32]; current approaches typically require domain-specific tuning rather than universal

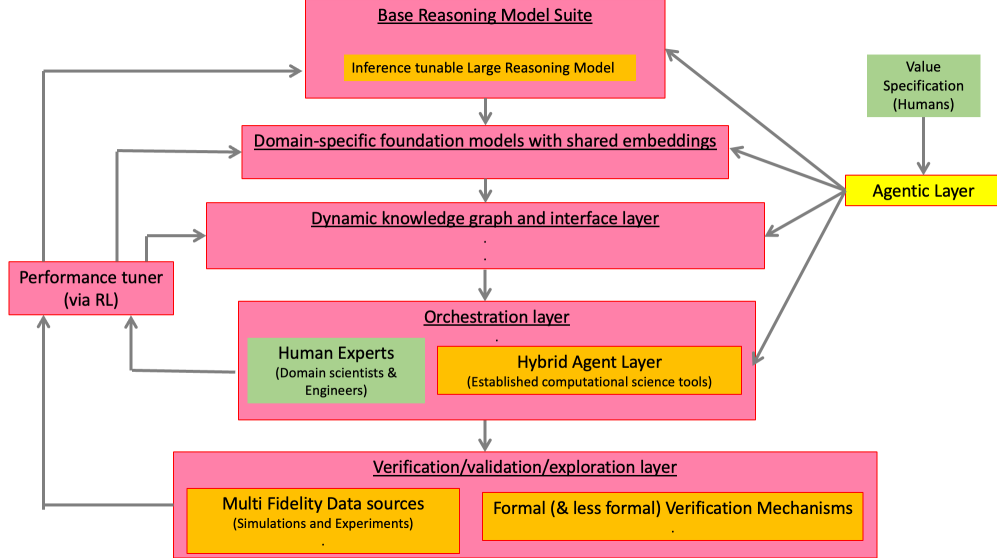


Figure 1: Exemplar architecture of an Active Inference AI system for scientific discovery.

solutions.

Component Interactions Figure 1 illustrates the flow of information between components. The base reasoning model receives queries from human users and orchestrates calls to domain foundation models, which return predictions with uncertainty estimates. These predictions update the dynamic knowledge graph, which in turn informs the verification layer’s decisions about whether claims require formal proof, computational simulation, or physical experimentation. The performance tuner (via reinforcement learning) adjusts exploration-exploitation tradeoffs based on reward signals derived from successful predictions, experimental confirmations, and human feedback on the utility of discoveries. Concretely: (i) reward signals for the RL tuner originate from prediction accuracy on held-out data, confirmation rates in downstream experiments, and explicit human ratings of hypothesis quality; (ii) the time horizons are domain-dependent. The architectural principles outlined here are intended as design guidance rather than a complete implementation specification; realizing them in practice will require substantial engineering and domain-specific adaptation.

The architectural principles outlined above find grounding in recent work on transformational scientific creativity. For instance, Schapiro et al. [82] formalize scientific conceptual spaces as directed acyclic graphs, where vertices represent generative rules and edges capture logical dependencies. This offers a concrete implementation pathway for the proposed dynamic knowledge graphs. Their distinction between modifying existing constraints versus fundamentally restructuring the space itself maps directly onto our architecture’s dual modes of reasoning (traversing established knowledge) and thinking (discovering new patterns that may violate existing assumptions). This convergence suggests that achieving transformational scientific discovery through AI systems requires systems capable of identifying and modifying the foundational axioms that constrain current scientific understanding—a capability the active inference framework aims to provide through its stacked architecture and integration of models, empirical feedback, and human guidance.

It is acknowledged that while the AI system can, in principle, be operated autonomously through well-defined interfaces between components, human interaction and decisions can be expected to play

a key role. The architectural principles outlined above find partial instantiation in contemporary systems, though none fully realize the complete vision of scientific intelligence. Appendix A examines some current implementations through the lens of our three-gap framework, and discusses both substantial progress and persistent limitations that illuminate the path forward. Appendix B gives high-level mathematical constructs for key components of the above system.

Challenges of Iterative Learning and Importance of Human Interactions While the aforementioned architecture presents a compelling vision of AI systems that learn from real-world interaction, incorporating feedback into iterative training poses fundamental challenges that cannot be overlooked. Scientific experiments produce sparse, noisy, and often contradictory signals. A single failed synthesis might stem from equipment miscalibration, modeling errors, or genuine chemical impossibility—yet the system must learn appropriately from each case. The tension between generalization and specificity becomes acute: overfitting to particular configurations may yield brittle models that fail to transfer across laboratories, while excessive generalization may miss critical context-dependent phenomena.

This inherent ambiguity in processing experimental feedback into actionable model refinements makes human judgment indispensable, not as a temporary scaffold but as a permanent architectural component. Thus, the challenge lies not merely in designing systems that can incorporate feedback, but in creating architectures that handle the full spectrum of empirical reality, including clear confirmations, ambiguous results, systematic biases and truly novel results. Effective human-AI collaboration must therefore go beyond simple oversightempirical studies demonstrate that human-AI teams can outperform either alone, but only when humans maintain genuine understanding rather than deferring uncritically to AI suggestions [5]. This partnership becomes especially critical when experiments and computations challenge fundamental assumptions.

5 Outlook

Active AI systems encompass external experience (empirical data) and internal experience (mental simulation). AI systems that can fluidly navigate between these modes will mark the transition from tools that find patterns to partners that discover principles. This perspective builds upon substantial progress in causal machine learning, active learning, and automated scientific discovery while addressing critical gaps. The causal machine learning community has made significant strides in developing methods for causal inference from observational data, with frameworks like Pearl’s causal hierarchy and recent advances in causal representation learning providing mathematical foundations for understanding interventions and counterfactuals [86, 52]. Similarly, active learning has evolved sophisticated strategies for optimal experimental design [88, 35], while automated discovery systems have demonstrated success in specific domains such as materials science and drug discovery [15, 91]. However, these communities have largely operated in isolation, with causal methods focusing primarily on statistical inference rather than physical mechanism discovery, active learning optimizing for narrow uncertainty reduction rather than conceptual breakthrough, and automated discovery systems excelling at interpolation within known spaces rather than extrapolation to genuinely novel phenomena. We note that efforts to bridge these communities are underway: the Acceleration Consortium [95] represents a major initiative integrating autonomous laboratories with AI-driven discovery, while recent work demonstrates promising integration of causal reasoning with materials optimization [44, 31, 16].

Current implementations prioritize task completion over understanding, optimization over exploration, and correlation over causation. The path forward requires AI systems that integrate causal reasoning not merely as a statistical tool but as the foundation for mental simulation and counterfactual experimentation, extending active learning beyond data efficiency to include the gen-

eration and testing of novel hypotheses that violate existing assumptions, and grounding automated discovery in continuous empirical feedback loops that prevent drift from physical reality. Most critically, while existing approaches excel within their prescribed domains, they lack the architectural foundation for the kind of open-ended, cross-domain reasoning that characterizes human scientific discovery—the ability to recognize when anomalous observations demand not just parameter updates but fundamental reconceptualization of the problem space itself.

Scientific discovery has always been a collaborative enterprise—across disciplines, institutions, and generations. AI systems represent new kinds of collaborative tools. The transition from static models to living systems marks a fundamental shift in how we conceive of AI systems that persist, that remember, that build intuition through repeated engagement with reality. Just as human scientists develop insight through years of experimentation, future AI systems will accumulate wisdom through continuous cycles of hypothesis, experiment, and revision. We thus call for the creation of new benchmarks and research programs centered around the proposed stacked architecture, moving evaluation beyond static datasets to interactive, discovery-oriented environments. Success should ultimately be judged not only by benchmark performance but by domain expert assessment of discovery quality. For example, whether AI-generated hypotheses lead to experiments that domain scientists find surprising and informative, whether the system identifies gaps in current understanding that experts recognize as important, and whether it proposes experimental designs that are both feasible and scientifically meaningful. Structured evaluation protocols, such as blinded assessment of AI-proposed versus human-proposed hypotheses by independent expert panels, could provide more rigorous qualitative evaluation [69].

Finally, it has to be emphasized that modern AI systems are already useful in their present form, and are being utilized effectively by scientific research groups across the world. However, even with future improvements, these tools bring many systemic hazards [66]: *a) false positives and false negatives*: spurious correlations can be mistaken for laws, while cautious priors may hide real effects, and thus rigorous uncertainty metrics and adversarial falsification must be built in; *b) epistemic overconfidence*: large models can exhibit poorly calibrated uncertainty estimates, particularly on inputs far from the training distribution, a phenomenon where predictive confidence remains high even when accuracy degrades [40]; ensemble disagreement among multiple independently trained models provides one diagnostic for such failures; *c) erosion of insight and rigor*: over time, there is significant risk of researchers losing key scientific skills; *d) Cost*: simulation-driven exploration can consume resources long after the expected reduction in uncertainty from additional experiments becomes negligible relative to their cost (i.e., marginal information gain saturates); resource allocation algorithms (schedulers) must weigh value against resources; *e) instrumental drift*: scientific instruments and sensors evolve over time, and their characteristics may shift due to wear, recalibration, or environmental changes; this can be addressed through standard calibration practices such as round-robin testing across laboratories [108], but without continual residual checks and rapid retraining, predictions may silently bias. These issues have to be continually acknowledged, recognizing and safeguards should be embedded into the scientific process.

Acknowledgment

This piece has benefitted directly or indirectly from many discussions with Jason Pruet (OpenAI), Venkat Raman, Venkat Viswanathan, Alex Gorodetsky (U. Michigan), Rick Stevens (Argonne National Laboratory/U. Chicago), Earl Lawrence (Los Alamos National Laboratory) and Brian Spears (Lawrence Livermore National Laboratory). This work was partly supported by Los Alamos National Laboratory under the grant #AWD026741 at the University of Michigan.

References

- [1] Kamal Acharya, Waleed Raza, Carlos Dourado, Alvaro Velasquez, and Houbing Herbert Song. Neurosymbolic reinforcement learning and planning: A survey. *IEEE Transactions on Artificial Intelligence*, 5(5), 2023.
- [2] Frank Arute, Kunal Arya, Ryan Babbush, et al. Quantum supremacy using a programmable superconducting processor. *Nature*, 574(7779), 2019.
- [3] Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15619–15629, 2023.
- [4] Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zhohus, et al. V-jepa 2: Self-supervised video models enable understanding, prediction and planning. *arXiv preprint arXiv:2506.09985*, 2025.
- [5] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. Does the whole exceed its parts? the effect of ai explanations on complementary team performance. pages 1–16, 2021.
- [6] Elias Bareinboim and Judea Pearl. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352, 2016.
- [7] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- [8] Yves Bertot and Pierre Castéran. *Interactive theorem proving and program development: Coq’Art: the calculus of inductive constructions*. Springer Science & Business Media, 2013.
- [9] Celeste Biever. Ai scientist ‘team’ joins the search for extraterrestrial life. *Nature*, 641(8063):568–569, 2025.
- [10] Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578, 2023.
- [11] Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. Augmenting large language models with chemistry tools. *Nature Machine Intelligence*, pages 1–11, 2024.
- [12] Markus J Buehler. Agentic deep graph reasoning yields self-organizing knowledge networks. *arXiv preprint arXiv:2502.13025*, 2025.
- [13] Markus J Buehler. In situ graph reasoning and knowledge expansion using graph-preflexor. *Advanced Intelligent Discovery*, 2025.
- [14] Benjamin Burger, Phillip M Maffettone, Vladimir V Gusev, Catherine M Aitchison, Yang Bai, Xiaoyan Wang, Xiaobo Li, Ben M Alber, Andrea Cognigni, Gavin Sherborne, et al. A mobile robotic chemist. *Nature*, 583(7815):237–241, 2020.
- [15] Keith T Butler, Daniel W Davies, Hugh Cartwright, Olexandr Isayev, and Aron Walsh. Machine learning for molecular and materials science. *Nature*, 559(7715):547–555, 2018.
- [16] Keith T Butler, Daniel W Davies, Hugh Cartwright, Olexandr Isayev, and Aron Walsh. Machine learning for molecular and materials science. *Nature*, 559(7715):547–555, 2018.
- [17] Ingrid Campo-Ruiz. Artificial intelligence may affect diversity: architecture and cultural context reflected through chatgpt, midjourney, and google maps. *Humanities and Social Sciences Communications*, 12(1):1–13, 2025.
- [18] Nancy Cartwright. Causation: One word, many things. *Philosophy of Science*, 71(5):805–819, 2004.

- [19] Krzysztof Chalupka, Tobias Bischoff, Pietro Perona, and Frederick Eberhardt. Unsupervised discovery of el nino using causal feature learning on microlevel climate data. In *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*, pages 72–81, 2016.
- [20] Krzysztof Chalupka, Pietro Perona, and Frederick Eberhardt. Visual causal feature learning. *arXiv preprint arXiv:1412.2309*, 2015.
- [21] Matthew Chang, Gunjan Chhablani, Alexander Clegg, Mikael Dallaire Cote, Ruta Desai, Michal Hlavac, Vladimir Karashchuk, Jacob Krantz, Roozbeh Mottaghi, Priyam Parashar, et al. Partnr: A benchmark for planning and reasoning in embodied multi-agent tasks. *arXiv preprint arXiv:2411.00081*, 2024.
- [22] Yuan Chiang, Elvis Hsieh, Chia-Hong Chou, and Janosh Riebesell. Llamp: Large language model made powerful for high-fidelity materials knowledge retrieval and distillation. *arXiv preprint arXiv:2401.17244*, 2024.
- [23] Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. Chemberta: Large-scale self-supervised pretraining for molecular property prediction. *CoRR*, 2020.
- [24] François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
- [25] Jae-Won Chung, Yile Gu, Insu Jang, Luoxi Meng, Nikhil Bansal, and Mosharaf Chowdhury. Reducing energy bloat in large model training. In *Proceedings of the ACM SIGOPS 30th Symposium on Operating Systems Principles*, pages 144–159, 2024.
- [26] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [27] Jerome Cornfield, William Haenszel, E Cuyler Hammond, Abraham M Lilienfeld, Michael B Shimkin, and Ernst L Wynder. Smoking and lung cancer: Recent evidence and a discussion of some questions. *Journal of the National Cancer Institute*, 22(1):173–203, 1959.
- [28] Kourosh Darvish, Marta Skreta, Yuchi Zhao, Naruki Yoshikawa, Sagnik Som, Miroslav Bogdanovic, Yang Cao, Han Hao, Haoping Xu, Alan Aspuru-Guzik, et al. Organa: A robotic assistant for automated chemistry experimentation and characterization. *arXiv preprint arXiv:2401.06949*, 2024.
- [29] Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishnamurthy, Alex Smola, and Andrew McCallum. Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning. In *International Conference on Learning Representations*, 2018.
- [30] Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham Chinya, Yongqiang Cao, Sri Harsha Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, et al. Loihi: A neuromorphic manycore processor with on-chip learning. *Ieee Micro*, 38(1):82–99, 2018.
- [31] Juan J de Pablo, Nicholas E Jackson, Michael A Webb, Long-Qing Chen, Jeffrey E Moore, Dane Morgan, Ryan Jacobs, Tresa Pollock, Darrell G Schlom, Eric S Toberer, et al. New frontiers for the materials genome initiative. *npj Computational Materials*, 5(1):41, 2019.
- [32] Marina Dubova, Arseny Moskvichev, and Kevin Zollman. Against theory-motivated experimentation in science. *MetaArXiv preprint*, 2022.
- [33] Event Horizon Telescope Collaboration. First m87 event horizon telescope results. i. the shadow of the supermassive black hole. *The Astrophysical Journal Letters*, 875(1):L1, 2019.
- [34] Ronald A Fisher. *The design of experiments*. Oliver and Boyd, 1935.
- [35] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *International Conference on Machine Learning*, pages 1183–1192, 2017.
- [36] Alireza Ghafarollahi and Markus J Buehler. Protagents: protein discovery via large language model multi-agent collaborations combining physics and machine learning. *Digital Discovery*, 2024.

- [37] Kurt Gödel. Über formal unentscheidbare sätze der principia mathematica und verwandter systeme i. *Monatshefte für Mathematik und Physik*, 38(1):173–198, 1931.
- [38] Mourad Gridach, Jay Nanavati, Khaldoun Zine El Abidine, Lenon Mendes, and Christina Mack. Agentic ai for scientific discovery: A survey of progress, challenges, and future directions. *International Conference on Learning Representations*, 2025.
- [39] Xuemei Gu and Mario Krenn. Interesting scientific idea generation using knowledge graphs and llms: Evaluations with 100 research group leaders. *arXiv preprint arXiv:2405.17044*, 2024.
- [40] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, pages 1321–1330, 2017.
- [41] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- [42] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pages 1861–1870, 2018.
- [43] Miguel A Hernán and James M Robins. *Causal inference: What if*. Chapman & Hall/CRC, 2020.
- [44] Hideki Hino, Takuya Obo, Genshiro Kitagawa, Koichi Niihara, and Takashi Taniguchi. Active learning in materials science with emphasis on adaptive sampling using uncertainties for targeted design. *npj Computational Materials*, 6(1):21, 2020.
- [45] John PA Ioannidis. Why most published research findings are false. *PLoS Medicine*, 2(8):e124, 2005.
- [46] Peter Jansen, Marc-Alexandre Côté, Tushar Khot, Erin Bransom, Bhavana Dalvi Mishra, Bodhisattwa Prasad Majumder, Oyvind Tafjord, and Peter Clark. Discoveryworld: A virtual environment for developing and evaluating automated scientific discovery agents. *Advances in Neural Information Processing Systems*, 37:10088–10116, 2024.
- [47] Philip Nicholas Johnson-Laird. *Mental models: Towards a cognitive science of language, inference, and consciousness*. Number 6. Harvard University Press, 1983.
- [48] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Zidek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [49] Daniel Kahneman. *Thinking, Fast and Slow*. Farrar, Straus and Giroux, 2011.
- [50] Narayanan Kasthuri, Kenneth Jeffrey Hayworth, Daniel Raimund Berger, Richard Lee Schalek, José Angel Conchello, Seymour Knowles-Barley, Dongil Lee, Amelio Vázquez-Reina, Verena Kaynig, Thouis Raymond Jones, et al. Saturated reconstruction of a volume of neocortex. *Cell*, 162(3):648–661, 2015.
- [51] Lance Kavalsky, Vinay I Hegde, Eric Muckley, Matthew S Johnson, Bryce Meredig, and Venkatasubramanian Viswanathan. By how much can closed-loop frameworks accelerate computational materials discovery? *Digital Discovery*, 2(4):1112–1125, 2023.
- [52] Nan Rosemary Ke, Philemon Langlais, Rui Shu, et al. Learning causal representations for robust domain adaptation. *arXiv preprint arXiv:2011.08857*, 2022.
- [53] Joo-Young Kim, Bongjin Kim, and Tony Tae-Hyoung Kim. Processing-in-memory for ai. 2022.
- [54] Ross D King, Jem Rowland, Stephen G Oliver, Michael Young, Wayne Aubrey, Emma Byrne, Maria Liakata, Magdalena Markham, Pinar Pir, Larisa N Soldatova, et al. The automation of science. *Science*, 324(5923):85–89, 2009.
- [55] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In *International conference on machine learning*, pages 5338–5348. PMLR, 2020.
- [56] Thomas S. Kuhn. *The Structure of Scientific Revolutions*. University of Chicago Press, 1962.

- [57] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40:e253, 2017.
- [58] Pat Langley. Integrated systems for computational scientific discovery. 38(20):22598–22606, 2024.
- [59] Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 62(1):1–62, 2022.
- [60] Hongwei Liu, Zilong Zheng, Yuxuan Qiao, Haodong Duan, Zhiwei Fei, Fengzhe Zhou, Wenwei Zhang, Songyang Zhang, Dahua Lin, and Kai Chen. Mathbench: Evaluating the theory and application proficiency of llms with a hierarchical mathematics benchmark. *arXiv preprint arXiv:2405.12209*, 2024.
- [61] Konrad Lorenz. *Behind the mirror: A search for a natural history of human knowledge*. Harcourt Brace Jovanovich, 1973.
- [62] Karttikeya Mangalam, Haoqi Fan, Yanghao Li, Chao-Yuan Wu, Bo Xiong, Christoph Feichtenhofer, and Jitendra Malik. Reversible vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [63] Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B Tenenbaum, and Jiajun Wu. The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision. In *International Conference on Learning Representations*, 2023.
- [64] Kim Martineau. How the von neumann bottleneck is impeding ai computing, 2024. IBM Research Blog, accessed 30 June 2025.
- [65] Juan Mateos-Garcia and Joel Klinger. Is there a narrowing of ai research? 2023.
- [66] Lisa Messeri and MJ Crockett. Artificial intelligence and illusions of understanding in scientific research. *Nature*, 627(8002):49–58, 2024.
- [67] Tom M Mitchell. The need for biases in learning generalizations. *CS Tech Report CBM-TR-117, Rutgers University*, 1980.
- [68] Leonardo de Moura and Sebastian Ullrich. The lean 4 theorem prover and programming language. In *Automated Deduction—CADE 28: 28th International Conference on Automated Deduction, Virtual Event, July 12–15, 2021, Proceedings 28*, pages 625–635. Springer, 2021.
- [69] Sebastian Musslick, Laura K Bartlett, Sreejan Humayun Chandramouli, Marina Dubova, Fernand Gobet, Thomas L Griffiths, Ralph Hertwig, William R Holmes, et al. Automating the practice of science: Opportunities, challenges, and implications. *Proceedings of the National Academy of Sciences*, 122(5):e2401238121, 2025.
- [70] Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work: Scratchpads for intermediate computation with language models. 2021.
- [71] Pierre-Yves Oudeyer, Frederic Kaplan, and Verena V Hafner. Intrinsic motivation systems for autonomous mental development. *IEEE Transactions on Evolutionary Computation*, 11(2), 2007.
- [72] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [73] Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *International Conference on Machine Learning*, pages 2778–2787, 2017.
- [74] Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashesh Chattopadhyay, Morteza Mardani, Thorsten Kurth, David Hall, Zongyi Li, Kamyar Azizzadenesheli, et al. Fourcastnet: A global data-driven high-resolution weather model using adaptive fourier neural operators. *arXiv preprint arXiv:2202.11214*, 2022.

- [75] Judea Pearl. *Causality*. Cambridge university press, 2009.
- [76] Roger Penrose. *The Emperor’s New Mind: Concerning Computers, Minds, and the Laws of Physics*. Oxford University Press, 1989.
- [77] Roger Penrose. *Shadows of the Mind: A Search for the Missing Science of Consciousness*. Oxford University Press, 1994.
- [78] Karl Popper. *The Logic of Scientific Discovery*. Hutchinson & Co., 1959.
- [79] Michael H Prince, Henry Chan, Aikaterini Vriza, Tao Zhou, Varuni K Sastry, Yanqi Luo, Matthew T Dearing, Ross J Harder, Rama K Vasudevan, and Mathew J Cherukara. Opportunities for retrieval and tool augmented large language models in scientific facilities. *npj Computational Materials*, 10(1):251, 2024.
- [80] Chuan Qin, Xin Chen, Chengrui Wang, Pengmin Wu, Xi Chen, Yihang Cheng, Jingyi Zhao, Meng Xiao, Xiangchao Dong, Qingqing Long, et al. Scihorizon: Benchmarking ai-for-science readiness from scientific data to large language models. *arXiv preprint arXiv:2503.13503*, 2025.
- [81] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8), 2019.
- [82] Samuel Schapiro, Jonah Black, and Lav R Varshney. Transformational creativity in science: A graphical theory. *arXiv preprint arXiv:2504.18687*, 2025.
- [83] Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and Emad Barsoum. Agent laboratory: Using llm agents as research assistants. *arXiv preprint arXiv:2501.04227*, 2025.
- [84] Jürgen Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (1990–2010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230–247, 2010.
- [85] Bernhard Schölkopf. Causality for machine learning. In *Probabilistic and causal inference: The works of Judea Pearl*, pages 765–804. 2022.
- [86] Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- [87] Kristof Schütt, Pieter-Jan Kindermans, Huziel Enoc Saucedo Felix, Stefan Chmiela, Alexandre Tkatchenko, and Klaus-Robert Müller. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. *Advances in neural information processing systems*, 30, 2017.
- [88] Burr Settles. Active learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 6(1):1–114, 2012.
- [89] Sabina J Sloman, Daniel M Oppenheimer, Stephen B Broomell, and Cosma R Shalizi. Robustness of bayesian adaptive experimental designs. *arXiv preprint arXiv:2205.13698*, 2022.
- [90] Richard Sutton. The bitter lesson. *Incomplete Ideas (blog)*, 13(1):38, 2019.
- [91] Kyle Swanson, Weiyi Wu, Nathaniel L Bulaong, John E Pak, and James Zou. The virtual lab of ai agents designs new sars-cov-2 nanobodies. *Nature*, pages 1–3, 2025.
- [92] Nathan J Szymanski, Bernardus Rendy, Yuxing Fei, Rishi E Kumar, Tanjin He, David Milsted, Matthew J McDermott, Max Gallber, Xiang Zeng, Junhui Jia, et al. An autonomous laboratory for the accelerated synthesis of novel materials. *Nature*, 624(7990):86–91, 2023.
- [93] Nathaniel Thomas, Tess Smidt, Steven Kearnes, Lusann Yang, Li Li, Kai Kohlhoff, and Patrick Riley. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv preprint arXiv:1802.08219*, 2018.
- [94] Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482, 2024.

- [95] University of Toronto. The acceleration consortium. <https://acceleration.utoronto.ca>, 2024.
- [96] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, 2017.
- [97] Peihao Wang, Rameswar Panda, Lucas Torroba Hennigen, Philip Greengard, Leonid Karlinsky, Rogerio Feris, David Daniel Cox, Zhangyang Wang, and Yoon Kim. Learning to grow pretrained models for efficient transformer training. In *The Eleventh International Conference on Learning Representations*.
- [98] Xuezhi Wang, Dale Wei, Yizhong Dong, Nan Bao, Michelle Yang, Denny Yu, Zijian Guo, Quoc V Le, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *Advances in Neural Information Processing Systems*, 2022.
- [99] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [100] Karen Williams, Elizabeth Bilsland, Andrew Sparkes, Wayne Aubrey, Michael Young, Larisa N Soldatova, Kurt De Grave, Jan Ramon, Michaela de Clare, Thierry Delaveau, et al. Cheaper faster drug development validated by the repositioning of drugs against neglected tropical diseases. *Journal of the Royal Society Interface*, 12(104):20141289, 2015.
- [101] Stephen Wolfram. *A New Kind of Science*. Wolfram Media, 2002.
- [102] Stephen Wolfram. Can ai solve science?, March 2024. <https://writings.stephenwolfram.com/2024/03/can-ai-solve-science/>.
- [103] Michael Wooldridge and Nicholas R Jennings. Intelligent agents: Theory and practice. *The Knowledge Engineering Review*, 10(2):115–152, 1995.
- [104] Yihang Xiao, Jinyi Liu, Yan Zheng, Xiaohan Xie, Jianye Hao, Mingzhi Li, Ruitao Wang, Fei Ni, Yuxiao Li, Jintian Luo, et al. Cellagent: An llm-driven multi-agent framework for automated single-cell data analysis. *bioRxiv*, pages 2024–05, 2024.
- [105] Qi Xin, Quyu Kong, Hongyi Ji, Yue Shen, Yuqi Liu, Yan Sun, Zhilin Zhang, Zhaorong Li, Xunlong Xia, Bing Deng, et al. Bioinformatics agent (bia): Unleashing the power of large language models to reshape bioinformatics workflow. *bioRxiv*, pages 2024–05, 2024.
- [106] Hector Zenil et al. The future of fundamental science led by generative closed-loop artificial intelligence. *arXiv preprint arXiv:2307.07522*, 2023.
- [107] Guoqiang Zhang, JP Lewis, and W Bastiaan Kleijn. On exact bit-level reversible transformers without changing architectures. *arXiv preprint arXiv:2407.09093*, 2024.
- [108] ZwickRoell. Round robin tests. <https://www.zwickroell.com/services/round-robin-tests/>, 2024.

Appendix A: Current Implementations of Agentic Systems

A comprehensive review of agentic systems for scientific discovery can be found in Ref. [38]; see also Zenil et al. [106] for a complementary perspective. Below, a few references that are related to abstraction, reasoning and reality gaps are provided.

Before the transformer era, several systems demonstrated that AI could make genuine scientific discoveries through closed-loop interaction with physical laboratories. Adam and Eve [54, 100] combined ontological reasoning with active learning to autonomously design and execute falsifiable experiments. More recently, the Robot Chemist [14] demonstrated autonomous exploration of reaction conditions, while A-Lab [92] achieved autonomous synthesis of novel materials. These systems embody key principles of our proposed architecture: persistent scientific memory (through ontologies and databases), closed-loop empirical validation, and hypothesis-driven experimentation.

However, they operate within well-defined search spaces with predetermined objectives, limiting their capacity for the kind of open-ended discovery that might restructure the problem space itself.

Recent systems demonstrate varying degrees of success in elevating from statistical patterns to scientific abstractions. ChemCrow [11] integrates eighteen expert-designed tools to bridge token-level operations with chemical reasoning, enabling tasks such as reaction prediction and molecular property analysis. ProtAgents [36] employs reinforcement learning to navigate the conceptual space of protein design, moving beyond sequence statistics to optimize for biochemical properties. Agent Laboratory’s [83] achieves high success rates in data preparation and experimentation phases while exhibiting notable failures during literature review.

The reasoning gap manifests most clearly in limited capacity for genuine causal inference. Co-scientist [10] represents the current frontier, successfully designing and optimizing cross-coupling reactions through iterative experimentation, though its reasoning remains fundamentally correlational. LLaMP [22] attempts to address this limitation by grounding material property predictions in atomistic simulations, effectively implementing a preliminary form of mental experimentation. These systems, while promising, cannot yet perform the counterfactual reasoning that distinguishes scientific understanding from mere pattern matching.

The reality gap presents both tangible progress and stark limitations. Systems such as Organa [28] demonstrate sophisticated integration with laboratory robotics, automating complex experimental protocols in electrochemistry and materials characterization. CALMS [79] extends this integration by providing context-aware assistance during experimental execution. However, these implementations reveal brittleness: when experimental outcomes deviate from expected patterns, current systems lack the adaptive capacity to reformulate hypotheses or recognize when their fundamental assumptions require revision.

Multi-agent architectures such as BioInformatics Agent [105] and CellAgent [104] represent attempts to address these limitations through specialized collaboration, with distinct agents handling data retrieval, analysis, and validation. While these systems demonstrate improved performance on well-structured tasks, they do not yet perform the open-ended exploration that characterizes genuine discovery. The coordination overhead and brittleness of inter-agent communication often negate the benefits of specialization when confronting novel phenomena.

These implementations and others are already accelerating science, but also collectively reveal a critical insight: current systems excel at automating well-defined scientific workflows but falter when required to navigate the uncertain terrain of genuine discovery. They can execute sophisticated experimental protocols, analyze complex datasets, and even generate plausible hypotheses, yet they lack the metacognitive capabilities to recognize when they are operating beyond their training domains. The contribution of this perspective is to identify the architectural requirements for moving beyond these limitations, not to claim that closed-loop discovery is impossible, but to argue that scaling to genuinely novel discovery requires integrated resolution of the abstraction, reasoning, and reality gaps.