

ORIGINAL RESEARCH

Open access

Algorithmic Curiosity as Scientific Method: A Conceptual Proposal for Exploration-Driven Materials AI

Thomas Andersen^{1*}, Lars Nielsen¹, Mette Sørensen²

Abstract

In the rapidly evolving domain of artificial intelligence for materials science, exploitation-driven approaches that prioritize the optimization of predefined objective functions—such as targeted properties, minimal prediction error, or high accuracy—have come to dominate the field, systematically constraining exploration to regions of chemical space that are already anticipated to yield incremental gains while leaving vast swaths of potentially transformative materials undiscovered. Yet, the history of scientific progress, from the serendipitous isolation of novel compounds to paradigm-shifting insights into structure-property relationships, demonstrates that genuine discovery is propelled not solely by goal-directed pursuit but by a deeper, intrinsic drive toward the unknown, a principle that has begun to find formal expression in artificial intelligence through concepts of intrinsic motivation, novelty-seeking, and curiosity-driven algorithms. This paper proposes algorithmic curiosity as a foundational design principle for materials AI: an AI system architecture that actively seeks novelty, uncertainty, surprising patterns, and unexplored regions of materials space, irrespective of immediate utility or alignment with pre-specified rewards. It articulates five core components—novelty detection, uncertainty seeking, surprise maximization, coverage maximization, and prediction error seeking—alongside four operational principles that translate this philosophical shift into practical system behavior. Finally, the proposal delineates a phased implementation path for the broader materials AI community, offering a blueprint that promises to rebalance the exploration-exploitation trade-off and restore the spirit of open-ended scientific inquiry at the heart of materials discovery. By elevating curiosity from a peripheral heuristic to a central architectural imperative, this framework aims to unlock scientific breakthroughs that goal-optimized systems are structurally incapable of anticipating.

Keywords Algorithmic curiosity, Exploration-driven materials AI, Intrinsic motivation, Novelty search, Surprise maximization, Curiosity-driven scientific method

*Correspondence:

Thomas Andersen

thomas.andersen@gmail.com

¹ Department of Materials Data Science and AI, University of Copenhagen, Copenhagen, Denmark

² Department of Smart Engineering Materials, Technical University of Denmark, Lyngby, Denmark

Introduction

The contemporary landscape of artificial intelligence applied to materials science is characterized by an overwhelming emphasis on exploitation—on the relentless refinement of models and searches toward clearly defined, high-utility objectives such as superior mechanical strength, optimal electronic band gaps, or minimal lattice thermal conductivity [1-3]. In this paradigm, every computational or

experimental step is evaluated against an objective function that quantifies “success” in terms of proximity to a pre-chosen target, rendering the AI system a sophisticated optimizer rather than an autonomous explorer. Materials that deviate from these targets, even if they exhibit entirely unanticipated and scientifically profound behaviors, are routinely deprioritized or discarded as outliers. This exploitation dominance, while undeniably productive for incremental engineering advances, imposes a structural

blind spot: the AI never ventures far from what it already “knows” or expects to be valuable, thereby systematically excluding the unexpected, the surprising, and the genuinely novel.

The present paper contends that this imbalance is not merely a technical limitation but a conceptual one, rooted in an epistemology that equates scientific value with immediate utility. In contrast, this work introduces “algorithmic curiosity” as a countervailing design principle for materials AI—one that reframes the system’s primary drive as the active pursuit of novelty, uncertainty, and surprise within the vast, high-dimensional chemical space of possible materials. Algorithmic curiosity does not abolish goal-directed optimization; rather, it complements and, at times, deliberately overrides it by allocating computational and experimental resources to regions where the system’s own expectations are most likely to be violated [4-9].

This proposal is timely. Recent surveys of machine learning in molecular and materials science underscore how dominant paradigms—ranging from supervised property prediction to generative inverse design—remain tethered to reward-maximizing frameworks [2, 5]. Yet parallel developments in reinforcement learning and autonomous agents have demonstrated that intrinsic motivation mechanisms can yield more robust and creative behaviors than pure extrinsic reward signals [7, 8]. By synthesizing these insights with the unique challenges of materials discovery—principally the combinatorial explosion of chemical space and the epistemic opacity of structure-property landscapes—this blueprint articulates a new architectural philosophy. The sections that follow first document the extent of exploitation dominance and its epistemic costs, then examine the role of curiosity in both human science and contemporary AI, before presenting the formal proposal, its five constituent components, and the operational and evaluative scaffolding required for its realization. In so doing, the manuscript offers not merely an incremental tweak to existing pipelines but a foundational reorientation of how artificial intelligence can serve as a genuine partner in scientific discovery.

The Exploitation Dominance

Goal-driven optimization has become the default operating system of materials AI. Virtually every major methodological advance of the past decade is framed around the construction and maximization of explicit objective functions

that map materials descriptors directly onto desired performance metrics. Whether through supervised learning pipelines that minimize prediction error on labeled datasets, Bayesian optimization routines that seek expected improvement in a scalar property, or generative models trained to produce molecules with user-specified attributes, the underlying logic remains identical: define the target, quantify deviation from it, and let gradient-based or sampling-based search converge toward it [2].

This exploitation-centric worldview is amply documented across the literature. Foundational reviews such as that by Butler *et al.* [2] emphasize how models are routinely trained to predict and optimize specific observables such as formation energies, band gaps, or ionic conductivities. Recent advances in solid-state materials science, as surveyed by Schmidt *et al.* [3], similarly highlight applications where AI accelerates the identification of compounds that satisfy narrow design criteria, often within well-characterized families of perovskites, oxides, or metal-organic frameworks. Inverse design strategies, articulated by Zunger [4], explicitly frame the discovery process as a search for materials that realize target functionalities, treating chemical space as a landscape to be navigated toward known optima rather than as a territory to be charted for its own sake. Data-driven continuous representations of molecules, pioneered by Gómez-Bombarelli *et al.* [5], further exemplify the pattern: variational autoencoders and generative adversarial networks are conditioned on desired properties, ensuring that sampled candidates lie close to regions already deemed promising by the objective function.

The same logic permeates active-learning and autonomous-experimentation workflows [6-15]. Generative active learning frameworks across polymer spaces, for instance, as developed by Jiang and Webb [10], iteratively propose candidates expected to deliver targeted rheological behaviors, while dynamic Bayesian optimizers in human-in-the-loop systems, described by Biswas *et al.* [11], prioritize experiments that promise maximal information gain with respect to a predefined property landscape. Active learning in scanning tunneling microscopy, reported by Narasimha *et al.* [12], focuses on uncovering structure-property correlations only within regions already hypothesized to be relevant, and autonomous laboratories, as demonstrated by Szymanski *et al.* [13], accelerate synthesis toward materials that satisfy pre-programmed performance thresholds. Even experimental discovery campaigns in ferroelectric materials

via active learning, conducted by Liu *et al.* [14], remain anchored to structure–property relationships that align with known application domains. Curiosity-aware molecular reinforcement learning further illustrates the same bias when intrinsic rewards are subordinated to property optimization targets [15–22].

Collectively, these approaches—spanning [2–5, 10–15] and related works—illustrate a field that has internalized exploitation as its epistemic default. Objective functions define what counts as “good”; acquisition functions in Bayesian optimization or reinforcement learning reward moves that reduce uncertainty only insofar as that reduction serves the target; and generative pipelines are rewarded exclusively for fidelity to user-specified constraints. What is missed in this regime is precisely what cannot be anticipated: materials whose properties lie outside the manifold of current expectations, compounds that reveal entirely new physical mechanisms, or structures whose scientific interest emerges only after their synthesis and characterization. The chemical space is astronomically large; yet, exploitation-driven AI traverses only narrow corridors already illuminated by prior human intuition or limited datasets [6]. The result is a self-reinforcing loop in which the “interesting” is equated with the “useful,” and the surprising is algorithmically invisible. Algorithmic curiosity is proposed precisely to break this loop.

As summarized in **Table 1**, algorithmic curiosity introduces a fundamental epistemic shift from utility-driven optimization toward discovery-oriented exploration.”

Table 1. Conceptual differentiation between exploitation-driven AI and algorithmic curiosity systems

Dimension	Exploitation-driven materials AI	Algorithmic curiosity systems
Primary objective	Optimize predefined target properties	Explore novelty, uncertainty, and surprise
Epistemic orientation	Utility-centered	Discovery-centered
Search behavior	Local optimization around known optima	Systematic traversal of unknown regions
Treatment of	Discarded or	Prioritized as

outliers	deprioritized	signals of novelty
Role of uncertainty	Reduced to improve prediction accuracy	Actively maximized as an exploration driver
Reward structure	External (task-specific objective functions)	Intrinsic (novelty, surprise, and error signals)
Exploration strategy	Constrained by acquisition functions	Open-ended and internally guided
Knowledge expansion	Incremental refinement	Structural expansion of chemical space
Failure interpretation	Error to be minimized	Signal of hidden structure
Scientific role	Optimizer of known goals	Generator of new scientific questions

Curiosity in Science and AI

Curiosity has long served as the animating force of human scientific practice. From the accidental observation of strange optical phenomena that led to quantum mechanics to the systematic exploration of unknown territories in the periodic table, breakthroughs frequently arise not from the pursuit of a narrowly defined goal but from an intrinsic delight in the novel, the anomalous, and the unexplained. Human scientists allocate time and resources to phenomena that violate expectations, even when immediate applications are absent, because such violations promise deeper insight into the underlying structure of nature. This epistemic virtue—prioritizing the surprising over the merely confirmatory—has been formalized in philosophy of science as the driver of paradigm shifts and in cognitive science as the mechanism that sustains long-term exploratory behavior [9].

Artificial intelligence research has increasingly recognized the power of analogous intrinsic motivation mechanisms. Rather than relying solely on extrinsic reward signals, certain reinforcement learning architectures derive internal rewards from the novelty of encountered states, the reduction of prediction error, or the maximization of information gain about the environment. Intrinsically motivated exploration of learned goal spaces, as formalized by Laversanne-Finot *et al.* [7], demonstrates how agents

can autonomously generate their own learning curricula, discovering competencies far beyond those encoded in any hand-crafted objective. Deep intrinsically motivated exploration in continuous control, advanced by Saglam and Kozat [8], further shows that novelty-seeking policies can outperform purely reward-maximizing baselines in complex, high-dimensional domains by avoiding premature convergence to local optima. At a foundational level, curiosity-driven exploration has been analyzed through the lenses of neuroscience and computational modeling by Modirshanechi *et al.* [9], revealing that dopamine-modulated surprise signals and uncertainty-seeking behaviors constitute a biologically plausible substrate for open-ended learning.

Within materials-specific contexts, early applications of intrinsic rewards to molecular reinforcement learning, explored by Thiede *et al.* [15], illustrate how curiosity can guide traversal of chemical space toward regions where model confidence is low or where structural motifs diverge sharply from training data. Parallel calls for curiosity-creativity elements in human-AI materials discovery, voiced by Ozin *et al.* [6], reinforce the same theme. These developments collectively suggest that algorithmic curiosity is not an alien imposition upon materials AI but an extension of principles already proving effective in adjacent domains. When transposed to materials discovery, such mechanisms promise to counteract the epistemic narrowing imposed by exploitation dominance [2, 3]. Instead of repeatedly refining predictions within well-characterized subspaces, a curiosity-driven system would actively seek out materials whose properties or structures challenge its current world model, thereby generating the very data that fuel genuine scientific advance.

The epistemological payoff is profound: curiosity reframes discovery as an open-ended dialogue between model and reality rather than a unidirectional optimization toward human-specified targets. By institutionalizing surprise maximization and novelty detection as first-class objectives, materials AI can recover the creative, serendipitous character that has historically defined scientific progress. The proposal developed in the following sections, therefore, seeks not to supplant existing tools but to embed curiosity as a core architectural primitive, ensuring that exploration is pursued for its own sake alongside—and at times in deliberate tension with—exploitation.

The Proposal

This paper advances algorithmic curiosity as a new design principle for materials AI. Definition 1: Algorithmic curiosity is an AI system design principle in which the system actively seeks out novelty, uncertainty, surprising patterns, and unexplored regions of material space, regardless of immediate utility or alignment with pre-specified reward functions. The goal is to discover what is unexpected rather than merely to optimize what is already expected.

Algorithmic curiosity, therefore, represents a deliberate philosophical and architectural shift. Where conventional materials AI treats chemical space as an optimization landscape whose contours are defined by human-specified objectives [4, 5], a curiosity-driven system treats the same space as an open territory whose intrinsic structure—its pockets of surprise, its gradients of uncertainty, and its latent regularities—becomes the primary object of inquiry. The system is rewarded internally for encountering configurations that deviate maximally from its current predictive model, for visiting underrepresented subspaces, and for generating data that force model revision. Utility remains relevant but is demoted from sole arbiter to one among several competing drives [6].

This proposal is distinct from existing paradigms in three critical respects. First, it is not equivalent to active learning, which, even when framed around uncertainty sampling, ultimately serves the goal of improving performance on a downstream task [10, 11]; algorithmic curiosity values uncertainty for its own sake, independent of any eventual property prediction. Second, it diverges from Bayesian optimization, whose acquisition functions balance exploration against exploitation only insofar as both serve a global objective function [14]; curiosity-driven systems may deliberately ignore that objective for extended periods to pursue pure novelty. Third, it is not random exploration, which lacks any internal structure or memory of past surprises; algorithmic curiosity maintains a dynamic model of its own knowledge gaps and prioritizes traversals that are likely to be maximally informative in an epistemic rather than instrumental sense [9, 15, 23-29].

Conceptually, one may envision a spectrum running from pure exploitation—where every action is scored solely by contribution to a fixed objective—to pure curiosity, where actions are scored solely by their capacity to violate expectations. Most current systems cluster near the exploitation pole [2, 3, 13]. The blueprint articulated here advocates a principled migration toward the curiosity pole, with hybrid regimes emerging naturally as intermediate

stations. By making curiosity an explicit, first-class design criterion rather than an ad-hoc add-on, materials AI can begin to operate as a true scientific instrument: one that not only answers posed questions but also poses its own.

Figure 1 illustrates the hierarchical architecture of algorithmic curiosity, highlighting its linear progression from materials space input to discovery-oriented knowledge expansion.

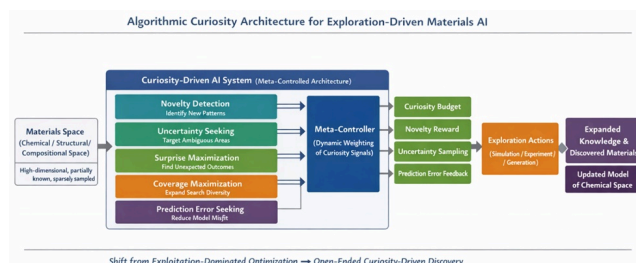


Figure 1. The hierarchical architecture of algorithmic curiosity

Components of Algorithmic Curiosity

Algorithmic curiosity is realized through five interlocking components, each operationalizing a distinct facet of the intrinsic drive toward the unknown.

Novelty detection

The system continuously monitors the similarity between newly encountered or proposed materials and the corpus of previously observed or simulated structures. Novelty is quantified not merely by Euclidean distance in descriptor space but by the degree to which a candidate lies outside the manifold spanned by the training distribution, using density estimation or reconstruction error in learned latent representations. This component ensures that the AI preferentially allocates resources to materials whose atomic arrangements, bonding motifs, or compositional ratios deviate qualitatively from historical experience [15], thereby surfacing candidates that conventional similarity-based screening would overlook.

Uncertainty seeking

Epistemic uncertainty—arising from insufficient data rather than inherent stochasticity—is actively maximized. Ensemble disagreement, dropout-based variance, or Gaussian process posterior width serve as proxies; the

system deliberately selects experiments or simulations in regions where its own confidence is lowest. Unlike uncertainty sampling in active learning [12], which terminates once a task-specific model reaches acceptable accuracy, uncertainty seeking here is open-ended: the system treats persistent ignorance as an intrinsically valuable state worth sustaining and exploring [9].

Surprise maximization

The system maintains an internal predictive model of expected outcomes (structure, spectra, properties) and seeks inputs that produce the largest divergence between prediction and reality. Surprise is operationalized as information gain or negative log-likelihood of observed data under the current model. By chasing high-surprise events, the AI forces rapid model evolution, uncovering hidden causal structures that no amount of targeted optimization would reveal [8].

Coverage maximization

Beyond local novelty, the system maintains a global map of explored versus unexplored regions of materials space—whether through adaptive tessellation, kernel density estimation, or graph-based connectivity measures—and prioritizes moves that expand the frontier of known territory. This component prevents the AI from cycling within already-familiar clusters and ensures broad, systematic coverage even in the absence of any external performance metric [7].

Prediction error seeking

The system deliberately probes inputs where prior prediction errors were largest, treating those errors not as failures to be minimized but as signals of latent structure worthy of deeper investigation. Error gradients become exploratory beacons, guiding the AI toward the very phenomena that expose the incompleteness of its current theory [6].

Table 2 consolidates the five components of algorithmic curiosity into a unified functional framework linking computational mechanisms to scientific discovery outcomes.

Table 2. Integrated framework of algorithmic curiosity: components, mechanisms, and scientific functions

Component	Operational	System	Scientific
-----------	-------------	--------	------------

	mechanism	function	contribution
Novelty detection	Density estimation/latent space distance	Identify out-of-distribution materials	Surfaces of structurally unseen compounds
Uncertainty seeking	Ensemble variance/Bayesian posterior width	Target epistemic gaps	Expands knowledge in poorly understood regions
Surprise maximization	Prediction–observation divergence (information gain)	Trigger model revision	Reveals hidden causal relationships
Coverage maximization	Space mapping / adaptive tessellation	Ensure global exploration	Prevents local trapping in known clusters
Prediction error seeking	Error-driven sampling policies	Focus on model failure regions	Converts failure into discovery signal

These five components do not operate in isolation; they interact dynamically through a meta-controller that weights their relative influence according to the current state of knowledge. Together they constitute a self-sustaining loop in which novelty begets uncertainty, uncertainty begets surprise, surprise begets coverage, and coverage begets new prediction errors—continuously propelling the system into uncharted regions of chemical space. When instantiated within materials AI architectures, they promise to transform discovery from a goal-constrained search into an open-ended scientific conversation with matter itself.

Operational Principles

To translate algorithmic curiosity from conceptual blueprint into functional architecture, four operational principles provide the concrete mechanisms by which a materials AI system can sustain intrinsic motivation without collapsing back into exploitation dominance. These principles are not optional heuristics but mandatory design constraints that govern resource allocation, reward shaping, experiment selection, and model updating.

Curiosity budget

Rather than permitting the system to expend its entire computational or experimental capacity on objective-function maximization, a fixed fraction—typically between 20% and 40%—of the total budget must be ring-fenced exclusively for curiosity-driven actions. This allocation is enforced at the meta-controller level and cannot be overridden by any downstream reward signal. As Butler *et al.* [2] have shown in their comprehensive survey of machine learning for molecular and materials science, current pipelines devote nearly 100% of resources to exploitation; the curiosity budget therefore acts as a structural safeguard, ensuring that novelty-seeking is not merely an afterthought but a guaranteed portion of every training or discovery cycle. Over time, the budget ratio itself can be adapted based on the rate of surprising discoveries, yet the principle demands that it never drops to zero.

Novelty reward

The system receives an internal scalar reward proportional to the degree of novelty detected in any newly evaluated material. This reward is computed independently of any external property target and is added directly to the loss landscape or policy gradient. Drawing on the intrinsic motivation frameworks reviewed by Laversanne-Finot *et al.* [7] and Saglam and Kozat [8], the novelty reward is derived from reconstruction error in a variational autoencoder or from negative log-density under a normalizing flow trained on the historical dataset. Critically, this signal is allowed to compete with—and occasionally dominate—any exploitation-oriented reward, creating deliberate tension that prevents premature convergence. In practice, the principle ensures that a material exhibiting an entirely new bonding motif or an unanticipated electronic structure receives a substantial internal bonus even if its targeted property lies far from current optima.

Uncertainty sampling

At each iteration, the next candidate for simulation or synthesis is selected not solely by expected improvement on a user-defined objective but by a composite score that includes the highest epistemic uncertainty across the ensemble of models. This principle operationalizes Component 2 (Uncertainty Seeking) by maintaining a continuously updated Gaussian process or dropout ensemble whose variance map directly informs acquisition. Modirshanechi *et al.* [9] demonstrate in their neuroscience-

informed modeling that such uncertainty-seeking behavior mirrors biological curiosity and leads to faster coverage of state spaces; in materials AI, the same mechanism forces the system to probe regions where its own predictions are least reliable, thereby generating precisely the data that expose hidden structure-property relationships.

Prediction error feedback

Every discrepancy between predicted and observed outcomes is logged and used to generate a secondary curiosity signal that biases future exploration toward similar discrepancies. This principle closes the loop between Component 5 (Prediction Error Seeking) and the global exploration policy: large errors are not treated as training failures to be minimized but as beacons that illuminate promising subspaces. Thiede *et al.* [15] illustrate how prediction-error-driven intrinsic rewards can steer molecular reinforcement learning away from well-explored islands; the present principle elevates that insight to a first-class operational rule, ensuring that the system actively returns to and deepens investigation of its own predictive failures.

Taken together, these four principles create a self-regulating dynamical system in which exploitation and curiosity coexist in deliberate, tunable tension. One may visualize the curiosity-exploitation spectrum as a two-dimensional phase diagram: the horizontal axis represents the fraction of budget allocated to pure exploitation (0–1), while the vertical axis represents the strength of the internal novelty reward (likewise 0–1). Current materials AI systems cluster near the lower-right corner (high exploitation, zero novelty reward). The proposed blueprint advocates trajectories that deliberately traverse the upper-left quadrant, where curiosity budget and novelty reward are both elevated, producing hybrid regimes that retain engineering utility while systematically uncovering the unexpected. By enforcing these principles at the architectural level—rather than as post-hoc add-ons—the materials AI community can move from a philosophy of constrained optimization to one of principled open-ended exploration.

Success Criteria

Evaluating whether a curiosity-driven materials AI system is operating as intended requires a shift away from conventional performance metrics such as mean absolute error or hit rate on known targets. Instead, five distinct

success criteria, each grounded in the intrinsic goals of algorithmic curiosity, provide a multidimensional scorecard.

Novel discovery rate

The system must demonstrate a sustained rate of identifying materials whose structural or compositional features lie outside any previously sampled manifold, quantified by a novelty score that exceeds a predefined threshold derived from the training distribution. This criterion directly measures the effectiveness of Component 1 and ensures that the AI is not merely rediscovering variants of known compounds.

Space coverage

Success is gauged by the fraction of the high-dimensional chemical space that has been visited at least once, tracked via adaptive binning or persistent homology. High coverage indicates that Principle 4 and Component 4 are functioning, preventing the system from collapsing into familiar clusters even when the curiosity budget is active.

Surprising discoveries

The system should regularly produce outcomes whose properties deviate sharply from internal model predictions, measured as information gain or Kullback-Leibler divergence between prior and posterior beliefs. Ozin *et al.* [6] emphasize that such surprises lie at the heart of human-AI collaborative creativity; this criterion formalizes that insight for autonomous systems.

Downstream utility

Although utility is deliberately deprioritized during exploration, a successful curiosity-driven system must eventually yield materials that, once characterized, prove valuable for unanticipated applications. This criterion is evaluated retrospectively: novel discoveries are later tested against a broad suite of external objectives, confirming that curiosity does not preclude long-term impact.

Human curiosity match

Expert materials scientists review a random sample of the AI-generated candidates and rate them on a scale of “scientifically intriguing.” Alignment between human judgment and the system’s internal curiosity signals validates that algorithmic curiosity mirrors the epistemic virtues that have historically driven scientific progress [9].

These criteria are not mutually exclusive; the most robust systems will excel across all five simultaneously. By adopting them, the field can move beyond narrow benchmark chasing and begin to quantify the very openness that defines genuine scientific inquiry.

Implementation Path

Realizing algorithmic curiosity within the materials AI community demands a staged rollout that respects both technical feasibility and cultural readiness. The five-phase path outlined below begins with low-risk pilots and culminates in community-wide adoption.

Pilot curiosity systems

Deploy small-scale prototypes in tightly constrained subspaces (for example, binary oxides or small organic molecules) where computational cost remains modest. These pilots test the four operational principles in isolation, collecting baseline statistics on novelty detection and surprise maximization before scaling.

Hybrid systems

Integrate curiosity modules into existing exploitation-driven pipelines, allowing users to dial the curiosity budget from 0% to 50%. This phase, informed by the hybrid regimes visible on the conceptual spectrum, lets practitioners observe how curiosity augments rather than disrupts familiar workflows [2, 13].

Fully curiosity-driven

Remove all external objective functions for extended periods, permitting the system to explore solely under the guidance of the internal curiosity signals. Success at this stage confirms that the five components can sustain autonomous discovery without any human-specified target.

Integration with human curators

Introduce a collaborative loop in which human experts periodically review curiosity-generated candidates and provide lightweight feedback that modulates the novelty reward landscape. This phase leverages the human curiosity match criterion to refine the system's epistemic taste.

Community adoption

Publish open-source reference implementations, standardized curiosity metrics, and benchmark suites so that curiosity-driven exploration becomes a routine design choice alongside Bayesian optimization and generative modeling. At this stage, algorithmic curiosity transitions from proposal to standard practice, rebalancing the exploration-exploitation trade-off across the entire field.

Each phase includes explicit checkpoints tied to the success criteria, ensuring measurable progress and allowing early termination if foundational assumptions prove untenable. The path is deliberately incremental yet philosophically consistent, preserving the core commitment to open-ended discovery at every scale.

Objections and Replies

Any proposal that challenges the dominant exploitation paradigm inevitably encounters skepticism. Three representative objections merit careful consideration.

“Curiosity wastes resources on useless materials.” The reply is that many of the most consequential breakthroughs in materials science—high-temperature superconductors, topological insulators, and halide perovskites—were initially regarded as useless or anomalous until their unexpected properties were recognized. By institutionalizing surprise maximization [8, 15], algorithmic curiosity deliberately surfaces candidates that lie outside current utility models; the “waste” is therefore an investment in the unknown whose downstream value cannot be predicted in advance.

“We have specific societal goals; why explore randomly?” Algorithmic curiosity is not random. It is governed by explicit, memory-dependent mechanisms—novelty detection, uncertainty seeking, and prediction error feedback—that systematically target epistemic gaps rather than stochastic diffusion. As demonstrated in intrinsically motivated reinforcement learning [7, 9], such directed exploration consistently outperforms both random search and pure exploitation when the objective landscape is non-stationary or partially observable, precisely the situation in materials discovery.

“This is just active learning with different acquisition functions.” The distinction is philosophical as well as operational. Active learning, even in its most uncertainty-aware forms [11, 12], ultimately serves a fixed downstream task; algorithmic curiosity treats uncertainty and novelty as ends in themselves, independent of any eventual

application. The curiosity budget and novelty reward ensure that exploration can diverge from any external objective for indefinite periods, a freedom no conventional active-learning pipeline permits.

These replies underscore that algorithmic curiosity is neither inefficient nor redundant; it is a necessary corrective to the epistemic narrowing that exploitation dominance has imposed.

Conclusion

This paper has articulated algorithmic curiosity as a foundational design principle for exploration-driven materials AI. By documenting the epistemic costs of exploitation dominance, reviewing the transformative role of intrinsic motivation in both human science and contemporary artificial intelligence, and presenting a detailed blueprint comprising five components, four operational principles, five success criteria, and a five-phase implementation path, the proposal offers the community a concrete route toward open-ended discovery. Algorithmic curiosity does not replace goal-directed optimization; it restores balance, ensuring that materials AI can once again serve as a genuine partner in the scientific enterprise rather than merely an accelerated optimizer of

pre-specified targets. The time has come for the field to move beyond narrow utility and embrace curiosity as a core architectural imperative. Only then can artificial intelligence help uncover the truly unexpected wonders still hidden within the vast chemical universe.

Acknowledgements

None

Conflict of interest

None

Financial support

None

Ethics statement

None

Received: 14 Nov 2024 Revised: 07 Jan 2025 Accepted: 21 Feb 2025
Published online: 18 July 2025

Rights and permissions

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

Von Lilienfeld OA. Introducing machine learning: Science and technology. *Mach Learn Sci Technol*. 2020;1(1):010201.

Butler KT, Davies DW, Cartwright H, Isayev O, Walsh A. Machine learning for molecular and materials science. *Nature*. 2018;559(7715):547-55.

Schmidt J, Marques MR, Botti S, Marques MA. Recent advances and applications of machine learning in solid-state materials science. *npj Comput Mater*. 2019;5(1):83.

Zunger A. Inverse design in search of materials with target functionalities. *Nat Rev Chem*. 2018;2(4):0121.

Gómez-Bombarelli R, Wei JN, Duvenaud D, Hernández-Lobato JM, Sánchez-Lengeling B, Sheberla D, et al. Automatic chemical design using a data-driven continuous representation of molecules. *ACS Cent Sci*. 2018;4(2):268-76.

Ozin G, Siler T, Qian C, Zhou W. The curiosity-creativity element in HI-AI materials discovery. *Matter*. 2024;7(3):718-22.

Laversanne-Finot A, Péré A, Oudeyer PY. Intrinsically motivated exploration of learned goal spaces. *Front*

Neurorobot. 2021;14:555271.

Saglam B, Kozat SS. Deep intrinsically motivated exploration in continuous control. *Mach Learn*. 2023;112(12):4959-93.

Modirshanechi A, Kondrakiewicz K, Gerstner W, Haesler S. Curiosity-driven exploration: Foundations in neuroscience and computational modeling. *Trends Neurosci*. 2023;46(12):1054-66.

Jiang S, Webb MA. Generative active learning across polymer architectures and solvophobicities for targeted rheological behavior. *npj Comput Mater*. 2025;12(1).

<https://doi.org/10.1038/s41524-025-01900-2>.

Biswas A, Liu Y, Creange N, Liu YC, Jesse S, Yang JC, et al. A dynamic Bayesian optimized active recommender system for curiosity-driven partially Human-in-the-loop automated experiments. *npj Comput Mater*. 2024;10(1):29.

Narasimha G, Kong D, Regmi P, Jin R, Gai Z, Vasudevan R, et al. Uncovering multiscale structure-property correlations via active learning in scanning tunneling microscopy. *npj Comput Mater*. 2025;11(1):189.

Szymanski NJ, Rendy B, Fei Y, Kumar RE, He T, Milsted D, et al. An autonomous laboratory for the accelerated synthesis of inorganic materials. *Nature*. 2023;624(7990):86.

Liu Y, Kelley KP, Vasudevan RK, Funakubo H, Ziatdinov MA, Kalinin SV. Experimental discovery of structure–property relationships in ferroelectric materials via active learning. *Nat Mach Intell*. 2022;4(4):341-50.

Thiede LA, Krenn M, Nigam A, Aspuru-Guzik A. Curiosity in exploring chemical spaces: Intrinsic rewards for molecular reinforcement learning. *Mach Learn Sci Technol*. 2022;3(3):035008.

Ma J, Cao B, Dong S, Tian Y, Wang M, Xiong J, et al. MLMD: A programming-free AI platform to predict and design materials. *npj Comput Mater*. 2024;10(1):59.

Brown KA. Model, guess, check: Wordle as a primer on active learning for materials research. *npj Comput Mater*. 2022;8(1):97.

Xie Y, Vandermause J, Sun L, Cepellotti A, Kozinsky B. Bayesian force fields from active learning for simulation of

inter-dimensional transformation of stanene. *npj Comput Mater*. 2021;7(1):40.

Jablonka KM, Jothiappan GM, Wang S, Smit B, Yoo B. Bias free multiobjective active learning for materials design and discovery. *Nat Commun*. 2021;12(1):2312.

De Breuck PP, Wang HC, Rignanese GM, Botti S, Marques MA. Generative AI for crystal structures: A review. *npj Comput Mater*. 2025;11:370.

Guo L, Liu Y, Chen Z, Yang H, Donadio D, Cao B. Generative deep learning for predicting ultrahigh lattice thermal conductivity materials. *npj Comput Mater*. 2025;11(1):97.

Burda Y, Edwards H, Pathak D, Storkey A, Darrell T, Efron AA. Large-scale study of curiosity-driven learning. *arXiv preprint arXiv:1808.04355*. 2018 Aug 13.

Dai R, Song L, Liu H, Liang Z, Yu D, Mi H, et al. Cde: Curiosity-driven exploration for efficient reinforcement learning in large language models. *arXiv preprint arXiv:2509.09675*. 2025 Sep 11.

Wei J, Chu X, Sun XY, Xu K, Deng HX, Chen J, et al. Machine learning in materials science. *InfoMat*. 2019;1(3):338-58.

Shoaib M, Švecová L, Scholleová H, Vozňáková I, Bargoni A. Unravelling green ambidexterity innovation with green strategic intent: Exploring mediation and moderation effects. *J Intellect Cap*. 2025:1-28.

Li PP, Liu H, Li Y, Wang H. Exploration–exploitation duality with both tradeoff and synergy: The curvilinear interaction effects of learning modes on innovation types. *Manag Organ Rev*. 2023;19(3):498-532.

Jiang M, Rocktäschel T, Grefenstette E. General intelligence requires rethinking exploration. *R Soc Open Sci*. 2023;10(6):230539.

Zhang C, Liu H, Li W. An exploration-driven framework for path planning in complex buildings using improved MADDPG. *J Build Eng*. 2025;107:112626.

Savov N, Kazemi N, Mahdi M, Paudel DP, Wang X, Van Gool L. Exploration-driven generative interactive environments. In: *Proceedings of the Computer Vision and Pattern Recognition Conference*; 2025. p. 27597-607.