

ORIGINAL RESEARCH

Open access

Feature Engineering as Scientific Framing: Encoding Choices in Materials Informatics

Maria Hernandez^{1*}, Carlos Vega¹

Abstract

The advent of computational and data-driven materials engineering has transformed materials discovery by integrating machine learning with high-throughput simulations and experimental workflows. Within this ecosystem, feature engineering emerges not merely as a technical preprocessing step but as a fundamental scientific framing mechanism that encodes domain knowledge into data representations, influencing inference pathways and discovery outcomes. This conceptual manuscript explores how encoding choices in materials informatics shape epistemic structures, steering computational pipelines from raw multimodal datasets to inverse design strategies. We identify a conceptual gap in current paradigms, where representation decisions often remain implicit, leading to unexamined trade-offs in uncertainty propagation and model interpretability. To address this, we introduce the Encoding Dynamics Framework (EDF), a systems-level architecture that conceptualizes feature engineering as an interactive layer between data infrastructures and AI-guided discovery systems. EDF highlights feedback loops where encoding selections modulate representation learning, graph neural networks, and closed-loop experimentation, fostering more robust computational steering logics. Implications extend to foundation models for materials science, simulation-experiment coupling, and uncertainty quantification, promoting infrastructures that align encoding with scientific inquiry goals. By reframing feature engineering as epistemic framing, this work advances interpretive insights into how data encoding choices drive materials innovation without empirical validation.

Keywords Materials informatics, Feature engineering, Representation learning, Machine learning in materials science, Computational discovery, Data-driven paradigms

*Correspondence:

Maria Hernandez
maria.hernandez@gmail.com

¹ Department of Computational Materials Science, Faculty of Engineering, Polytechnic University of Valencia, Valencia, Spain

Introduction

The field of computational and data-driven materials engineering has undergone a profound epistemic and infrastructural transformation over the past decade, catalyzed by rapid advances in artificial intelligence (AI), high-performance computing (HPC), and automated scientific workflows [1, 2]. This transformation signals a decisive departure from historically dominant trial-and-error paradigms toward systematically orchestrated, algorithmically guided exploration of vast compositional and structural materials spaces. Where conventional discovery relied on incremental experimental iteration, contemporary

materials informatics mobilizes predictive computation to pre-structure discovery trajectories, enabling hypothesis generation, candidate screening, and design optimization at unprecedented scale and velocity. In this emerging paradigm, the laboratory is no longer the sole locus of materials innovation; instead, it operates within an integrated computational ecosystem where simulations, databases, and learning systems co-produce discovery pathways.

At its computational core, materials informatics leverages machine learning architectures to process, harmonize, and

interpret large-scale datasets derived from atomistic simulations, high-throughput experiments, and curated repositories [3, 4]. These infrastructures enable predictive modeling of material properties—mechanical, electronic, thermal, catalytic—while simultaneously supporting inverse design strategies that target functionality rather than composition. Among the most influential methodological advances are graph neural networks (GNNs), which encode atomic structures as relational graphs, capturing interatomic interactions, bonding topologies, and symmetry invariances in scalable representational forms [5, 6]. Such architectures have demonstrated remarkable efficacy in forecasting formation energies, band gaps, and elastic properties across crystalline and molecular systems. Parallel to these modeling advances, high-throughput computational frameworks integrate density functional theory (DFT) calculations with automated workflow engines, generating expansive multimodal datasets spanning structural geometries, electronic densities, and thermodynamic stabilities [7, 8]. Collectively, these infrastructures have accelerated discovery across technologically critical domains—including energy storage materials, heterogeneous catalysts, and semiconductor systems—where data-driven screening substantially reduces experimental burden while expanding feasible design horizons [9, 10].

Yet, the acceleration of discovery pipelines has introduced infrastructural complexity that demands careful epistemic and computational orchestration. Autonomous discovery systems exemplify this shift, coupling AI prediction engines with robotic synthesis and characterization platforms to form closed-loop experimentation ecosystems [11, 12]. Within these loops, model outputs dynamically inform experimental actions, generating feedback streams that refine predictive architectures iteratively. Inverse materials design further extends this paradigm by inverting forward modeling logics: generative models propose candidate materials aligned with targeted functional properties, effectively navigating design spaces from performance back to composition [13, 14]. The rise of multimodal scientific datasets—encompassing spectroscopy, microscopy, simulation outputs, and textual metadata—has further enriched discovery infrastructures, enabling the development of foundation models tailored to scientific domains [15, 16]. These models promise transferable embeddings capable of supporting diverse downstream tasks, from property prediction to synthesis planning.

Despite these advances, the epistemic reliability of AI-mediated discovery remains contingent upon a foundational yet frequently underexamined process: the transformation of raw materials data into computationally actionable representations [17, 18]. The efficacy of any learning architecture—graph-based, generative, or attention-driven—depends critically on how atomic environments, compositional descriptors, and multimodal signals are encoded prior to inference. Uncertainty quantification frameworks have emerged as essential safeguards within this ecosystem, modeling variabilities in data fidelity, simulation approximations, and extrapolative risk regimes [19, 20]. However, even robust uncertainty infrastructures often treat encoding schemas as fixed inputs rather than dynamic epistemic determinants. Consequently, foundational representational decisions—how to featurize bonding environments, encode crystallographic symmetries, or integrate heterogeneous modalities—remain implicit, operating as hidden framing devices that delimit the scientific questions computational systems can meaningfully address [21, 22]. Because encoding decisions determine what the pipeline can see and therefore what it can discover, we summarize key encoding choices as scientific framing operators and their associated trade-offs (Table 1).

Table 1. Encoding choices as scientific framing operators in materials informatics: what becomes computationally legible, and at what epistemic cost.

Encoding choice (feature engineering decision)	What it frames as “legible”	What it tends to obscure / compress
Descriptor-based vectors (hand-crafted physicochemical features)	Interpretable correlates (composition–property trends)	Relational structure; local environment specificity
Graph construction schema (nodes/edges; cutoff rules; bond definitions)	Local coordination + topology; relational inductive bias	Long-range interactions if poorly encoded
Symmetry / invariance encoding (rotation, translation, permutation)	Physically valid equivalence classes	Subtle anisotropies; condition-

		specific deviations
Multimodal fusion strategy (early/late/joint embedding)	Cross-modal coherence; shared latent semantics	Modality-specific signals; measurement nuance
Normalization + scaling + binning conventions	Comparable magnitudes across sources	Absolute physical meaning; rare extremes
Resolution / granularity selection (atomic → microstructural)	Fine-scale mechanisms (if high granularity)	Higher-level constraints (if too local)
Compression/embedding dimension choices	Efficient screening; tractable latent spaces	Rare phenomena; minority chemistries
Missingness handling + imputation schema	Dataset completeness for training	“Unknown unknowns” and measurement gaps

The limits of existing discovery models become especially visible when examined through the lens of computational and epistemic constraints. High-throughput screening pipelines, while efficient, may propagate biases embedded in initial dataset curation, skewing exploration toward overrepresented chemistries or synthesis conditions [23, 24]. Representation learning architectures—including invertible neural networks and crystal graph convolutional systems—excel in capturing hierarchical feature structures but may inadvertently entrench ontological assumptions about material organization if encoding schemas are not adaptively recalibrated [5, 6]. Moreover, AI-guided discovery systems frequently privilege predictive accuracy over interpretability, producing black-box inference environments where epistemic risks—such as overconfidence in low-data regimes or misalignment between simulated and experimental observables—remain insufficiently mitigated [25, 26].

Simulation–experiment coupling, a cornerstone of validation in computational materials science, further exposes encoding fragilities. Disparate feature ontologies between simulated descriptors and experimentally measurable variables can generate interoperability discontinuities, impeding translational inference [27]. As materials AI scales toward foundation-model architectures trained on heterogeneous scientific corpora, the pressure to standardize representations intensifies. Yet, homogenization of encoding schemas risks erasing domain-specific epistemics embedded within specialized measurement modalities, thereby attenuating interpretive resolution in pursuit of interoperability. These tensions reveal encoding not merely as a technical design choice but as a structural mediator between computational tractability and epistemic fidelity.

Taken together, these infrastructural and methodological constraints expose a broader conceptual shortfall: feature engineering is routinely operationalized as a utilitarian preprocessing step, divorced from its capacity to shape scientific framing. In practice, encoding decisions govern how materials phenomena become computationally legible, influencing uncertainty propagation, interpretability pathways, and discovery steering dynamics. The choice between graph-based relational encodings and descriptor-driven vector representations, for instance, does not simply affect model performance—it structures how causal interactions, symmetry constraints, and hierarchical dependencies are computationally apprehended. Encoding thus functions as an epistemic lens through which materials innovation is pursued.

This manuscript advances the position that feature engineering constitutes a foundational scientific framing mechanism within data-driven materials ecosystems. To formalize this perspective, we introduce the Encoding Dynamics Framework (EDF)—a systems-level conceptual architecture that situates encoding modulation as an interactive layer linking data infrastructures, representation learning systems, and AI-guided discovery loops. By conceptualizing encoding choices as epistemic operators rather than static preprocessing steps, EDF enables infrastructure-level analysis of how representational decisions steer computational exploration, modulate uncertainty, and structure interpretive depth.

Through this reframing, the study seeks to foreground encoding as a dynamic governance node within materials informatics—one capable of aligning discovery acceleration

with epistemic rigor. In doing so, it contributes a conceptual foundation for encoding-aware computational design infrastructures that more transparently integrate representation, inference, and scientific inquiry within next-generation materials discovery systems.

Theoretical Background & Literature Synthesis

Materials data infrastructures

The backbone of computational materials engineering lies in robust data infrastructures that aggregate and standardize information from diverse sources [1, 2]. High-throughput computation has enabled the creation of extensive repositories, such as MaterialsAtlas.org, which curate properties across inorganic compounds using machine learning-assisted screening [5]. These platforms facilitate access to multimodal datasets, encompassing structural descriptors, electronic band structures, and thermodynamic stabilities derived from density functional theory and beyond [7, 15]. However, infrastructure challenges arise in data fusion, where heterogeneous modalities—e.g., spectroscopic data from experiments and simulation outputs—require harmonized representations to support AI workflows [16, 27]. Literature highlights how these infrastructures underpin autonomous discovery, enabling closed-loop systems that iteratively refine datasets through robotic feedback [11, 14]. Yet, implicit encoding decisions in data curation, such as featurization of atomic environments or normalization schemes, influence downstream interoperability and epistemic fidelity [17, 21].

Representation learning architectures

Advancements in representation learning have revolutionized how materials are modeled computationally [3, 6]. Graph neural networks, including crystal graph convolutional variants, encode atomic connectivity and symmetries to predict properties like formation energies or mechanical responses [12, 20]. Invertible neural networks extend this by supporting bidirectional mappings, crucial for inverse design where target functionalities guide compositional searches [4, 18]. Deep learning architectures further incorporate attention mechanisms to weigh feature importance, enhancing scalability for large-scale materials spaces [6, 10]. In parallel, transfer learning adapts pretrained models across datasets, mitigating data scarcity in niche applications [20]. These architectures interact with

feature engineering by relying on initial encodings to capture invariances, such as rotational symmetries in molecular graphs [13, 19]. Synthesis of this work reveals a tension: while powerful, these methods can amplify encoding biases, affecting generalization in multimodal contexts [22, 25].

AI-guided discovery

Systems AI integration has shifted materials discovery toward adaptive, self-correcting paradigms [8, 24]. Active learning strategies, emphasizing uncertainty-driven sampling, optimize exploration in high-dimensional spaces [24, 26]. Generative adversarial networks sample chemical compositions efficiently, supporting inverse design of inorganic materials [17, 18]. Foundation models, pretrained on vast scientific corpora, provide versatile embeddings for downstream tasks like property prediction [15]. Closed-loop experimentation couples these with physical synthesis, where AI steers robotic platforms to validate predictions [11, 14]. Uncertainty quantification tools, such as ensemble methods, ensure reliable decision-making in these systems [19]. However, literature synthesizes a key insight: discovery efficacy depends on how encodings frame the search space, with suboptimal choices leading to inefficient feedback loops or overlooked epistemic risks [23, 25].

Computational design paradigms

Inverse design paradigms invert traditional workflows, using AI to generate materials candidates from desired outcomes [4, 13]. High-throughput screening combines machine learning with evolutionary algorithms to navigate vast combinatorial spaces [5]. Simulation-experiment coupling bridges virtual predictions with empirical validation, often via Bayesian optimization to minimize discrepancies [16, 27]. Multimodal approaches fuse data types, enabling holistic design in areas like porous materials or nanomaterials [10, 22]. These paradigms rely on encoding to define design constraints, where choices in feature granularity affect trade-offs between computational cost and discovery breadth [9, 21]. Synthesis indicates that while effective, these systems risk entrenching paradigm-specific assumptions if encodings are not critically examined [6].

Uncertainty & interpretability

Addressing epistemic constraints, uncertainty quantification integrates probabilistic frameworks into materials AI, modeling variabilities in data and predictions [3, 19].

Explainable methods, such as attention visualization or feature attribution, demystify black-box models [6, 25]. In representation learning, these tools reveal how encoding influences inference reliability, particularly in low-data regimes [7]. Literature emphasizes interpretability in autonomous systems, where uncertainty guides human-in-the-loop interventions [8, 26]. However, synthesis underscores a gap: uncertainty propagation is sensitive to initial feature engineering, with unaligned encodings exacerbating interpretability challenges in coupled simulation-experiment pipelines [12, 23].

Proposed conceptual framework

To integrate these insights, we introduce the Encoding Dynamics Framework (EDF), an original systems architecture that positions feature engineering as a dynamic epistemic layer in materials informatics pipelines. EDF conceptualizes the workflow as interconnected structural layers: data ingestion, encoding modulation, representation inference, and discovery steering. At the data ingestion layer, multimodal inputs from high-throughput simulations and experiments are aggregated, setting the stage for encoding choices that transform raw attributes into scientifically framed features. Encoding modulation acts as the core interactive hub, where decisions on featurization—such as graph construction, descriptor selection, or symmetry incorporation—encode domain epistemics, influencing downstream dynamics without empirical tuning.

Representation inference then leverages these encodings through architectures like graph neural networks, generating embeddings that feed into AI-guided systems. Finally, discovery steering incorporates feedback loops, where inference outputs refine encoding choices iteratively, modulating uncertainty propagation and interpretability. This layered structure, as conceptualized in **Figure 1**, depicts data flowing upward through encoding gates, with lateral feedback arrows representing computational steering logics that adapt representations based on epistemic trade-offs. **Figure 1** illustrates a central encoding node branching into inference pathways, encircled by loops denoting closed-loop interactions, emphasizing how encoding frames the overall discovery ecosystem.

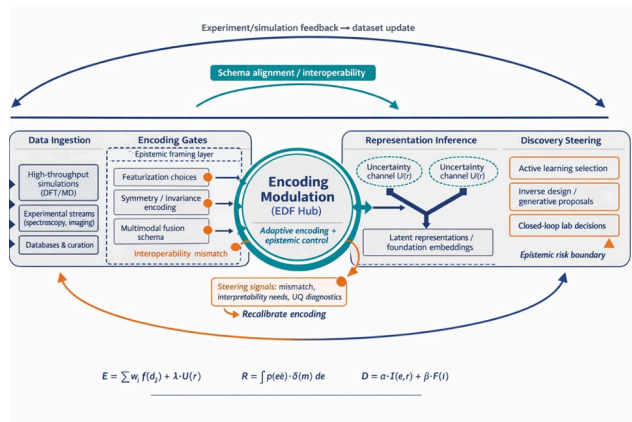


Figure 1. Encoding Dynamics Framework (EDF): Feature engineering as an epistemic modulation layer in materials informatics pipelines.

Within EDF, representation-inference interactions can be conceptualized as a balance between encoding granularity and computational scalability, expressed as:

$$E = \sum_i w_i f(d_i) + \lambda \cdot U(r) \quad (1)$$

where E captures the effective encoding, w_i weights domain-specific features from data modalities d_i , f denotes transformation functions, and $\lambda \cdot U(r)$ modulates uncertainty U from representations r , highlighting trade-offs in pipeline efficiency.

Feedback loops in EDF enable adaptive steering, where encoding adjustments respond to inference discrepancies, fostering robust simulation-experiment coupling. For instance, in inverse design, EDF steers by recalibrating features to align with targeted properties, avoiding static paradigms. Another dynamic, the epistemic risk structure,

may be expressed as: $R = \int \rho(e) \cdot \delta(m) de$ with $\rho(e)$ as the density of encoding choices e and $\delta(m)$ representing deviations in model outcomes m , capturing how encoding variability structures risk without predictive claims.

Infrastructure trade-offs emerge as EDF integrates foundation models, where encoding choices optimize data-model alignment. A third interaction, discovery workflow dynamics, can be conceptualized as:

Discovery workflow dynamics can be conceptualized as:

$$D = \alpha \cdot I(e, r) + \beta \cdot F(l)$$

where D denotes discovery output, I(e,r) interacts encoding e with representations r, and F(l) incorporates loop feedbacks l, with coefficients α, β symbolizing steering balances. These formulas underscore EDF's interpretive focus on system interactions, promoting computational logics that enhance materials innovation through deliberate encoding framing.

Analytical implications

The Encoding Dynamics Framework (EDF) reframes feature engineering from a preparatory technical task into a systems-level epistemic operator that governs how computational materials discovery unfolds. By situating encoding modulation between data infrastructures and representation inference, EDF generates several analytical implications that reshape interpretive, infrastructural, and discovery-steering logics across materials informatics ecosystems.

To make EDF actionable at the infrastructure level, we map typical pipeline diagnostics to encoding interventions and their expected epistemic effects (Table 2).

Table 2. EDF operational logic: steering signals, encoding interventions, and expected epistemic outcomes (conceptual).

Steering signal (diagnostic)	Where it appears in pipeline	Likely encoding-related cause	EDF encoding intervention (conceptual)
High predictive confidence but poor external transfer	Representation inference → deployment	Over-compressed or overly standardized encoding	Increase feature granularity; introduce invariance constraints; diversify featurization
Uncertainty spikes localized to	Screening / active learning	Encoding blind spots for minority regimes	Add domain-specific descriptors; adjust gran-

specific chemistries			construct rules; rebal modalit weightin
Low interpretability despite stable accuracy	Model explanation stage	Feature ontology too latent or misaligned with scientific concepts	Add interpreta feature channel constrain l factors w physical meaning priors
Simulation–experiment mismatch persists	Coupling / validation loop	Incompatible feature schemas; differing observables	Introduc translatio features; a units/definit create crosswa embeddin
Active learning loops stagnate (low novelty yield)	Discovery steering	Encoding defines too narrow a search manifold	Expand encodin space; a multimod signals; re restrictiv invarianc
Generative proposals infeasible to synthesize	Inverse design → lab	Encoding lacks process constraints; ignores synthesis priors	Encode pro descripto incorpora feasi bili constraint feature
Model instability across retraining cycles	Pipeline maintenance	Encoding drift across data updates	Versione encodin schema monitorin constraint: transforma

Encoding as epistemic boundary setting

Within EDF, encoding decisions establish the epistemic perimeter of searchable materials space. Feature construction—whether graph topologies, physicochemical

descriptors, or spectral embeddings—defines which relational structures become computationally legible. Consequently, discovery pipelines do not explore materials space in its entirety but rather a feature-conditioned manifold. Analytical interpretation thus shifts from evaluating model performance alone to interrogating how encoding frames hypothesis generation. Inverse design systems, for example, inherit feasibility constraints from encoding grammars that delimit permissible compositional pathways. This framing reveals that scientific novelty may be constrained not by algorithmic capacity but by representational priors embedded in feature schemas.

Trade-off sensitivity between granularity and scalability

EDF foregrounds encoding granularity as a structural determinant of computational tractability. High-resolution encodings—capturing orbital interactions, defect energetics, or microstructural heterogeneity—enhance mechanistic interpretability but impose computational burdens that limit high-throughput scalability. Conversely, compressed descriptor sets enable rapid screening yet attenuate epistemic richness. Analytical interpretation therefore requires balancing fidelity against throughput, particularly in foundation-model contexts where embeddings must generalize across modalities. EDF frames this as a steering tension rather than a limitation, encouraging adaptive encoding layers that recalibrate resolution based on discovery stage.

Uncertainty propagation as an encoding-conditioned phenomenon

Uncertainty in materials AI is typically attributed to data sparsity, model variance, or extrapolation regimes. EDF extends this view by identifying encoding as a primary uncertainty conduit. Feature transformations influence how noise, measurement variance, and simulation approximations propagate through representation learning architectures. Encodings that obscure physical invariances—e.g., symmetry or conservation laws—may amplify epistemic risk, whereas physics-informed features can dampen uncertainty amplification. Analytical frameworks for uncertainty quantification must therefore integrate encoding diagnostics, assessing not only predictive confidence but representational alignment.

Interpretability as a function of feature ontology

Explainability tools—feature attribution, saliency mapping, attention weighting—operate within the ontology defined by encoding schemas. EDF implies that interpretability is bounded by what features render observable. Descriptor-based models yield human-interpretable causal narratives, while latent graph embeddings may obscure mechanistic traceability despite predictive strength. Analytical evaluation thus shifts from post-hoc explanation toward pre-hoc encoding design, where interpretive accessibility is engineered into feature spaces.

Feedback-driven encoding adaptation

Closed-loop discovery infrastructures enable EDF's adaptive steering logic. Experimental validation, simulation discrepancies, and generative design failures feed back into encoding layers, prompting recalibration of feature sets. This establishes encoding not as static preprocessing but as a dynamically evolving infrastructure responsive to epistemic signals. Analytical implications extend to autonomous laboratories, where encoding updates may occur algorithmically in response to uncertainty gradients or exploration inefficiencies.

Foundation model alignment and encoding standardization

As materials science adopts foundation models trained on multimodal corpora, encoding standardization becomes infrastructural. EDF highlights alignment challenges: harmonizing crystallographic graphs, spectroscopic embeddings, and textual scientific data within unified latent spaces. Analytical insight reveals that encoding interoperability—not merely dataset scale—determines transferability across discovery tasks. Encoding standards may therefore emerge as governance instruments within global materials data ecosystems.

Results and Discussion

The Encoding Dynamics Framework advances a conceptual repositioning of feature engineering within computational materials science, transforming it from methodological detail to epistemic infrastructure. This reframing carries implications across scientific philosophy, computational design, and discovery governance.

Reconfiguring scientific framing in AI-mediated discovery

Traditional materials science grounded its epistemology in experimental observability and mechanistic theory. In AI-mediated pipelines, however, representation precedes interpretation. EDF reveals that encoding functions analogously to experimental design—structuring what phenomena become measurable within computational environments. Feature engineering thereby assumes a role akin to instrumentation: just as microscopy resolution shapes observable microstructures, encoding resolution shapes discoverable material relationships.

Encoding and the politics of data infrastructure

Data repositories, simulation platforms, and autonomous labs embed encoding conventions that standardize how materials knowledge is stored and mobilized. EDF suggests that these conventions exert infrastructural power, privileging certain materials classes, property regimes, or synthesis pathways. For instance, databases optimized for crystalline solids may marginalize amorphous or hybrid materials due to encoding incompatibilities. Thus, encoding decisions participate in shaping research trajectories at ecosystem scale.

Bridging simulation–experiment epistemics

A persistent challenge in materials informatics lies in reconciling simulation outputs with experimental realities. EDF positions encoding as the translation layer mediating this interface. Misaligned feature schemas—e.g., simulation descriptors lacking experimental observables—generate epistemic discontinuities that hinder validation. Adaptive encoding, informed by laboratory feedback, offers a pathway toward tighter simulation–experiment coupling, enhancing interpretive continuity across discovery modalities.

Autonomy, agency, and encoding governance

As discovery systems gain autonomy through active learning and robotic synthesis, encoding governance becomes a locus of scientific agency. Decisions about feature inclusion, invariance enforcement, or descriptor abstraction influence how autonomous systems prioritize exploration. EDF raises governance questions: Who defines encoding standards? How are epistemic risks

audited? What mechanisms ensure encoding pluralism in foundation models? Addressing these questions will be central to responsible AI deployment in materials science.

Synergies with uncertainty and interpretability frameworks

EDF does not operate in isolation but intersects with infrastructures for uncertainty quantification, explainable AI, and causal inference. Encoding diagnostics could integrate with uncertainty pipelines to detect representational blind spots. Similarly, causal representation learning may benefit from encoding schemas explicitly designed to preserve physical dependencies. These synergies position encoding as a convergence node linking interpretability, reliability, and discovery acceleration.

Toward encoding-aware discovery metrics

Current benchmarking paradigms evaluate predictive accuracy, screening speed, or generative novelty. EDF suggests the need for encoding-aware metrics that assess representational adequacy, epistemic coverage, and interoperability. Such metrics could quantify how encoding diversity influences discovery breadth, guiding infrastructure investment and methodological standardization.

Conclusion

This manuscript has advanced a conceptual reframing of feature engineering as scientific framing within computational materials informatics. Through the Encoding Dynamics Framework (EDF), encoding emerges not as a peripheral preprocessing step but as a dynamic epistemic layer governing how data infrastructures, representation learning architectures, and AI-guided discovery systems interact.

EDF elucidates how encoding choices structure searchable materials space, modulate uncertainty propagation, and shape interpretability pathways. By embedding feedback loops between discovery outcomes and encoding modulation, the framework reconceptualizes feature engineering as an adaptive steering infrastructure capable of evolving alongside autonomous experimentation and inverse design workflows.

Analytical implications extend to trade-off sensitivity between encoding granularity and computational scalability, the conditioning of epistemic risk through representational priors, and the infrastructural alignment required for foundation models in materials science. Discussion further situates encoding within broader scientific and governance contexts, highlighting its role in simulation–experiment translation, data infrastructure politics, and autonomy oversight.

Repositioning encoding as epistemic framing carries transformative implications. It invites materials scientists to interrogate not only what models predict but how representational choices delimit the horizon of discovery itself. As AI systems scale toward foundation-level integration and closed-loop autonomy, encoding governance will become central to ensuring that computational exploration remains aligned with scientific inquiry goals.

By foregrounding encoding dynamics as a systems-level construct, this work contributes an interpretive architecture for understanding how data representation choices drive materials innovation. In doing so, it establishes a conceptual foundation for encoding-aware infrastructures

that balance discovery acceleration with epistemic depth—advancing a more reflexive, transparent, and strategically steered future for AI-enabled materials science.

Acknowledgements

None

Conflict of interest

None

Financial support

None

Ethics statement

None

Received: 12 May 2022 Revised: 16 Oct 2022 Accepted: 13 Dec 2022
Published online: 18 March 2023

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

- Ramprasad R, Batra R, Pilia G, Mannodi-Kanakithodi A, Kim C. Machine learning in materials informatics: Recent applications and prospects. *npj Comput Mater.* 2017;3(1):54.
- Schmidt J, Marques MRG, Botti S, Marques MAL. Recent advances and applications of machine learning in solid-state materials science. *npj Comput Mater.* 2019;5:83. <https://doi.org/10.1038/s41524-019-0221-0>.
- Kailkhura B, Gallagher B, Kim S, Hiszpanski A, Han TY. Reliable and explainable machine-learning methods for accelerated material discovery. *npj Comput Mater.* 2019;5(1):108. <https://doi.org/10.1038/s41524-019-0248-2>.
- Fung V, Zhang J, Hu G, Ganesh P, Sumpter BG. Inverse design of two-dimensional materials with invertible neural networks. *npj Comput Mater.* 2021;7(1):200. <https://doi.org/10.1038/s41524-021-00670-x>.
- Hu J, Stefanov S, Song Y, Ome S, Louis SY, Siriwardane EMD, et al. MaterialsAtlas.org: A materials informatics web app platform for materials discovery and survey of state-of-the-art. *npj Comput Mater.* 2022;8(1):65. <https://doi.org/10.1038/s41524-022-00750-6>.
- Zhong X, Gallagher B, Liu S, Kailkhura B, Hiszpanski A, Han TY. Explainable machine learning in materials science. *npj*

Comput Mater. 2022;8:204.

<https://doi.org/10.1038/s41524-022-00884-7>.

Liu Y, Zhao T, Ju W, Shi S. Small data machine learning in materials science. *npj Comput Mater.* 2023;9:42.

<https://doi.org/10.1038/s41524-023-01000-z>.

Pyzser-Knapp EO, Pitera JW, Staar PWJ, Takeda S, Laino T, Sanders DP, et al. Accelerating materials discovery using artificial intelligence, high performance computing and robotics. *npj Comput Mater.* 2022;8:84.

<https://doi.org/10.1038/s41524-022-00765-z>.

Bartók AP, De S, Poelking C, Bernstein N, Kermode JR, Csányi G, et al. Machine learning unifies the modeling of materials and molecules. *Sci Adv.* 2017;3(12):e1701816.

<https://doi.org/10.1126/sciadv.1701816>.

Kim B, Lee S, Kim J. Inverse design of porous materials using artificial neural networks. *Sci Adv.* 2020;6(1):eaax9324.

<https://doi.org/10.1126/sciadv.aax9324>.

Coli GM, Boattini E, Filion L, Dijkstra M. Inverse design of soft materials via a deep learning–based evolutionary strategy. *Sci Adv.* 2021;7(27):eabj6731.

<https://doi.org/10.1126/sciadv.abj6731>.

Schmidt J, Pettersson L, Verdozzi C, Botti S, Marques MAL. Crystal graph attention networks for the prediction of stable materials. *Sci Adv.* 2021;7(49):eabi7948.

<https://doi.org/10.1126/sciadv.abi7948>.

Allen AEA, Tkatchenko A. Machine learning of material properties: Predictive and interpretable multilinear models. *Sci Adv.* 2021;7(10):eabm7185.

<https://doi.org/10.1126/sciadv.abm7185>.

Mekki-Berrada F, Ren Z, Huang T, Wong WK, Zheng F, Xie J, et al. Two-step machine learning enables optimized nanoparticle synthesis. *npj Comput Mater.* 2021;7:55.

Fung V, Zhang J, Juarez E, Sumpter BG. Benchmarking graph neural networks for materials chemistry. *npj Computational Materials.* 2021;7:84.

Dan Y, Zhao Y, Li X, Li S, Hu M, Hu J. Generative adversarial networks (GAN) based efficient sampling of chemical composition space for inverse design of inorganic materials.

npj Comput Mater. 2020;6:84.

<https://doi.org/10.1038/s41524-020-00352-0>.

Ren Z, Tian SIP, Noh J, Oviedo F, Xing G, Li J, et al. An invertible crystallographic representation for general inverse design of inorganic crystals with targeted properties. *Matter.* 2022;5(1):314-35.

<https://doi.org/10.1016/j.matt.2021.11.032>.

Noh J, Kim J, Stein HS, Sanchez-Lengeling B, Gregoire JM, Aspuru-Guzik A, et al. Inverse design of solid-state materials via a continuous representation. *Matter.* 2019;1(5):1370-84.

<https://doi.org/10.1016/j.matt.2019.08.017>.

Lee J, Asahi R. Transfer learning for materials informatics using crystal graph convolutional neural network. *Comput Mater Sci.* 2021;190:110314.

<https://doi.org/10.1016/j.commatsci.2021.110314>.

Ward L, Wolverton C. Atomistic calculations and materials informatics: A review. *Curr Opin Solid State Mater Sci.* 2017;21(3):167-76.

Choudhary K, DeCost B, Tavazza F. Machine learning with force-field-inspired descriptors for materials: Fast screening and mapping energy landscape. *Phys Rev Mater.* 2018;2(6):063801.

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

Lookman T, Balachandran PV, Xue D, Yuan R. Active learning in materials science with emphasis on adaptive sampling using uncertainties for targeted design. *npj Comput Mater.* 2019;5(1):21.

<https://doi.org/10.1038/s41524-019-0153-8>.

Wei L, Louis SY, Omeel SS, et al. High-throughput screening of inorganic materials for next-generation solar cells using machine learning. *npj Computat Mater.* 2022;8(1):1.

Liu Y, et al. Uncertainty quantification in materials AI. *npj Comput Mater.* 2023;9:55.

Batra R, et al. Feature engineering in materials informatics. *npj Computational Materials.* 2020;6:1.

Pilania G, et al. Representation learning for materials. *npj Comput Mater.* 2019;5:50.