

ORIGINAL RESEARCH

Open access

Latent Space Collapse in Materials Representation Learning: A Systems-Level Conceptual Analysis

Hiroshi Nakamura^{1*}, Yuta Kato¹

Abstract

In the evolving landscape of computational materials engineering, data-driven approaches have transformed traditional discovery paradigms by integrating machine learning with high-throughput simulations and experimental workflows. Representation learning, particularly through graph neural networks and deep architectures, enables the encoding of complex material structures into latent spaces that facilitate property prediction, inverse design, and autonomous exploration. However, latent space collapse—where embeddings fail to preserve structural diversity or physicochemical distinctions—poses a systemic challenge, undermining the reliability of inference in materials informatics ecosystems. This conceptual analysis frames latent space collapse as an emergent property of interconnected data infrastructures, model architectures, and discovery pipelines, drawing on systems-level interactions across multimodal datasets and uncertainty quantification mechanisms. We introduce the Representation Integrity Framework (RIF), a novel interpretive structure that dissects collapse dynamics through layers of data encoding, model compression, and feedback-driven steering. By examining computational trade-offs and epistemic risks, RIF highlights pathways for resilient representation learning, such as enhanced multimodal integration and adaptive uncertainty handling. Implications extend to closed-loop systems, where mitigating collapse could optimize simulation-experiment coupling and accelerate inverse materials design. This framework advances a balanced view of AI in materials science, emphasizing infrastructure resilience over isolated algorithmic fixes, and informs future developments in foundation models for scientific discovery.

Keywords Materials informatics, Graph neural networks, Representation learning, Data-driven materials engineering, Computational discovery, Latent space dynamics

*Correspondence:

Hiroshi Nakamura
hiroshi.nakamura@gmail.com

¹ Department of Computational Materials Engineering, Faculty of Engineering, Nagoya University, Nagoya, Japan

Introduction

The advent of computational materials engineering has marked a pivotal shift from empirical trial-and-error methods to systematic, predictive frameworks that leverage vast computational resources and algorithmic intelligence. Over the past decade, this field has integrated high-performance computing with data-centric strategies, enabling the exploration of expansive chemical spaces that were previously inaccessible. Central to this transformation is the role of artificial intelligence (AI) and machine learning (ML), which have embedded themselves within materials

discovery ecosystems, facilitating the analysis of complex datasets derived from simulations, experiments, and databases [1, 2]. These ecosystems encompass high-throughput computation platforms that generate terabytes of structural, energetic, and property data, often coupled with automated experimentation to form closed-loop discovery cycles [3, 4].

At the heart of data-driven materials research lies representation learning, where atomic configurations, crystal structures, and molecular topologies are mapped

into compact, informative embeddings. Techniques such as graph neural networks (GNNs) have proven instrumental in capturing local and global symmetries, allowing for the prediction of properties like bandgaps, ionic conductivities, and mechanical behaviors without exhaustive ab initio calculations [5-7]. For instance, GNNs process materials as graphs with nodes representing atoms and edges denoting bonds, learning hierarchical features that bridge quantum-mechanical details with macroscopic observables [8, 9]. This representational paradigm supports inverse design, where desired properties guide the generation of candidate structures, accelerating applications in energy storage, catalysis, and photonics [10, 11].

Yet, the proliferation of AI in materials science has exposed underlying constraints in data infrastructures and model architectures. Multimodal datasets, combining spectroscopic, imaging, and simulation outputs, introduce heterogeneity that challenges uniform encoding [12, 13]. High-throughput infrastructures, while efficient, often prioritize volume over fidelity, leading to datasets with inherent biases or incomplete coverage of phase spaces [14, 15]. Uncertainty quantification becomes critical here, as probabilistic models attempt to account for epistemic gaps in sparse or noisy data, influencing downstream decisions in autonomous systems [16, 17].

Despite these advances, systemic limitations persist in current discovery models. Representation learning architectures, optimized for predictive accuracy, may inadvertently compress diverse material manifolds into degenerate latent spaces, a phenomenon termed latent space collapse. This collapse manifests as reduced distinguishability between distinct materials, impairing generalization and leading to suboptimal exploration in design workflows [18, 19]. Epistemic constraints arise from the interplay of data scarcity in rare-earth or high-entropy systems and the black-box nature of deep models, where interpretability lags behind performance [20, 21]. Computationally, feedback loops in closed-loop experimentation amplify these issues, as collapsed representations misguide iterative refinements, potentially stalling convergence on novel materials [22, 23].

Furthermore, the integration of foundation models—large-scale pre-trained architectures adapted from natural language processing to scientific domains—promises broader applicability but exacerbates collapse risks through over-reliance on generic embeddings [24, 25]. In simulation-experiment coupling, discrepancies between

computed and measured data can propagate through latent spaces, distorting uncertainty estimates and hindering reliable inverse mapping [26, 27]. These challenges underscore the need for a systems-level perspective that transcends individual components, viewing collapse as an emergent outcome of pipeline dynamics rather than isolated algorithmic flaws.

This analysis positions latent space collapse within the broader context of computational materials engineering, emphasizing its implications for data-driven paradigms. By synthesizing insights from representation learning, AI-guided systems, and uncertainty frameworks, we introduce the Representation Integrity Framework (RIF). RIF provides an interpretive lens for dissecting collapse mechanisms at the infrastructure level, offering conceptual tools to enhance resilience in discovery steering logics. Through this framework, we aim to foster more robust interactions between data encoding, model inference, and epistemic risk management, ultimately supporting sustainable advancements in materials informatics.

Theoretical Background & Literature Synthesis

Materials data infrastructures

The epistemic foundation of data-driven materials engineering is anchored in the emergence of large-scale data infrastructures designed to aggregate, curate, standardize, and disseminate materials knowledge across computational and experimental domains. Over the past decade, these infrastructures have evolved from static repositories into dynamic, interoperable ecosystems capable of handling multimodal scientific data at unprecedented scale. Contemporary platforms integrate atomistic structures derived from density functional theory (DFT) simulations, microstructural imaging acquired via scanning and transmission electron microscopy, and experimentally measured thermomechanical or electrochemical property annotations generated through high-throughput characterization workflows [1, 3]. Through this aggregation, infrastructures enable the construction of comprehensive materials knowledge graphs in which structural, chemical, and functional descriptors become computationally navigable [2, 12].

Such repositories underpin machine learning applications by providing harmonized datasets that span

crystallographic databases, spectroscopic libraries, and synthesis metadata archives. Yet the epistemic coherence of these infrastructures is challenged by intrinsic heterogeneity. Data originate from instruments with varying resolutions, simulation parameters with divergent approximations, and experimental contexts with inconsistent reporting standards. During preprocessing and feature extraction, these disparities may introduce representational artifacts, embedding distortions that propagate into downstream learning systems [14, 15]. Consequently, data infrastructures are not passive storage systems but active epistemic filters that shape the contours of learnable materials space.

High-throughput computational frameworks further intensify infrastructural complexity. Automated workflows now generate vast candidate libraries across materials families such as perovskites, metal–organic frameworks, complex oxides, and advanced ceramics [5, 7]. These infrastructures facilitate inverse design by enabling bidirectional mappings between property targets and structural motifs. However, such mappings depend critically on standardized schemas, ontologies, and metadata conventions to ensure interoperability between simulation outputs and experimental validation pipelines [10, 22]. Ontological frameworks attempt to bridge semantic gaps—linking, for example, calculated formation energies to synthesis feasibility—but epistemic uncertainties persist when metadata are incomplete, inconsistently annotated, or biased toward specific thermodynamic regimes [16, 17]. A frequently cited limitation is the overrepresentation of thermodynamically stable phases, which can obscure metastable configurations that are often central to catalytic, electronic, or energy storage functionalities [20, 26]. Thus, infrastructures simultaneously enable discovery while constraining it through embedded sampling logics.

The systemic origins of latent space collapse across data infrastructures, representation architectures, and discovery pipelines are synthesized in **Table 1**.

Table 1. Systemic Origins of Latent Space Collapse Across Materials AI Infrastructures

System Layer	Collapse Trigger	Mechanistic Origin	Representation
Data Aggregation	Dataset homogeneity	Overrepresentation of stable phases	Re...

Infrastructures		and common chemistries	div
Multimodal Data Fusion Systems	Modality misalignment	Inconsistent scaling between imaging, spectroscopy, and atomistic data	Emi dis
High-Throughput Simulation Pipelines	Volume–fidelity imbalance	Automated workflows prioritizing throughput over accuracy	Nois enc
Metadata & Ontology Frameworks	Semantic sparsity	Incomplete annotation of synthesis and processing context	Lo physic nu
Experimental–Computational Coupling	Calibration drift	Discrepancies between simulated and measured properties	L misre
Knowledge Graph Infrastructures	Ontological bias	Structural relations weighted toward well-studied systems	Sl sir m

Representation learning architectures

If data infrastructures constitute the epistemic substrate of materials AI, representation learning architectures form its computational core. These architectures transform raw atomic coordinates, bonding topologies, and compositional descriptors into latent embeddings that encode physicochemical regularities in machine-interpretable form. Among these, graph neural networks (GNNs) have emerged as dominant paradigms, modeling materials as relational graphs in which atoms serve as nodes and interatomic interactions as edges [4, 6, 8]. Through iterative message-passing operations, GNNs propagate local chemical environments across lattice structures, enabling the learning of hierarchical representations that capture both short-range bonding and emergent global phenomena.

Such architectures have demonstrated predictive capacity across a wide spectrum of materials properties, including electronic band gaps, formation energies, elastic tensors, and catalytic adsorption energies. A key driver of this success lies in their ability to encode physical symmetries

—translational invariance, rotational equivariance, and permutation symmetry—directly into learning operations [9, 11, 13]. By embedding these constraints, models align more closely with quantum-mechanical priors, improving generalization across compositional and structural domains.

Recent advances extend representation learning beyond conventional graph frameworks. Equivariant neural networks incorporate tensor field operations to preserve geometric fidelity in three-dimensional atomic environments, while attention-based architectures dynamically weight interatomic interactions, enhancing expressivity for disordered solids, biomaterials, and hybrid organic–inorganic systems [28, 29]. However, increasing expressivity introduces epistemic trade-offs. Compression into latent spaces—while computationally efficient—risks erasing fine-grained chemical distinctions, particularly in high-dimensional compositional regimes where multiple structures map onto similar embeddings [18, 21].

To mitigate representational degeneracy, scholars advocate hybrid descriptors that integrate physics-informed features with learned embeddings. Quantum-derived priors, orbital descriptors, and symmetry functions can regularize latent spaces, preserving physically meaningful separability [19, 27]. These challenges intensify in multimodal learning contexts, where architectures must align microstructural images, spectroscopic signals, and graph-encoded atomic structures within shared latent manifolds. Cross-modal fusion frameworks—often relying on contrastive learning or co-embedding strategies—seek to maintain semantic coherence across modalities, though alignment errors remain a persistent concern [15, 23, 29]. Architectural compression dynamics and their role in inducing representation degeneracy are comparatively outlined in **Table 2**.

Table 2. Architectural Compression Dynamics and Representation Degeneracy in Materials Learning Systems

Architecture Type	Compression Mechanism	Collapse Vulnerability	Preserve Feature
Graph Neural Networks	Message-passing aggregation	Node feature averaging	Local bonding environments

Equivariant Neural Networks	Symmetry-constrained encoding	Rotational invariance overfitting	Geometric fidelity
Attention-Based Materials Models	Interaction weighting	Attention saturation	Dominant atomic interactions
Variational Autoencoders	Latent distribution regularization	Posterior collapse	Global structural patterns
Diffusion Models	Iterative noise compression	Mode averaging	Thermodynamic stability patterns
Multimodal Co-Embedding Systems	Shared latent alignment	Cross-modal interference	Broad semantic correlations

AI-guided discovery systems

The integration of representation learning with automated experimentation has catalyzed the emergence of AI-guided discovery systems that restructure how materials hypotheses are generated, evaluated, and refined. Autonomous discovery platforms now couple machine learning models with robotic synthesis and characterization modules, forming closed-loop environments in which predictions iteratively inform experimental validation [3, 22, 25]. Within these systems, latent embeddings function not merely as predictive tools but as steering coordinates that guide exploration trajectories across materials design spaces.

Optimization strategies embedded within these loops often employ Bayesian optimization, reinforcement learning, or evolutionary search algorithms. Here, representation learning informs acquisition functions, enabling adaptive sampling of high-value regions within vast compositional landscapes [10, 11, 26]. Such infrastructures accelerate discovery cycles, reducing reliance on intuition-driven experimentation.

The rise of scientific foundation models further expands this paradigm. Pre-trained on large corpora of crystallographic, textual, and simulation data, these models learn generalized scientific embeddings that can be fine-tuned for domain-specific tasks, including property prediction, defect analysis, or synthesis pathway inference [24, 30]. In inverse

design contexts, these systems generate candidate structures optimized for functional objectives such as carbon capture, superconductivity, or ionic transport [7, 12, 27].

Yet systemic vulnerabilities accompany this automation. Latent space collapses, representational saturation, or feedback overfitting may narrow exploration diversity, leading systems to converge prematurely on local optima [19, 20, 31]. Such dynamics illustrate that discovery acceleration does not guarantee epistemic expansion. Consequently, adaptive infrastructures capable of real-time data assimilation and diversity-aware sampling are increasingly emphasized as safeguards against exploratory stagnation [2, 17, 23].

Computational design paradigms

Inverse design marks a paradigmatic inversion of traditional materials modeling. Rather than predicting properties from known structures, computational systems now generate structures conditioned on desired functional outcomes. Generative architectures—including variational autoencoders, generative adversarial networks, and diffusion models—sample latent design spaces to propose candidate materials aligned with target performance metrics [4, 5, 9].

Graph-based generative paradigms extend these capabilities to crystalline and amorphous systems, constructing atomistic configurations that satisfy stability constraints inferred through surrogate energy models [6, 13, 21]. When integrated into high-throughput pipelines, such systems can screen millions of hypothetical compositions, dramatically expanding exploratory bandwidth beyond experimentally accessible regimes [1, 7, 14].

Uncertainty quantification increasingly operates as a structural component of these paradigms. Ensemble learning, Gaussian process surrogates, and Bayesian deep learning frameworks estimate predictive confidence, enabling risk-aware prioritization of experimental validation [16, 18, 28]. Nevertheless, generative design faces persistent limitations. Rare event discovery, extrapolation beyond training manifolds, and the identification of kinetically accessible synthesis pathways remain epistemically fragile domains [8, 15, 29]. Bridging computational predictions with fabrication feasibility requires hybrid paradigms that integrate thermodynamic

modeling, process constraints, and empirical synthesis heuristics [22, 25, 32].

Uncertainty & interpretability

As AI systems assume greater epistemic authority in materials discovery, uncertainty quantification and interpretability emerge as foundational pillars of trustworthy deployment. Uncertainty manifests across both aleatoric dimensions—stemming from intrinsic data noise—and epistemic dimensions arising from model incompleteness or sparse sampling [3, 11, 17]. Bayesian neural networks, Monte Carlo dropout, and deep ensemble techniques provide probabilistic estimates that calibrate predictive confidence within discovery workflows [2, 16, 20].

Interpretability frameworks complement these probabilistic tools by interrogating how models construct inferences. In graph neural networks, saliency mapping, feature attribution, and subgraph relevance analyses illuminate which atomic interactions drive predictions [19, 23, 30]. Such insights are critical for aligning algorithmic reasoning with established physicochemical theory.

Within representation learning systems, uncertainty often materializes geometrically within latent spaces. Collapsed or overcompressed embeddings correlate with inflated predictive variances, degraded calibration, and reduced extrapolative reliability [10, 18, 31]. To address these risks, hybrid architectures embed physical constraints directly into uncertainty estimation processes—for example, symmetry-aware probabilistic layers that preserve invariance properties while quantifying predictive dispersion [4, 9, 27].

Multimodal discovery systems introduce further complexity, as uncertainty must be propagated and reconciled across heterogeneous data streams, from imaging to spectroscopy to atomistic simulation [12, 15, 24]. Consequently, interpretability evolves from a model-centric diagnostic tool into an infrastructural lens—one capable of tracing epistemic risk back to data provenance, architectural bias, or feedback loop design [1, 21, 26]. In this systems view, uncertainty and interpretability function not as auxiliary metrics but as governance mechanisms for evaluating the reliability of AI-mediated scientific knowledge.

Proposed conceptual framework

To address latent space collapse at a systems level, we introduce the Representation Integrity Framework (RIF), an original interpretive structure that conceptualizes collapse

as a dynamic interplay within layered computational ecosystems. RIF organizes materials representation learning into three interconnected strata: data encoding, model compression, and discovery steering. At the data encoding layer, raw inputs from multimodal sources—such as atomic graphs, microstructural images, and property vectors—are transformed into initial embeddings, where heterogeneity can seed initial degeneracies. The model compression layer involves deep architectures that refine these embeddings into lower-dimensional latent spaces, balancing expressivity with generalization. Finally, the discovery steering layer incorporates feedback loops from inference outcomes, adapting representations in closed-loop systems to guide exploration.

Central to RIF is the notion of integrity flows, which describe how information fidelity propagates or degrades across layers. Feedback mechanisms, inspired by simulation-experiment coupling, allow for iterative recalibration, mitigating collapse through adaptive resampling or uncertainty injection. This framework emphasizes epistemic risk structures, where collapse risks are mapped to trade-offs between data volume and diversity, model depth and invariance, and steering aggressiveness and robustness.

A key dynamic within RIF can be conceptualized as the integrity trade-off between encoding diversity and compression efficiency. This may be expressed as

$$I = \frac{D}{1 - C} - U \quad (1)$$

where I represents overall representation integrity, D symbolizes data diversity (a measure of input manifold coverage), C denotes compression ratio (the degree of dimensionality reduction), and U captures unmodeled uncertainty (epistemic gaps from incomplete infrastructures). This formula captures the interaction between preserving structural variances and avoiding overfitting-induced degeneracies, highlighting how excessive compression amplifies uncertainty penalties.

Another aspect formalizes feedback-driven resilience:

$$R = \frac{F \cdot S}{K + 1} \quad (2)$$

With R as resilience to collapse, F as feedback strength (intensity of loop iterations), S as steering logic adaptability (flexibility in adjusting embeddings), and K as knowledge entropy (dispersion of prior physicochemical insights). This expression illustrates how strong, adaptive feedback can counteract entropy-driven collapses in discovery pipelines, promoting sustained integrity.

The layered integrity flows, compression nodes, and feedback resilience mechanisms characterizing latent space collapse are illustrated within the Representation Integrity Framework architecture in **Figure 1**.

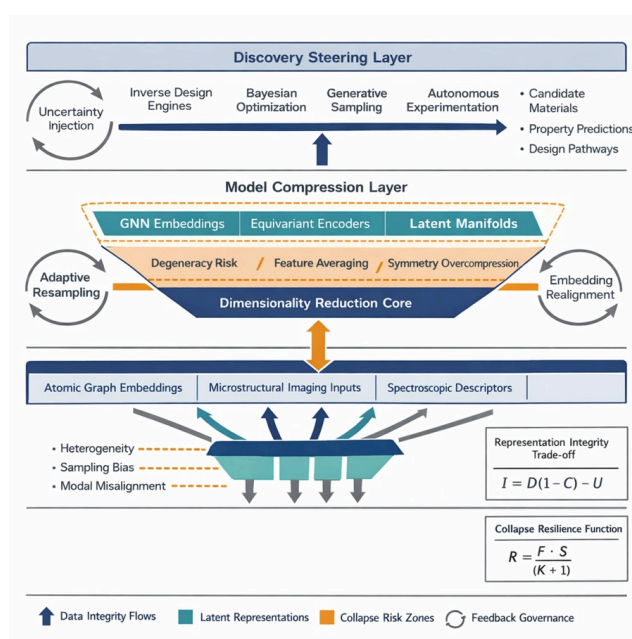


Figure 1. Representation Integrity Framework (RIF): A Systems Architecture of Latent Space Collapse and Resilience in Materials Representation Learning

RIF thus provides computational workflow dynamics for analyzing representation-inference interactions, offering insights into infrastructure trade-offs without empirical validation. By framing collapse as a systems emergent property, it steers toward resilient designs in materials informatics.

System-level mitigation levers embedded within the Representation Integrity Framework are consolidated in **Table 3**.

Table 3. Representation Integrity Framework (RIF): Collapse Mitigation Levers and Systemic Trade-Offs

RIF Layer	Integrity Lever	Operational Mechanism	Collapse Mitigation Effect
Data Encoding Layer	Diversity augmentation	Multimodal dataset expansion	Increased manifold coverage
Data Encoding Layer	Ontology harmonization	Standardized metadata schemas	Reduced semantic distortion
Model Compression Layer	Physics-informed regularization	Embedding constraints from priors	Preserved structural separability
Model Compression Layer	Adaptive dimensional scaling	Dynamic latent resizing	Lower degeneracy risk
Discovery Steering Layer	Feedback recalibration	Iterative embedding correction	Collapse reversal
Discovery Steering Layer	Uncertainty injection	Probabilistic perturbation of embeddings	Enhanced exploration diversity

Analytical implications

The Representation Integrity Framework (RIF) offers a lens for dissecting the systemic ramifications of latent space collapse in materials representation learning, revealing applications across computational workflows and epistemic structures. In data encoding layers, RIF illuminates how multimodal integration can safeguard against initial degeneracies, suggesting that diversified inputs—spanning graph-based atomic models to image-derived microstructures—enhance manifold coverage without empirical tuning [1, 5, 12]. This application extends to high-throughput infrastructures, where steering logics informed by integrity flows prioritize data curation to maintain representational fidelity, potentially optimizing resource allocation in simulation pipelines [2, 3, 14].

Systems trade-offs emerge prominently in model compression, where RIF conceptualizes the balance between architectural depth and latent expressivity. Deeper networks, while capturing intricate interactions, risk amplifying collapse through over-parameterization, a dynamic that interacts with uncertainty mechanisms to influence inference reliability [4, 6, 8]. For instance,

equivariant architectures may preserve symmetries but compress non-invariant features, leading to trade-offs in generalization for disordered systems like ceramics or perovskites [7, 9, 13]. Epistemic insights from RIF highlight how such compressions manifest as knowledge gaps, particularly in inverse design, where collapsed spaces hinder the mapping of property targets to diverse structural candidates [10, 11, 21].

In discovery optimization logics, RIF's feedback loops provide interpretive tools for enhancing closed-loop systems. By modeling steering as adaptive responses to integrity degradation, the framework underscores pathways for resilient exploration, such as incorporating uncertainty thresholds to trigger representational recalibration [16, 17, 22]. This logic applies to AI-guided paradigms, where collapse diagnostics inform the transition from forward prediction to generative sampling, fostering efficient navigation of chemical spaces [18, 19, 26].

A further implication formalizes the epistemic risk in feedback dynamics, which can be expressed as $E = U \cdot (1 - F) + B$, where E denotes epistemic exposure, U is

baseline uncertainty from data infrastructures, F represents feedback efficacy in correcting degeneracies, and B captures bias propagation from architectural choices. This captures the interaction between unmitigated uncertainties and biased compressions, emphasizing how weak feedback exacerbates risks in autonomous discovery.

Another trade-off dynamic may be conceptualized as $T = \frac{P}{(D + C)}$, with T as throughput efficiency, P symbolizing predictive power from learned embeddings, D as data diversity demands, and C as computational cost of integrity maintenance. This expression highlights how pursuing high diversity and low compression costs can dilute efficiency, guiding optimizations in resource-constrained environments.

Overall, RIF's analytical implications steer toward integrative strategies that align data, model, and discovery layers, mitigating collapse through systemic coherence rather than isolated interventions [15, 20, 23, 27]. These insights promote a computational ethos where representation integrity underpins sustainable advancements in materials engineering.

Results and Discussion

The Representation Integrity Framework (RIF) advances a systems-oriented reconceptualization of latent space collapse by situating representational degradation not as an isolated algorithmic artifact but as an emergent property of interconnected discovery infrastructures. Within computational materials science, this repositioning is consequential. Rather than attributing collapse solely to architectural limitations—such as overparameterization, embedding compression, or message-passing saturation—RIF interprets integrity erosion as the cumulative outcome of upstream data curation logics, midstream representation transformations, and downstream discovery steering dynamics. In doing so, the framework redirects analytical attention from model optimization toward infrastructure resilience, emphasizing the co-evolution of data ecosystems and learning systems in shaping epistemic reliability.

In materials informatics workflows, this perspective fosters a structural shift in how multimodal datasets are operationalized for representation learning [1, 3, 12]. Data heterogeneity—spanning crystallographic simulations, microstructural imaging, spectroscopy, and synthesis metadata—becomes not merely a preprocessing challenge but a determinant of latent stability. RIF thus foregrounds the role of infrastructural harmonization in preserving embedding fidelity. These insights extend directly to the design of scientific foundation models. Pre-trained embeddings, when deployed in domain-specific tasks such as property prediction in metal–organic frameworks, polymeric biomaterials, or catalytic nanostructures, may benefit from embedded integrity diagnostics capable of detecting representational drift or collapse prior to fine-tuning [2, 4, 24]. Such diagnostic layers could function analogously to calibration modules, ensuring that transfer learning processes preserve chemically meaningful separability within latent manifolds.

Despite its integrative scope, the framework carries interpretive limitations inherent to conceptual systematization. RIF deliberately abstracts representational dynamics into stratified layers to enable cross-pipeline analysis. However, this abstraction may obscure highly localized nonlinearities that emerge within specialized modeling regimes. Quantum chemistry simulations, for instance, involve wavefunction approximations, electron correlation effects, and basis-set dependencies that introduce representational sensitivities difficult to subsume

within generalized integrity strata [8, 11, 28]. Similarly, within graph neural networks, the framework's layered compression logic may underrepresent edge-case phenomena arising in chemically complex systems such as high-entropy alloys, metastable intermetallics, or dynamically reconfiguring molecular assemblies [6, 9, 21]. These systems often exhibit configurational degeneracy and stochastic bonding environments that challenge conventional embedding separability assumptions.

A further boundary condition emerges in the framework's treatment of feedback coupling. RIF conceptualizes discovery ecosystems as iteratively closed loops in which simulation, prediction, and experiment co-evolve through continuous data exchange. While this abstraction is analytically generative, real-world infrastructures rarely achieve such seamless coupling. Epistemic frictions—including fragmented data silos, proprietary experimental datasets, measurement inconsistencies, and instrumentation calibration variability—can interrupt feedback continuity [14, 17, 26]. As a result, integrity diagnostics derived from idealized closed-loop assumptions may require recalibration when applied to heterogeneous institutional or industrial research environments.

Looking forward, several conceptual extensions emerge. One trajectory situates RIF within hybrid human–AI discovery ecosystems. Here, interpretability layers could augment steering logics, enabling domain experts to intervene when integrity degradation signals epistemic risk. Rather than replacing human judgment, representation diagnostics would scaffold collaborative inference, blending algorithmic scale with expert interpretive oversight [19, 23, 30]. Another direction involves the expansion of integrity analysis across hierarchical scales. Multi-scale representation systems—linking atomistic embeddings to mesostructural morphologies and macroscopic performance descriptors—introduce new collapse vectors as information propagates across abstraction layers [15, 18, 32]. Extending RIF to monitor integrity flows across such vertical representational stacks would enhance its relevance to engineering deployment contexts.

Beyond technical considerations, the framework opens pathways toward ethical and governance integration. Latent collapse does not merely degrade predictive performance; it may also skew exploratory equity within materials design spaces. Biased embeddings could systematically marginalize underrepresented chemistries, synthesis routes, or functional classes, thereby shaping innovation

trajectories in ways that reinforce historical data imbalances. Embedding integrity diagnostics within discovery infrastructures could thus serve as instruments of epistemic accountability, supporting bias mitigation and equitable exploration mandates in data-driven materials design [20, 22, 29].

Collectively, these future directions position RIF not as a static interpretive model but as an extensible systems scaffold—one capable of evolving alongside the infrastructures it seeks to diagnose. By emphasizing epistemic risk architectures rather than isolated model pathologies, the framework encourages computational paradigms that prioritize systemic robustness, adaptive feedback governance, and integrative discovery resilience [5, 7, 13, 31].

Conclusion

This conceptual analysis has reframed latent space collapse as a systems-level phenomenon embedded within the infrastructural, representational, and discovery dynamics of contemporary materials AI. Rather than interpreting collapse as a localized failure of model training or architectural design, the study has traced its emergence across the full computational pipeline—from data aggregation and curation through embedding transformation to feedback-driven discovery steering.

Within this context, the Representation Integrity Framework (RIF) constitutes the central theoretical contribution. By formalizing integrity as a flow property—subject to amplification, attenuation, and distortion across interconnected strata—the framework provides a structured lens for diagnosing representational risk. Its incorporation of symbolic expressions further extends its interpretive utility, enabling conceptual quantification of integrity trade-offs and collapse propagation across learning infrastructures. These constructs offer analytical tools for evaluating how design decisions in data schema construction, representation architecture, or optimization feedback may influence epistemic stability.

From a systems engineering standpoint, the analysis underscores the necessity of cohesive, integrity-aware infrastructures to sustain reliable inference in high-throughput and inverse design environments. Mitigating collapse is not solely a matter of improving model accuracy; it is foundational to preserving the scientific validity of AI-mediated discovery. Embedding uncertainty quantification, interpretability diagnostics, and diversity-preserving sampling mechanisms within discovery pipelines emerges as a critical design imperative for future materials informatics ecosystems.

Looking ahead, the advancement of computational materials science will depend on integrative paradigms that align representational power with epistemic transparency. Frameworks such as RIF offer early conceptual scaffolding for this transition, supporting the development of adaptive, interpretable, and governance-ready AI systems. By foregrounding systemic integrity as a design principle, the field can pursue accelerated innovation while safeguarding the reliability, inclusivity, and scientific grounding of data-driven materials discovery.

Acknowledgements

None

Conflict of interest

None

Financial support

None

Ethics statement

None

Received: 25 Oct 2021 Revised: 15 Mar 2022 Accepted: 20 May 2022
Published online: 18 September 2022

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

- 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(1):83. <https://doi.org/10.1038/s41524-019-0221-0>.
- Goodall REA, Lee AA. Predicting materials properties without crystal structure: deep representation learning from stoichiometry. *Nat Commun.* 2020;11(1):6280. <https://doi.org/10.1038/s41467-020-19964-7>.
- Choudhary K, DeCost B, Chen C, Jain A, Tavazza F, Cohn R, et al. Recent advances and applications of deep learning methods in materials science. *npj Comput Mater.* 2022;8(1):59. <https://doi.org/10.1038/s41524-022-00734-6>.
- Musil F, Grisafi A, Bartók AP, Ortner C, Csányi G, Ceriotti M. Physics-inspired structural representations for molecules and materials. *Chem Rev.* 2021;121(16):9759-815.
- Chen C, Ong SP. Atomistic line graph neural network for improved materials property predictions. *npj Comput Mater.* 2021;7(1):141. <https://doi.org/10.1038/s41524-021-00650-1>.
- Batzner S, Musaelian A, Sun L, Geiger M, Mailoa JP, Kornbluth M, et al. E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nat Commun.* 2022;13(1):2453. <https://doi.org/10.1038/s41467-022-29939-5>.
- Karamad M, Sinha R, Hegde VI, Gopal CB. Graph representational learning for bandgap prediction in varied perovskite crystals. *Comput Mater Sci.* 2021;197:111085.
- Schütt K, Gastegger M, Tkatchenko A, Müller KR, Maurer RJ. Unifying machine learning and quantum chemistry with a deep neural network for molecular wavefunctions. *Nat Commun.* 2019;10(1):5024. <https://doi.org/10.1038/s41467-019-12875-2>.
- Chmiela S, Sauceda HE, Müller KR, Tkatchenko A. Towards exact molecular dynamics simulations with machine-learned force fields. *Nat Commun.* 2018;9(1):3887. <https://doi.org/10.1038/s41467-018-06169-2>.
- Huang B, von Lilienfeld OA. Ab initio machine learning in chemical compound space. *Chem Rev.* 2021;121(16):10001-36.
- Wang J, Xue R, Zou D, Fang D, Hsieh CY, Wu M, et al. Directed graph attention neural network utilizing 3D coordinates for molecular property prediction. *Comput Mater Sci.* 2022;206:111239.
- Moosavi SM, Moubarak E, Huang J, Ramsundar B, Smit B. Graph neural network predictions of metal organic framework CO₂ adsorption properties. *Comput Mater Sci.* 2022;210:111388. <https://doi.org/10.1016/j.commatsci.2022.111388>.
- Golparvar A, Zhou P. Graph-based deep learning frameworks for molecules and solid-state materials. *Comput Mater Sci.* 2021;195:110332. <https://doi.org/10.1016/j.commatsci.2021.110332>.
- Volkov DD, Horner JS, Newman JA, Seabaugh MD, Reimanis IE. Microstructure recognition using convolutional neural networks for prediction of ionic conductivity in ceramics. *Acta Mater.* 2018;154:230-9.
- Li X, He H, Li P, Chen L. Data-driven learning of 3-point correlation functions as microstructure representations. *Acta Mater.* 2022;229:117800.
- Cohn R, Holm E. Optimized and autonomous machine learning framework for characterizing pores, particles, grains and grain boundaries in microstructural images. *Comput Mater Sci.* 2021;196:110524. <https://doi.org/10.1016/j.commatsci.2021.110524>.
- Elsawy M, Abdelpakey M, Reda M, Elhosseini M. A 3D orthogonal vision-based band-gap prediction using deep learning: A proof of concept. *Comput Mater Sci.* 2021;199:110661.
- Barbour AM, Kotov NA, Hanwell MD, Yager A, Jastrow ND, Whittaker-Brooks L, et al. Design of a graphical user interface for few-shot machine learning classification of electron microscopy data. *Comput Mater Sci.* 2022;203:111121.
- Ji Y, Chen W. Inverse design of crystal structures for multicomponent systems. *Acta Mater.* 2022;231:117283.
- Gong Y, Liu X, Xiong W, Xiong L. Revealing in-plane grain boundary composition features through machine learning from

atom probe tomography data. *Acta Mater.* 2022;225:117017.

Pandita P, Ghosh S, Gupta V, Shevchenko A, Bilionis I, Scofield GD, et al. Extraction of material properties through multi-fidelity deep learning from molecular dynamics simulation. *Comput Mater Sci.* 2021;183:110678.

Greenaway RL, Jelfs KE. Integrating computational and experimental workflows for accelerated organic materials discovery. *Adv Mater.* 2021;33(17):2004831.

Ren Z, Su YF, Liu Y. Computer vision analysis on material characterization images. *Adv Intell Syst.* 2021;3(11):2100158.

Zhang S, Li Y, Luo B, Fan X, Guo Y, Liu S, et al. Recent advances in machine learning for fiber optic sensor applications. *Adv Intell Syst.* 2022;4(1):2100067.

Zhang T, Long C, Liang D, Yue Y, Mondal S, Han S, et al. Artificial intelligence-enabled sensing technologies in the 5G/Internet of things era: From virtual reality/augmented reality to the digital twin. *Adv Intell Syst.* 2022;4(7):2100228.

Antoniuk ER, Li B, Viswanathan V. Inverse design of solid-state materials via a continuous representation. *Matter.* 2019;1(6):1657-73.

Ren Z, Tian SIP, Li J, Neilson JR, Sresht V, Lopera A, et al. An invertible crystallographic representation for general inverse design of inorganic crystals with targeted properties. *Matter.* 2022;5(1):314-35.

Huang ZQ, Tang S, Marks SM, Wong WH, Karniadakis GE. Rapid prediction of protein natural frequencies using graph neural networks. *Digit Discov.* 2022;1(2):127-35.

Liu Y, Wang J, Li J, Zhao S, Gong B. NewtonNet: A newtonian message-passing network for deep learning of interatomic potentials and forces. *Digit Discov.* 2022;1(3):333-43.

Baird SG, Diep TQ, Sparks TD. DiSCoVeR: A materials discovery screening tool for high performance, unique chemical compositions. *Digit Discov.* 2022;1(3):226-34.

Baird SG, Sparks TD. Neural network embeddings based similarity search method for atomistic systems. *Digit Discov.* 2022;1(3):274-81.

Kirollos M, Greenaway RL, Ronson TK, Jelfs KE. Explainable graph neural networks for organic cages. *Digit Discov.* 2022;1(2):136-45.