

Digital Discovery

Accepted Manuscript

This article can be cited before page numbers have been issued, to do this please use: M. Q. Ha, D. Le, V. Nguyen, H. Kino, S. Curtarolo and H. Dam, *Digital Discovery*, 2025, DOI: 10.1039/D5DD00400D.



This is an Accepted Manuscript, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this Accepted Manuscript with the edited and formatted Advance Article as soon as it is available.

You can find more information about Accepted Manuscripts in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this Accepted Manuscript or any consequences arising from the use of any information it contains.

Beyond Interpolation: Integration of Data and AI-Extracted Knowledge for High-Entropy Alloy Discovery

Minh-Quyet Ha,¹ Dinh-Khiet Le,¹ Viet-Cuong Nguyen,² Hiori Kino,³ Stefano Curtarolo,^{4,5} and Hieu-Chi Dam^{1,6, a)}

¹Japan Advanced Institute of Science and Technology, 1-1 Asahidai, Nomi, Ishikawa 923-1292, Japan

²HPC SYSTEMS Inc., 3-9-15 Kaigan, Minato, Tokyo 108-0022, Japan

³Research Center for Materials Informatics, Department of Advanced Data Science, The Institute of Statistical Mathematics, 10-3 Midori-cho, Tachikawa, Tokyo 190-8562, Japan

⁴Department of Mechanical Engineering and Materials Science, Duke University, Durham, NC 27708, USA

⁵Center for Extreme Materials, Duke University, Durham, NC 27708, USA

⁶International Center for Synchrotron Radiation Innovation Smart (SRIS), Tohoku University, 2-1-1 Katahira, Aoba-ku, Sendai 980-8577, Japan

(Dated: 14 December 2025)

Discovering novel high-entropy alloys (HEAs) with desirable properties is challenged by the vast compositional space and the complexity of phase formation mechanisms. Several inductive screening methods that excel at interpolation have been developed; however, they struggle with extrapolating to novel alloy systems. This study introduces a framework that addresses the extrapolation limitation by systematically integrating knowledge extracted from material datasets with expert knowledge derived from scientific literature using large language models (LLMs). Central to our framework is the elemental substitution principle, which identifies chemically similar elements that can be interchanged while preserving desired properties. To model and combine evidence from these multiple sources of knowledge, we employ the Dempster–Shafer theory, which provides a mathematical foundation for reasoning under uncertainty. Our framework consistently outperforms conventional phase selection models that rely on single-source knowledge across all experiments, showing notable advantages in predicting phase stability for compositions containing elements absent from training data. Importantly, the framework effectively complements the strengths of the existing methods. Moreover, it provides interpretable reasoning that elucidates element substitutability patterns critical to alloy stability in HEA formation. These results highlight the framework’s potential for knowledge integration, offering an efficient approach to exploring the vast compositional space of HEAs with enhanced generalizability and interpretability.

I. INTRODUCTION

High-entropy alloys (HEAs), also known as multi-principal element alloys (MPEAs), have garnered significant attention owing to their exceptional mechanical properties, thermal stability, and corrosion resistance^{1–3}. Typically consisting of five or more principal elements in near-equiatomic ratios, these alloys utilize high-configurational entropy to stabilize single-phase solid solutions^{4–6}. However, identifying stable compositions remains a significant challenge due to the vast compositional space and the complex interplay of factors such as mixing entropy, enthalpy, atomic size differences, and electronic structure. These challenges, including exploring expansive design spaces, handling sparse data, and managing uncertainty, represent broader issues in combinatorial materials research, where efficient navigation strategies of compositional possibilities are essential.

A useful framework for understanding this challenge is a decision-making model in which researchers must bal-

ance *exploitation* and *exploration*^{7,8}, as illustrated in Figure 1. Exploitation focuses on well-characterized regions of the design space, having sufficient data for reliable property predictions. This approach supports steady, incremental improvements to existing alloys. In these data-rich regions, uncertainty is primarily *aleatoric*, arising from irreducible variability within the system. Conversely, exploration targets novel regions where data is insufficient for reliable property predictions. These regions introduce higher *epistemic* uncertainty that can be decreased as we collect more data through systematic experimentation. Although exploration bears greater risk, it offers the exciting potential to uncover groundbreaking and fundamentally new alloys with exceptional properties. Achieving an optimal balance between these two strategies is crucial for advancing HEA development.

Data-driven methods have emerged as transformative tools for guiding these exploitation-exploration decisions, enabling the processing of large datasets and streamlining the search for promising HEAs^{9–13}. High-throughput approaches, such as CALPHAD^{3,14,15}, AFLOW^{16–18}, and Hamiltonian models^{19,20}, alongside machine learning (ML)²¹, have significantly reduced the time and cost associated with evaluating candidate compositions. While

^{a)}Electronic mail: dam@jaist.ac.jp

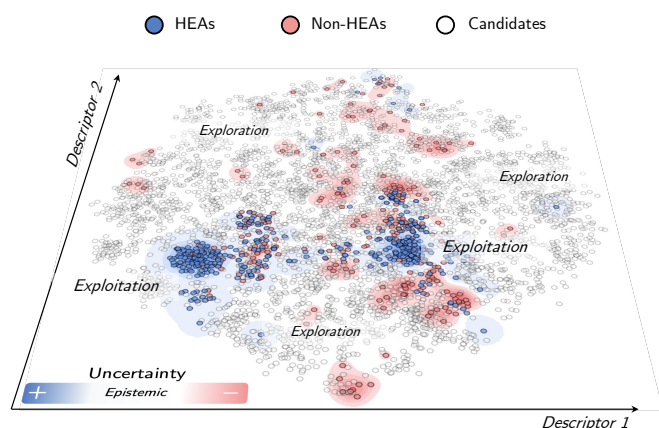


FIG. 1. Illustration of decision-making scenarios in high-entropy alloy (HEA) discovery. Colored regions represent well-established areas of the HEA compositional space, characterized by sufficient data suitable for effective exploitation. In contrast, white regions depict unexplored areas with sparse or no existing data, highlighting opportunities for risky yet potentially transformative exploration that could lead to discovering groundbreaking alloys with fundamentally new and exceptional properties. HEAs and Non-HEAs denote alloys that respectively form or do not form a stable high-entropy phase.

conventional ML models excel at *interpolation*, accurately predicting outcomes for compositions similar to those in the training sets (supporting exploitation), they struggle with *extrapolation* to novel systems, limiting exploration capability²². Although careful feature engineering can partially address extrapolation challenges²³, designing features that generalize across vast compositional spaces remains practically difficult^{22,24}. This *interpolation–extrapolation* dichotomy needs to be overcome as HEA discovery obviously requires venturing into uncharted territory.

A critical aspect of managing exploration–exploitation balance is uncertainty quantification, which falls into two categories. Epistemic uncertainty arises from incomplete or sparse data and is reducible through targeted information gathering, while aleatoric uncertainty corresponds to intrinsic variability within the system and is irreducible regardless of data volume²⁵. Traditional methods, such as Bayesian neural networks, Gaussian processes, and Monte Carlo dropout, are commonly employed to quantify these uncertainties^{26,27}. However, they often falter in early-stage materials discovery, where data is sparse or conflicting^{28–30}.

An alternative framework, the Dempster–Shafer theory^{31–33}, also known as evidence theory, offers a more flexible means of representing uncertainty. Unlike Bayesian methods, which assign probabilities to individual elements within a set of possibilities (denoted as Ω), evidence theory assigns non-negative weights (summing to one) to subsets of Ω . This enables the explicit representation of ignorance rather than requiring

an assumption about a prior probability distribution²⁵, allowing for nuanced characterization of both epistemic and aleatoric uncertainties. Thus, this framework can guide researchers to specific regions of the compositional space for either efficient exploitation or effective exploration^{22,34,35}.

However, collecting additional data to reduce epistemic uncertainty is often impractical due to high costs and experimental constraints. Expert knowledge offers a valuable alternative for mitigating this uncertainty. Domain specialists bring insights accumulated across multiple studies and contexts, providing heuristics that extend beyond any single dataset^{36–38}. Physics-informed neural networks (PINNs) exemplify one approach to incorporating domain knowledge by embedding a priori physical laws, enabling inference of governing equations from limited observations when those laws are explicit and well-defined³⁹. Yet their performance degrades when the underlying physics is only partially understood or key constraints remain unknown. More broadly, expert knowledge often resides in unstructured forms, such as laboratory notebooks, informal rules of thumb, or tacit experience, making its integration with structured, data-driven models a significant challenge.

To bridge this gap, this study introduces a framework that integrates knowledge from material datasets with expert domain knowledge accessed through AI systems—in this implementation, large language models (LLMs) extracting insights from scientific literature—while accounting for inherent uncertainties in each source. This uncertainty-aware integration enables systematic predictions beyond the interpolative boundaries of conventional data-driven methods. Central to our methodology is the *elemental substitution* principle^{40,41}, a well-established concept in alloy design wherein chemically similar elements can be interchanged while preserving target properties. We treat observed alloy pairs as evidence for substitutability patterns, then consolidate this empirical data with AI-derived insights obtained through state-of-the-art LLMs, including GPT-4o, GPT-4.5, Claude Opus 4, and Grok3. These LLMs leverage documented knowledge from related scientific domains through *knowledge integration* to assess elemental substitutability beyond the training dataset, not by generating information beyond their training corpus. Through Dempster–Shafer theory, the framework systematically models and combines these diverse evidence sources while quantifying both epistemic and aleatoric uncertainties. By providing accurate predictions in well-characterized regions alongside uncertainty-aware guidance for data-sparse spaces, this framework demonstrates—using HEAs as a proof of concept—the viability of materials discovery through uncertainty-aware AI integration.



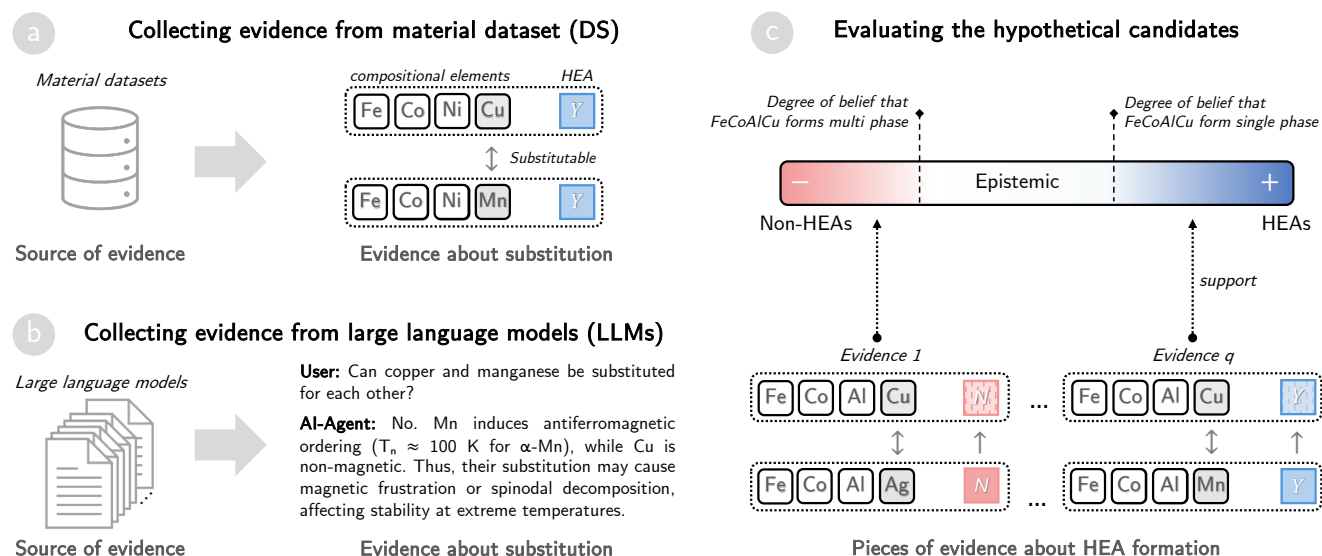


FIG. 2. **Hybrid framework integrating Data and AI-extracted Knowledge for high-entropy alloy (HEA) discovery.** (a–b) Schematic depicting the collection of substitutability evidence from a single material dataset (DS) and large language models (LLMs). (c) Schematic illustrating the assessment of hypothetical candidate properties using aggregated evidence derived from substitution-based methods.

II. METHODOLOGY

Each alloy A in the dataset \mathcal{D} is represented by its constituent elements. The property of interest y_A , for any alloy A , can be either HEA or \overline{HEA} . Here, HEA denotes alloys that form a stable high-entropy phase (single-phase solid solution), while \overline{HEA} (or Non- HEA) denotes alloys that do not form a stable high-entropy phase (multi-phase structures). To determine elemental substitutability, we assess the similarity between different element combinations by adapting *evidence theory*, which models and aggregates diverse pieces of evidence obtained from \mathcal{D} . Similarities between objects can manifest in various forms⁴²; e.g., pairwise ratings, object sorting, communal associations, substitutability, and correlation. In this study, we specifically focus on the *solid-solution formability* of element combinations and quantify their similarities based on elemental substitutability.

Our approach is intuitively illustrated using the example of element substitutability between Mn and Cu in Figure 2. Suppose we observe from materials datasets that two alloys, FeCoNiCu and FeCoNiMn, both form HEAs. This provides evidence that Cu can substitute for Mn in this context. Meanwhile, consulting domain knowledge through LLMs might reveal that metallurgists consider Cu-Mn pairs as non-substitutable, contributing additional conflicting evidence. Our proposed framework models and combines these independent pieces of evidence using evidence theory, potentially resulting in stronger belief in their substitutability than either source alone would provide. When predicting whether a new alloy, such as FeCoAlCu, forms an HEA, the framework can leverage existing data about FeCoAlMn and the es-

tablished Cu-Mn substitutability to make informed predictions.

A. Transforming Materials Data to Substitutability Evidence

Consider two alloys, A_i and A_j in \mathcal{D} , that share at least one common element. This non-disjoint pair of alloys provides evidence regarding the substitutability between the element combinations:

$$C_t = A_i \setminus (A_i \cap A_j) \quad \text{and} \quad C_v = A_j \setminus (A_i \cap A_j).$$

The intersection $A_i \cap A_j$ serves as the *context* for measuring similarity. If y_{A_i} and y_{A_j} agree (i.e., both are classified as HEA or both as \overline{HEA}), we infer that C_t and C_v are substitutable; otherwise, they are non-substitutable, as shown in Figure 2a.

The symmetric substitutability assumption ($C_t \rightarrow C_v$ and $C_v \rightarrow C_t$ are the same) used in this work represents a context-averaged approximation. While empirically validated for near-equiatomic HEAs, this assumption may limit accuracy for systems with strong directional substitution preferences. However, this symmetric treatment is justified in this study by two factors: first, the limited training data in our data-sparse scenarios makes learning separate directional patterns statistically infeasible; second, for near-equiatomic multi-principal element HEAs characterized by disordered random solid solutions, elements occupy statistically similar local environments, rendering symmetric substitution a physically reasonable first-order approximation.

Evidence for similarity is captured by defining a *frame of discernment*³² $\Omega_{sim} = \{\text{similar, dissimilar}\}$, encom-



1 passing all possible outcomes. The evidence from A_i and
 2 A_j is then represented by a *mass function* (or *basic prob-*
 3 *ability assignment*) $m_{A_i, A_j}^{C_t, C_v}$. This mass function assigns
 4 non-zero probability to the non-empty subsets of Ω_{sim} ,
 5 as:

$$m_{A_i, A_j}^{C_t, C_v}(\{\text{similar}\}) = \begin{cases} \alpha, & \text{if } y_{A_i} = y_{A_j}, \\ 0, & \text{otherwise,} \end{cases} \quad (1)$$

$$m_{A_i, A_j}^{C_t, C_v}(\{\text{dissimilar}\}) = \begin{cases} \alpha, & \text{if } y_{A_i} \neq y_{A_j}, \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

$$m_{A_i, A_j}^{C_t, C_v}(\Omega_{sim}) = 1 - \alpha. \quad (3)$$

6 Here, the parameter $0 < \alpha < 1$ is determined through
 7 an exhaustive search for optimal cross-validation per-
 8 formance, as shown in Supplementary Section 1. Intu-
 9 itively, $m_{A_i, A_j}^{C_t, C_v}(\{\text{similar}\})$ and $m_{A_i, A_j}^{C_t, C_v}(\{\text{dissimilar}\})$ rep-
 10 resent the extent to which alloys A_i and A_j support sub-
 11 stitutability or non-substitutability of C_t and C_v . Fur-
 12 ther, $m_{A_i, A_j}^{C_t, C_v}(\Omega_{sim})$ encodes epistemic uncertainty (i.e.,
 13 lack of definitive information). The probabilities assigned
 14 to these three subsets of Ω_{sim} must sum to 1.

15 Assuming that we collect q pieces of evidence from \mathcal{D}
 16 to compare C_t and C_v , each piece of evidence corresponds
 17 to a pair of alloys that generates a mass function $m_i^{C_t, C_v}$.
 18 These q mass functions are combined via *Dempster's rule*
 19 of combination³¹ to obtain a joint mass function $m_{\mathcal{D}}^{C_t, C_v}$:

$$m_{\mathcal{D}}^{C_t, C_v}(\omega) = \left(m_1^{C_t, C_v} \oplus m_2^{C_t, C_v} \oplus \dots \oplus m_q^{C_t, C_v} \right)(\omega), \quad (4)$$

20 where $\omega \subseteq \Omega_{sim}$, $\omega \neq \emptyset$ and \oplus denotes the Dempster's
 21 rule of combinations, as described in Supplementary Sec-
 22 tion 2. When no relevant evidence is available, $m_{\mathcal{D}}^{C_t, C_v}$
 23 is initialized with a mass of 1 on $\{\text{similar}, \text{dissimilar}\}$,
 24 indicating total uncertainty.

25 B. Transforming Domain Knowledge to Substitutability 26 Evidence

27 In addition to evidence collected from material
 28 datasets (DS), we focus on evidence derived from do-
 29 main knowledge, utilizing LLMs to extract insights from
 30 a vast corpus of scientific literature. Specifically, we use
 31 a set of state-of-the-art LLMs including GPT-4o, GPT-
 32 4.5, Claude Opus 4, and Grok3 to assess element sub-
 33 stitutability based on expert perspectives within a given
 34 domain, as illustrated in Figure 2b. The proposed model
 35 evaluates the substitutability of element pairs from the
 36 perspective of a domain expert, ensuring that the anal-
 37 ysis aligns with established scientific reasoning. To en-
 38 hance result reliability, we implement a two-step prompt-
 39 ing procedure:

- **Question 1:** Do you possess sufficient knowledge or data to evaluate the substitutability of elements C_t and C_v within the context of [domain knowledge]?
- **Question 2:** If the answer to the first question is Yes, the LLM further rates element substitutability as High, Medium, or Low, based on insights distilled from relevant scientific literature in the given domain.

49 Detailed prompts used for each LLM are provided in
 50 Supplementary File 1. This approach is based on the as-
 51 sumption that, when given clear and structured prompts,
 52 these LLMs can simulate expert reasoning across multi-
 53 ple scientific domains. This capability stems from their
 54 extensive training on scientific literature, which enables
 55 them to provide contextually relevant, domain-specific
 56 feedback tailored to the challenges of HEA discovery.

57 Elemental substitutability is not universal and is
 58 property-specific, strongly associated with functionality
 59 and applications. For example, substitution for struc-
 60 tural stability differs from substitution targeting the
 61 magnetic, optical, or mechanical properties. Recogniz-
 62 ing this property-specific nature, our framework requires
 63 careful domain selection tailored to the target property
 64 to ensure accurate predictions. To facilitate the extrac-
 65 tion of domain knowledge, we focus on five key scientific
 66 domains, including *corrosion science*, *materials mechan-*
 67 *ics*, *metallurgy*, *solid-state physics*, and *materials science*.
 68 These domains are selected due to their critical roles in
 69 understanding and optimizing HEAs, specifically tailored
 70 for phase stability prediction⁵. Each domain contributes
 71 essential insights into different aspects of alloy design.

- **Corrosion science:** This domain examines chemical degradation mechanisms and protective strategies, essential for ensuring long-term durability.
- **Materials mechanics:** This domain investigates mechanical properties such as strength, ductility, and toughness, crucial for structural performance.
- **Metallurgy:** This domain analyzes phase formation, phase diagrams, and microstructure control, offering insights into alloy stability and processing methods.
- **Solid-state physics:** This domain explores atomic-scale interactions, electronic structure, and thermal behavior, all of which influence phase stability and material performance.
- **Materials science:** This domain serves as an integrative field that synthesizes perspectives from the other domains, emphasizing the relationships between composition, structure, properties, and performance to optimize alloy design strategies.

91 The evidence collected from the LLM for each do-
 92 main is categorized into one of four outcomes: High,



TABLE I. Possible outcomes generated by an LLM for each domain-specific criterion, along with the corresponding mass functions $m_{LLMs}^{C_t, C_v}(\{\text{similar}\})$, $m_{LLMs}^{C_t, C_v}(\{\text{dissimilar}\})$, and $m_{LLMs}^{C_t, C_v}(\Omega_{sim})$. Here, $0 < \beta < 1$ indicates our confidence in LLM's response, with determination details provided in Supplementary Section 1.

Q1	Q2	$m_{LLMs}^{C_t, C_v}(\{\text{similar}\})$	$m_{LLMs}^{C_t, C_v}(\{\text{dissimilar}\})$	$m_{LLMs}^{C_t, C_v}(\Omega_{sim})$	Interpretation
No	—	0	0	1	LLM does not provide sufficient domain knowledge
Yes	High	β	0	$1 - \beta$	C_t and C_v are considered <i>highly</i> substitutable
Yes	Medium	$\beta/2$	$\beta/2$	$1 - \beta$	C_t and C_v are considered <i>moderately</i> substitutable
Yes	Low	0	β	$1 - \beta$	C_t and C_v are considered <i>poorly</i> substitutable

1 Medium, Low, or No Knowledge. Further, these outcomes
2 are mapped to a corresponding mass function denoted
3 as $m_{LLMs}^{C_t, C_v}$, as shown in Table I. If the LLM indicates No
4 Knowledge, then the entire mass is assigned to the set
5 $\{\text{similar}, \text{dissimilar}\}$, reflecting complete epistemic uncer-
6 tainty. Conversely, if the LLM provides a specific substi-
7 tutability rating (High, Medium, and Low), then a portion
8 of the mass is allocated to either $\{\text{similar}\}$ or $\{\text{dissimilar}\}$,
9 while the remaining mass is assigned to Ω_{sim} to account
10 for residual uncertainty in the prediction.

11 Notably, all LLMs (GPT-4o, GPT-4.5, Claude Opus
12 4, and Grok3) are used as pre-trained models *out-of-*
13 *the-box* without any fine-tuning, retraining, or in-context
14 literature provision. These models are queried directly
15 through their respective API interfaces using the two-
16 step prompting procedure described above and detailed
17 in Supplementary File 1. The LLMs leverage knowledge
18 from scientific literature encountered during their origi-
19 nal pre-training by the respective model developers; we
20 do not modify these models in any way. Each LLM pro-
21 vides independent assessments that are later combined
22 using Dempster-Shafer theory (Section II.C).

23 C. Combining Evidence from Multiple Sources

24 In this study, a *source* S refers to an independent
25 knowledge provider that generates evidence about ele-
26 mental substitutability. Our multi-source framework in-
27 tegrates two kinds of independent sources:

- 28 • **DS-source:** A material dataset \mathcal{D} provides em-
29 pirical evidence by analyzing alloy pairs that differ
30 by element substitution (Section II A). This dataset
31 contains factual observations about the target do-
32 main (e.g., which alloy compositions form HEAs).
- 33 • **LLM sources:** We query 4 state-of-the-art LLMs
34 (GPT-4o, GPT-4.5, Claude Opus 4, Grok3) across
35 5 scientific domains (corrosion science, materials
36 mechanics, metallurgy, solid-state physics, materi-
37 als science), creating $4 \times 5 = 20$ independent knowl-
38 edge sources (Section II B). Each combination of an
39 LLM and a domain provides documented scientific
40 knowledge from related or similar domains to the
41 target domain.

42 To integrate substitutability evidence collected from
43 multiple sources, Dempster's rule of combination with a

44 *reliability-aware discounting* step is used^{32,43}. Recogniz-
45 ing that substitutability is property-specific and differ-
46 ent sources capture different aspects of elemental substi-
47 tutability, our framework implements an adaptive mech-
48 anism that evaluates each source's relevance to the target
49 property. This reliability-aware discounting automati-
50 cally assigns higher weights to sources that align well with
51 the specific property being predicted while suppressing
52 sources that capture irrelevant substitutability criteria,
53 thereby preventing inappropriate knowledge integration.
54 For each source S , we compute a dataset-specific dis-
55 count factor as:

$$\gamma_S = \text{disc}(m_S^{C_t, C_v}, \mathcal{D}) \in [0, 1], \quad (5)$$

56 where $\text{disc}(\cdot)$ quantifies how well the substitutability evi-
57 dence collected from source S generalizes to the alloy
58 properties in \mathcal{D} . The reliability of each source is assessed
59 using the macro-averaged F1 score with 10-fold cross-
60 validation. For instance, if a source S has historically
61 demonstrated accurate predictions on alloys similar to
62 those in \mathcal{D} , we assign γ_S a value closer to 1. Conversely,
63 if S performs poorly or unpredictably for alloys in \mathcal{D} , γ_S
64 is reduced accordingly.

The original mass function $m_S^{C_t, C_v}$ for source S is then
modified by incorporating the discount factor γ_S , leading
to an adjusted function $\gamma_S m_S^{C_t, C_v}$:

$$\begin{aligned} \gamma_S m_S^{C_t, C_v}(\{\text{similar}\}) &= \gamma_S \times m_S^{C_t, C_v}(\{\text{similar}\}), \\ \gamma_S m_S^{C_t, C_v}(\{\text{dissimilar}\}) &= \gamma_S \times m_S^{C_t, C_v}(\{\text{dissimilar}\}), \\ \gamma_S m_S^{C_t, C_v}(\Omega_{sim}) &= 1 - \gamma_S + \gamma_S \times m_S^{C_t, C_v}(\Omega_{sim}). \end{aligned} \quad (6)$$

65 This redistribution shifts mass from definitive conclu-
66 sions $\{\text{similar}\}$ and $\{\text{dissimilar}\}$ to the ambiguous set
67 $\{\text{similar}, \text{dissimilar}\}$, thereby encoding epistemic uncer-
68 tainty for less reliable sources. Therefore, when all mass
69 functions are subsequently merged using Dempster's rule,
70 less credible sources exert a weaker influence on the final
71 decision.

72 Assuming p sources $\{S_1, S_2, \dots, S_p\}$, the substitutabil-
73 ity evidence gathered from them is aggregated using
74 Dempster's rule of combination:

$$m^{C_t, C_v}(\omega) = \left(\gamma_{S_1} m_{S_1}^{C_t, C_v} \oplus \gamma_{S_2} m_{S_2}^{C_t, C_v} \oplus \dots \oplus \gamma_{S_p} m_{S_p}^{C_t, C_v} \right)(\omega), \quad (7)$$



where ω denotes non-empty subsets of Ω_{sim} . The rule iteratively integrates evidence while normalizing conflicts (such as empty-set intersections arising from contradictory sources). This approach preserves diverse insights, from data-driven correlations to LLM-derived domain knowledge, while mitigating the influence of unreliable sources. Critically, when evidence about substitutability is insufficient or conflicting, Dempster's rule of combination assigns high mass to $m^{C_t, C_v}(\Omega_{sim})$, explicitly signaling uncertainty rather than forcing confident predictions. This naturally prevents overfitting in data-sparse scenarios common in materials discovery.

Similar analyses are conducted for all pairs of element combinations, resulting in a symmetric matrix M , where $M[t, v] = M[v, t] = m^{C_t, C_v}(\{\text{similar}\})$.

D. Evaluating Hypothetical Candidates by Analogy-Based Inference

To predict whether a *new* alloy A_{new} is likely to form an HEA, we employ a substitution-based inference approach utilizing the similarity matrix M . The process begins with a known alloy A_k , labeled y_{A_k} , and identifies the subset $C_t \subset A_k$ that, when replaced by C_v , generates A_{new} (Figure 2 c). If C_t and C_v are deemed substitutable, then $y_{A_{new}}$ is more likely to match y_{A_k} ; conversely, if they are dissimilar, $y_{A_{new}}$ may differ.

We formalize this inference using a frame of discernment³² $\Omega_{HEA} = \{\text{HEA}, \overline{\text{HEA}}\}$ and define a mass function $m_{A_k, C_t \leftarrow C_v}^{A_{new}}$ to model the evidence collected from A_k and the substitution of C_t for C_v , denoted as $C_t \leftarrow C_v$. This mass function distributes belief among $\{\text{HEA}\}$, $\{\overline{\text{HEA}}\}$, or $\{\text{HEA}, \overline{\text{HEA}}\}$ according to the similarity $M[t, v]$ and the label of A_k as:

$$m_{A_k, C_t \leftarrow C_v}^{A_{new}}(\{\text{HEA}\}) = \begin{cases} M[t, v], & \text{if } y_{A_k} = \text{HEA}, \\ 0, & \text{otherwise,} \end{cases} \quad (8)$$

$$m_{A_k, C_t \leftarrow C_v}^{A_{new}}(\{\overline{\text{HEA}}\}) = \begin{cases} M[t, v], & \text{if } y_{A_k} = \overline{\text{HEA}}, \\ 0, & \text{otherwise,} \end{cases} \quad (9)$$

$$m_{A_k, C_t \leftarrow C_v}^{A_{new}}(\Omega_{HEA}) = 1 - M[t, v]. \quad (10)$$

Here, the probability mass assigned to $\{\text{HEA}\}$ and $\{\overline{\text{HEA}}\}$ reflects the confidence levels with which A_k and the substitution of C_v for C_t support the probabilities that A_{new} is or is not an HEA, respectively. The mass assigned to subset $\{\text{HEA}, \overline{\text{HEA}}\}$ represents epistemic uncertainty, signifying cases where the available evidence does not provide definitive information regarding the properties of A_{new} . The total probability mass assigned to all three non-empty subsets of Ω_{HEA} is constrained to sum to 1, ensuring a consistent probabilistic framework. An illustrative example employing the Dempster-Shafer

theory for the evaluation of hypothetical candidates is provided in Supplementary Section 3.

We assume that multiple pieces of evidence can be collected, each derived from a distinct pair of host alloy A_{host} and substitution pair $C_t \leftarrow C_v$, for a new alloy candidate A_{new} . These individual pieces of evidence are systematically combined using Dempster's rule of combination to generate a final mass function $m^{A_{new}}$. This function integrates all available analogies, resolving potential inconsistencies and contradictions among the sources. The resulting combined evidence offers a coherent assessment, aiding in informed decision-making regarding whether further resource-intensive experiments are necessary to validate the HEA formation ability of A_{new} .

III. EXPERIMENTAL SETTING

In this section, we present the design of experiments, which assess both the *predictive capability* and *interpretability* of our proposed method. Additionally, we provide comparisons against alternative approaches, including single-source evidential methods and other data-driven classifiers.

A. Datasets

Experiments are conducted considering four computational datasets of quaternary alloys, one experimental dataset of quaternary alloys, and one experimental dataset of quinary high-entropy borides (HEB), summarized in Table II. HEBs are single-phase ceramics containing multiple transition metal cations randomly distributed on the metal sublattice of a boride structure, offering unique combinations of metallic and ceramic properties⁴⁴. Despite different bonding mechanisms, HEBs exhibit similarly high elemental selectivity as HEAs—boron's restrictive bonding requirements create stringent constraints on metal selection, analogous to the selective substitutability patterns in metallic HEAs, making them suitable for testing our framework's core principle of managing uncertainty in highly selective multi-component systems.

- $\mathcal{D}_{0.9T_m}$ and \mathcal{D}_{1350K} : These computational datasets include *all possible quaternary* alloys generated from a set of 26 elements: Fe, Co, Ir, Cu, Ni, Pt, Pd, Rh, Au, Ag, Ru, Os, Si, As, Al, Re, Mn, Ta, Ti, W, Mo, Cr, V, Hf, Nb, and Zr. The stability of these alloys is predicted using methods proposed by Chen *et al.*⁴⁵ at two different temperatures: $0.9T_m$ (approximately 90% of the melting temperature T_m of the alloy) and 1350 (K). These predictions are obtained via a high-throughput computational workflow, which employs a regular-solution model^{46,47} using binary interaction param-



TABLE II. Summary of alloy datasets used in evaluation experiments. No. alloys: Total number of alloys present in each dataset. No. positive label: Number of alloys classified as forming HEA phases in datasets $\mathcal{D}_{0.9T_m}$ and \mathcal{D}_{1350K} , the number of alloys exhibiting non-zero magnetization in \mathcal{D}_{Mag} , and the number of alloys with a non-zero Curie temperature in \mathcal{D}_{T_C} . The percentage values in parentheses represent the proportion of positive labels within each dataset.

Dataset	No. alloys	Physical properties	Positive label	No. positive label
$\mathcal{D}_{0.9T_m}$	14,950 quaternary alloys	Stability	HEA	4,218 (28%)
\mathcal{D}_{1350K}	14,950 quaternary alloys	Stability	HEA	1,402 (9%)
\mathcal{D}_{Mag}	5,968 quaternary alloys	Magnetization (T)	Magnetic	2,428 (41%)
\mathcal{D}_{T_C}	5,968 quaternary alloys	Curie temperature (K)	Non-zero Curie Temperature	2,355 (39%)
\mathcal{D}_{HEA}^{exp}	55 quaternary alloys	Stability	HEA	40 (73%)
\mathcal{D}_{HEB}^{exp}	19 quinary alloys-borides	Stability	HEB	15 (79%)

eters derived from *ab initio* density functional theory (DFT) to compute and compare Gibbs free energies of solid solutions against competing inter-metallic phases^{16–18}.

- \mathcal{D}_{Mag} and \mathcal{D}_{T_C} : These computational datasets comprise 5,968 quaternary high-entropy alloys (HEAs)³⁵, each formed by selecting four elements from a set of 21 transition metals: Fe, Co, Ir, Cu, Ni, Pt, Pd, Rh, Au, Ag, Ru, Os, Tc, Re, Mn, Ta, W, Mo, Cr, V, and Nb. Their magnetizations (\mathcal{D}_{Mag}) and Curie temperatures (\mathcal{D}_{T_C}) in the body-centered cubic (BCC) phase are computed using the Korringa–Kohn–Rostoker coherent approximation method⁴⁸. These datasets are derived from an original pool of 147,630 equiatomic quaternary HEAs.
- \mathcal{D}_{HEA}^{exp} : The experimental dataset includes 55 experimentally verified quaternary HEAs from peer-reviewed publications^{45,49,50}. The dataset includes both HEA (40 alloys) and non-HEA (15 alloys) compositions, providing balanced representation for validation.
- \mathcal{D}_{HEB}^{exp} : The experimental dataset includes 19 experimentally verified quinary HEBs from peer-reviewed publications⁴⁴. The dataset includes 15 quinary systems forming HEB.

B. Design of experiments

We begin by verifying the reliability of the elemental substitutability knowledge queried from large language models (LLMs). Specifically, we compare the LLM-derived substitutability knowledge with the well-established Hume–Rothery criteria for elemental substitution.

With that reliability confirmed, we turn to predictive capability. Two experiments on four computational datasets serve as the framework’s proving ground to evaluate predictive capability of our proposed framework: (1) Cross-validation on quaternary alloys, assessing performance with randomly partitioned training sets (1%-30%

of data) to determine how effectively LLM-derived knowledge aligns with material-specific relationships across different data availability scenarios, with particular focus on data-limited conditions; and (2) Extrapolation on quaternary alloys, simulating real discovery scenarios by excluding alloys containing a specific element from training and evaluating performance on compositions that incorporate this previously unseen element. These computational datasets, free from experimental bias and large enough for robust statistics, provide the controlled environment needed for *framework development*.

To benchmark our multi-source method, we compare its predictive performance against two baseline approaches.

- **Single-source methods:** These methods rely exclusively on one source of evidence, either a material dataset or domain knowledge derived from only one LLM from the set of state-of-the-art models under investigation.
- **Traditional classification method:** We employ logistic regression (LR)⁵¹.

Hyper-parameters of these methods are tuned via systematic grid search, as detailed in Supplementary Section 1. Hereinafter, we define models employing the evidential method (based on the Dempster–Shafer theory) as follows: models trained solely on material datasets are termed DS-source models; those leveraging evidence from LLMs are termed LLM-source models; and those integrating both sources are termed multi-source models. Notably, the LLM-source models are obtained by combining 20 independent sources—each of the 4 LLMs (GPT-4o, GPT-4.5, Claude Opus 4, Grok3) queried across 5 scientific domains—through Dempster–Shafer theory (Section II C). The multi-source model further integrates this combined LLM-source with the DS-source using the same framework. Models utilizing logistic regression and support vector machines are referred to as LR-based model.

To assess the real-world applicability of our framework, we next validate its predictive performance on experimentally verified alloys. This validation examines whether the proposed framework can accurately predict phase stability for experimentally synthesized alloys. Our framework integrates LLM-derived knowledge with

TABLE III. Confusion matrix comparing LLM consensus predictions with Hume–Rothery rules for 351 element pairs considered in this study.

LLMs	Hume–Rothery rules		
	Substitutable	Non-substitutable	Total
	33 pairs (<i>True positive</i>)	45 pairs (<i>False positive</i>)	78 pairs
	4 pairs (<i>False negative</i>)	269 pairs (<i>True negative</i>)	273 pairs
	Total 37 pairs	314 pairs	351 pairs

substitutability patterns extracted from computational datasets. This reflects real-world scenarios where researchers must consider all available knowledge to fill the gaps raised by limited experimental data before selecting candidates for expensive synthesis. Finally, after evaluating the predictive performance across all settings, we analyze the element substitutability patterns captured using the multi-source approach to gain deeper insights into the underlying HEA formation mechanisms of quaternary alloys.

C. Materials descriptors

Descriptors, which are the representation of alloys, play a crucial role in building a recommender system to explore potential new HEAs. In this research, the raw data of alloys is represented in the form of element combinations. Several descriptors have been studied in materials informatics to represent the compounds⁵². To employ the data-driven approaches for this work, we applied compositional descriptor⁵³ and binary elemental descriptor.

Compositional descriptors represent each alloy through 135 features derived from 15 atomic properties of constituent elements. These properties include structural parameters (*atomic number, mass, period, and group*), electronic characteristics (*first ionization energy, second ionization energy, Pauling electronegativity and Allen electronegativity*), size factors (*van der Waals, covalent, and atomic radii*), and thermophysical properties (*melting point, boiling point, density, specific heat*). For each atomic property, we calculate statistical numbers, including mean, standard deviation, and pairwise covariances across the alloy's elements, to represent the alloy. The compositional descriptors can be applied not only to crystalline systems but also to molecular systems. However, the descriptors cannot easily distinguish alloys with different numbers of constituent elements, because they treat the atomic properties as statistical distributions. Therefore, the descriptors cannot be applied when extrapolating to alloys with a different number of components.

Binary elemental descriptors use binary encoding to indicate element presence (1) or absence (0) in an alloy. The number of binary elemental descriptors corresponds to the number of element types included in the train-

ing data. In this study, the binary elemental descriptors are used to represent the alloys in the DS-source, LLM-source, and multi-source models. In contrast, the compositional descriptors are applied for the LR-based model.

IV. RESULTS AND DISCUSSIONS

A. Reliability Assessment of LLM-Based Elemental Substitutability Knowledge

Verifying the reliability of large language model (LLM) responses is a prerequisite for trusting downstream predictions. We therefore validate element-substitutability knowledge extracted from LLM queries against the empirical Hume–Rothery rules⁵⁴, which are a set of basic rules for predicting elemental substitution. These rules stipulate that elements readily substitute in solid solutions when: (i) atomic radius mismatch is lower than 15%, (ii) they share similar crystal structures and valence states, and (iii) they have similar electronegativity. When electronegativity differences exceed critical thresholds, metals typically form intermetallic compounds rather than solid solutions. For this validation, we use an electronegativity difference threshold of 0.55. For valency comparison in metallic alloy systems, we consider the effective valency⁵⁵ (number of electrons effectively contributing to metallic cohesion). While most metals exhibit a single characteristic valency, certain transition metals (e.g., Fe, Co, Mn, Cr) can exhibit multiple effective valencies in different alloy environments. In our analysis, two elements are considered to have similar valency if they share at least one common valence state.

We aggregated substitutability assessments from four LLMs, including GPT-3, Claude Opus 4, GPT-4o, and GPT-4.5, for 351 element pairs using our DST framework. Each pair is classified as substitutable if the combined belief for substitutability exceeds that for non-substitutability. Comparison against Hume–Rothery predictions reveals strong alignment: 86% of element pairs show identical classifications with high recall rates for substitutable labels and high precision for non-substitutable labels, as shown in Table III. Specifically, 33 of 37 pairs (89%) deemed substitutable by Hume–Rothery rules are correctly identified by LLMs, while



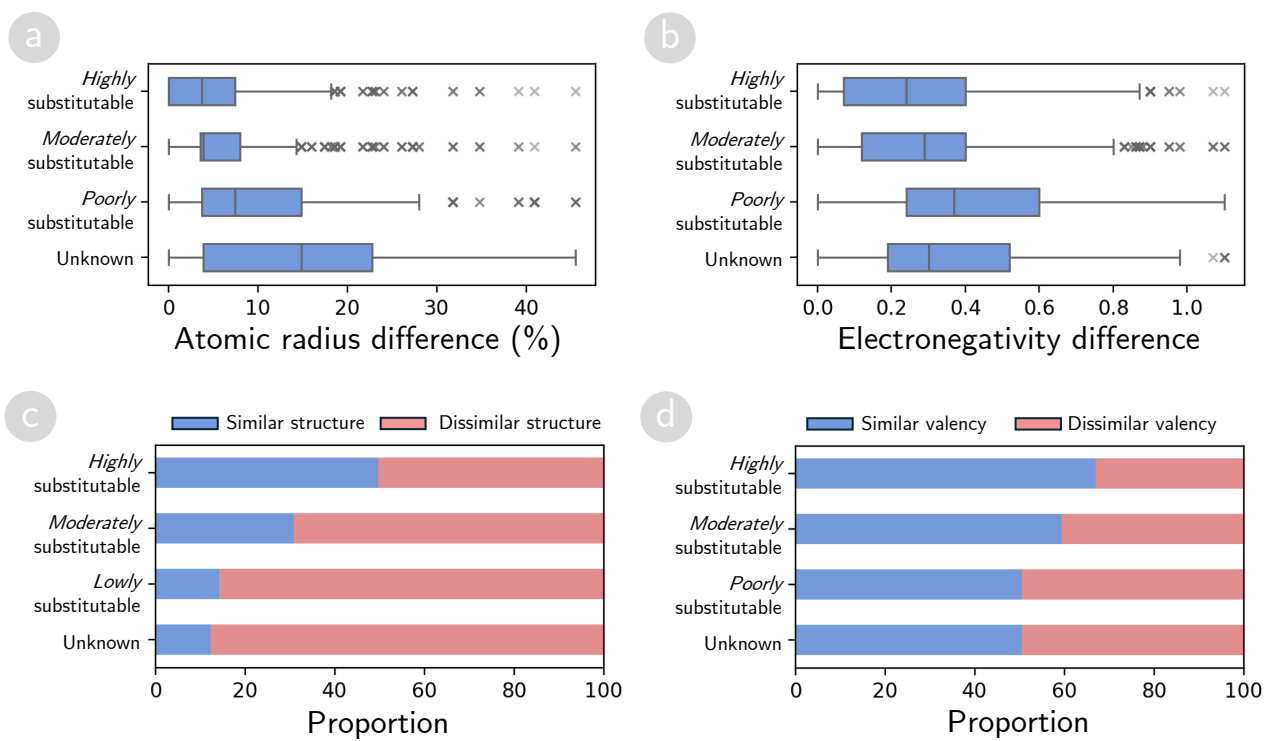


FIG. 3. Validation of LLM-extracted substitutability against Hume–Rothery rules. (a, b) Distribution of atomic radius differences (a) and electronegativity differences (b) for element pairs categorized by LLM-predicted substitutability levels (highly, moderately, and poorly substitutable, plus unknown). Box plots show median, interquartile range, and outliers. (c, d) Proportions of element pairs with similar versus dissimilar crystal structures and valency, grouped by substitutability levels.

Open Access Article. Published on 19 December 2025. Downloaded on 1/14/2026 8:34:16 PM.
This article is licensed under a Creative Commons Attribution 3.0 Unported Licence.



269 of 273 pairs classified as non-substitutable by LLMs matched Hume–Rothery rules, achieving a precision of 99%.

The 14% misalignment consists entirely of cases where LLMs identify additional substitutable pairs beyond the traditional Hume–Rothery criteria. Among the 45 misaligned pairs, most satisfy the size and electronegativity requirements but exceed traditional thresholds for valency or crystal structure differences. Remarkably, experimental validation supports these context-specific predictions: 14 of these pairs have been confirmed to form single-phase binary systems⁵⁶, as shown in Supplementary Table 3. Additionally, Cr and Nb differ in valence electron counts (Cr: 6, Nb: 5), placing them outside general substitutability criteria. However, when incorporated into quaternary systems, they demonstrate successful substitution—Cr in quaternary system Cr–Al–Ti–V can be replaced by Nb (forming Nb–Al–Ti–V), and similarly in Cr–Ta–Ti–V and Nb–Ta–Ti–V systems, both form stable single-phase BCC structures.

This asymmetric difference reflects a fundamental distinction between general rules and context-specific knowledge. The Hume–Rothery rules, developed through careful empirical observation, provide general guidelines with well-defined thresholds (e.g., 15% for radius difference) that have successfully guided alloy design for

decades. These universal criteria ensure high reliability across diverse alloy systems. In contrast, LLMs capture context-dependent substitutability documented in materials literature⁵⁷, in which specific processing conditions, alloy compositions, or applications enable successful substitution despite exceeding general thresholds. LLMs integrate knowledge from documented experimental systems across material families for general substitutability assessment, explaining why they complement conservative Hume–Rothery rules with context-specific insights. Detailed analysis of all 45 pairs with experimental validation status is provided in Supplementary Table 3.

Figure 3 analyzes in detail the alignment of LLM’s response with each criterion of substitutability from Hume–Rothery rules. Element pairs that LLMs identified as highly substitutable exhibit significantly lower atomic radius differences and electronegativity differences compared to pairs identified as poorly substitutable, as shown in Figure 3(a–b). Additionally, highly substitutable pairs predominantly share similar crystal structures and valencies, while poorly substitutable pairs rarely do as shown in Figure 3(c–d).

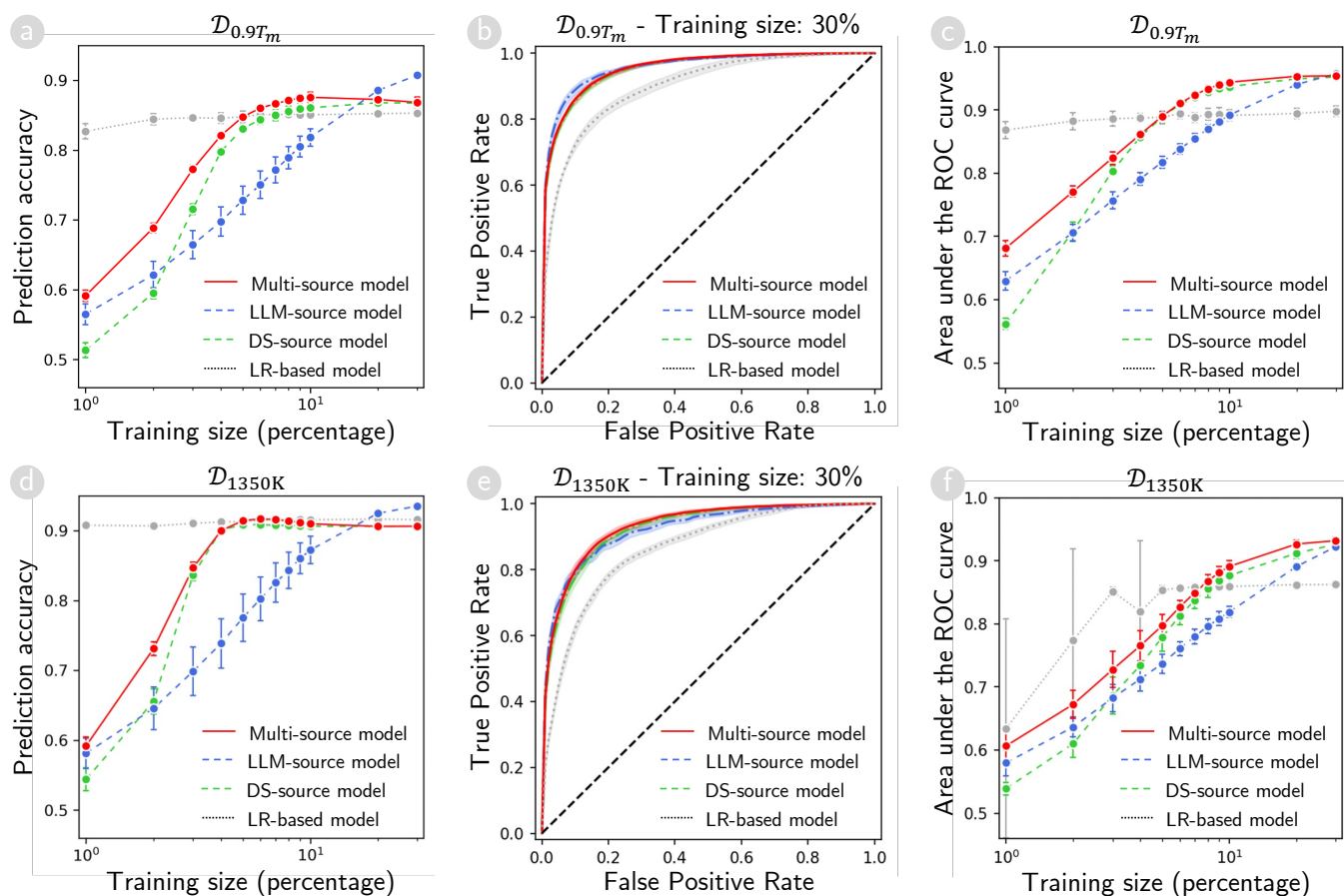


FIG. 4. **Predictive capability evaluation via cross-validation on quaternary-alloy datasets $\mathcal{D}_{0.9T_m}$ and \mathcal{D}_{1350K} .** (a, d) Classification accuracy of the multi-source, single-source, and LR-based models on two quaternary alloy datasets $\mathcal{D}_{0.9T_m}$ and \mathcal{D}_{1350K} . (b, e) Receiver operating characteristic (ROC) curves for the same models at a 30% training-set size on these datasets. (c, f) Area under the ROC curves (AUC) for each model across different training-set sizes, providing an overall measure of discriminative performance. In all subplots, red lines indicate the multi-source model (using both DS and LLM sources), green and blue lines represent single-source models (using either DS or LLM sources), and gray lines represent the LR-based model.

B. Cross-Validation Analysis of Multi-Source Knowledge Integration

For the experiment, we systematically vary the training set size from 1% to 30% of each quaternary-alloy dataset, incrementing by 1% up to 10%, followed by steps of 20% and 30%. The variation enables the assessment of how different methods handle data scarcity versus model availability.

Figures 4(a,d) and 5(a,d) show the classification accuracy of the single-source, multi-source, and LR-based models on the four datasets. At smaller training sizes (approximately 1%–10%), the LR-based model achieves the highest overall accuracy, outperforming evidential models, which explicitly model element substitutability to predict alloy properties. Among the evidential models, single-source LLM models initially outperform DS-source models, attributed to LLM-derived domain-specific insights that assist in mitigating data limita-

tions. However, multi-source models remain competitive and sometimes achieve the highest accuracy among evidential models, even with limited data. As the training size exceeds 10%, DS-source models exhibit superior performance on the magnetization and Curie temperature datasets while achieving comparable accuracy to LLM-source models on alloy stability datasets. Conversely, the accuracy of LR-based models plateaus and is eventually outperformed by evidential models. These findings underscore the importance of incorporating LLM-based, DS-source, or multi-source knowledge to improve quaternary-alloy property predictions.

Although prediction accuracy provides a convenient single-metric overview, it relies on a fixed classification threshold (typically 0.5), which may not be optimal for imbalanced datasets, where HEAs (positive class) are relatively rare. Under these conditions, LR-based models may serve effectively at extremely small training sizes when they effectively predict the dominant (Non-HEA) class by default, thereby inflating accuracy. However,



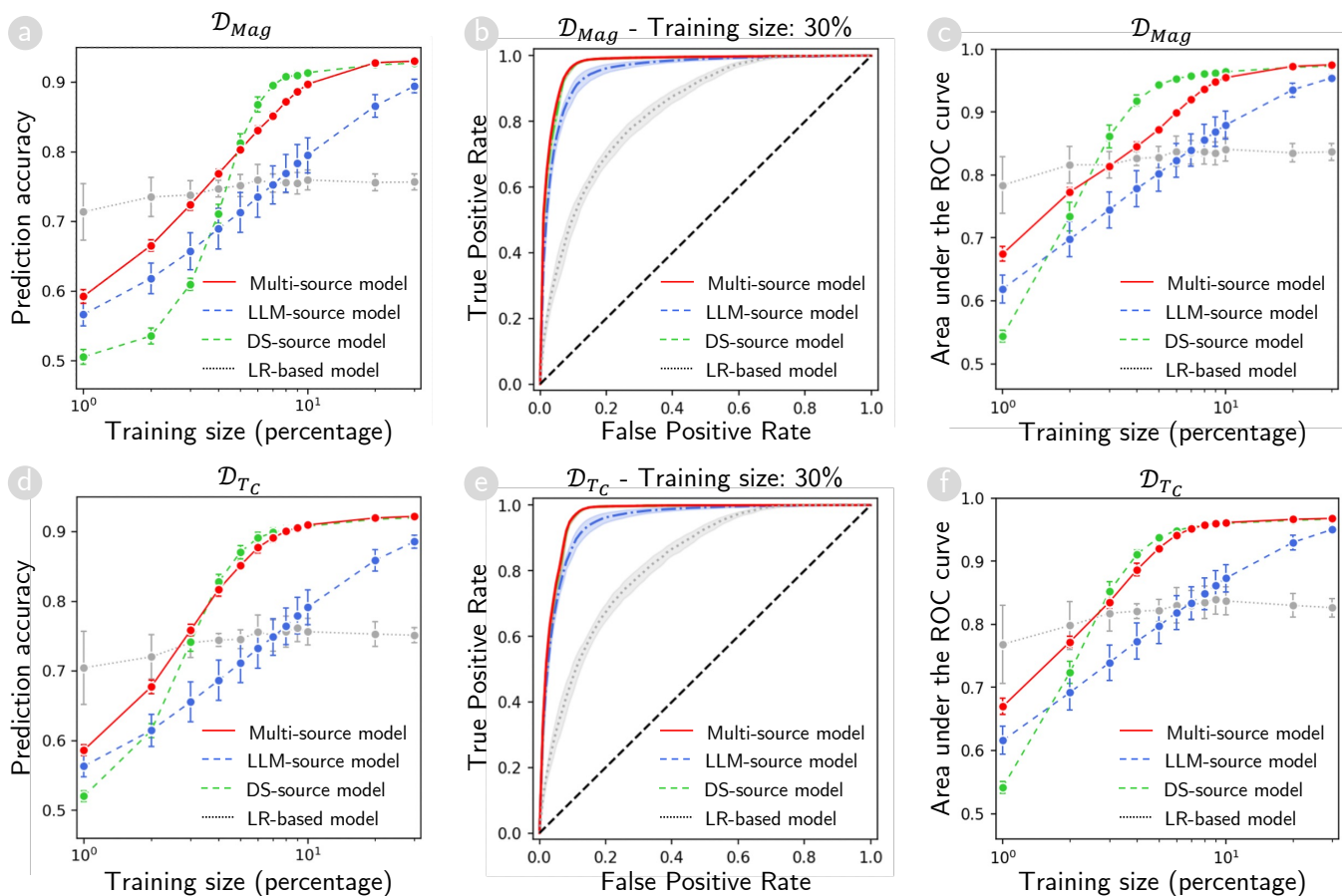


FIG. 5. **Predictive capability evaluation via cross-validation on quaternary-alloy datasets \mathcal{D}_{Mag} and \mathcal{D}_{T_C} .** (a, d) Classification accuracy of the multi-source, single-source, and LR-based models on two quaternary alloy datasets \mathcal{D}_{Mag} and \mathcal{D}_{T_C} . (b, e) Receiver operating characteristic (ROC) curves for the same models at a 30% training-set size on these datasets. (c, f) Area under the ROC curves (AUC) for each model across different training-set sizes, providing an overall measure of discriminative performance. In all subplots, red lines indicate the multi-source model (using both DS and LLM sources), green and blue lines represent single-source models (using either DS or LLM sources), and gray lines represent the LR-based model.

1 this approach fails to address scenarios where different
2 types of misclassifications (false positives versus false neg-
3 atives) incur different costs.

4 To effectively capture these trade-offs under dynamic
5 thresholds, we analyze receiver operating characteristic
6 (ROC) curves across the four datasets, which illustrate
7 variations in true positive rate (TPR) and false positive
8 rate (FPR) of each model across all possible decision
9 boundaries. Figures 4(b,e) and 5(b,e) depict the ROC
10 curves for the multi-source models, LLM-source models,
11 DS-source models, and LR-based models at a 30% train-
12 ing size. Overall, the multi-source and DS-source mod-
13 els exhibit comparable ROC performance and outper-
14 form the other models. The LLM-source models achieve
15 results comparable to the best ones on the alloy sta-
16 bility datasets $\mathcal{D}_{0.9T_m}$ and \mathcal{D}_{1350K} but lag behind DS-
17 source models on the magnetization and Curie tempera-
18 ture datasets \mathcal{D}_{Mag} and \mathcal{D}_{T_C} . Therefore, knowledge col-
19 lected from the five considered research domains may
20 not fully capture the magnetic and thermal properties

21 reflected in those datasets. Meanwhile, the LR-based
22 models consistently show the lowest performance across
23 all four datasets.

24 To further assess the ROC performance of each model
25 at different training sizes, we analyze the AUC distri-
26 bution from 1% to 30% training data, as shown in Fig-
27 ures 4(c,f) and 5(c,f). When the training set is extremely
28 small, LLM-based models generally attain an early ad-
29 vantage, presumably because domain insights compen-
30 sate for limited alloy observations. However, as data ac-
31 cumulates, DS-source models typically outperform LLM-
32 source models, suggesting that direct data-driven cues
33 from quaternary-alloy datasets become increasingly de-
34 cisive. In contrast, multi-source models maintain robust
35 performance across all training sizes, benefitting from
36 their ability to merge domain-specific substitutability in-
37 sights with empirical data. Multi-source models leverage
38 complementary evidence, enabling an effective balance
39 between TPR and FPR. On stability datasets $\mathcal{D}_{0.9T_m}$ and
40 \mathcal{D}_{1350K} , DS-source and multi-source models achieve com-



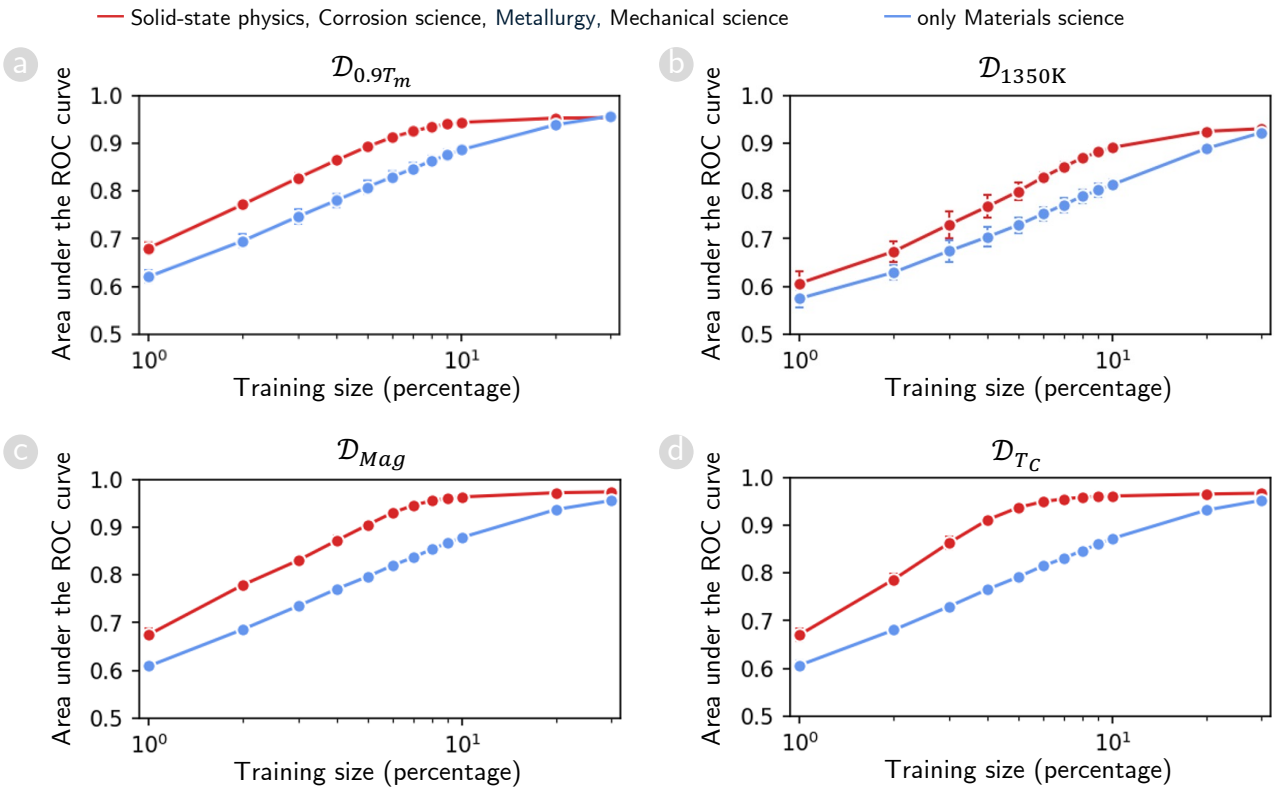


FIG. 6. **Performance comparison of explicit versus implicit domain integration.** Area under ROC curves for predicting HEA stability ($\mathcal{D}_{0.9T_m}$, \mathcal{D}_{1350K}) and magnetic properties (\mathcal{D}_{Mag} , \mathcal{D}_{T_C}) using two domain integration strategies: (i) systematic combination of four specialized domains (solid-state physics, corrosion science, metallurgy, materials mechanics) shown in red, (ii) only using materials science, which serves as an integrative field that synthesizes perspectives from four specialized domains, shown in blue.

1 parable AUC early on and remain highly competitive as
2 training data accumulates. For magnetization and Curie-
3 temperature datasets, DS-source models briefly outper-
4 form multi-source models at moderate training sizes (ap-
5 proximately 6–20%), but this gap diminishes at larger
6 training sizes.

7 We note that the LLM-derived substitutability ma-
8 trix \mathbf{M} remains fixed across all training sizes (LLMs
9 are used out-of-the-box without retraining); improved
10 performance with larger training sets results from hav-
11 ing more host compositions available to apply this fixed
12 knowledge through substitution-based inference (Sec-
13 tion II.D). This explains why LLM-source and multi-
14 source models benefit from increased training data de-
15 spite the LLM knowledge itself remaining unchanged.

16 Figure 6 provides compelling evidence for the effective-
17 ness of our systematic evidence combination approach
18 compared to relying on materials science as an integra-
19 tive domain that synthesizes perspectives from the other
20 four domains. Significantly, using only materials science
21 knowledge yields substantially lower performance by 10-
22 20% across all datasets than our multi-source framework,
23 which systematically combines evidence from the four

24 specialized domains, across different prediction tasks.
25 This performance gap demonstrates the fundamental ad-
26 vantage of our Dempster–Shafer-based approach: while
27 materials science provides a static, pre-integrated per-
28 spective that may obscure domain-specific nuances, our
29 framework preserves distinct domain insights and adap-
30 tively weights them based on their alignment with target
31 properties. The superior performance of our systematic
32 combination method validates that explicit, property-
33 aware evidence synthesis outperforms implicit knowledge
34 fusion, particularly when different domains contribute
35 varying degrees of relevant information for specific mate-
36 rial properties such as stability, magnetization, or Curie
37 temperature.

38 While LLM-source models generally perform well, our
39 results reveal two scenarios where they potentially un-
40 derperform compared to data-driven approaches.

1. *Property-specific predictions with weak domain alignment:* For magnetic property datasets (\mathcal{D}_{Mag} , \mathcal{D}_{T_C}), DS-source substantially outperforms LLM-source, showing a larger performance gap than observed for phase stability datasets (Figures 4 and 5). The five selected domains (corrosion



TABLE IV. Prediction accuracy and Areas under the receiver operating characteristic (ROC) curves of various methods on quaternary-alloy datasets in extrapolation experiments. For each dataset, alloys containing a specific element e are systematically excluded from the training set and used exclusively for testing. Results are reported as mean accuracy and mean AUC, averaged across all elements e within each dataset, with standard deviations reflecting variability across elements.

Evaluation criteria	Methods	$\mathcal{D}_{0.9T_m}$	\mathcal{D}_{1350K}	\mathcal{D}_{Mag}	\mathcal{D}_{T_C}
Prediction accuracy	Multi-source model	0.86 ± 0.06	0.92 ± 0.04	0.86 ± 0.19	0.86 ± 0.18
	LLM-source model	0.84 ± 0.09	0.90 ± 0.09	0.81 ± 0.21	0.86 ± 0.18
	DS-source model	0.50 ± 0.04	0.51 ± 0.05	0.48 ± 0.07	0.50 ± 0.10
	LR-based model	0.83 ± 0.05	0.91 ± 0.04	0.67 ± 0.15	0.68 ± 0.13
Area under ROC curves	Multi-source model	0.93 ± 0.06	0.92 ± 0.08	0.95 ± 0.06	0.94 ± 0.07
	LLM-source model	0.91 ± 0.11	0.90 ± 0.12	0.95 ± 0.06	0.94 ± 0.07
	DS-source model	0.50 ± 0.00	0.50 ± 0.00	0.50 ± 0.00	0.50 ± 0.00
	LR-based model	0.85 ± 0.11	0.82 ± 0.10	0.84 ± 0.06	0.84 ± 0.06

science, materials mechanics, metallurgy, solid-state physics, materials science) were optimized for structural stability and do not adequately capture magnetic exchange interactions or spin configurations.

2. *Data-rich regimes:* At large training sizes ($>20\%$, Figures 4 and 5), DS-source matches or exceeds LLM-source performance across all datasets. When sufficient data exists, empirical patterns extracted directly from the dataset provide adequate information, and general domain knowledge offers minimal additional value.

In conclusion, LLM-source models excel in data-scarce scenarios by leveraging domain-specific insights to mitigate sparsity-related challenges. As data availability increases, DS-source models outperform LLM-source models, particularly where DS-derived evidence provides sufficient information for a purely data-driven learning approach. Multi-source models, which integrate insights derived from LLM and DS-sources, demonstrate robust and consistent performance across various training sizes.

C. Extrapolation Analysis of Multi-Source Knowledge Integration

Having assessed the proposed framework via cross-validation (Section IV B), we examine its *extrapolation* performance on quaternary alloys containing an element e , which is excluded during training. Unlike the cross-validation experiments, the training set size is not varied for this set of experiments. Instead, for each element e , we remove all e -containing alloys from the dataset and train each model on the remaining alloys that do not contain e . Further, we evaluate the ability of the models to predict the properties of e -containing alloys. This procedure tests whether the learned models can generalize to compositions containing unseen elements in their training datasets.

Table IV reveals distinct performance patterns across model types. DS-source models fail in this scenario,

achieving ~ 0.50 accuracy (random guessing) across all datasets because they cannot extract substitutability patterns for absent element e from training data. In contrast, LLM-source models achieve substantially higher accuracies across all datasets. Multi-source models mostly outperform LLM-source on phase stability datasets ($\mathcal{D}_{0.9T_m}$ and \mathcal{D}_{1350K}) but achieve nearly identical performance on magnetic property datasets (\mathcal{D}_{Mag} and \mathcal{D}_{T_C}).

This convergence of multi-source and LLM-source performance on magnetic datasets reflects proper uncertainty handling rather than a limitation. When element e is absent from training, DS-source has no observed substitutability patterns involving e . Following the principle established in Section II A, DS-source assigns unit mass to the uncertainty set, explicitly representing total ignorance about e -containing compositions. When this total uncertainty combines with confident LLM evidence through Dempster's rule (Equation 7), the final multi-source prediction is naturally dominated by informative LLM knowledge. The framework thus explicitly represents *unknown* rather than forcing unreliable predictions from insufficient data, demonstrating principled uncertainty quantification in extrapolation scenarios.

Figure 7 illustrates the ROC curves, showing that the multi-source and LLM-source models consistently exhibit higher TPR at comparable FPR across all datasets. Conversely, DS-source models exhibit near-random discrimination, as evidenced by their diagonal ROC curves, while LR-based models yield moderate performance between these extremes. To quantify these visual differences, Table IV also lists AUC for each dataset. Multi-source models achieve the highest AUC scores (0.92–0.95), followed closely by LLM-source models (0.90–0.95), while LR-based models peak at approximately 0.85, and DS-source models hover at approximately 0.50.

Figure 8a–c illustrates knowledge integration in extrapolation simulations for Os-based alloys using the $\mathcal{D}_{0.9T_m}$ dataset. Specifically, Figures 8a and 8b present maps reconstructed from element substitutability patterns derived from the DS-source and multi-source models, respectively, both trained on $\mathcal{D}_{0.9T_m}$ dataset excluding Os-based alloys. Details of the visualization method



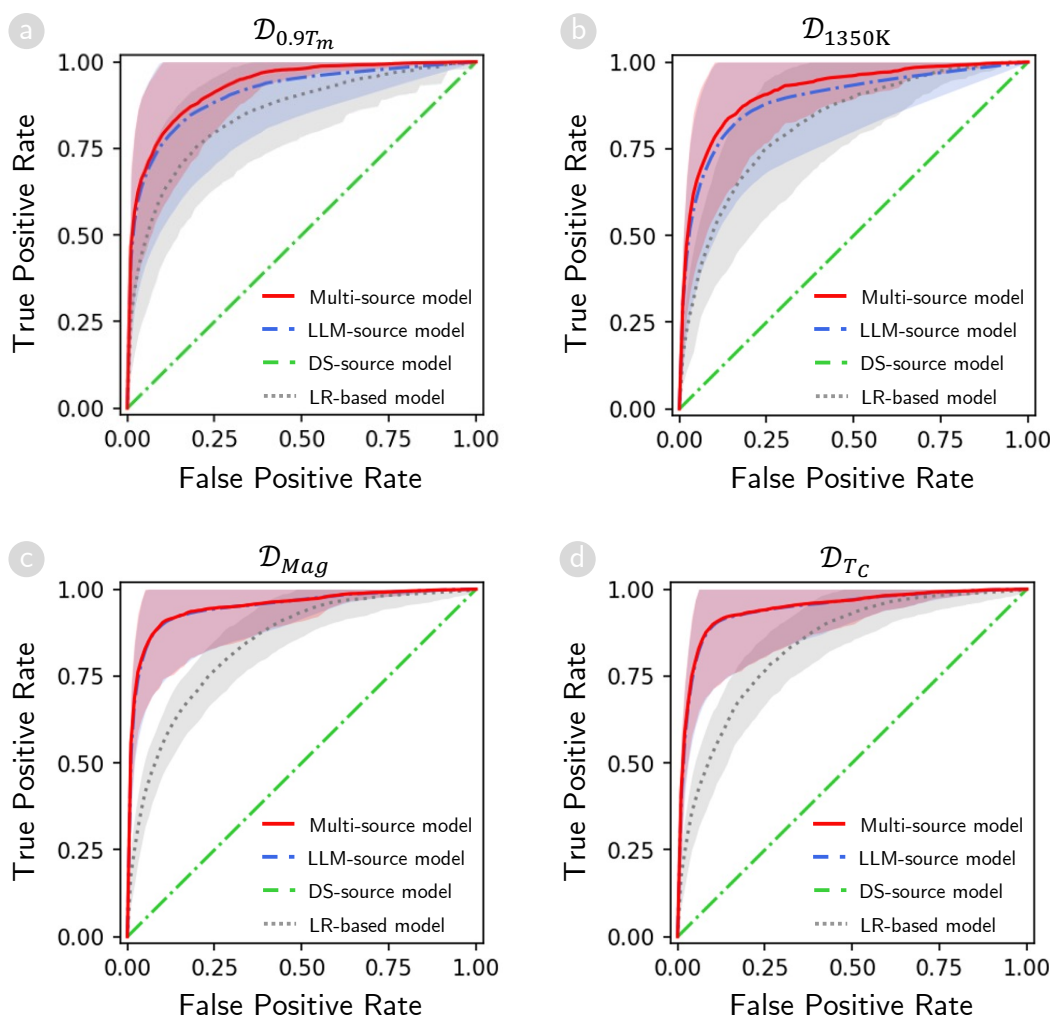


FIG. 7. **Predictive capability evaluation via extrapolation on quaternary-alloy datasets.** For each dataset, alloys containing a specific element e are systematically excluded from the training set and used exclusively for testing. (a–d) Area under the receiver operating characteristic (ROC) curves (AUC) is plotted for each model on their respective test sets in the extrapolation experiments. In all subplots, red lines represent the multi-source model (integrating both DS and LLM sources), green and blue lines represent single-source models (using either DS or LLM sources), and gray lines represent the LR-based model.

are shown in Supplementary Section 4. In these visualizations, the observed alloys are well-structured into sub-clusters according to their phase formation behavior, with blue markers indicating HEA-forming alloys and red markers representing non-HEA alloys. The Os-based candidate alloys, depicted as white circular markers, consistently form a distinct sub-cluster in the upper region of each map. In these visualizations, the background coloration indicates the predicted probability of HEA formation, with deeper blue regions suggesting higher probability of forming stable HEAs.

The limitations of the DS-source model become evident in Figure 8a, where the phase behavior of Os-based alloys remains undetermined due to the absence of Os-

containing alloys in the training dataset. This knowledge gap leaves researchers with no guidance when exploring the uncharted territory of Os-based alloys, forcing them to rely on random selection. In contrast, our multi-source approach addresses this limitation by integrating expert insights distilled from scientific literature using LLMs, as illustrated in Figure 8b. The effectiveness of this approach is visually confirmed in Figure 8c, where the multi-source model's predictions closely align with the actual phase behavior of the candidates. This qualitative assessment is complemented by quantitative evaluation in Supplementary Table 4, which reports that the multi-source model achieves an impressive 88% prediction accuracy for Os-based alloys, validating our approach's ca-



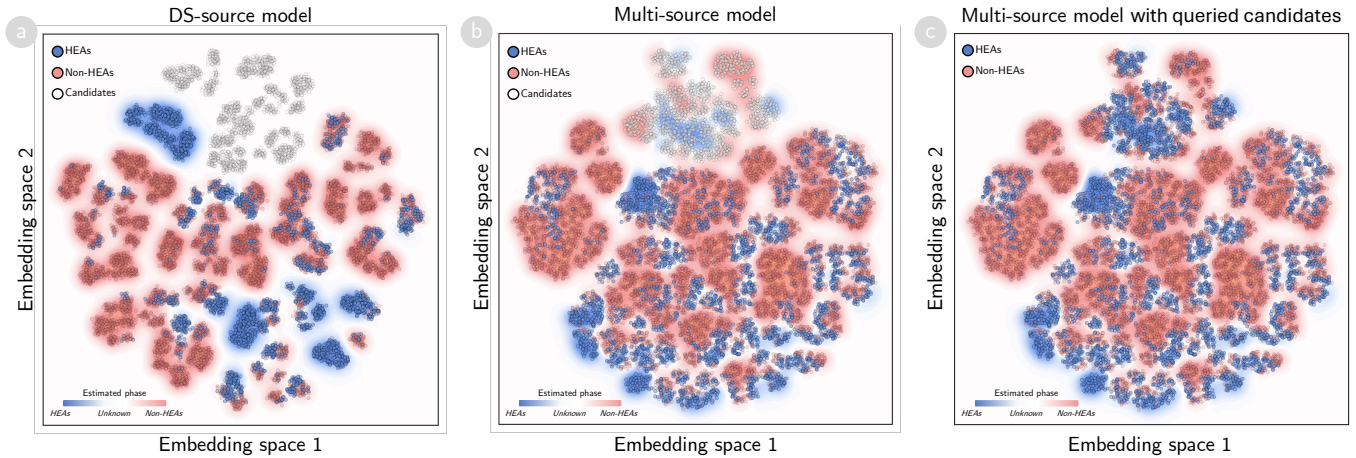


FIG. 8. **Visualization of Os-based alloy extrapolation in dataset $D_{0.9T_m}$.** (a) Alloy map generated from element substitutability patterns extracted using the DS-source model after excluding Os-based alloys from training. (b-c) Alloy maps generated from element substitutability patterns extracted using the multi-source model after excluding Os-based alloys from training. The map in (c) incorporates queried labels for Os-based candidate alloys. Marker colors represent phase formation: blue for HEA alloys, red for non-HEA alloys, and white for Os-based candidate alloys. Background coloration indicates the predicted phase formation probability according to the DS-source model (a) and multi-source model (b-c), with deeper blue shades suggesting higher probability of HEA formation.

1 pability to effectively extrapolate to unexplored compo-
2 sitional spaces. In summary, these results confirm that
3 leveraging multi-source or LLM-based evidence signifi-
4 cantly enhances discriminative power in the extrapola-
5 tion scenario.

6 D. Effectiveness Assessment on Experimental High-Entropy 7 Alloy Data

8 To assess the real-world applicability of our frame-
9 work, we validated its performance on experimentally
10 verified alloys from the literature. This validation ex-
11 amines whether the proposed framework, developed pri-
12 marily using computational datasets, can accurately pre-
13 dict phase stability for experimentally synthesized al-
14 loys. Our framework integrates LLM-derived knowledge
15 with substitutability patterns extracted from computa-
16 tional databases using the methodology described in Sec-
17 tion II A. This reflects real-world scenarios where re-
18 searchers must consider all available knowledge before
19 selecting candidates for expensive synthesis.

20 We performed 5-fold cross-validation on experimental
21 datasets: $\mathcal{D}_{\text{HEA}}^{\text{exp}}$ of 55 experimentally confirmed alloys.
22 For the HEA dataset $\mathcal{D}_{\text{HEA}}^{\text{exp}}$, we integrated LLM knowl-
23 edge with substitutability patterns extracted from com-
24 putational datasets $\mathcal{D}_{1350\text{K}}$, $\mathcal{D}_{\text{AFLOW}}$, $\mathcal{D}_{\text{CALPHAD}}$, and
25 $\mathcal{D}_{\text{LTVC}}$. Details of the computational datasets are in-
26 troduced in the Supplementary Section 6. Notably, the
27 predictions from these computational methods for the 55
28 experimentally confirmed alloys are not utilized in our
29 framework training, ensuring unbiased validation.

30 For benchmarking on the HEA dataset, we compared

31 our framework against four empirical rules (ERs)^{58–61},
32 two free-energy models (FEM)^{3,62}, and a valence-electron
33 concentration (VEC) model⁶³. Supplementary Table 2
34 provides details of these baseline models. Addition-
35 ally, we compared our framework with results obtained
36 from computational datasets $\mathcal{D}_{\text{AFLOW}}$ ¹⁵, $\mathcal{D}_{\text{LTVC}}$ ¹⁹, and
37 $\mathcal{D}_{1350\text{K}}$ ⁴⁵. These computational datasets are collected
38 by using high-throughput approaches and Hamiltonian
39 models.

40 Figure 9a presents ROC curves demonstrating that
41 our multi-source integration framework consistently out-
42 performs empirical phase selection models such as ERs,
43 FEMs, and VEC, while achieving performance compa-
44 rable to costly computational methods. These results
45 confirm that systematically integrating diverse evidence
46 sources through our DST framework enhances prediction
47 accuracy across different material classes. The frame-
48 work's value does not lie in replacing established meth-
49 ods but in effectively combining their complementary
50 strengths, creating a unified platform that enhances prac-
51 tical decision-making in materials discovery.

52 To investigate the underlying mechanisms of forming
53 HEAs, we analyzed the elemental substitutability pat-
54 terns extracted by our framework from multiple evidence
55 sources. Specifically, we integrated substitutability infor-
56 mation from the experimental dataset $\mathcal{D}_{\text{HEA}}^{\text{exp}}$, computa-
57 tional datasets ($\mathcal{D}_{1350\text{K}}$, $\mathcal{D}_{\text{AFLOW}}$, $\mathcal{D}_{\text{CALPHAD}}$, $\mathcal{D}_{\text{LTVC}}$),
58 and LLM-derived knowledge.

59 Figure 9b presents the substitutability matrix for 26
60 elements relevant to HEA stability, along with their hi-
61 erarchical clustering structure. The dendrogram is gen-
62 erated via hierarchical agglomerative clustering (HAC)
63 with the complete linkage criterion, grouping elements



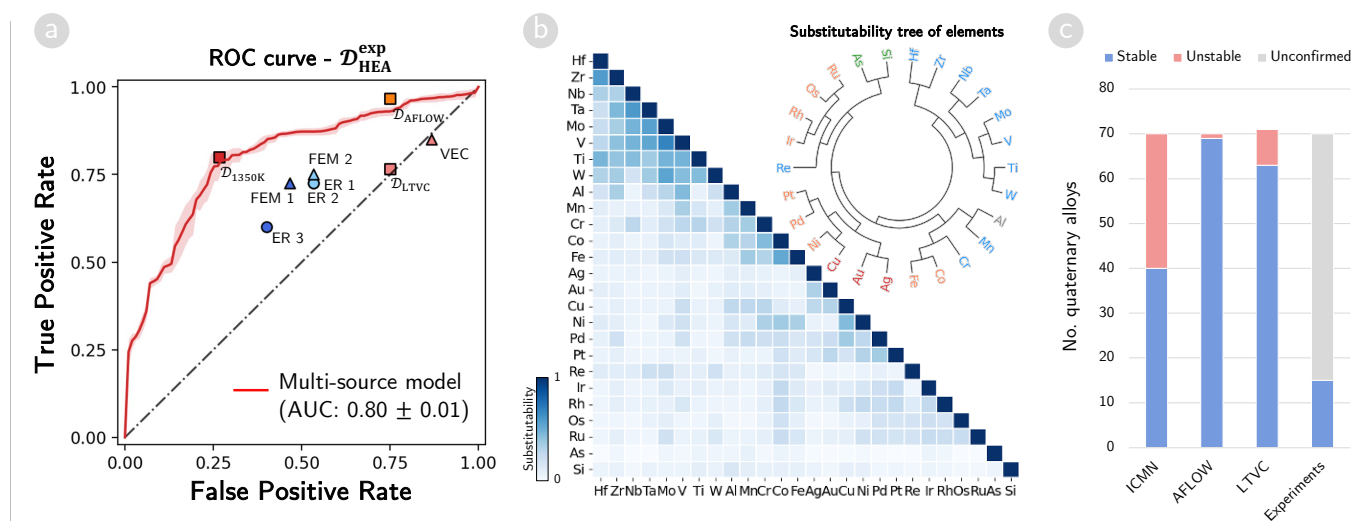


FIG. 9. **Effectiveness Assessment of Multi-Source Knowledge Integration for High-Entropy Alloy Formation.** (a) Receiver operating characteristic (ROC) curves for the phase estimation task on experimental dataset $\mathcal{D}_{\text{HEA}}^{\text{exp}}$. Red line represent the multi-source model (integrating both DS and LLM sources) and gray dashed line represent the random selection. Coloured scatter points represent results of ERs, FEMs, VEC, and computational methods that return only a single stable/unstable estimation. (b) Substitutability matrix and substitutability tree for 26 elements. Matrix values represent substitutability scores derived from integrated computational datasets, experimental dataset and LLM sources. The substitutability tree is generated using hierarchical agglomerative clustering with complete linkage criterion. Element colors: blue (early transition metals), orange (intermediate transition metals), gray (post-transition elements). (c) Predicted phase stability for 70 possible quaternary alloys from Group 1 elements (Hf, Zr, Nb, Ta, Mo, V, Ti, W). Bars show number of alloys predicted as single-phase obtained from computational datasets ($\mathcal{D}_{\text{AFLOW}}^{15}$, $\mathcal{D}_{\text{LTVC}}^{19}$, and $\mathcal{D}_{1350\text{K}}^{45}$) and experimentally verified single-phase HEAs^{45,49,50}.

1 based on similar substitutability patterns. The substi-
2 tutability analysis reveals three distinct element groups
3 with strong intra-group substitutability. Group 1 com-
4 prises eight early transition metals from periodic groups
5 4–6: Ti, Zr, Hf (group 4); V, Nb, Ta (group 5); and Mo,
6 W (group 6). Cr, while belonging to group 6, exhibits
7 unique behavior, showing moderate substitutability with
8 Group 1 elements but high substitutability with Fe, Co,
9 Mn, and Al, which together form Group 2. Group 3
10 contains primarily late transition metals from periodic
11 groups 9–11, including Rh, Ir, Pd, Pt, Ni, Cu, Au, Ag.
12 Notably, Groups 1 and 3 show weak inter-group substi-
13 tutability but moderate substitutability with the bridg-
14 ing Group 2.

15 The exceptional intra-group substitutability of Group
16 1 elements (Ti, Zr, Hf, V, Nb, Ta, Mo, W), exhibiting
17 notably higher scores than Groups 2 and 3, suggests a
18 design principle: quaternary combinations should read-
19 ily form stable single-phase HEAs. Critically, this sub-
20 stitutability matrix (Figure 9b) is derived by fusing evi-
21 dence from multiple independent sources—experimental
22 HEA dataset ($\mathcal{D}_{\text{HEA}}^{\text{exp}}$), computational databases ($\mathcal{D}_{1350\text{K}}$,
23 $\mathcal{D}_{\text{AFLOW}}$, $\mathcal{D}_{\text{LTVC}}$), and 20 LLM-domain sources—through
24 Dempster–Shafer integration; such high mutual sub-
25 stitutability indicates unanimous agreement across all
26 sources regarding these patterns. Figure 9c validates
27 this prediction: all three computational datasets unani-
28 mously predict single-phase formation for all 70 possible

29 Group 1 quaternaries, and all 15 experimentally synthe-
30 sized compositions form single-phase HEAs (100% suc-
31 cess rate). This agreement is consistent with established
32 principles for refractory high-entropy alloys^{41,64}: early
33 transition metals (groups 4–6) preferentially form stable
34 BCC solid solutions due to similar atomic sizes and com-
35 patible electronic structures, with single-phase stability
36 thermodynamically reinforced by configurational entropy
37 that lowers Gibbs free energy at elevated temperatures⁶⁵.

38 E. Effectiveness Assessment on Experimental High-Entropy 39 Boride Data

40 We extend our analysis to high-entropy borides
41 (HEBs), where boron's restrictive bonding requirements
42 create similarly high elemental selectivity as observed in
43 HEAs⁶⁶. Despite different underlying mechanisms, both
44 systems share the key challenge of identifying rare viable
45 combinations within vast compositional spaces, making
46 HEBs suitable for demonstrating our framework's appli-
47 cability to diverse multi-component materials with strin-
48 gent compatibility constraints.

49 In this experiment, we applied our framework to a
50 dataset of 19 experimentally confirmed quinary borides
51 collected from previous studies. Using these validated
52 compositions as training data, our framework was then
53 employed to rank 314 potential quinary boride candi-



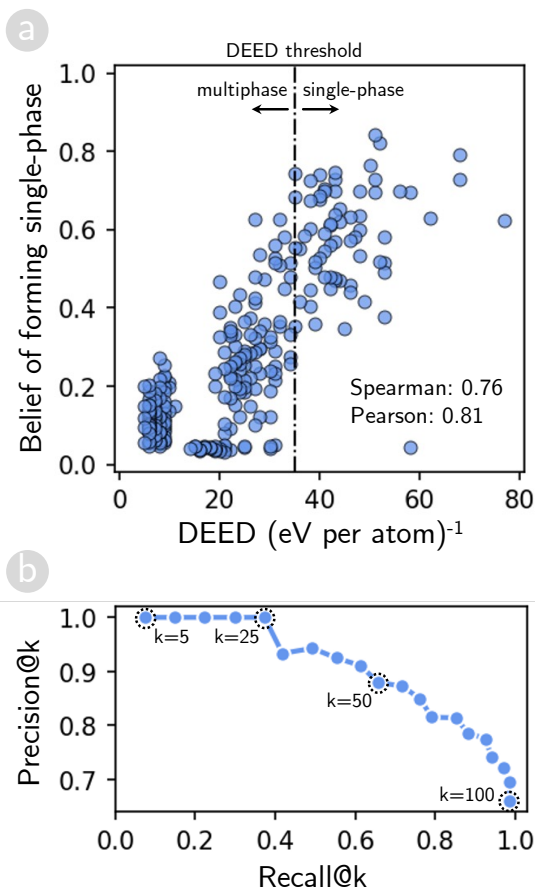


FIG. 10. **Effectiveness Assessment of Multi-Source Knowledge Integration for High-Entropy Borides Formation.** (a) Correlation analysis between our framework's single-phase formation belief and the disordered enthalpy-entropy descriptors (DEED) for 275 quinary boride candidates. The dashed line indicates the DEED threshold of 35 (eV per atom)⁻¹ for single-phase prediction. (b) Precision@k and Recall@k performance metrics evaluated at k values from 5 to 100 with increments of 5

1 dates formed by boron as the anion and the following
2 metals: Cr, Hf, Ir, Mn, Mo, Nb, Ta, Ti, V, W, Y, Zr. To
3 benchmark our framework, we compared the rankings ob-
4 tained by our framework with those derived using the dis-
5 ordered enthalpy-entropy descriptors (DEED)⁴⁴, which
6 represents the state-of-the-art descriptor based on ab-
7 initio calculations for guiding experimental discovery of
8 new single-phase high-entropy carbonitrides and borides.
9 Figure 10a illustrates the correlation between DEED
10 values and the belief of forming single-phase structures
11 for 275 of the 314 quinary boride candidates. For the
12 remaining 39 candidates, our framework could not pro-
13 vide reliable predictions due to insufficient training data
14 coverage, resulting in maximum uncertainty values that
15 rendered these predictions uninformative for comparison

16 purposes. The results demonstrate a strong positive lin-
17 ear correlation between the single-phase formation belief
18 derived from our framework and the DEED values, with
19 Pearson and Spearman correlation coefficients of 0.81 and
20 0.76, respectively. The previous DEED study established
21 a threshold of 35 (eV per atom)⁻¹ to distinguish be-
22 tween single-phase and multiphase candidates, where val-
23 ues above this threshold indicate predicted single-phase
24 formation.
25 The strong correlation for the 275 confident predic-
26 tions, combined with explicit uncertainty flagging for 39
27 candidates, demonstrates effective uncertainty quantifi-
28 cation. To further validate this mechanism, we analyzed
29 prediction accuracy at varying uncertainty thresholds, as
30 shown in Supplementary Figure 8. The results reveal
31 a systematic trade-off: as the uncertainty threshold de-
32 creases (accepting more uncertain predictions as confi-
33 dent), prediction accuracy degrades accordingly. This
34 behavior confirms that high uncertainty values success-
35 fully flag regions where evidence is insufficient, prevent-
36 ing overconfident extrapolation beyond the training data.
37 The explicit uncertainty quantification thus serves as a
38 critical safeguard against overfitting in data-sparse sce-
39 narios, distinguishing our approach from conventional
40 machine learning methods that would force predictions
41 regardless of data sufficiency.
42 To evaluate our framework's practical utility as a ma-
43 terials discovery tool, we analyzed how well it ranks
44 promising candidates compared to the established DEED
45 method. We measured this using standard ranking met-
46 rics: Precision@k (what percentage of our top k recom-
47 mendations are actually good) and Recall@k (what per-
48 centage of all good candidates we capture in our top k
49 recommendations). The results show impressive perfor-
50 mance: when we look at our top 25 recommendations
51 (k=25), all of them were also predicted to form single-
52 phase structures by the DEED method, giving us perfect
53 precision, as shown in Figure 10b. More broadly, to cap-
54 ture 50% of all the promising candidates identified by
55 DEED, our method requires selecting approximately the
56 top 35-40 candidates and maintains over 90% precision,
57 meaning that more than 90% of these top-ranked can-
58 didates are correctly identified as single-phase according
59 to DEED. Even when capturing 75% of the promising
60 candidates, our precision remains above 85%. These re-
61 sults demonstrate that our framework effectively priori-
62 tizes the most promising compositions for experimental
63 synthesis.
64 The strong performance on high-entropy borides, com-
65 bined with the previous results on high-entropy alloys,
66 establishes the framework's capability to handle uncer-
67 tainty in compositionally selective multi-component ma-
68 terial systems. Notably, while computational databases
69 such as AFLOW and CALPHAD carry inherent uncer-
70 tainties from DFT approximations and thermodynamic
71 extrapolations¹⁸, the Dempster-Shafer theory explicitly
72 models these through mass assignments to ignorance,
73 enabling robust integration with experimental data and



mitigating risks of systematic errors in guiding alloy synthesis. The discount factor mechanism (Equations 5–7) automatically downweights unreliable sources based on cross-validation performance, preventing error propagation by allowing high-quality evidence to dominate when computational predictions conflict with experimental observations.

F. Limitations and Future Extensions

Previous sections have demonstrated the framework's effectiveness across computational and experimental datasets. We now examine its current limitations and corresponding opportunities for future development.

Context-Independent Evidence Weighting: The current implementation employs fixed weighting parameters for each source without considering the specific context of elemental substitution. For instance, metallurgical knowledge may be more reliable for refractory elements, while solid-state physics insights may better inform noble metal substitutability. Future extensions could implement context-dependent weighting, wherein discount factors vary based on the element pair under consideration. This could be achieved by conditioning discount factors on elemental properties such as atomic radius, electronegativity, or periodic group membership, enabling the framework to recognize element-specific reliability patterns across different knowledge sources.

From Uncertainty Quantification to Discovery Navigation: This study proposes a framework to integrate multi-source knowledge and quantify uncertainty for candidate materials. However, a subsequent challenge remains: how to effectively utilize these uncertainty measures to select candidates for experimental validation under limited resources. This candidate selection problem inherently involves balancing exploration (investigating compositions with high uncertainty that may reveal novel alloys) and exploitation (refining predictions in promising regions with moderate uncertainty). Active learning provides a principled approach to this challenge by identifying experiments that maximally reduce epistemic uncertainty, prioritizing candidates where additional data would most improve model reliability. Reinforcement learning complements this by learning optimal selection policies through iterative experimental feedback, dynamically adjusting the exploration–exploitation balance as the discovery campaign progresses. Together, these techniques could transform the current prediction framework into a comprehensive decision-support system for accelerated materials discovery.

Symmetric Substitutability Assumption: The symmetric substitutability assumption ($A \rightarrow B$ and $B \leftarrow A$ are equivalent) represents a context-averaged approximation

that may limit accuracy for systems with strong directional substitution preferences. This symmetric treatment is justified in this study by two factors: first, the limited training data in our data-sparse scenarios makes learning separate directional patterns statistically infeasible; second, for near-equiatomic multi-principal element HEAs characterized by disordered random solid solutions, elements occupy statistically similar local environments, rendering symmetric substitution a physically reasonable first-order approximation. Future extensions could incorporate asymmetric substitutability by maintaining separate $A \rightarrow B$ and $B \leftarrow A$ matrices and collecting directional evidence from LLMs through modified prompts.

Broaden Scope Beyond Phase Stability: To serve the purpose of screening the element combinations forming HEA phases, the proposed framework focuses on the fundamental question of whether the HEA phase exists. We design a frame of discernment $\Omega_{HEA} = \{HEA, \overline{HEA}\}$ to model the existence of HEA phases with mass functions. Consequently, our framework has not answered essential questions regarding the structure and other properties of the HEAs. However, by redesigning the frame of discernment to reflect the additional properties of interest, we can also construct a model that can recommend potential alloys forming HEA phases with desirable properties. Extending to mechanical, electronic, or catalytic properties represents another promising direction as sufficient property-specific data becomes available⁶⁷.

Scalability to Higher-Order Systems: The current validation focuses primarily on quaternary alloy systems, with limited exploration of higher-order compositions. Extension to quinary and higher-order alloys could be achieved through hierarchical decomposition, wherein quaternary systems serve as baseline evidence augmented by pairwise substitutability relationships. However, more complex systems may require sparse approximation techniques and substantially larger materials databases to maintain predictive reliability.

V. CONCLUSIONS

The central contribution of this work lies in demonstrating that the interpolation–extrapolation dichotomy inherent to conventional data-driven materials discovery can be systematically addressed through principled integration of multi-source knowledge. Crucially, the framework does not indiscriminately combine all available evidence; rather, it evaluates the reliability of each source based on its alignment with the target property, ensuring that only relevant domain knowledge contributes meaningfully to predictions. By employing elemental substitutability as a unifying concept and leveraging Dempster–Shafer theory to combine empirical observations with insights extracted from scientific literature via LLMs, the framework effectively bridges data-rich and data-sparse regions in materials exploration. Our frame-



work demonstrates superior performance compared to traditional data-driven approaches and empirical phase selection rules, while achieving accuracy comparable to computationally expensive methods, particularly when predicting phase stability for compositions containing previously unseen elements. These results highlight that the significance of the framework does not reside in superseding established methods, but rather in effectively synthesizing their complementary strengths while representing epistemic limitations transparently.

Beyond HEAs, this framework could accelerate discovery in several materials classes facing similar challenges of vast compositional spaces and sparse data, including functional ceramics⁴⁴, and catalytic materials³⁴. Through successful validation on diverse alloy systems, this study demonstrates that uncertainty-aware AI integration provides a viable path forward for accelerated materials discovery. The element substitutability patterns extracted using this framework may also inform synthetic strategies for targeted property optimization across diverse material applications.

AUTHOR CONTRIBUTIONS

M.-Q. H.: Conceptualization, Methodology, Software, Formal analysis, Validation, Investigation, Writing - Original Draft, Writing - Review & Editing, Visualization. D.-K. L.: Software, Investigation, Data Curation. V.-C. N.: Software, Formal analysis, Data Curation. H. K.: Investigation, Validation, Writing - Review & Editing. S. C.: Investigation, Validation, Writing - Review & Editing. H.-C. D.: Conceptualization, Methodology, Validation, Investigation, Writing - Original Draft, Writing - Review & Editing, Visualization, Supervision, Project administration, Funding acquisition.

CONFLICT OF INTERESTS

The authors report there are no competing interests to declare.

DATA AVAILABILITY

Publicly Available Datasets: Data for this article, including experimental and computational datasets supporting high-entropy alloy phase prediction, are available at Zenodo at <https://doi.org/10.5281/zenodo.17074832>.

Code Availability: Code for the uncertainty-aware AI integration framework is available at GitHub at <https://github.com/minhquyet2308/Uncertainty-Aware-AI-Intergration>, with an archived version available at Zenodo at <https://doi.org/10.5281/zenodo.17744151>.

ACKNOWLEDGEMENTS

This work is supported by the JST-CREST Program (Innovative Measurement and Analysis), under Grant number JPMJCR2235; and the JSPS KAKENHI Grant Numbers 20K05301, JP19H05815, 20K05068, 23KJ1035, 23K03950, and JP23H05403. S.C. acknowledges support by US-DoD (ONR MURI program N00014-21-1-251). H. K. acknowledges support by the Japan Science and Technology Agency (JST) ASPIRE Program under the project "International Collaborative Research Network for Advanced Atomic Layer Processes". The authors thank Dr. Huan Tran and Dr. Xiomara Campilongo for fruitful discussions.

REFERENCES

1 J.-W. Yeh, S.-K. Chen, S.-J. Lin, J.-Y. Gan, T.-S. Chin, T.-T. Shun, C.-H. Tsau, and S.-Y. Chang, *Advanced Engineering Materials* **6**, 299 (2004).
2 B. Cantor, I. Chang, P. Knight, and A. Vincent, *Materials Science and Engineering: A* **375-377**, 213 (2004).
3 O. Senkov and D. Miracle, *Journal of Alloys and Compounds* **658** (2015), 10.1016/j.jallcom.2015.10.279.
4 J. M. Rickman, H. M. Chan, M. P. Harmer, J. A. Smeltzer, C. J. Marvel, A. Roy, and G. Balasubramanian, *Nature Communications* **10**, 2618 (2019).
5 M.-H. Tsai and J.-W. Yeh, *Materials Research Letters* **2**, 107 (2014).
6 C. Toher, C. Oses, D. Hicks, and S. Curtarolo, *npj Comput. Mater.* **5**, 69 (2019).
7 G. Deshmukh, N. J. Wichrowski, N. Evangelou, P. G. Ghanekar, S. Deshpande, I. G. Kevrekidis, and J. Greeley, *npj Computational Materials* **10**, 116 (2024).
8 M. Ghorbani, M. Boley, P. N. H. Nakashima, and N. Birbilis, *Scientific Reports* **14**, 8299 (2024).
9 M.-H. Tsai, *Entropy* **18**, 252 (2016).
10 M.-H. Tsai, R.-C. Tsai, T. Chang, and W.-F. Huang, *Metals* **9**, 247 (2019).
11 W. Huang, P. Martin, and H. L. Zhuang, *Acta Mater.* **169**, 225 (2019).
12 Z. Rao, P.-Y. Tung, R. Xie, Y. Wei, H. Zhang, A. Ferrari, T. Klaver, F. Körmann, P. T. Sukumar, A. K. da Silva, Y. Chen, Z. Li, D. Ponge, J. Neugebauer, O. Gutfleisch, S. Bauer, and D. Raabe, *Science* **378**, 78 (2022).
13 J. Roberts, B. Rijal, S. Divilov, J.-P. Maria, W. G. Fahrenholtz, D. E. Wolfe, D. W. Brenner, S. Curtarolo, and E. Zurek, *npj Computational Materials* **10**, 142 (2024).
14 D. Alman, *Entropy* **15**, 4504 (2013).
15 F. Zhang, C. Zhang, S. Chen, J. Zhu, W. Cao, and U. Kattner, *Calphad* **45**, 1 (2014).
16 M. Esters, C. Oses, S. Divilov, H. Eckert, R. Friedrich, D. Hicks, M. J. Mehl, F. Rose, A. Smolyanyuk, A. Calzolari, X. Campilongo, C. Toher, and S. Curtarolo, *Comp. Mat. Sci.* **216**, 111808 (2023).
17 C. Oses, M. Esters, D. Hicks, S. Divilov, H. Eckert, R. Friedrich, M. J. Mehl, A. Smolyanyuk, X. Campilongo, A. van de Walle, J. Schroers, A. G. Kusne, I. Takeuchi, E. Zurek, M. Buongiorno Nardelli, M. Fornari, Y. Lederer, O. Levy, C. Toher, and S. Curtarolo, *Comp. Mat. Sci.* **217**, 111889 (2023).
18 C. Toher and S. Curtarolo, *Journal of Phase Equilibria and Diffusion* **45**, 219 (2024).
19 Y. Lederer, C. Toher, K. S. Vecchio, and S. Curtarolo, *Acta Materialia* **159**, 364 (2018).



- 1 ²⁰V. Stanev, C. Oses, A. G. Kusne, E. Rodriguez, J. Paglione,
2 S. Curtarolo, and I. Takeuchi, *npj Comput. Mater.* **4**, 29 (2018).
- 3 ²¹G. L. W. Hart, T. Mueller, C. Toher, and S. Curtarolo, *Nature*
4 *Reviews Materials* **6**, 730 (2021).
- 5 ²²M.-Q. Ha, D.-N. Nguyen, V.-C. Nguyen, T. Nagata, T. Chikyow,
6 H. Kino, T. Miyake, T. Denœux, V.-N. Huynh, and H.-C. Dam,
7 *Nature Computational Science* **1**, 470 (2021).
- 8 ²³J. He, R. Yin, C. Wang, C. Liu, D. Xue, Y. Su, L. Qiao, T. Look-
9 man, and Y. Bai, *Journal of Materiomics* **11**, 100913 (2025).
- 10 ²⁴T. L. Pham, H. Kino, K. Terakura, T. Miyake, K. Tsuda, I. Taki-
11 gawa, and H. C. Dam, *Science and Technology of Advanced*
12 *Materials* **18**, 756 (2017), pMID: 29152012.
- 13 ²⁵E. Hüllermeier and W. Waegeman, *Machine Learning* **110**, 457
14 (2021).
- 15 ²⁶E. Brochu, V. M. Cora, and N. de Freitas, “A tutorial on bayesian
16 optimization of expensive cost functions, with application to
17 active user modeling and hierarchical reinforcement learning,”
18 (2010), arXiv:1012.2599 [cs.LG].
- 19 ²⁷J. Snoek, H. Larochelle, and R. P. Adams, in *Advances in Neural*
20 *Information Processing Systems*, Vol. 25, edited by F. Pereira,
21 C. Burges, L. Bottou, and K. Weinberger (Curran Associates,
22 Inc., 2012).
- 23 ²⁸E. Hüllermeier and K. Brinker, *Fuzzy Sets and Systems* **159**,
24 2337 (2008), theme: Information Processing.
- 25 ²⁹E. P. George, D. Raabe, and R. O. Ritchie, *Nat. Rev. Mater.* **4**,
26 515 (2019).
- 27 ³⁰T. Konno, H. Kurokawa, F. Nabeshima, Y. Sakishita, R. Ogawa,
28 I. Hosako, and A. Maeda, *Phys. Rev. B* **103**, 014509 (2021).
- 29 ³¹A. P. Dempster, *Journal of the Royal Statistical Society: Series*
30 *B (Methodological)* **30**, 205 (1968).
- 31 ³²G. Shafer, *A Mathematical Theory of Evidence* (Princeton Uni-
32 versity Press, 1976).
- 33 ³³T. Denœux, D. Dubois, and H. Prade, in *A Guided Tour of Arti-*
34 *ficial Intelligence Research*, Vol. 1, edited by P. Marquis, O. Pap-
35 ini, and H. Prade (Springer Verlag, 2020) Chap. 4, pp. 119–150.
- 36 ³⁴N. Nu Thanh Ton, M.-Q. Ha, T. Ikenaga, A. Thakur, H.-C. Dam,
37 and T. Taniike, *2D Materials* **8**, 015019 (2020).
- 38 ³⁵M.-Q. Ha, D.-N. Nguyen, V.-C. Nguyen, H. Kino, Y. Ando,
39 T. Miyake, T. Denœux, V.-N. Huynh, and H.-C. Dam, *Journal*
40 *of Applied Physics* **133**, 053904 (2023).
- 41 ³⁶E. O. Pyzer-Knapp, J. W. Pitera, P. W. J. Staar, S. Takeda,
42 T. Laino, D. P. Sanders, J. Sexton, J. R. Smith, and A. Curioni,
43 *npj Computational Materials* **8**, 84 (2022).
- 44 ³⁷D. H. Cook, P. Kumar, M. I. Payne, C. H. Belcher, P. Borges,
45 W. Wang, F. Walsh, Z. Li, A. Devaraj, M. Zhang, M. Asta, A. M.
46 Minor, E. J. Lavernia, D. Apelian, and R. O. Ritchie, *Science*
47 **384**, 178 (2024).
- 48 ³⁸S. Liu, T. Wen, A. S. Pattamatta, and D. J. Srolovitz, *Materials*
49 *Today* **80**, 240 (2024).
- 50 ³⁹Z. Chen, Y. Liu, and H. Sun, *Nature Communications* **12**, 6136
51 (2021).
- 52 ⁴⁰B. Cantor, K. Kim, and P. J. Warren, in *Metastable, Mechan-*
53 *ically Alloyed and Nanocrystalline Materials (2001)*, *Journal of*
54 *Metastable and Nanocrystalline Materials*, Vol. 13 (Trans Tech
55 Publications Ltd, 2002) pp. 27–32.
- 56 ⁴¹D. Miracle and O. Senkov, *Acta Materialia* **122**, 448 (2017).
- 57 ⁴²A. Tversky, *Psychological Review* **84**, 327 (1977).
- 58 ⁴³P. Smets, *International Journal of Approximate Reasoning* **9**, 1
59 (1993).
- 60 ⁴⁴S. Divilov, H. Eckert, D. Hicks, C. Oses, C. Toher, R. Friedrich,
61 M. Esters, M. J. Mehl, A. C. Zettl, Y. Lederer, E. Zurek, J.-
62 P. Maria, D. W. Brenner, X. Campilongo, S. Filipović, W. G.
63 Fahrenholtz, C. J. Ryan, C. M. DeSalle, R. J. Creales, D. E.
64 Wolfe, A. Calzolari, and S. Curtarolo, *Nature* **625**, 66 (2024).
- 65 ⁴⁵W. Chen, A. Hilhorst, G. Bokas, S. Gorsse, P. J. Jacques, and
66 G. Hautier, *Nature Communications* **14**, 2856 (2023).
- 67 ⁴⁶A. Takeuchi and A. Inoue, *MATERIALS TRANSACTIONS* **46**,
68 2817 (2005).
- 69 ⁴⁷A. Takeuchi and A. Inoue, *Intermetallics* **18**, 1779 (2010).
- 70 ⁴⁸T. Fukushima, H. Akai, T. Chikyow, and H. Kino, *Phys. Rev.*
71 *Materials* **6**, 023802 (2022).
- 72 ⁴⁹C. K. H. Borg, C. Frey, J. Moh, T. M. Pollock, S. Gorsse, D. B.
73 Miracle, O. N. Senkov, B. Meredig, and J. E. Saal, *Scientific*
74 *Data* **7**, 430 (2020).
- 75 ⁵⁰G. Khanna R, M. K. Singh, D. K. Rai, and S. Samal, *Materials*
76 *Letters* **365**, 136404 (2024).
- 77 ⁵¹M. P. LaValley, *Circulation* **117**, 2395 (2008).
- 78 ⁵²A. Seko, A. Togo, and I. Tanaka, “Descriptors for machine learn-
79 ing of materials data,” in *Nanoinformatics* (Springer Singapore,
80 Singapore, 2018) pp. 3–23.
- 81 ⁵³A. Seko, H. Hayashi, K. Nakayama, A. Takahashi, and I. Tanaka,
82 *Phys. Rev. B* **95**, 144110 (2017).
- 83 ⁵⁴F. C. T., *Nature* **138**, 7 (1936).
- 84 ⁵⁵U. Mizutani, *MRS Bulletin* **37**, 169 (2012).
- 85 ⁵⁶H. O. M. S. E. Mueller, in *Alloy Phase Diagrams* (ASM Inter-
86 national, 2016).
- 87 ⁵⁷Z. Pei, J. Yin, P. K. Liaw, and D. Raabe, *Nature Communica-*
88 *tions* **14**, 54 (2023).
- 89 ⁵⁸X. Yang and Y. Zhang, *Materials Chemistry and Physics* **132**,
90 233 (2012).
- 91 ⁵⁹S. Guo, Q. Hu, C. Ng, and C. Liu, *Intermetallics* **41**, 96 (2013).
- 92 ⁶⁰W. Zhijun, Y. Huang, Y. Yang, J. Wang, and C. Liu, *Scripta*
93 *Materialia* **94** (2015), 10.1016/j.scriptamat.2014.09.010.
- 94 ⁶¹A. K. Singh, N. Kumar, A. Dwivedi, and A. Subramaniam,
95 *Intermetallics* **53**, 112 (2014).
- 96 ⁶²M. C. Tropicovsky, J. R. Morris, P. R. C. Kent, A. R. Lupini,
97 and G. M. Stocks, *Phys. Rev. X* **5**, 011041 (2015).
- 98 ⁶³S. Guo, C. Ng, J. Lu, and C. T. Liu, *Journal of Applied Physics*
99 **109**, 103505 (2011).
- 100 ⁶⁴O. Senkov, G. Wilks, D. Miracle, C. Chuang, and P. Liaw, *Inter-*
101 *metallics* **18**, 1758 (2010).
- 102 ⁶⁵B. S. Murty, J.-W. Yeh, S. Ranganathan, and P. P. Bhattachar-
103 jee, *High-Entropy Alloys*, 2nd ed. (Elsevier, Amsterdam, 2019).
- 104 ⁶⁶J. Gild, Y. Zhang, T. Harrington, S. Jiang, T. Hu, M. C. Quinn,
105 W. M. Mellor, N. Zhou, K. Vecchio, and J. Luo, *Scientific Re-*
106 *ports* **6**, 37946 (2016).
- 107 ⁶⁷S. Nakanowatari, K. Takahashi, H. C. Dam, and T. Taniike,
108 *ACS Catalysis* **15**, 8691 (2025).



Data Availability Statement (DAS)

Publicly Available Datasets: Data for this article, including experimental and computational datasets supporting high-entropy alloy phase prediction, are available at Zenodo at <https://doi.org/10.5281/zenodo.17074832>.

Code Availability: Code for the uncertainty-aware AI integration framework is available at GitHub at <https://github.com/minhquyet2308/Uncertainty-Aware-AI-Intergration>, with an archived version available at Zenodo at <https://doi.org/10.5281/zenodo.17744151>.

