LLaMP: Large Language Model Made Powerful for High-fidelity Materials Knowledge Retrieval

Yuan Chiang^{1,2*}, Elvis Hsieh^{1*}, Chia-Hong Chou³, Janosh Riebesell^{2,4}
¹University of California, Berkeley, ²Lawrence Berkeley National Laboratory,
³Foothill College, ⁴Cavendish Laboratory, University of Cambridge, UK
{cyrusyc, htelvis92}@berkeley.edu

Abstract

Materials science research requires multi-step reasoning and precise material informatics retrieval, where minor errors can propagate into significant failures in downstream experiments. Despite their general success, Large Language Models (LLMs) often struggle with hallucinations, handling domain-specific data effectively (e.g., crystal structures), and integrating experimental workflows. To address these challenges, we introduce LLaMP, a hierarchical multi-agent framework designed to emulate the materials science research workflow. The high-level supervisor agent decomposes user requests into sub-tasks and coordinates with specialized assistant agents. These assistant agents handle domain-specific tasks, such as retrieving and processing data from the Materials Project (MP) or conducting simulations as needed. This pipeline facilitates iterative refinement of material property retrieval and enables the simulation of real-world research workflows. To ensure reliability, we propose a novel metric combining uncertainty and confidence estimate to evaluate the selfconsistency of responses from LLaMP and baseline methods. Our experiments demonstrate LLaMP's superior performance in material property retrieval, crystal structure editing, and annealing molecular dynamics simulations using pre-trained interatomic potentials. Unlike prior work focused solely on material property prediction or discovery, LLaMP serves as a foundation for autonomous materials research by combining grounded informatics and enabling iterative experimental processes. Code and live demo are available at https://github.com/chiang-yuan/llamp.

1 Introduction

The generation of convincing yet unreliable information poses a pressing challenge to large

language models (LLMs), particularly for scientific applications. LLMs are prone to hallucination—providing false information with high confidence (Xu et al., 2024; Bang et al., 2023). This issue is especially concerning for knowledge-intensive tasks, where users rely on AI systems for accurate guidance (Lewis et al., 2020). LLMs often lack up-to-date factual knowledge beyond their training data, requiring verification against trusted sources (Mallen et al., 2023). In science, the proliferation of generative models may exacerbate misinformation risks, accentuating the importance of ensuring reliable information sources.

Current approaches to enhance LLM accuracy in domain-specific knowledge often involve finetuning pre-trained models (Gupta et al., 2022; Dagdelen et al., 2024) or prompt engineering (Yang et al., 2023; Zheng et al., 2023). While easy to deploy, these models suffer from diminished reproducibility, lack of a memory base, and untraceable fine-tuning history. Even though fine-tuning can encode a certain amount of domain-specific knowledge into LLMs, it is constrained by scalability and intrinsic memory capacity (Morris et al., 2025; Carlini et al., 2018; Schwarzschild et al., 2024; Kandpal et al., 2022). Prompt engineering, while effective, also compromises the generalizability, limiting the overall flexibility. Therefore, a more sensible approach involves equipping LLMs with external data sources, allowing them to generate holistic responses via few-shot adaptation to factual information (Lewis et al., 2021) that can reliably support real-world scientific research and decisionmaking.

In this work, we propose LLaMP, a hierarchical multi-agent framework that leverages Materials Project (MP), arXiv, Wikipedia, and atomistic simulation tools. The framework serves as a safeguard against LLM hallucination by grounding them in high-fidelity material informatics from large-scale material database, including computational data

^{*}Equal Contribution.

from quantum-mechanical first-principles calculations and expert-curated material synthesis recipes, and further enables the capabilities of complex downstream tasks. The hierarchical planning of supervisor and assistant agents improves tool-usage performance and enhances consistency in final responses with self-correcting mechanism. LLaMP supports a wide range of advanced capabilities, including multi-modal searching, tensor and 3D crystal structure retrieval and operation, and languagedriven simulation. The framework not only can accurately retrieve high-fidelity, higher-order materials data but also can combine different modalities to perform complex, knowledge-intensive inferences and operations essential for real-world materials science applications.

Our contributions are as follows: (1) we introduce a hierarchical agentic framework that enables LLMs to access high-fidelity materials informatics and perform complex simulation; (2) we propose a statistical metric to assess the self-consistency of LLM responses in high-precision, reproducibilitycritical settings; (3) we evaluate the performance of LLaMP, standard and specialized LLMs in predicting material properties, including bulk moduli, electronic bandgaps, formation energies, and magnetic orderings; (4) we demonstrate real-world applications in materials science, including inorganic synthesis, crystal structure generation and editing, and high-throughput atomistic simulations using pre-trained force fields, highlighting the potential of language-driven experimentation for automating and accelerating scientific discovery.

2 Related Work

Materials science databases and benchmarks.

The Materials Project is a multi-institution effort to explore and compute the properties of all known inorganic materials (Jain et al., 2013) and molecules (Spotte-Smith et al., 2023). The initiative leverages high-throughput electronic structure calculations (Kresse and Furthmüller, 1996; Shao et al., 2015) based on density functional theory (DFT), providing large-scale open-source database and analysis algorithms, with the goal to drastically reduce the time and cost required for materials discovery by focusing experiments on the promising candidates from computational screening. Most of the atomic structures are selected from the Inorganic Crystal Structure Database (ICSD) (Zagorac et al., 2019) and undergo standardized relaxation

procedures, followed by post-processing or additional calculations for higher-order material properties such as electron and phonon bandgaps, elastic tensors, and dielectric tensors. Complementing MP, the Automatic FLOW for Materials Discovery (AFLOW) provides a database of over 3.5 million compounds and extensive computational tools (Curtarolo et al., 2012b), while the Open Quantum Materials Database (OQMD) focuses on thermodynamic and structural properties (Kirklin et al., 2015a), and the Novel Materials Discovery (NO-MAD) Laboratory offers a FAIR-compliant platform for managing and sharing materials science data (Scheidgen et al., 2023b). Despite the value of these resources, MP distinguishes itself with its curated, high-fidelity dataset, which is continuously updated and vetted by the community, ensuring reliability and scientific validity for downstream applications.

Prompting and fine-tuning LLM in science do-

main. Prompt-based methods have been used as effective tools for automating the data extraction process from the literature. Polak and Morgan (2023) employ a prompt workflow to extract the cooling rates of metallic glasses and yield strengths of high entropy alloys. Zheng et al. (2023) implement a ChatGPT metal-organic framework (MOF) synthesis assistant through embedding and searching on preselected papers. StructChem (Ouyang et al., 2024) leverages step-by-step reasoning and iteratively refines results to solve college-level chemistry questions. Yang et al. (2023) use GPT-4 to extract experimentally measured bandgaps to train a graph neural network for accurate bandgap prediction from crystal structures. Other works address the challenges of extracting complex materials informatics from diverse formats such as tables and unstructured texts (Hira et al., 2024; Schilling-Wilhelmi et al., 2024). Despite the success in the specific data extraction tasks, prompt-based methods face challenges in reproducibility when the used prompts are fine-grained to work for specific edge cases. They are also still prone to hallucination and less generalizable to combine different data sources due to the deliberately designed prompt.

Several other knowledge-grounded, domainspecific language models lean on the fine-tuning approach against pre-selected data and literature. For instance, ChemGPT (Frey et al., 2022) involves fine-tuning GPT-neo on self-referencing embedded strings (SELFIES) representations of small molecules. Jablonka et al. (2024) demonstrated GPT-3 fine-tuned against online corpora could outperform purpose-trained models on classification, regression, and inverse design of high-entropy alloys and molecules. Dagdelen et al. (2024) finetuned GPT-3 on ~500 prompt-completion pairs to enhance LLM's capability to extract useful information on materials chemistry from text paragraphs. Cha et al. (2024) further curated instruction data to fine-tune Llama for material sciencespecific tasks. Xie et al. (2023) combines questionanswering fine-tuning with multi-task learning. However, the fine-tuned models without augmentation inherently lack awareness of the up-to-date results, and any data is only available after their training. Furthermore, these works focus on general (undergraduate-level) question answering (Zaki et al., 2023; Song et al., 2023; Wang et al., 2023) instead of factual grounding on expert-curated database and downstream experimental workflow.

LLM agent and tool usage. An emerging class of LLM applications take advantage of LLM text completion and instruction following capability for function calling. This approach extends LLMs with expert-curated tools to improve the quality of control for downstream applications (Chase, 2022; Qin et al., 2023; Wang et al., 2024; Shinn et al., 2023; Du et al., 2024; Lu et al., 2024). In the science domain, such an approach is widely adopted. Coscientist (Boiko et al., 2023) combines tools such as search engines, Python, and document index for autonomous chemical research. ChemCrow (M. Bran et al., 2024) gathers multiple molecule and safety tools to enhance organic chemistry experiments and molecule design. Ghafarollahi and Buehler (2024) propose AtomAgents for alloy design and analysis. Zhang et al. (2024) develop a retrieval-based agent on their curated dataset. However, most prior works adopt the flat planning strategy, where a single agent accesses all the available tools, resulting in a lack of self-correcting and iterative refinement capabilities. We mitigate this through hierarchical structure of multi-agents (see Section 3).

3 Method

3.1 Hierarchical Multi-Agent Planning

Overview. Flat planning, where a single agent sees all the available tools and related API schemas, quickly exceeds the context window and incurs a huge cost for multimodal data processing in

material science. To manage heterogeneous data sources and diverse queries, we introduce hierarchical multi-agent planning, featuring a supervisor agent overseeing multiple specialized assistant agents (Figure 1). This design offers three major advantages over flat planning commonly implemented in previous works (Boiko et al., 2023; M. Bran et al., 2024): (1) modularity of the system ensures that each assistant agent can focus on domain-specific queries while the supervisor agent handles higher-level reasoning and task allocation; (2) the hierarchical structure improves the overall accuracy and efficiency by minimizing the context window consumption and schema parsing; (3) the multi-agent collaboration allows the iterative refinement on multimodal data (Figure A.1) and long-horizon simulation (Appendix C).

Supervisor agent. The supervisor agent acts as a router and decision-maker, handling abstract logic between user requests and assistant agents. Here, we follow Yao et al. (2023) to augment the agent's action space A with a language space L to create an expanded action space of $\hat{A} = A \cup \mathcal{L}$ on GPT-4. Instead of executing a direct function call, the supervisor dynamically delegates queries to assistant agents via natural language, ensuring a humaninterpretable reasoning trace. This approach enhances contextual grounding, ensuring that factual information and numerically accurate simulations. By modularizing factual grounding and simulation tasks, the supervisor effectively reduces error propagation where reasoning errors in materials properties compound across multiple steps.

Assistant agent. One of the critical failures in the scientific domain (Miret and Krishnan, 2024) is the LLM-generated hallucinations and incorrect outputs that mislead downstream experiments. To mitigate this, each MP assistant agent directly retrieves materials informatics from the Materials Project, ensuring grounded reasoning with expert-verified data.

Under the modular architecture, we assign a specialized agent for each specific tool or API endpoint. It reduces context window consumption, as each agent handles only the relevant schema for its task, avoiding unnecessary schema parsing. Additionally, the assistant agents can refine their API calls based on feedback, significantly improving task completion rates through the iterative self-correcting mechanism.

The full list of agents and tools are defined in A.2.

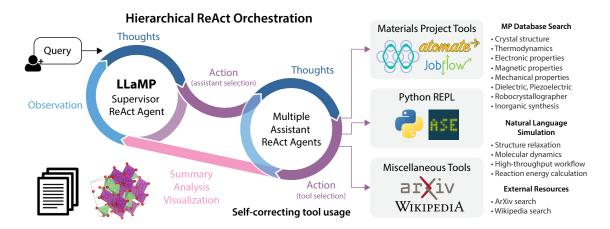


Figure 1: Hierarchical multi-agent planning in LLaMP. Designed to emulate the materials science research workflow, LLaMP employs a two-level agentic structure, Supervisor ReAct agent decomposes user queries into sub-tasks and oversees multiple specialized assistant agents. These assistants retrieve high-fidelity data from the Materials Project (MP) and perform domain-specific tasks like crystal structure manipulation and molecular dynamics simulations. For a detailed example, refer to Figure A.1.

Each MP assistant agent employs a self-correcting mechanism, enabling agents to refine their API calls and improve task completion rates. The framework's modularity enables seamless integration of new assistant agents, allowing for extensibility to various materials discovery methods and experimental techniques (Zeni et al., 2024; Luo et al., 2023; Pilania et al., 2017; Wen et al., 2024, 2023).

3.2 Self-Consistency of Response (SCoR)

When LLMs are integrated into scientific work-flows and deployed in high-stakes settings (*i.e.* self-driving labs), it is important for these models to have consistent and predictable behaviors (Liang et al., 2023). For numeric knowledge retrieval tasks, we define the following metrics:

Precision. (sample standard deviation) measures the uncertainty in the model's responses where n is the number of valid responses from N trials and $\hat{\sigma}$ is the standard deviation of a valid response:

$$\text{Precision} = \frac{\hat{\sigma}}{\sqrt{n}} \ge 0$$

Coefficient of Precision (CoP). maps the precision to (0,1]:

$$\operatorname{CoP} = \exp\left(-\operatorname{Precision}\right) = \exp\left(-\frac{\hat{\sigma}}{\sqrt{n}}\right) \in (0,1] \,.$$

Confidence. measures the ratio of generating n valid responses in N trials:

Confidence
$$=\frac{n}{N}$$
.

Self-consistency of Response (SCoR). is then defined as

$$SCoR = CoP \times Confidence \in [0, 1]$$
.

The limit of SCoR = 1 is reached when the model yields the same response to a given query every time. At the limit of SCoR = 0, the model is either very inconsistent (with large variance across the responses) or very reluctant (with low confidence) to answer the query. Despite the simplicity in definition, SCoR effectively reflects the reproducibility and practical usability of the method, which is important when the method is incorporated into broader systems where the stable and expected behaviors are prioritized. Detailed metric calculation can be found in Appendix A.3.

4 Experiments

4.1 Material Properties Retrieval

Response quality and consistency. The prior benchmarks in materials science often focus on graduate school-level question-answering datasets (Wang et al., 2023; Zaki et al., 2023) or are limited to a single data modality (Rubungo et al., 2024). These approaches overlook the complexity of materials science, which involves reasoning across multiple properties and modalities. To bridge this gap, we benchmark material property retrieval on bulk modulus, formation energy, and bandgap for both common and multi-element materials. We consider four baselines: StructChem with GPT-4 (prompting-based) (Ouyang et al., 2024),

		Bulk	Modulus K (0	GPa)			Formati	on Energy ΔI	H_f (eV)	
	Precision↓	CoP	Confidence	SCoR↑	MAE↓	Precision↓	CoP	Confidence	SCoR↑	MAE↓
LLaMP (GPT-4)	2.698	0.900	1.000	0.900	14.574	0.006	0.994	0.940	0.934	0.007
LLaMP (Sonnet)	1.816	0.562	1.000	0.562	<u>15.104</u>	0.000	1.000	1.000	1.000	0.000
LLaMP (Gemini)	5.178	0.053	1.000	0.053	16.251	0.076	0.932	0.620	0.576	0.166
LLaMP (Llama3)	12.993	0.036	0.800	0.029	50.308	0.000	1.000	0.250	0.250	1.377
HoneyBee	54.105	0.010	0.900	0.010	110.151	0.132	0.881	0.440	0.430	1.518
StructChem	0.000	1.000	0.200	0.200	41.017	0.000	1.000	0.200	0.200	3.146
Darwin	0.001	0.999	0.500	0.499	156.266	0.003	0.997	1.000	0.997	2.245
GPT-4+Serp	2.221	0.833	0.300	0.433	29.937	0.025	0.977	0.560	0.791	11.669
GPT-4	0.186	0.910	1.000	0.910	41.225	0.000	1.000	0.180	0.200	1.680
Sonnet	0.009	0.992	1.000	0.992	41.033	0.022	0.979	1.000	0.979	294.360
Gemini-Pro	6.065	0.169	1.000	0.169	43.429	0.467	0.657	1.000	0.657	1.412
Llama 3	11.222	0.010	1.000	0.010	41.874	2.346	0.139	0.960	0.137	4.657
	Electronic Bandgap E_g - Common (eV)			Electronic Bandgap E_g - Multi-element (eV)						
	Biee			o(c	•)	Liccio	ine Danc	$L_g = L_g = W(u)$	ti-ciciliciit	(ev)
	Precision ↓	СоР	Confidence	SCoR↑	MAE↓	Precision ↓	CoP	Confidence	SCoR†	MAE↓
LLaMP (GPT-4)										
LLaMP (GPT-4) LLaMP (Sonnet)	Precision↓	CoP	Confidence	SCoR↑	MAE↓	Precision↓	CoP	Confidence	SCoR↑	MAE↓
, ,	Precision↓ 0.000	CoP 1.000	Confidence 0.800	SCoR↑ 0.800	MAE↓ 0.000	Precision↓ 0.047	CoP 0.958	Confidence 0.960	SCoR↑ 0.918	MAE↓ 0.167
LLaMP (Sonnet)	Precision↓ 0.000 0.145	CoP 1.000 0.870	O.800 0.600	SCoR↑ 0.800 0.522	MAE↓ 0.000 <u>0.298</u>	Precision↓ 0.047 0.046	CoP 0.958 0.962	Confidence 0.960 1.000	SCoR↑ 0.918 0.962	MAE↓ 0.167 <u>0.304</u>
LLaMP (Sonnet) LLaMP (Gemini)	0.000 0.145 0.627	CoP 1.000 0.870 0.571	0.800 0.600 0.600	SCoR↑ 0.800 0.522 0.343	MAE↓ 0.000 0.298 1.327	Precision↓ 0.047 0.046 0.003	CoP 0.958 0.962 0.997	0.960 1.000 0.500	SCoR↑ 0.918 0.962 0.997	MAE↓ 0.167 0.304 0.637
LLaMP (Sonnet) LLaMP (Gemini) LLaMP (Llama3)	0.000 0.145 0.627 0.051	CoP 1.000 0.870 0.571 0.952	0.800 0.600 0.600 0.800	SCoR↑ 0.800 0.522 0.343 0.761	MAE↓ 0.000 0.298 1.327 1.038	0.047 0.046 0.003 0.169	CoP 0.958 0.962 0.997 0.848	0.960 1.000 0.500 0.800	SCoR↑ 0.918 0.962 0.997 0.678	MAE↓ 0.167 0.304 0.637 1.094
LLaMP (Sonnet) LLaMP (Gemini) LLaMP (Llama3) HoneyBee	Precision↓ 0.000 0.145 0.627 0.051 0.249	CoP 1.000 0.870 0.571 0.952 0.800	0.800 0.600 0.600 0.800 1.000	SCoR↑ 0.800 0.522 0.343 0.761 0.800	MAE↓ 0.000 0.298 1.327 1.038 1.242	Precision↓ 0.047 0.046 0.003 0.169 0.299	CoP 0.958 0.962 0.997 0.848 0.779	Confidence 0.960 1.000 0.500 0.800 0.740	SCoR↑ 0.918 0.962 0.997 0.678 0.601	MAE↓ 0.167 0.304 0.637 1.094 1.640
LLaMP (Sonnet) LLaMP (Gemini) LLaMP (Llama3) HoneyBee StructChem	Precision↓ 0.000 0.145 0.627 0.051 0.249 0.017	CoP 1.000 0.870 0.571 0.952 0.800 0.984	0.800 0.600 0.600 0.800 1.000 1.000	0.800 0.522 0.343 0.761 0.800 0.984	MAE↓ 0.000 0.298 1.327 1.038 1.242 0.986	Precision↓ 0.047 0.046 0.003 0.169 0.299 0.000	CoP 0.958 0.962 0.997 0.848 0.779 1.000	0.960 1.000 0.500 0.800 0.740 0.200	SCoR↑ 0.918 0.962 0.997 0.678 0.601 0.200	MAE↓ 0.167 0.304 0.637 1.094 1.640 0.973
LLaMP (Sonnet) LLaMP (Gemini) LLaMP (Llama3) HoneyBee StructChem Darwin	Precision↓ 0.000 0.145 0.627 0.051 0.249 0.017 0.002	CoP 1.000 0.870 0.571 0.952 0.800 0.984 0.998	0.800 0.600 0.600 0.800 1.000 1.000	SCoR↑ 0.800 0.522 0.343 0.761 0.800 0.984 0.998	MAE↓ 0.000 0.298 1.327 1.038 1.242 0.986 1.224	0.047 0.046 0.003 0.169 0.299 0.000 0.000	CoP 0.958 0.962 0.997 0.848 0.779 1.000 1.000	Confidence 0.960 1.000 0.500 0.800 0.740 0.200 1.000	SCoR↑ 0.918 0.962 0.997 0.678 0.601 0.200 1.000	MAE↓ 0.167 0.304 0.637 1.094 1.640 0.973 1.951
LLaMP (Sonnet) LLaMP (Gemini) LLaMP (Llama3) HoneyBee StructChem Darwin GPT-4+Serp	Precision↓ 0.000 0.145 0.627 0.051 0.249 0.017 0.002 0.040	CoP 1.000 0.870 0.571 0.952 0.800 0.984 0.998 0.963	Confidence 0.800 0.600 0.600 0.800 1.000 1.000 1.000 1.000	SCoR↑ 0.800 0.522 0.343 0.761 0.800 0.984 0.998 0.963	MAE↓ 0.000 0.298 1.327 1.038 1.242 0.986 1.224 1.012	0.047 0.046 0.003 0.169 0.299 0.000 0.000	CoP 0.958 0.962 0.997 0.848 0.779 1.000 1.000	Confidence 0.960 1.000 0.500 0.800 0.740 0.200 1.000 0.660	SCoR↑ 0.918 0.962 0.997 0.678 0.601 0.200 1.000 0.660	MAE↓ 0.167 0.304 0.637 1.094 1.640 0.973 1.951
LLaMP (Sonnet) LLaMP (Gemini) LLaMP (Llama3) HoneyBee StructChem Darwin GPT-4+Serp GPT-4	Precision↓ 0.000 0.145 0.627 0.051 0.249 0.017 0.002 0.040 0.032	CoP 1.000 0.870 0.571 0.952 0.800 0.984 0.998 0.963 0.970	Confidence 0.800 0.600 0.600 0.800 1.000 1.000 1.000 1.000 1.000	SCoR↑ 0.800 0.522 0.343 0.761 0.800 0.984 0.998 0.963 0.970	MAE↓ 0.000 0.298 1.327 1.038 1.242 0.986 1.224 1.012 0.959	Precision↓ 0.047 0.046 0.003 0.169 0.299 0.000 0.000 -	CoP 0.958 0.962 0.997 0.848 0.779 1.000 1.000	Confidence 0.960 1.000 0.500 0.800 0.740 0.200 1.000 0.660 0.000	SCoR↑ 0.918 0.962 0.997 0.678 0.601 0.200 1.000 0.660 0.000	MAE↓ 0.167 0.304 0.637 1.094 1.640 0.973 1.951 0.576

Table 1: Performance metrics of LLaMP and baselines on material properties prediction tasks. The metrics from left to right are precision (sample standard deviation), coefficient of precision (CoP), confidence, self-consistency of response (SCoR), and mean absolute error (MAE), with theoretical values computed in the Materials Project taken as the ground truth. All the values are the average metrics over five runs and the sampled materials. Better method has high SCoR and MAE simultaneously. Full values across five runs are provided in Figure A.2

HoneyBee (fine-tuned) (Cha et al., 2024), Darwin (fine-tuned) (Xie et al., 2023), GPT-4 with search engine tool (RAG), and vanilla LLMs (gpt-4, llama3-8b, gemini-1.0-pro). Performance is assessed through Precision, CoP, SCoR, and MAE metrics, as defined in Section 3.2. We argue that any useful LLM agents to be included in the scientific workflow should have high SCoR and low error on the material's properties. Notably, LLaMP consistently outperforms other models, achieving the highest SCoR and the lowest errors across material properties, making it highly suitable for scientific workflows. Despite employing extensive prompting strategies, StructChem struggles due to its lack of domain-specific knowledge, leading to frequent refusals when it cannot verify outputs. Similarly, GPT-4 with SerpAPI suffers from noisy and inconsistent web sources, which degrade precision.

For bulk modulus prediction, vanilla LLMs, particularly Llama 3-8b, frequently rely on low-

fidelity online data, leading to significant deviations for elements like Cr, Mn, and Fe, compared to MP theoretical values. Interestingly, Llama 3-8b usually cites spurious references in the responses despite the largest response variance but occasionally agrees with MP values. In contrast, LLaMP outperforms vanilla LLMs and reduces the MAE from around 40 to 14.57 GPa.

For formation energy prediction, vanilla LLMs suffer from low SCoR and high MAE ranging from 1.5 to 5.5 eV, which is impractical for material discovery requiring meV-level precision. This is not unexpected, since accurate formation energy prediction requires the computation of multiple energetics (energies of the compound itself and its elemental constituents).

For bandgaps prediction, we query 10 common compounds and 10 multi-element materials that are less commonly encountered in the literature. Vanilla LLMs perform surprisingly well on the bandgaps of common semiconductors (Table 1),

with expected systematic deviation from MP values retrieved by LLaMP¹. This is likely due to the extensive literature on experimental semiconductor bandgaps. On the contrary, vanilla LLMs lack intrinsic knowledge of the bandgaps for the queried multi-element materials and exhibit low confidence or refuse to make predictions (Table 1, Table B5.8), whereas LLaMP retrieves accurate data with high SCoR and low MAE simultaneously and correctly identifies the stable polymorph's bandgap when multiple forms are present.

Ablation study. Our ablation study evaluates the impact of factual grounding and function-calling capabilities on LLaMP's performance. In Table 1, we examine three variants: (1) LLaMP; (2) GPT-4+ReAct with SerpAPI for internet browsing; (3) vanilla GPT-4. LLaMP achieved the best performance when using the complete set of MP tools, highlighting the importance of grounding in upto-date, high-fidelity materials databases. In Section 3, we mentioned the importance of hierarchical planning for robust function calls. Evaluating several backbone models on bulk moduli, formation energy prediction, and bandgap, we found LLaMP's grounding performance correlates with the function-calling capability of backbone LLM: Claude-3.5-Sonnet (#1) > Gemini-1.5-Flash (#24) > and Llama3-8B (#46). The number following each model refers to its ranking on the Berkeley Function-Calling Leaderboard at the time of the experiment (Yan et al., 2024). These results highlight LLaMP's robust adaptability across different LLM backbones, demonstrating that it's hierarchical multi-agent planning consistently enhances performance and scale with the function-calling capability.

High-fidelity and higher-order data retrieval.

To assess LLaMP's effectiveness in handling large-scale, complex material properties, we randomly selected 800 unary, binary, and ternary compounds from the Materials Project (MP). We then evaluated LLaMP, GPT-3.5, and GPT-4 on magnetic ordering classification and total magnetization prediction, tasks that are inherently challenging due to the subtle distinctions between magnetic orderings and magnetization units. Our

result indicates that vanilla LLMs suffer from hallucinations and misclassify the magnetic orderings of materials. LLaMP with GPT-4 as backend counteracts the intrinsic bias of GPT models, increasing the classification accuracy to 0.98 and R^2 of magnetization prediction to 0.992 (Table 2). We note that GPT-3.5 as backend struggles to distinguish total_magnetization from magnetization_per_formula_unit and often forgets to normalize the values. In the magnetic orderings queries, LLaMP with GPT-3.5 as backend fails to distinguish ferromagnetic (FM) and ferrimagnetic (FiM) orderings, while LLaMP with GPT-4 as backend gracefully separates the two classes (Figure 2a, d).

Table 2: Retrieval performance of LLaMP, GPT-3.5, and GPT-4 on magnetic orderings and magnetization. LLaMP with GPT-4 and GPT-3.5 as backend LLM are compared.

	Magnetic C	Ordering	Magne	tization
	Accuracy	F1	MAE	R^2
LLaMP (GPT-4)	0.98	0.89	0.045	0.992
GPT-4	0.48	0.26	1.611	-0.201
LLaMP (GPT-3.5)	0.96	0.88	1.896	0.407
GPT-3.5	0.23	0.18	1.988	-0.024

We further test the capability of LLaMP and LLMs for higher-order data (such as tensors, 3D crystal structures, curves). As shown in Table B5.2, GPT-3.5 hallucinates the values for the components in the elastic tensor of NaCl, with serious erroneous values such as $C_{11}=289.2\,\mathrm{GPa}$ —a significant deviation from DFT-calculated values (76 GPa). It also omits the values for $C_{22},C_{33},C_{55},C_{66}$ and fails to represent the full elastic tensor in a matrix format, despite the query explicitly requesting the full elastic tensor. This highlights the limitation of intrinsic knowledge in LLMs to recall higher-order, more complex data.

4.2 Real-World Applications

Inorganic synthesis recipes. Equipped with the MP synthesis endpoint (Kononova et al., 2019), LLaMP can extract synthesis recipes and summarize detailed step-by-step procedures grounded on real experimental papers with associated DOI references, as demonstrated in the example queries (Table B5.9 and B5.10).

Vanilla LLMs often give seemingly correct and verbose synthesis procedures but pull irrelevant compounds into the recipes and overlook more

¹Bandgaps calculated from generalized gradient approximation (GGA) functional are known to underestimate the experimental values by 40-50% (Borlido et al., 2020). Strategies to improve bandgap prediction at moderate or low computational cost will be included in MP in the future.

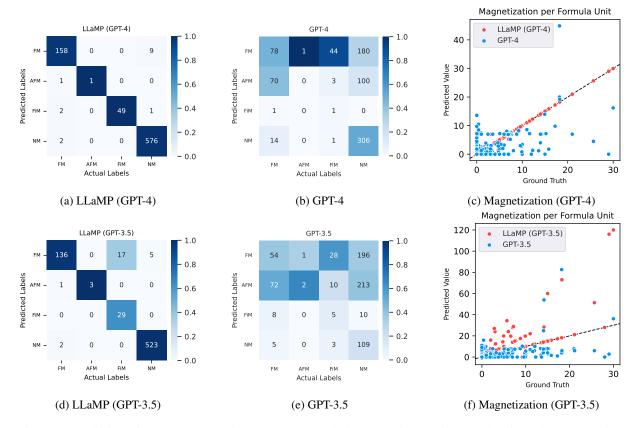


Figure 2: Prediction of LLaMP, GPT-3.5, and GPT-4 on (a,b,d,e) magnetic orderings and (c,f) total magnetization per formula unit of randomly selected materials. Confusion matrix presents the number of entries in each class. Colormap represents the percentage of correct classification. We define FM as ferromagnetic, AFM as antiferromagnetic, FiM as ferrimagnetic, and NM as nonmagnetic.

Table 3: Positive-unlabeled (PU) classification of LLaMP and baseline methods on inorganic material synthesizablity. (*) Evaluations on 352,236 positive and 40,817 unlabeled compounds by Kim et al. (2024).

	Accuracy	F1	Precision	Recall
LLaMP (GPT-4)	0.800	0.773	0.895	0.680
LLaMP (Sonnet)	0.818	0.812	0.848	0.780
GPT-4	0.600	0.649	0.578	0.740
Sonnet	0.530	0.230	0.636	0.140
Llama3	0.480	0.623	0.489	0.860
Gemini	0.590	0.388	0.765	0.260
GPT-4*	-	_	0.151	0.620
GPT-3.5 (FT)*	-	-	0.558	0.951
stoi-CGNF*	-	-	0.541	0.942

optimal or efficient reactions. In the example of YMnO₃ (Table B5.9), GPT-3.5 suggests the possible reaction pathways from two common oxide precursors (Y₂O₃ and MnO₂). However, it pulls irrelevant lithium compounds (Li₂CO₃ and LiOH) into the recipe and overlooks the fact that metathesis reactions (Li et al., 2015; Todd et al., 2021) require less applied energy than high-temperature sintering,

which relies on solid-state diffusion (Maximenko and Olevsky, 2004). Vanilla LLMs also exhibit uncertainty about specific synthesis details, such as heating temperature, duration, cooling rate, *etc*. In some edge cases such as LiFePO₄ presented in Table B5.10, the cited references are associated with the real papers but the contents are dissociated from the cited title and hallucinated from the pre-training corpus.

We further compare the performance of LLaMP on synthesizability prediction with stoichiometric convolutional graph neural fingerprint (stoi-CGNF) (Jang et al., 2024) and fine-tuned LLMs (Kim et al., 2024). We follow the positive-unlabeled (PU) classification task proposed in (Kim et al., 2024) by randomly selecting a subset of positive (probable) and unlabeled (unlikely) inorganic compounds and compare the classification performances of LLaMP with different backend LLMs and baselines. As presented in Table 3, LLaMP effectively enhances the performances of backbone GPT-4 and Sonnet LLMs by a significant margin of 20%, with the classification precision of LLaMP (GPT-4) up to

0.895.



Figure 3: Generation and manipulation of crystal structures using LLMs to insert an additional lithium atom at the interstitial site in diamond cubic silicon structure. Blue: Si. Green: Li. Question-answer pairs are listed in Table B5.11. Additional atoms extended through bonds are visualized.

Table 4: Structural parameters of generated crystals compared with diamond cubic (DC) silicon. Key metrics include fractional coordinates of inserted Li atom $(x,y,z)_{\rm Li}$, average Si—Si bond length $\ell_{\rm SiSi}$, cell volume V, and angular parameters (Si—Si and Si—Li—Si). GPT-4 refuse to respond due to safeguards against incomplete atomic structure information. Values in parentheses indicate the error rate relative to the ground truth.

Method	$(x, y, z)_{Li}$	ℓ_{SiSi} (Å)	$V(\mathring{\text{A}}^3)$	Angles (°)
				Si-Si-Si / Si-Li-Si
LLaMP	(0.5, 0.5, 0.5)	2.36 (0.0%)	40.33 (0.0%)	109.47 (0.0%) / 62.96
GPT-3.5	(0.5, 0.5, 0.5)	2.71 (+15.0%)	67.05 (+66.3%)	98.28 (-10.2%) / 67.69
GPT-4	-	-	-	-
DC Si (mp-149)	-	2.36	40.33	109.47 / -

RAG-assisted crystal editing. Fine-tuned LLMs for text-encoded atomistic information have shown the capability to edit stable crystals under the constraints of atomic positions and charges (Gruver et al., 2023). In this context, we delve into the examination and comparison of the crystal editing capabilities between LLaMP and GPT-3.5, without resorting to fine-tuning or tailored prompt messages in previous work. Figure 3 showcases the structures generated by LLaMP and vanilla GPT-3.5 without RAG, both instructed to insert one lithium atom at the tetrahedral site of the diamond cubic silicon structure (Table B5.11). Notably, both LLaMP and GPT-3.5 place an additional Li atom at fractional coordinate (0.5, 0.5, 0.5). However, the Si structure retrieved by LLaMP adheres to the MP convention, positioning two Si bases at (0.125, 0.125, 0.125) and (0.875, 0.875, 0.875). This causes the inserted Li atom to be hexagonal interstitial instead of tetrahedral interstitial.

GPT-3.5 locates the Li atom at the tetrahedral site given the "luckily chosen" Si bases at (0,0,0) and (0.25,0.25,0.25); however, the resulting cell

volume and shape are highly distorted, and the Si-Si bond length and Si-Si-Si angle deviate significantly from the ground truth (Table 4), highlighting the limitations in the intrinsic encoding of LLMs for atomistic information and the challenges associated with zero-shot generation of crystal structures. In contrast, the LLaMP-retrieved MP structure serves as a robust prior, anchoring the lattice parameters of the generated structure to the correct values.

Language-driven simulation. LLaMP equipped with Python REPL and atomistic simulation workflow package atomate2 performs well out of the box for complex multi-step simulations using pretrained universal interatomic potential MACE-MP-0 (Batatia et al., 2023) through language instruction. As demonstrated in Appendix C, LLaMP is able to follow multi-step instructions to fetch stable crystal structure from MP, generate a supercell of atomic structure, and run annealing molecular dynamics simulation with varying temperatures from 300K to 800K and back to 300K. After the simulation is finished, LLaMP can read the simulation trajectories and plot the temperature profile over time (Appendix C.1).

Failure and safeguard modes. In scientific research, refusing ambiguous or unsupported queries is critical to ensuring reliability. LLaMP demonstrates this capability by identifying limitations in data availability or query specificity and providing appropriate responses. For instance, when asked for the bulk modulus of "stainless steel," LLaMP points out the ambiguity of the query due to variable compositions. Similarly, it flags unavailable synthesis recipes or computed properties for specific materials. Examples of these safeguard modes can be found in Appendix A.5.

5 Conclusion

We introduce LLaMP, a hierarchical multi-agent framework that integrates LLMs with high-fidelity material informatics and domain-specific workflows for scientific research. By leveraging modular agents and robust hierarchical planning, LLaMP bridges the gap between general-purpose LLMs and complex materials science experimental workflows. Our results show that LLaMP consistently outperforms prior methods and demonstrates broad capabilities across real-world research tasks, including material property retrieval, crystal structure

editing, and multimodal data analysis. Through grounding in expert-verified databases and iterative self-correction, LLaMP mitigates hallucination, enhances response consistency, and enables complex experimental simulations. This framework highlights multi-agent knowledge grounding as a promising path toward scalable and autonomous scientific workflow

Limitation We recognize the effectiveness of LLaMP's framework relies on backbone LLM's function calling and reasoning capabilities. Sometimes LLMs misunderstand the description of schemas and therefore yield unexpected behaviors. The correctness of LLaMP is also subject to the quality of theoretical prediction and the comprehensiveness of the data in MP. Other than the underpredicted bandgaps by GGA functional, MP's ongoing effort to search all possible spin configurations has also not been completed. Most of the existing calculations in MP start from high-spin ferromagnetic configurations, which may overlook many antiferromagnetic ground states below the current energy convex hull. While MP is one of the most comprehensive materials databases, the crystal structures available in MP are not exhaustive but continuously expanding (Merchant et al., 2023), and would benefit from additional inter-metallic compounds and high-entropy materials from other databases such as AFLOW, OQMD, NOMAD, etc. (Curtarolo et al., 2012a; Kirklin et al., 2015b; Scheidgen et al., 2023a).

References

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity. *arXiv preprint*. ArXiv:2302.04023 [cs].

Ilyes Batatia, Philipp Benner, Yuan Chiang, Alin M. Elena, Dávid P. Kovács, Janosh Riebesell, Xavier R. Advincula, Mark Asta, William J. Baldwin, Noam Bernstein, Arghya Bhowmik, Samuel M. Blau, Vlad Cărare, James P. Darby, Sandip De, Flaviano Della Pia, Volker L. Deringer, Rokas Elijošius, Zakariya El-Machachi, and 49 others. 2023. A foundation model for atomistic materials chemistry. arXiv preprint. ArXiv:2401.00096 [cond-mat, physics:physics].

Daniil A. Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. 2023. Autonomous chemical research

with large language models. *Nature*, 624(7992):570–578. Number: 7992 Publisher: Nature Publishing Group.

Pedro Borlido, Jonathan Schmidt, Ahmad W. Huran, Fabien Tran, Miguel A. L. Marques, and Silvana Botti. 2020. Exchange-correlation functionals for band gaps of solids: benchmark, reparametrization and machine learning. *npj Computational Materials*, 6(1):1–17. Number: 1 Publisher: Nature Publishing Group.

Nicholas Carlini, Chang Liu, Jernej Kos, Úlfar Erlingsson, and Dawn Song. 2018. The secret sharer: Measuring unintended neural network memorization & extracting secrets. *arXiv preprint arXiv:1802.08232*, 5.

Junbum Cha, Wooyoung Kang, Jonghwan Mun, and Byungseok Roh. 2024. Honeybee: Locality-enhanced Projector for Multimodal LLM. *arXiv preprint*. ArXiv:2312.06742.

Harrison Chase. 2022. LangChain.

Stefano Curtarolo, Wahyu Setyawan, Gus L. W. Hart, Michal Jahnatek, Roman V. Chepulskii, Richard H. Taylor, Shidong Wang, Junkai Xue, Kesong Yang, Ohad Levy, Michael J. Mehl, Harold T. Stokes, Denis O. Demchenko, and Dane Morgan. 2012a. AFLOW: An automatic framework for high-throughput materials discovery. *Computational Materials Science*, 58:218–226.

Stefano Curtarolo, Wahyu Setyawan, Gus LW Hart, Michal Jahnatek, Roman V Chepulskii, Richard H Taylor, Shidong Wang, Junkai Xue, Kesong Yang, Ohad Levy, and 1 others. 2012b. Aflow: An automatic framework for high-throughput materials discovery. *Computational Materials Science*, 58:218–226.

John Dagdelen, Alexander Dunn, Sanghoon Lee,
Nicholas Walker, Andrew S. Rosen, Gerbrand Ceder,
Kristin A. Persson, and Anubhav Jain. 2024. Structured information extraction from scientific text with large language models. *Nature Communications*, 15(1):1418. Publisher: Nature Publishing Group.

Yu Du, Fangyun Wei, and Hongyang Zhang. 2024. Anytool: Self-reflective, hierarchical agents for large-scale api calls. *arXiv* preprint arXiv:2402.04253.

Nathan Frey, Ryan Soklaski, Simon Axelrod, Siddharth Samsi, Rafael Gomez-Bombarelli, Connor Coley, and Vijay Gadepally. 2022. Neural Scaling of Deep Chemical Models.

Alireza Ghafarollahi and Markus J. Buehler. 2024. AtomAgents: Alloy design and discovery through physics-aware multi-modal multi-agent artificial intelligence. *arXiv preprint*. ArXiv:2407.10022.

Nate Gruver, Anuroop Sriram, Andrea Madotto, Andrew Wilson, C. Lawrence Zitnick, and Zachary Ulissi. 2023. Fine-Tuned Language Models Generate Stable Inorganic Materials as Text.

- Tanishq Gupta, Mohd Zaki, N. M. Anoop Krishnan, and Mausam. 2022. MatSciBERT: A materials domain language model for text mining and information extraction. *npj Computational Materials*, 8(1):1–11. Number: 1 Publisher: Nature Publishing Group.
- Kausik Hira, Mohd Zaki, Dhruvil Sheth, Mausam, and N. M. Anoop Krishnan. 2024. Reconstructing the materials tetrahedron: challenges in materials information extraction. *Digital Discovery*, 3(5):1021–1037.
- Kevin Maik Jablonka, Philippe Schwaller, Andres Ortega-Guerrero, and Berend Smit. 2024. Leveraging large language models for predictive chemistry. *Nature Machine Intelligence*, 6(2):161–169. Publisher: Nature Publishing Group.
- Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards, Stephen Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, and Kristin A. Persson. 2013. Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Materials*, 1(1):011002.
- Jidon Jang, Juhwan Noh, Lan Zhou, Geun Ho Gu, John M. Gregoire, and Yousung Jung. 2024. Synthesizability of materials stoichiometry using semisupervised learning. *Matter*, 7(6):2294–2312. Publisher: Elsevier.
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2022. Large language models struggle to learn long-tail knowledge. arxiv. arXiv preprint arXiv:2211.08411.
- Seongmin Kim, Yousung Jung, and Joshua Schrier. 2024. Large Language Models for Inorganic Synthesis Predictions. *Journal of the American Chemical Society*, 146(29):19654–19659. Publisher: American Chemical Society.
- Scott Kirklin, James Saal, Bryce Meredig, Alex Thompson, Jeff Doak, Muratahan Aykol, Stephan Rühl, and Chris Wolverton. 2015a. The open quantum materials database (oqmd): Assessing the accuracy of dft formation energies. *npj Computational Materials*, 1:15010.
- Scott Kirklin, James E. Saal, Bryce Meredig, Alex Thompson, Jeff W. Doak, Muratahan Aykol, Stephan Rühl, and Chris Wolverton. 2015b. The Open Quantum Materials Database (OQMD): assessing the accuracy of DFT formation energies. *npj Computational Materials*, 1(1):1–15. Publisher: Nature Publishing Group.
- Olga Kononova, Haoyan Huo, Tanjin He, Ziqin Rong, Tiago Botari, Wenhao Sun, Vahe Tshitoyan, and Gerbrand Ceder. 2019. Text-mined dataset of inorganic materials synthesis recipes. *Scientific Data*, 6(1):203. Number: 1 Publisher: Nature Publishing Group.
- G. Kresse and J. Furthmüller. 1996. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. *Physical Review B*,

- 54(16):11169–11186. Publisher: American Physical Society.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *arXiv preprint*. ArXiv:2005.11401 [cs] version: 4.
- Jun Li, Andrew E. Smith, Peng Jiang, Judith K. Stalick, Arthur W. Sleight, and M. A. Subramanian. 2015. True Composition and Structure of Hexagonal "YAlO3", Actually Y3Al3O8CO3. *Inorganic Chemistry*, 54(3):837–844. Publisher: American Chemical Society.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, and 31 others. 2023. Holistic Evaluation of Language Models. *arXiv preprint*. ArXiv:2211.09110 [cs].
- Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. 2024. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*.
- Youzhi Luo, Chengkai Liu, and Shuiwang Ji. 2023. Towards Symmetry-Aware Generation of Periodic Materials.
- Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D. White, and Philippe Schwaller. 2024. Augmenting large language models with chemistry tools. *Nature Machine Intelligence*, pages 1–11. Publisher: Nature Publishing Group.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories. *arXiv preprint*. ArXiv:2212.10511 [cs].
- Andrey L. Maximenko and Eugene A. Olevsky. 2004. Effective diffusion coefficients in solid-state sintering. *Acta Materialia*, 52(10):2953–2963.
- Amil Merchant, Simon Batzner, Samuel S. Schoenholz, Muratahan Aykol, Gowoon Cheon, and Ekin Dogus Cubuk. 2023. Scaling deep learning for materials discovery. *Nature*, pages 1–6. Publisher: Nature Publishing Group.

- Santiago Miret and N. M. Anoop Krishnan. 2024. Are LLMs Ready for Real-World Materials Discovery? *arXiv preprint*. ArXiv:2402.05200 [cond-mat].
- John X Morris, Chawin Sitawarin, Chuan Guo, Narine Kokhlikyan, G Edward Suh, Alexander M Rush, Kamalika Chaudhuri, and Saeed Mahloujifar. 2025. How much do language models memorize? *arXiv* preprint arXiv:2505.24832.
- Siru Ouyang, Zhuosheng Zhang, Bing Yan, Xuan Liu, Yejin Choi, Jiawei Han, and Lianhui Qin. 2024. Structured Chemistry Reasoning with Large Language Models. *arXiv preprint*. ArXiv:2311.09656 [cs].
- G. Pilania, J. E. Gubernatis, and T. Lookman. 2017. Multi-fidelity machine learning models for accurate bandgap predictions of solids. *Computational Materials Science*, 129:156–163.
- Maciej P. Polak and Dane Morgan. 2023. Extracting Accurate Materials Data from Research Papers with Conversational Language Models and Prompt Engineering. *arXiv preprint*. ArXiv:2303.05352 [condmat].
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. ToolLLM: Facilitating Large Language Models to Master 16000+ Real-world APIs. arXiv preprint. ArXiv:2307.16789.
- Andre Niyongabo Rubungo, Kangming Li, Jason Hattrick-Simpers, and Adji Bousso Dieng. 2024. Llm4mat-bench: Benchmarking large language models for materials property prediction. *arXiv preprint arXiv:2411.00177*.
- Markus Scheidgen, Lauri Himanen, Alvin Noe Ladines, David Sikter, Mohammad Nakhaee, Ádám Fekete, Theodore Chang, Amir Golparvar, José A. Márquez, Sandor Brockhauser, Sebastian Brückner, Luca M. Ghiringhelli, Felix Dietrich, Daniel Lehmberg, Thea Denell, Andrea Albino, Hampus Näsström, Sherjeel Shabih, Florian Dobener, and 10 others. 2023a. NO-MAD: A distributed web-based platform for managing materials science research data. *Journal of Open Source Software*, 8(90):5388.
- Markus Scheidgen, Lauri Himanen, Alvin Noe Ladines, David Sikter, Mohammad Nakhaee, Ádám Fekete, Theodore Chang, Amir Golparvar, José A. Márquez, Sandor Brockhauser, Sebastian Brückner, Luca M. Ghiringhelli, Felix Dietrich, Daniel Lehmberg, Thea Denell, Andrea Albino, Hampus Näsström, Sherjeel Shabih, Florian Dobener, and 10 others. 2023b. Nomad: A distributed web-based platform for managing materials science research data. *Journal of Open Source Software*, 8(90):5388.

- Mara Schilling-Wilhelmi, Martiño Ríos-García, Sherjeel Shabih, María Victoria Gil, Santiago Miret, Christoph T. Koch, José A. Márquez, and Kevin Maik Jablonka. 2024. From Text to Insight: Large Language Models for Materials Science Data Extraction. *arXiv preprint*. ArXiv:2407.16867.
- Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary Lipton, and J Zico Kolter. 2024. Rethinking llm memorization through the lens of adversarial compression. *Advances in Neural Information Processing Systems*, 37:56244–56267.
- Yihan Shao, Zhengting Gan, Evgeny Epifanovsky, Andrew T.B. Gilbert, Michael Wormit, Joerg Kussmann, Adrian W. Lange, Andrew Behn, Jia Deng, Xintian Feng, Debashree Ghosh, Matthew Goldey, Paul R. Horn, Leif D. Jacobson, Ilya Kaliman, Rustam Z. Khaliullin, Tomasz Kuś, Arie Landau, Jie Liu, and 138 others. 2015. Advances in molecular quantum chemistry contained in the Q-Chem 4 program package. *Molecular Physics*, 113(2):184–215. Publisher: Taylor & Francis.
- Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language Agents with Verbal Reinforcement Learning. *arXiv preprint*. ArXiv:2303.11366 [cs].
- Yu Song, Santiago Miret, and Bang Liu. 2023. MatSci-NLP: Evaluating Scientific Language Models on Materials Science Language Tasks Using Text-to-Schema Modeling. *arXiv preprint*. ArXiv:2305.08264.
- Evan Walter Clark Spotte-Smith, Orion Archer Cohen, Samuel Blau, Jason Munro, Ruoxi Yang, Rishabh Guha, Hetal Patel, Sudarshan Vijay, Patrick Huck, Ryan Kingsbury, Matthew Horton, and Kristin Persson. 2023. A database of molecular properties integrated in the Materials Project.
- Paul K. Todd, Matthew J. McDermott, Christopher L.
 Rom, Adam A. Corrao, Jonathan J. Denney, Shyam S.
 Dwaraknath, Peter G. Khalifah, Kristin A. Persson, and James R. Neilson. 2021. Selectivity in Yttrium Manganese Oxide Synthesis via Local Chemical Potentials in Hyperdimensional Phase Space. *Journal of the American Chemical Society*, 143(37):15185–15194. Publisher: American Chemical Society.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, and 1 others. 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345.
- Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2023. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. arXiv preprint arXiv:2307.10635.

- Mingjian Wen, Matthew K. Horton, Jason M. Munro, Patrick Huck, and Kristin A. Persson. 2024. An equivariant graph neural network for the elasticity tensors of all seven crystal systems. *Digital Discovery*, 3(5):869–882. Publisher: RSC.
- Mingjian Wen, Evan Walter Clark Spotte-Smith, Samuel M. Blau, Matthew J. McDermott, Aditi S. Krishnapriyan, and Kristin A. Persson. 2023. Chemical reaction networks and opportunities for machine learning. *Nature Computational Science*, 3(1):12–24. Publisher: Nature Publishing Group.
- Tong Xie, Yuwei Wan, Wei Huang, Zhenyu Yin, Yixuan Liu, Shaozhou Wang, Qingyuan Linghu, Chunyu Kit, Clara Grazian, Wenjie Zhang, Imran Razzak, and Bram Hoex. 2023. DARWIN Series: Domain Specific Large Language Models for Natural Science. *arXiv preprint*. ArXiv:2308.13565 [cond-mat, physics:physics].
- Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. 2024. Hallucination is Inevitable: An Innate Limitation of Large Language Models. *arXiv preprint*. ArXiv:2401.11817 [cs].
- Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. 2024. Berkeley Function Calling Leaderboard.
- Samuel Yang, Shutong Li, Subhashini Venugopalan, Vahe Tshitoyan, Muratahan Aykol, Amil Merchant, Ekin Cubuk, and Gowoon Cheon. 2023. Accurate Prediction of Experimental Band Gaps from Large Language Model-Based Data Extraction.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. *arXiv preprint*. ArXiv:2210.03629 [cs].
- D. Zagorac, H. Müller, S. Ruehl, J. Zagorac, and S. Rehme. 2019. Recent developments in the Inorganic Crystal Structure Database: theoretical crystal structure data and related features. *Journal of Applied Crystallography*, 52(5):918–925. Publisher: International Union of Crystallography.
- Mohd Zaki, Jayadeva, Mausam, and N. M. Anoop Krishnan. 2023. MaScQA: A Question Answering Dataset for Investigating Materials Science Knowledge of Large Language Models. *arXiv preprint*. ArXiv:2308.09115.
- Claudio Zeni, Robert Pinsler, Daniel Zügner, Andrew Fowler, Matthew Horton, Xiang Fu, Sasha Shysheya, Jonathan Crabbé, Lixin Sun, Jake Smith, Bichlien Nguyen, Hannes Schulz, Sarah Lewis, Chin-Wei Huang, Ziheng Lu, Yichi Zhou, Han Yang, Hongxia Hao, Jielan Li, and 2 others. 2024. MatterGen: a generative model for inorganic materials design. *arXiv* preprint. ArXiv:2312.03687 [cond-mat].

- Huan Zhang, Yu Song, Ziyu Hou, Santiago Miret, and Bang Liu. 2024. HoneyComb: A Flexible LLM-Based Agent System for Materials Science. arXiv preprint. ArXiv:2409.00135.
- Zhiling Zheng, Oufan Zhang, Christian Borgs, Jennifer T. Chayes, and Omar M. Yaghi. 2023. Chat-GPT Chemistry Assistant for Text Mining and Prediction of MOF Synthesis. *Journal of the American Chemical Society*, 145(32):18048–18062. ArXiv:2306.11296 [cond-mat, physics:physics].

A Supplementary Information

A.1 Example Workflow

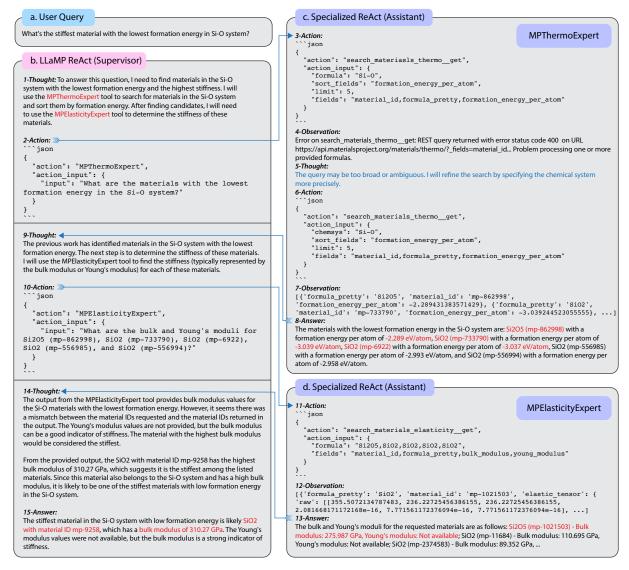


Figure A.1: Multimodal retrieval-augmented generation for materials informatics. (a) User query. (b) Supervisor ReAct agent capable of handling multiple assistant agents and high-level reasoning. (c-d) Assistant ReAct agents executing function calling and summarization. (c) MPThermoExpert and (d) MPElasticityExpert have access to the API schemas of thermo and elasticity endpoints on Materials Project, respectively. The selected details are highlighted in red, demonstrating the capabilities of RAG and ReAct implemented in LLaMP. The blue texts show LLaMP assistant ReAct agent can handle API calling errors and self-correct the input query accordingly.

A.2 List of Implemented Assistant Agents and Tools

Here we provide the comprehensive list of implemented **assistant agents** and tools. Notably, **MP Assistants** are designed with a highly modular architecture, making it straightforward to extend support for additional API endpoints from https://api.materialsproject.org/docs.

- **MPSummaryExpert**: summary provides amalgamated data for a material by combining subsets of data from many of the other API endpoints.
- **MPThermoExpert**: thermo provides computed thermodynamic data for a material such as formation energy and energy above hull.

- **MPElasticityExpert**: elasticity provides bulk, shear, and Young's modulus, poisson ratio, and universal anisotropy index.
- MPMagnetismExpert: magnetism provides computed magnetic ordering related data.
- **MPDielectricExpert**: dielectric provides computed dielectric data from density functional perturbation theory.
- **MPPiezoelectricExpert**: piezoelectric provides computed piezoelectric data from density functional perturbation theory.
- **MPElectronicExpert**: electronic_structure provides computed electronic structure related data for a material such as band gap and fermi level. Python objects for line-mode band structures, density of states, and fermi surfaces are also available.
- **MPSynthesisExpert**: synthesis provides a synthesis recipes for materials extracted from literature using text mining and natural language processing techniques.
- MPStructureRetriever: MaterialsStructureText fetches and saves pymatgen Structure objects to local JSON files.
- MLFFAgent: MLFFMD runs molecular dynamics simulations using pre-trained machine learning force fields; MLFFElastic calculates the elastic constants of a given material using pre-trained machine learning force fields.
- PythonREPLTool: Python REPL that LLMs could run the generated script.
- ArxivQueryRun: LangChain built-in tool that LLMs can use to send API request to ArXiv.
- WikipediaQueryRun: LangChain built-in tool that LLMs can use to send API request to Wikipedia.

A.3 Metric calculation in Table 1 and Figure A.2

The following procedures are adopted to calculate the metrics for material property regression benchmarks presented in Table 1 and Figure A.2:

- 1. Each method was presented with the same query asking for the property of multiple materials. Here we ask each method for bulk modulus, formation energy, and electronic bandgap of ten materials. For example, "What are the bulk moduli of the following metals: Sc, Ti, V, Cr, Mn, Fe, Co, Ni, Cu, Zn?".
- 2. Repeat the same query for five times and collect the responses.
- 3. The numerical values are extracted and formatted into 2D arrays with the aid by LLMs. If the range is provided in the response (*e.g.* Llama 3), the median value was used.
- 4. Calculate Precision, CoP, Confidence, and SCoR for each method across five trials on all materials. The code is provided and the pseudocode can be written as follows:

Input:arr (2D array of five responses for different materials) $n \leftarrow \text{count of valid responses (non-NaN values)}$ in each column of arr $\text{prec} \leftarrow \frac{\text{nanstd(arr,axis=0)}}{\sqrt{n}}$ $\text{cop} \leftarrow \text{mean}(\exp(-\text{prec}))$ $\text{conf} \leftarrow \text{mean}\left(\frac{n}{\# \text{ of trials }N}\right)$ $\text{scor} \leftarrow \begin{cases} 0 & \text{if conf} = 0 \text{ for all columns} \\ \text{cop} \times \text{conf} & \text{otherwise} \end{cases}$ $\text{prec} \leftarrow \text{mean}(\text{prec})$ Output:prec, cop, conf, scor

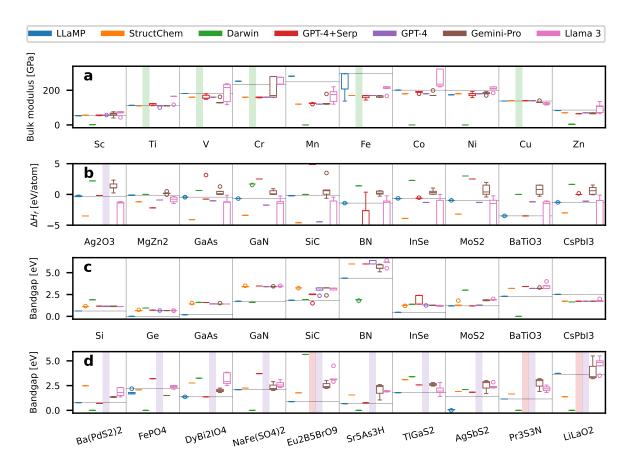


Figure A.2: Boxplot of LLaMP RAG responses, baseline methods, and LLM intrinsic knowledge on material properties. (a) Bulk moduli, K, of 3d transition metals. (b) Formation energies, ΔH_f , of common compounds. (c) Electronic bandgaps, E_g , of common intrinsic semiconductors. (d) Electronic bandgaps of multi-element (ternary or quaternary) materials. Missing predictions are marked by shaded areas. Fliers (Outliers) are marked in circles. Horizontal lines represent the MP reference data. All LLaMP results use GPT-4 as backend language provider. Method with higher SCoR has narrower distribution. LLaMP is effectively grounded on MP reference across different tasks and materials.

A.4 Material Property Retrieval Result

A.5 Safeguard or failure modes

- Question: What is the bulk modulus of stainless steel? ambiguous chemical formula LLaMP: Stainless steel is an alloy with a composition that can vary widely, including elements like iron, chromium, nickel, and molybdenum. Because of this variability, it does not have a single chemical formula or a unique Materials Project ID, and its bulk modulus can significantly vary based on its exact composition and heat treatment. The Materials Project database, focusing on inorganic crystalline materials, may not provide specific data on complex alloys like stainless steel. For precise values of the bulk modulus, consulting specific datasheets or literature on the type of stainless steel in question is recommended.
- Question: How to synthesize Y2Be2Fe(SiO5)2? data unavailable LLaMP: Currently, there are no synthesis recipes available in the Materials Project database for Y2Be2Fe(SiO5)2.
- Question: What's the phonon bandgap of Ac2HgGe? computed property unavailable LLaMP: The material Ac2HgGe, with a Materials Project ID of mp-861872, does not have its phonon bandgap information directly available through the Materials Project database. For detailed phonon properties, further analysis using specialized phonon calculation tools or databases would be necessary.

A.6 Prompt Template

We use the ReAct template hwchase17/react-multi-input-json from LangChain Hub (https://smith.langchain.com/hub/hwchase17/react-json) as follows:

```
Answer the following questions as best you can. You have access to the following tools:
{tools}
The way you use the tools is by specifying a JSON blob.
Specifically, this JSON should have an `action` key (with the name of the tool to use) and an `
    action_input` key (with the input to the tool going here).
The only values that should be in the "action" field are: {tool_names}
The $JSON_BLOB should only contain a SINGLE action, do NOT return a list of multiple actions.
    Here is an example of a valid $JSON_BLOB:
{{
    "action": $TOOL_NAME,
  "action_input": $INPUT
}}
ALWAYS use the following format:
Question: the input question you must answer
Thought: you should always think about what to do
Action:
$JSON_BLOB
Observation: the result of the action
 .. (this Thought/Action/Observation can repeat N times)
Thought: I now know the final answer
Final Answer: the final answer to the original input question
Begin! Reminder to always use the exact characters `Final Answer` when responding.
```

The above system prompt message is prepended with the following solicitation to guide the supervisor agent:

You name is LLaMP and you are a helpful agent that can consult materials-related data through Materials Project (MP) database, arXiv, Wikipedia, and run calculations/simulations via a python REPL and MP atomate2 computational workflow.

If you get an error, debug your code and try again. Only use the output of your code to answer the question. Ask user to clarify their queries if needed. Please note that you don't have direct control over MP but through multiple assistant agents to help you. You need to provide complete context in the input for assistants to do their job. REFINE and CRITQUE the output of the assistants if needed. You can end the conversation by saying "Bye".

B Sample Questions

Table 5: QA comparison between LLaMP and GPT-3.5.

	Overtion		Answer
	Question	LLaMP	GPT-3.5
B5.1	What are the bulk moduli of the following metals: Sc, Ti, V, Cr, Mn, Fe, Co, Ni, Cu, Zn?	The bulk moduli of the metals you asked for are as follows (in GPa): - Scandium (Sc): Voigt = 45.715, Reuss = 45.34, VRH = 45.528 - Titanium (Ti): Voigt = 112.955, Reuss = 112.807, VRH = 112.881 - Vanadium (V): Voigt = 181.232, Reuss = 181.232, VRH = 181.232 - Chromium (Cr): Voigt = 252.155, Reuss = 252.155, VRH = 252.155 - Manganese (Mn): Voigt = 280.959, Reuss = 280.959, VRH = 280.959 - Iron (Fe): Voigt = 296.003, Reuss = 295.207, VRH = 295.605 - Cobalt (Co): Voigt = 200.587, Reuss = 200.587, VRH = 200.587 - Nickel (Ni): Voigt = 197.399, Reuss = 197.29, VRH = 197.345 - Copper (Cu): Voigt = 151.394, Reuss = 151.394, VRH = 151.394 - Zinc (Zn): Voigt = 76.283, Reuss = 95.46, VRH = 85.872	The bulk modulus (K) is a measure of a material's resistance to uniform compression. It is defined as the ratio of the change in pressure to the fractional volume compression. The values for the bulk modulus can vary somewhat depending on the specific conditions and crystallographic structures of the materials. The values provided here are approximate and given in gigapascals (GPa): - Scandium (Sc): 55 GPa - Titanium (Ti): 110 GPa - Vanadium (V): 160 GPa - Chromium (Cr): 160 GPa - Manganese (Mn): 120 GPa - Iron (Fe): 170 GPa - Cobalt (Co): 180 GPa - Nickel (Ni): 180 GPa - Copper (Cu): 140 GPa - Zinc (Zn): 70 GPa These values represent typical bulk moduli for these metals, and actual values may vary depending on factors such as temperature and pressure.
			Continued on next page

Table 5 – continued from previous page

	Overtio-	Table 3 – continued	Answer
	Question	LLaMP	GPT-3.5
B5.2	What is	The elastic tensor for	The elastic tensor describes the relationship be-
	the full	NaCl is available for two	tween stress and strain in a material. It is a 4th-
	elastic	different entries on	rank tensor that characterizes the elastic prop-
	tensor of	Materials Project. For	erties of a crystal. In the case of NaCl (sodium
	NaCl?	material_id 'mp-22851',	chloride), the elastic tensor is typically repre-
		the elastic tensor in IEEE	sented by a 6x6 matrix.
		format is [[76.0, 1.0, 1.0,	The elastic constants of NaCl can vary slightly
		-0.0, 0.0, -0.0], [1.0, 76.0,	depending on the specific conditions and crys-
		1.0, -0.0, -0.0, 0.0], [1.0,	tallographic orientation. Here is an approximate
		1.0, 76.0, 0.0, -0.0, -0.0],	form of the elastic tensor for NaCl:
		[-0.0, -0.0, 0.0, -3.0, -0.0,	C11 = 289.2 GPa C12 = 97.8 GPa C44 = 81.6
		-0.0], [0.0, -0.0, -0.0, -0.0,	GPa C13 = C14 = C15 = C16 = C23 = C24 =
		-3.0, 0.0], [-0.0, 0.0, -0.0,	C25 = C26 = 0 GPa
		-0.0, 0.0, -3.0]]. For	Note: The values are given in gigapascals (GPa).
		material_id 'mp-22862',	The elastic constants Cij represent the compo-
		the elastic tensor in IEEE	nents of the elastic tensor. In this case, C11,
		format is [[47.0, 12.0,	C12, and C44 are the non-zero elastic constants.
		12.0, 0.0, -0.0, -0.0], [12.0,	The off-diagonal terms C13, C14, C15, C16,
		47.0, 12.0, 0.0, -0.0, -0.0],	C23, C24, C25, and C26 are typically zero due
		[12.0, 12.0, 47.0, 0.0, 0.0,	to cubic symmetry.
		-0.0], [0.0, 0.0, 0.0, 12.0,	Keep in mind that these values are approximate,
		-0.0, -0.0], [-0.0, -0.0, 0.0,	and there may be slight variations depending
		-0.0, 12.0, 0.0], [-0.0, -0.0,	on the source and the specific conditions under
		-0.0, -0.0, 0.0, 12.0]]	which the measurements were made. Experi-
			mental techniques, such as Brillouin scattering
			or ultrasonic measurements, are often used to
			determine the elastic constants of materials.
			Continued on next page

Table 5 – continued from previous page

	Answer				
ļ ļ	Question	LLaMP	GPT-3.5		
B5.3	What	The ground state	As of my last knowledge update in January		
	are the	bandgaps for the	2022, I can provide you with approximate val-		
ļ	bandgaps	requested materials are as	ues for the ground state bandgaps of the materi-		
ļ	of the	follows: Si: 1.12 eV, Ge:	als you mentioned. Keep in mind that these val-		
ļ	following	0.66 eV, GaAs: 1.42 eV,	ues can vary depending on factors like temper-		
ļ	materials:	GaN: 3.28 eV, SiC: 2.36	ature and crystal structure. Additionally, more		
ļ	Si, Ge,	eV, BN: 6.4 eV, InSe:	recent research may have refined these values.		
	and GaAs,	1.26 eV, MoS2: 1.29 eV,	Here are the approximate ground state bandgaps		
	GaN, SiC,	BaTiO3: 3.2 eV, CsPbI3:	in electron volts (eV):		
ļ	BN, InSe,	1.73 eV. Please note that	Silicon (Si): 1.1 eV Germanium (Ge): 0.67 eV		
ļ	MoS2,	these values are typical	Gallium Arsenide (GaAs): 1.42 eV Gallium Ni-		
ļ	BaTiO3,	for bulk materials and can	tride (GaN): 3.4 eV Silicon Carbide (SiC): 3.2		
ļ	CsPbI3?	vary depending on the	eV (for 4H-SiC) Boron Nitride (BN): 6 eV In-		
ļ	Only	specific structure and	dium Selenide (InSe): 1.25 eV Molybdenum		
	look for	calculation method.	Disulfide (MoS2): 1.3 eV (indirect bandgap)		
	ground-		Barium Titanate (BaTiO3): 3.2 eV Cesium		
	state		Lead Iodide (CsPbI3): 1.7 eV Please note that		
	bandgaps.		these values are general approximations and		
ļ ļ	8.1		may not represent the specific conditions or vari-		
			ations in different forms of these materials. For		
ļ			the most accurate and up-to-date information, I		
ļ ļ			recommend checking recent scientific literature		
ļ ļ			or databases.		
B5.4	Give	The maximum	The decomposition enthalpies (in kJ/mol) of the		
ļ ļ	me the	decomposition enthalpies	following materials are:		
ļ ļ	decom-	for the requested materials	Ag2O3: -335		
ļ ļ	position	are as follows: - Ag2O3:	MgZn2: -82		
ļ ļ	enthalpies	0.2565836414999989	GaAs: -317		
ļ ļ	of the	eV/atom - MgZn2: 0.1356	GaN: -505		
ļ	following	eV/atom - GaAs: 0.075	SiC: -40		
ļ	materials:	eV/atom - GaN: 0.07162	BN: -248		
	Ag2O3,	eV/atom - SiC: -0.172	InSe: -125		
	MgZn2,	eV/atom - BN: -0.358	MoS2: -393		
ļ	GaAs,	eV/atom - InSe: -0.01816	BaTiO3: -112		
ļ	GaN, SiC,	eV/atom - MoS2: 0.267	CsPbI3: -153		
ļ ļ	BN, InSe,	eV/atom - BaTiO3: 3.055			
	MoS2,	eV/atom - CsPbI3:			
	BaTiO3,	0.017536019000000902			
	CsPbI3. If	eV/atom Please note that			
	there are	the values for SiC, BN,			
	multiple	and InSe are negative,			
	values,	indicating that the			
	please	decomposition is			
	give me	exothermic. All other			
	the largest	values are positive,			
	value.	indicating endothermic			
	1	dagamnagitions	1		
		decompositions.	Continued on next page		

Table 5 – continued from previous page

		Table 3 – continued	Answer
	Question	LLaMP	GPT-3.5
B5.5	Give	The smallest formation	Ag2O3: -2.63 eV/atom MgZn2: -0.27 eV/atom
	me the	energies of the materials	GaAs: -4.13 eV/atom GaN: -3.51 eV/atom SiC:
	formation	you asked for are as	-4.98 eV/atom BN: -3.8 eV/atom InSe: -0.89
	energies	follows: Ag2O3: -0.3601	eV/atom MoS2: -3.31 eV/atom BaTiO3: -4.45
	of the	eV/atom, MgZn2:	eV/atom CsPbI3: -2.32 eV/atom
	following	-0.1306 eV/atom, GaAs:	
	materials:	-0.446 eV/atom, GaN:	
	Ag2O3,	-0.052 eV/atom, SiC:	
	MgZn2,	-0.204 eV/atom, BN: -1.41	
	GaAs,	eV/atom, InSe:	
	GaN, SiC,	-0.742582154921875	
	BN, InSe,	eV/atom, MoS2: -1.202	
	MoS2,	eV/atom, BaTiO3: -3.493	
	BaTiO3,	eV/atom, CsPbI3: -1.296	
	CsPbI3. If	eV/atom.	
	there are		
	multiple		
	values,		
	please		
	give		
	me the		
	smallest		
B5.6	value. Give	The smallest	The decomposition enthalpies of the materi-
B 3.0	me the	decomposition enthalpy	als are as follows: - Ag2O3: -1138.5 kJ/mol
	decom-	values for the requested	- MgZn2: -54.2 kJ/mol - GaAs: -452.2 kJ/mol -
	position	materials are as follows:	GaN: -510.9 kJ/mol - SiC: -2133.6 kJ/mol - BN:
	enthalpies	Ag2O3:	-2551.0 kJ/mol - InSe: -209.3 kJ/mol - MoS2: -
	of the	0.019373648083333528	739.3 kJ/mol - BaTiO3: -393.5 kJ/mol - CsPbI3:
	following	eV/atom, MgZn2: -0.0197	-65.3 kJ/mol
	materials:	eV/atom, GaAs: -0.446	
	Ag2O3,	eV/atom, GaN: -0.05231	
	MgZn2,	eV/atom, SiC: -0.204	
	GaAs,	eV/atom, BN: -1.41	
	GaN, SiC,	eV/atom, InSe: -0.06336	
	BN, InSe,	eV/atom, MoS2: -0.669	
	MoS2,	eV/atom, BaTiO3: -0.063	
	BaTiO3,	eV/atom, CsPbI3:	
	CsPbI3. If	-0.018602 eV/atom.	
	there are		
	multiple		
	values, please		
	give		
	me the		
	smallest		
	value.		
	varac.		Continued on next page
			Continued on next page

Table 5 – continued from previous page

	Owartion		Answer
	Question	LLaMP	GPT-3.5
B5.7	Give me	Here are the lowest	The energy above hull values for the given ma-
	the energy	energy above hull values	terials are as follows: - Ba2MnSe3: 0.0000
	above hull	for the requested	eV/atom - KAlSiO4: 0.0000 eV/atom - LiCuO:
	values	materials: Ba2MnSe3: 0.0	0.0000 eV/atom - H8C3N2O: 0.0000 eV/atom
	of the	eV/atom KAlSiO4: 0.0	- Dy2ZnS4: 0.0000 eV/atom - Rb2BeF4:
	following	eV/atom LiCuO: 0.0	0.0000 eV/atom - K9Fe2S7: 0.0000 eV/atom -
	materials:	eV/atom H8C3N2O:	Tb2FeC4: 0.0000 eV/atom - ScAlAu: 0.0000
	Ba2MnSe3,	0.0735 eV/atom	eV/atom - CaCrO4: 0.0000 eV/atom
	KAlSiO4,	Dy2ZnS4: 0.00966	
	LiCuO,	eV/atom Rb2BeF4: 0.0	
	H8C3N2O,	eV/atom K9Fe2S7: 0.0	
	Dy2ZnS4,	eV/atom Tb2FeC4: 0.0	
	Rb2BeF4,	eV/atom ScAlAu: 0.0	
	K9Fe2S7,	eV/atom CaCrO4: 0.0	
	Tb2FeC4,	eV/atom. Please note that	
	ScAlAu,	these are the lowest values	
	CaCrO4.	found in the Materials	
	If there	Project database and	
	are multi-	actual values may vary	
	ple values,	based on different entries.	
	please		
	give me		
	the lowest		
	one.		
			Continued on next page

Table 5 – continued from previous page

	Answer				
	Question	LLaMP	GPT-3.5		
D5 0	What				
B5.8		The bandgaps for the materials are:	I don't have real-time data access, and my train-		
			ing only includes information available up to		
	bandgaps of the	- Ba(PdS2)2: 0.779 eV	January 2022. As a result, I may not have spe-		
		- Most stable FePO4:	cific bandgap values for the materials you men-		
	following	2.1932 eV	tioned. However, I can provide you with general		
	materials:	- DyBi2IO4: 1.367 eV	information about some of these materials and		
	Ba(PdS2)2,	` ′	their known properties up to my last update:		
	FePO4,	2.094 eV	Ba(PdS2)2: I don't have specific bandgap infor-		
	DyBi2IO4,	- Eu2B5BrO9: 0.88 eV	mation for this compound. The bandgap may		
	` ′	2,- Sr5As3H: 0.676 eV	depend on factors such as crystal structure and		
		,- TlGaS2: 1.7885 eV	synthesis conditions.		
	Sr5As3H,	- Most stable AgSbS2:	FePO4: Iron(II) phosphate, known as ferrous		
	TlGaS2,	1.372 eV	phosphate, typically has a bandgap in the range		
	AgSbS2,	- Pr3S3N: 1.15 eV	of 2.5-3.0 eV.		
	Pr3S3N,	- Most stable LiLaO2:	DyBi2IO4: Dysprosium bismuth iodate. I don't		
	LiLaO2.	3.7336 eV	have specific bandgap information for this com-		
	If there	Please note that these	pound.		
	are poly-	values are obtained from	NaFe(SO4)2: Sodium iron sulfate. I don't		
	morphs,	the Materials Project	have specific bandgap information for this com-		
	give me	database and might be	pound.		
	only the	different in experimental	Eu2B5BrO9: Europium borate bromate. I don't		
	bandgap	results due to various	have specific bandgap information for this com-		
	of the	factors.	pound.		
	stable		Sr5As3H: Strontium arsenide hydride. I don't		
	material.		have specific bandgap information for this com-		
			pound.		
			TlGaS2: Thallium gallium sulfide. I don't		
			have specific bandgap information for this com-		
			pound.		
			AgSbS2: Silver antimony sulfide. This material		
			is a semiconductor with a bandgap typically		
			around 1.5-2.0 eV.		
			Pr3S3N: Praseodymium sulfide nitride. I don't		
			have specific bandgap information for this com-		
			pound.		
			LiLaO2: Lithium lanthanum oxide. I don't		
			have specific bandgap information for this com-		
			pound. The most stable structure for each material can		
			vary based on synthesis conditions and other fac-		
			tors. For accurate and up-to-date bandgap values and structural information, you may want to		
			ues and structural information, you may want to consult recent scientific literature or databases		
			that specialize in material properties.		
			Continued on next page		

Table 5 – continued from previous page

stion		Answer		
	LILAMD	GPT-3.5		
	LLaMP			
syn- ze	To synthesize YMnO ₃ , a solid-state reaction method can be employed. Here is a summarized procedure based on the synthesis details provided in the literature with DOI: 10.1021/ic502027k: 1. Dissolve Y(NO3)3·6H2O, Mn(NO3)2·4H2O, and a source of [OH-] in water to create a homogeneous solution. 2. Add citric acid to the solution to act as a mixing media. 3. Heat the mixture at 250°C for 2 hours to achieve pre-calcination. 4. Finally, heat the pre-calcined mixture at a temperature range of 700°C to 930°C for 10 to 18 hours in an air atmosphere to form YMnO ₃ .	Synthesizing YMnO ₃ (yttrium manganite) typically involves the solid-state reaction method. Please note that the details may vary depending on the specific conditions and equipment available in your laboratory. Here's a general procedure: Materials: Yttrium oxide (Y ₂ O ₃) Manganese oxide (MnO ₂ or Mn ₂ O ₃) High-purity manganese carbonate (MnCO ₃) Lithium carbonate (Li ₂ CO ₃) or lithium hydroxide (LiOH) Crucible and pestle High-temperature furnace Mortar and pestle Platinum crucible or alumina crucible Procedure: Step 1: Weighing the Precursors Weigh the appropriate amounts of yttrium oxide (Y ₂ O ₃), manganese oxide (MnO ₂ or Mn ₂ O ₃), and lithium carbonate (Li2CO ₃) or lithium hydroxide (LiOH) according to the stoichiometric ratio of YMnO ₃ . Step 2: Mixing the Powders Note: GPT-3.5 starts to fail numbering correctly in this response. Grind the powders thoroughly in a mortar and pestle to ensure a homogenous mixture. Step 3: Heating 3. Place the mixed powder in a crucible or alumina crucible. Heat the mixture in a high-temperature furnace. The specific temperature and heating rate depend on the exact synthesis conditions, but a typical range is 800-1100°C. Step 4: Sintering 5. Hold the mixture at the target temperature for a specific duration (several hours) to allow for sintering and the formation of YMnO ₃ . Step 5: Cooling 6. Slowly cool the sample to room temperature to avoid thermal stress. Step 6: Characterization 7. Characterize the synthesized YMnO ₃ using techniques such as X-ray diffraction (XRD), scanning electron microscopy (SEM), and energy-dispersive X-ray spectroscopy (EDS) to confirm the desired phase and composition. Notes: The choice of manganese precursor (MnO ₂ or Mn ₂ O ₃) and lithium source (carbonate or hydroxide) can impact the synthesis process and may require optimization. It's crucial to use high-purity starting materials to avoid impurities in the final product. The specific temperatures and times mentioned are general guidelines and may need to be adjuste		
	syn- ze O ₃ ?	syn- ze method can be employed. Here is a summarized procedure based on the synthesis details provided in the literature with DOI: 10.1021/ic502027k: 1. Dissolve Y(NO3)3·6H2O, Mn(NO3)2·4H2O, and a source of [OH-] in water to create a homogeneous solution. 2. Add citric acid to the solution to act as a mixing media. 3. Heat the mixture at 250°C for 2 hours to achieve pre-calcination. 4. Finally, heat the pre-calcined mixture at a temperature range of 700°C to 930°C for 10 to 18 hours in an air atmosphere to form		

Table 5 – continued from previous page

		Tuble 3 Continued	Answer
	Question	LLaMP	GPT-3.5
B5.10	How	LiFePO ₄ can be	Synthesizing LiFePO ₄ (lithium iron phosphate)
	to syn-	synthesized using various	typically involves a solid-state reaction method.
	thesize	methods. One method	Below is a step-by-step procedure for synthesiz-
	LiFePO ₄ ?	involves a	ing LiFePO ₄ :
	Please	high-temperature	Materials: - Lithium carbonate (Li2CO3) - Iron
	provide a	solid-state synthesis using	(II) acetate (Fe(CH ₃ COO) ₂) - Ammonium dihy-
	detailed	precursors such as Fe ₂ O ₃ ,	drogen phosphate (NH4H2PO4) - Phosphoric
	step-	LiOH, and NH ₄ H ₂ PO ₄ .	acid (H3PO4) - Deionized water - Ethanol -
	by-step	The precursors are mixed	Mortar and pestle - Crucible - Furnace
	procedure	and heated at	Procedure: 1. Weigh out the appropriate
	and the	temperatures up to 850°C	amounts of lithium carbonate (Li2CO3) and
	reference.	for several days. Another	iron (II) acetate (Fe(CH ₃ COO) ₂) in a 1:1 mo-
		method uses a sol-gel	lar ratio. For example, if you want to synthe-
		combustion approach with	size 1 mole of LiFePO ₄ , you would weigh out
		precursors like FePO ₄ and	74.55 grams of Li ₂ CO ₃ and 173.93 grams of
		Li ₂ CO ₃ , where the	Fe(CH ₃ COO) ₂ .
		mixture is calcined at	2. Grind the lithium carbonate and iron (II)
		700°C in a nitrogen	acetate together using a mortar and pestle to
		atmosphere. For detailed	ensure a homogenous mixture.
		procedures, please refer to	3. Weigh out the appropriate amount of ammo-
		the following scientific	nium dihydrogen phosphate (NH4H2PO4) and
		literature:	add it to the lithium carbonate and iron (II) ac-
		1. High-temperature solid-state synthesis -	etate mixture in a 1:1 molar ratio. For example,
		DOI: 10.1021/jp407510h	if you are synthesizing 1 mole of LiFePO ₄ , you would weigh out 122.00 grams of NH4H2PO4.
		2. Ion-exchange method -	4. Add a few drops of phosphoric acid (H3PO4)
		DOI: 10.1021/cm402617b	to the mixture to adjust the pH to around 7. This
		3. Solution method - DOI:	will help in the formation of LiFePO ₄ .
		10.1021/jp306936t	5. Add a small amount of deionized water to the
		4. Solid-state reactions -	mixture and mix it thoroughly to form a slurry.
		DOI: 10.1021/cm303259j	6. Slowly evaporate the water from the slurry
		5. Sol-gel combustion	by placing it on a hot plate or in a furnace at a
		method - DOI:	low temperature.
		10.1007/s10853-011-	7. Once the water has evaporated, heat the mix-
		6139-7	ture in a crucible in a furnace at around 600-
			700°C for several hours to initiate the solid-state
			reaction and form LiFePO ₄ .
			8. After the reaction is complete, allow the cru-
			cible to cool and collect the LiFePO ₄ powder.
			9. Wash the LiFePO ₄ powder with ethanol to
			remove any impurities and then dry it in a des-
			iccator.
			Reference: P. G. Bruce, B. Scrosati, JM. Taras-
			con, "Nanomaterials for Rechargeable Lithium
			Batteries," Angew. Chem. Int. Ed. 2008, 47,
			2930-2946. DOI: 10.1002/anie.200702505.
			Continued on next page

Table 5 – continued from previous page

Could you insert one Li atom at the tetrahedral interstitial site in the diamond cubic Si structure is as follows:
insert one Li atom at tetrahedral interstitial site in the diamond cubic Si structure is as follows: interstitial site in diamond cubic Si; C'emodule: 'pymatgen.core. structure', 'eclass': Structure', 'eleason': Structure',
hedral interstitial site in diamond cubic Si? Pull out cubic Si? Pull out the most stable DC Si structure from MP directly in JSON string, modify co. 895. 6.0.895. 8.75.

Table 5 – continued from previous page

	Question	Answer	
		LLaMP	GPT-3.5
B5.12	What	The stable FePO ₄ has a	The magnetic ordering of stable FePO ₄ is anti-
	is the	ferromagnetic (FM)	ferromagnetic. Its material id is mp-5405, and
	magnetic	magnetic ordering. Its	its space group is Pnma. The magnetization per
	ordering	material ID is	formula unit is 4.14 μB.
	of stable	'mp-753395', it belongs to	
	FePO ₄ ?	the Triclinic crystal	
	Give	system with space group	
	me its	symbol 'P-1' and number	
	material	2. The total magnetization	
	id, space	per formula unit is	
	group, and	5.00000015 Bohr	
	magneti-	magnetons.	
	zation per		
	formula		
	unit.		

C Other Capabilities

C.1 Automated Computational Workflow

The following example shows how LLaMP uses computational workflow atomate2 and job manager jobflow to run automated molecular dynamics simulation using foundation machine learning interatomic potential MACE-MP-0 (Batatia et al., 2023).

```
I would like you to help me run MD simulation with the follwing steps:
1. Retrieve the stable LiTaO3 structure from Materials Project
2. Create a supercell around the size of 100 atoms
3. Use atomate2 with foundation ML forcefield MACE to run Nose-Hoover NVT simulation
   for 2 ps with annealing temperature scheduling 300 to 800 back to 300 Kelvin.
  Use 2 fs timestep and save the trajectory to a ASE traj file.
4. Read the trajectory file and plot the system temperature vs time
> Entering new AgentExecutor chain...
Thought: Do I need to use a tool? Yes
Action: MPStructureRetriever
Action Input: LiTaO3
> Entering new AgentExecutor chain...
Action:
  `json
 "action": "search_materials_structure__get",
  "action_input": {
   "formula": "LiTaO3",
   "limit": 5,
   "fields": "material_id,structure"
}
```

```
All retrieved structures are saved as Pymatgen Structure JSON files to the following paths: mp -3666.json, mp-1105280.json, mp-754345.json, mp-1105216.json, mp-1105326.json

> Finished chain.
{'input': 'LiTaO3', 'output': 'All retrieved structures are saved as Pymatgen Structure JSON files to the following paths: mp-3666.json, mp-1105280.json, mp-754345.json, mp-1105216.json
```

```
, mp-1105326.json'}
Do I need to use a tool? Yes
Action: MPSummaryExpert
Action Input: mp-3666
> Entering new AgentExecutor chain...
Action:
```json
 "action": "search_materials_summary__get",
 "action_input": {
 "material_ids": "mp-3666"
}
```Error on search_materials_summary__get: `fields` must be specified in the query. Please
    revise arguments or try smaller request by specifying 'limit' in request.
 "action": "search_materials_summary__get",
  "action_input": {
   "material_ids": "mp-3666",
   "fields": "material_id, formula_pretty, composition, nsites, symmetry"
 }
}
```

```
[{'nsites': 10, 'composition': {'Li': 2.0, 'Ta': 2.0, '0': 6.0}, 'formula_pretty': 'LiTaO3', 'symmetry': {'crystal_system': 'Trigonal', 'symbol': 'R3c', 'number': 161, 'point_group': '3m', 'symprec': 0.1, 'version': '2.0.2'}, 'material_id': 'mp-3666'}]
Action:
```json
 "action": "Final Answer"
 "action_input": "Material ID mp-3666 corresponds to Lithium Tantalate (LiTaO3). It has a
 trigonal crystal system with space group R3c (number 161) and point group 3m. The
 structure consists of 10 sites, composed of Li, Ta, and O."
}
Finished chain.
{'input': 'mp-3666', 'output': 'Material ID mp-3666 corresponds to Lithium Tantalate (LiTaO3).
 It has a trigonal crystal system with space group R3c (number 161) and point group 3m. The
 structure consists of 10 sites, composed of Li, Ta, and O.'}
Do I need to use a tool? Yes
Action: MPStructureRetriever
Action Input: mp-3666
Entering new AgentExecutor chain...
Action:
```json
  "action": "search_materials_structure__get",
  "action_input": {
    "material_ids": "mp-3666",
    "limit": 1,
"fields": "material_id,structure"
  }
}
```

```
All retrieved structures are saved as Pymatgen Structure json files to the following paths: mp -3666.json

Do I need to use a tool? Yes Action: Python_REPL Action Input:
```

```
from pymatgen.core import Structure

# Load the structure from the JSON file
structure = Structure.from_file('mp-3666.json')

# Create a supercell to approximately reach 100 atoms
num_atoms = len(structure)
multiplier = int((100 / num_atoms) ** (1/3)) + 1 # Calculate the multiplier for each dimension
supercell = structure * (multiplier, multiplier, multiplier)

# Save the supercell structure to a file and print the number of atoms
supercell.to(fmt="poscar", filename="supercell.vasp")
print({"supercell_path": "supercell.vasp", "num_atoms": len(supercell)})
```

```
{'supercell_path': 'supercell.vasp', 'num_atoms': 270}
```

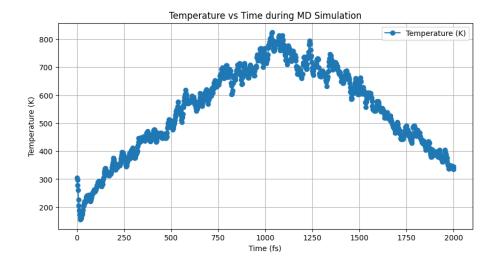
```
Do I need to use a tool? Yes
Action: MLFFMDAgent
Action Input: {"structure_path": "supercell.vasp", "forcefield": "MACE", "simulation_type": "NVT
       "temperature_schedule": [300, 800, 300], "simulation_time": 2, "timestep": 2,
     trajectory_file": "trajectory.traj"}Do I need to use a tool? Yes
> Entering new AgentExecutor chain...
Action:
 ``json
{
  "action": "MLFF MD",
  "action_input": {
    "run_mode": "local",
"project": "llamp-atomate2",
    "atom_path_or_dict": "supercell.vasp",
"force_field_name": "MLFF.MACE",
    "time_step": 2,
    "n_steps": 1000,
"ensemble": "NVT",
"dynamics": "langevin",
    "temperature": [300, 800, 300],
    "traj_file": "trajectory.traj",
     "traj_file_fmt": "ase",
    "traj_interval": 1
  }
}
...Action:
```

```
2024-05-12 14:42:31,797 INFO Started executing jobs locally
2024-05-12 14:42:31,837 INFO Starting job - Forcefield MD (040866d1-fbcb-4e0b-9044-0ed3cf9d0eef)
Using Materials Project MACE for \texttt{MACECalculator} with /global/homes/c/cyrusyc/.cache/
mace/5yyxdm76
Using \texttt{float32} for \texttt{MACECalculator}, which is faster but less accurate.
Recommended for MD. Use float64 for geometry optimization.

Default dtype float32 does not match model dtype float64, converting models to float32.
2024-05-12 14:43:58,683 INFO Finished job - Forcefield MD (040866d1-fbcb-4e0b-9044-0ed3cf9d0eef)
INFO:jobflow.core.job:Finished job - Forcefield MD (040866d1-fbcb-4e0b-9044-0ed3cf9d0eef)
2024-05-12 14:43:58,710 INFO Finished executing jobs locally
INFO:jobflow.managers.local:Finished executing jobs locally
{'energy': -2262.60595703125, 'n_steps': 1001}
```

```
Action:
   `json
  "action": "Final Answer",
  "action_input": "The molecular dynamics simulation was successfully completed. The final energy
         of the system after 1001 steps is -2262.60595703125 eV."
...Action:
···json
  "action": "Final Answer",
  "action_input": "The molecular dynamics simulation was successfully completed. The final energy
         of the system after 1001 steps is -2262.60595703125 eV."
}
> Finished chain.
{'input': '{"structure_path": "supercell.vasp", "forcefield": "MACE", "simulation_type": "NVT", "temperature_schedule": [300, 800, 300], "simulation_time": 2, "timestep": 2, "trajectory_file": "trajectory.traj"}', 'output': 'The molecular dynamics simulation was
     successfully completed. The final energy of the system after 1001 steps is -2262.60595703125
      eV.'}Do I need to use a tool? Yes
Action: Python_REPL
Action Input:
```

```
import matplotlib.pyplot as plt
from ase.io import read
# Load the trajectory
trajectory = read('trajectory.traj', index=':')
# Extract temperatures and time steps
temperatures = [step.get_temperature() for step in trajectory]
time_steps = [i * 2 for i in range(len(temperatures))] # Calculate time steps manually since
    each step is 2 fs
# Plotting
plt.figure(figsize=(10, 5))
plt.plot(time_steps, temperatures, label='Temperature_(K)', marker='o')
plt.title('Temperature_vs_Time_during_MD_Simulation')
plt.xlabel('Time_(fs)')
plt.ylabel('Temperature_(K)')
plt.grid(True)
plt.legend()
plt.savefig('temperature_vs_time.png')
plt.show()
print({"temperature_vs_time_plot": "temperature_vs_time.png"})
```



C.2 Simulation Code Generation and Execution

Below we present how LLaMP can generate, self-correct, and execute ASE simulation code to calculate the bulk modulus of aluminum.

Generate a script and run ASE EMT calculator to relax 2x2x2 Al supercell and calculate its bulk modulus

Entering new AgentExecutor chain...
To generate a script and run the ASE (Atomic Simulation Environment) EMT (Effective Medium Theory) calculator to relax a 2x2x2 Al supercell and calculate its bulk modulus, we need to write a Python script that uses ASE's functionalities. Here's a step-by-step breakdown of what the script will do:

1. Import necessary modules from ASE.
2. Create an aluminum (Al) supercell.
3. Set up the EMT calculator.
4. Relax the structure using the EMT calculator.
5. Calculate the bulk modulus.

Let's write the script:

```
from ase import Atoms
from ase.build import bulk
from ase.calculators.emt import EMT
from ase.constraints import StrainFilter, ExpCellFilter
from ase.optimize import BFGS
from ase.eos import EquationOfState
# Create an Al supercell
al = bulk('Al', 'fcc', a=4.05, cubic=True)
al = al.repeat((2, 2, 2)) # Repeat the unit cell to create a 2x2x2 supercell
# Set up the EMT calculator
al.set_calculator(EMT())
# Relax the structure
strain_filter = StrainFilter(al)
optimizer = BFGS(strain_filter)
optimizer.run(fmax=0.01)
# Calculate the bulk modulus
volumes = []
energies = []
```

```
for x in [0.98, 0.99, 1.00, 1.01, 1.02]:
    al_copy = al.copy()
    al_copy.set_calculator(EMT())
    cell = al_copy.get_cell()
    al_copy.set_cell(cell * x, scale_atoms=True)
    volumes.append(al_copy.get_volume())
    energies.append(al_copy.get_potential_energy())

eos = EquationOfState(volumes, energies)
v0, e0, B = eos.fit()
eos.plot('al-eos.png')

print(f'Bulk_modulus:_{B_/_kJ}_GPa')
```

```
Now, let's execute this script using the Python\_REPL tool.

Action: Python\_REPL
Action Input:
```

```
from ase import Atoms
from ase.build import bulk
from ase.calculators.emt import EMT
from ase.constraints import StrainFilter, ExpCellFilter
from ase.optimize import BFGS
from ase.eos import EquationOfState
from ase.units import kJ
# Create an Al supercell
al = bulk('Al', 'fcc', a=4.05, cubic=True)
al = al.repeat((2, 2, 2)) # Repeat the unit cell to create a 2x2x2 supercell
# Set up the EMT calculator
al.set_calculator(EMT())
# Relax the structure
strain_filter = StrainFilter(al)
optimizer = BFGS(strain_filter)
optimizer.run(fmax=0.01)
# Calculate the bulk modulus
volumes = []
energies = []
for x in [0.98, 0.99, 1.00, 1.01, 1.02]:
   al_copy = al.copy()
   al_copy.set_calculator(EMT())
   cell = al_copy.get_cell()
   al_copy.set_cell(cell * x, scale_atoms=True)
    volumes.append(al_copy.get_volume())
    energies.append(al_copy.get_potential_energy())
eos = EquationOfState(volumes, energies)
v0, e0, B = eos.fit()
print(f'Bulk_modulus:_{B_/_kJ}_GPa')
```

```
Energy
     Step Time
                                       fmax
BFGS: 0 01:37:31
                      -0.048066
                                     8.967521
BFGS:
      1 01:37:31
                       1.542253
                                    33.456077
                                     1.250352
      2 01:37:31
BFGS:
                      -0.154186
BFGS:
       3 01:37:31
                      -0.156157
                                     0.259261
       4 01:37:31
                      -0.156246
                                     0.000522
BFGS:
```

The output of the bulk modulus calculation seems to be incorrect because the value is extremely small and not in the expected range for aluminum. This is likely due to a unit conversion error. The ASE package uses eV and Angstroms for energy and length, respectively, and the bulk modulus should be in units of eV/Angstrom^3, which can be converted to GPa. The conversion factor from eV/Angstrom^3 to GPa is approximately 160.21766208. Let's correct the script to include the proper conversion factor and rerun the calculation.

Action: Python_REPL
Action Input:

```
from ase import Atoms
from ase.build import bulk
from ase.calculators.emt import EMT
from ase.constraints import StrainFilter, ExpCellFilter
from ase.optimize import BFGS
from ase.eos import EquationOfState
from ase.units import kJ
# Create an Al supercell
al = bulk('Al', 'fcc', a=4.05, cubic=True)
al = al.repeat((2, 2, 2)) # Repeat the unit cell to create a 2x2x2 supercell
# Set up the EMT calculator
al.set_calculator(EMT())
# Relax the structure
strain_filter = StrainFilter(al)
optimizer = BFGS(strain_filter)
optimizer.run(fmax=0.01)
# Calculate the bulk modulus
volumes = []
energies = []
for x in [0.98, 0.99, 1.00, 1.01, 1.02]:
   al_copy = al.copy()
   al_copy.set_calculator(EMT())
   cell = al_copy.get_cell()
   al_copy.set_cell(cell * x, scale_atoms=True)
   volumes.append(al_copy.get_volume())
   energies.append(al_copy.get_potential_energy())
eos = EquationOfState(volumes, energies)
v0, e0, B = eos.fit()
# Convert the bulk modulus from eV/Angstrom^3 to GPa
B GPa = B * 160.21766208
print(f'Bulk_modulus:_{B_GPa}_GPa')
```

```
Step
            Time
                        Energy
                                      fmax
BFGS: 0 01:37:48
                      -0.048066
                                     8.967521
BFGS: 1 01:37:48
                      1.542253
                                    33.456077
BFGS:
      2 01:37:48
                      -0.154186
                                     1.250352
BFGS:
       3 01:37:48
                      -0.156157
                                     0.259261
       4 01:37:48
                      -0.156246
                                     0.000522
Bulk modulus: 39.23492594026644 GPa
```

Do I need to use a tool? No
Final Answer: The corrected bulk modulus of the relaxed 2x2x2 Al supercell calculated using the
ASE EMT calculator is approximately 39.23 GPa.
Finished chain.

D Experiment Setup and Code Release

All the experiments were running on 1 Intel(R) Xeon(R) Gold 5120 CPU @ 2.20GHz of with average of 15 to 25 seconds inference time per prompt. Additionally, the total expenditure on API credits for running these experiments, including those that are not detailed in the paper, amounted to \$1,800, supported by OpenAI.

We make the code and data for our analysis available at https://github.com/chiang-yuan/llamp. We release both under the BSD 3-Clause License.

E Disclosure of AI Usage

The authors acknowledge the use of artificial intelligence (AI) tools in the preparation of this manuscript. Specifically, Microsoft Copilot, OpenAI ChatGPT, and Google Gemini Pro were utilized for general editing and code generation / completion purposes. All generated code and text was verified for correctness by one or more of the authors.