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

Yuan Chiang<sup>1,2\*†</sup> Elvis Hsieh<sup>1\*</sup> Chia-Hong Chou<sup>3</sup> Janosh Riebesell<sup>2,4</sup>

<sup>1</sup>University of California, Berkeley, <sup>2</sup>Lawrence Berkeley National Laboratory,

<sup>3</sup>Foothill College, <sup>4</sup>Cavendish Laboratory, University of Cambridge, UK

#### **Abstract**

Reducing hallucination of Large Language Models (LLMs) is imperative for use in the sciences, where reliability and reproducibility are crucial. However, LLMs inherently lack long-term memory, making it a nontrivial, ad hoc, and often biased task to fine-tune them on domain-specific literature and data. Here we introduce **LLaMP**, a multimodal retrieval-augmented generation (RAG) framework of hierarchical reasoning-and-acting (ReAct) agents that can dynamically and recursively interact with computational and experimental data from the Materials Project (MP) and run atomistic simulations via high-throughput workflow interface. Without fine-tuning, LLaMP demonstrates strong tool-usage ability to comprehend and integrate various modalities of materials science concepts, fetch relevant data stores on the fly, process higher-order data (such as crystal structure and elastic tensor), and streamline complex tasks in computational materials and chemistry. We propose a metric combining uncertainty and confidence estimates to evaluate the self-consistency of responses by LLaMP and vanilla LLMs. Our benchmark shows that LLaMP effectively mitigates the intrinsic bias in LLMs, counteracting the errors on bulk moduli, electronic bandgaps, and formation energies that seem to derive from mixed data sources. We also demonstrate LLaMP's capability to edit crystal structures and run annealing molecular dynamics simulations using pre-trained machine-learning interatomic potentials. The framework offers an intuitive and nearly hallucination-free approach to exploring and scaling materials informatics, and paves the way for future agentic scientific workflows and knowledge-grounded LLMs. 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 model (LLMs), particularly to their application in the sciences. LLMs are prone to hallucination—providing outright false information with high confidence [1, 39]. This issue is particularly concerning for knowledge-intensive tasks, where users rely on chatbots and other AI systems to provide accurate guidance [16]. LLMs often lack up-to-date factual knowledge on topics outside their training data, requiring rigorous verification against trusted external sources [22]. In the scientific community, where the integration of insights and data accuracy is already complex, the proliferation of generative models may exacerbate the risk of misinformation. This trend accentuates the importance of scrutinizing and ensuring the reliability of information sources.

Current approaches to enhance LLM accuracy in domain-specific knowledge often involve fine-tuning pre-trained models [6, 10] or tailored prompt engineering techniques [41, 45]. While these

<sup>\*</sup>Primary contributors

<sup>&</sup>lt;sup>†</sup>Correspondence to: cyrusyc@berkeley.edu

models are easy to deploy, they suffer from diminished reproducibility and data adherence due to the absence of a memory base, untraceable fine-tuning history, or opaque extraction processes. Even though fine-tuning can encode a certain amount of domain-specific knowledge into LLMs, it is constrained by scalability and intrinsic memory capacity. Fine-tuned LLMs struggle to retain in the long term the knowledge they were trained on as the training progresses, nor can they be aware of the up-to-date events and data beyond pre-training. Prompt engineering, while effective, also compromise the generalizability, thus limiting the overall power and flexibility of LLMs. 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 [17] that can reliably support real-world scientific research and decision-making.

In this work, we propose LLaMP, a multimodal retrieval-augmented generation (RAG) framework leveraging hierarchical reasoning-and-acting (ReAct) agents to interact with Materials Project (MP), arXiv, Wikipedia, and atomistic simulation tools. The framework serves as a safeguard against LLM hallucination and grounds them on high-fidelity material informatics derived from various sources, including computational data from quantum-mechanical first-principles calculations and expert-curated material synthesis recipes [14], and further enables the possibility of language-driven simulations. Through hierarchical planning of multiple ReAct agents [42], we demonstrate that LLaMP not only can correctly retrieve higher-order materials data such as tensors and 3D crystal structures but also can combine different modalities to perform complex, knowledge-intensive inferences and operations essential for real-world research applications.

Our contributions are as follows: (1) we introduce a multimodal RAG framework employing hierarchical ReAct agents that dynamically interact with the Materials Project, enabling LLMs to access high-fidelity materials informatics; (2) we propose a statistical metric to assess the self-consistency of LLM responses in high-precision, reproducibility-critical settings; (3) we evaluate the performance of LLaMP and standard LLMs in predicting key material properties, including bulk moduli, electronic bandgaps, formation energies, and magnetic orderings; (4) we showcase real-world applications in materials science, such as inorganic synthesis and crystal structure generation and editing; (5) we enhance LLaMP with high-throughput atomistic simulation workflows and pre-trained universal ML force fields, lowering the entry barriers to computational materials and chemistry.

#### 2 Background

Materials Project (MP) The Materials Project is a multi-institution effort to explore and compute the properties of all known inorganic materials [13] and molecules [31]. The initiative leverages high-throughput electronic structure calculations [15, 30] based on density functional theory (DFT), providing large-scale open-source database and analysis algorithms, with the ultimate 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) [43] 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, dielectric tensors, and more. MP provides these calculated material properties through API endpoints.

Natural language processing (NLP) in science NLP has found extensive application in extracting valuable information from scientific publications, with notable instances involving text-to-text or more recent image-to-text summarization techniques [10, 29, 34]. For summarizing crystal structures in textual form, Ganose and Jain [8] introduced the *robocrystallographer*, a toolkit designed for the analysis and generation of descriptions for crystalline materials. Their method condenses atomic structures into descriptive JSON representations that encompass coordination statistics, connectivity motifs, geometric features, and dimensionality. MP leverages robocrystallographer to generate human-level descriptions for 130K compounds which are accessible through MP website and API.

### 3 Related Work

**Prompting and fine-tuning in domain science** Prompt-based methods have been used as effective tools for automating data extraction process from the literature. Polak and Morgan [28] employ a prompt workflow to extract the cooling rates of metallic glasses and yield strengths of high entropy

alloys. Zheng et al. [45] implement a ChatGPT metal-organic framework (MOF) synthesis assistant through embedding and searching on preselected papers. StructChem [26] leverages step-by-step reasoning, and iteratively refines results to solve college-level chemistry questions. Yang et al. [41] use GPT-4 to extract experimentally measured bandgaps to train a graph neural network for accurate bandgap prediction from crystal structures. 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, domain-specific language models lean on the fine-tuning approach against pre-selected data and literature. For instance, ChemGPT [7] involves fine-tuning GPT-neo on self-referencing embedded strings (SELFIES) representations of small molecules. Jablonka et al. [12] 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. [6] fine-tuned GPT-3 on  $\sim\!500$  prompt-completion pairs to enhance LLM's capability to extract useful information on materials chemistry from text paragraphs. However, the fine-tuned models without augmentation inherently lack awareness of the up-to-date results and any data only available after their training. Moreover, fine-tuned LLMs still suffer from limited memory retention and are prone to forget during continual training [35].

**LLM function calling and tool usage** An emerging class of LLM applications, including this work, 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. Coscientist [3] combines tools such as search engines, Python, and document index for autonomous chemical research. ChemCrow [21] gathers multiple molecule and safety tools to enhance organic chemistry experiment and molecule design.

However, most prior works adopt *flat planning* strategy, where a single agent accesses all the available tools, resulting in a lack of self-correcting tool usage capabilities. This often leads to premature reasoning stop and summarization when the agent encounters tool usage errors. We mitigate this through *hierarchical planning* of multiple ReAct agents (see Section 4.1).

#### 4 Method

#### 4.1 Hierarchical orchestration

**Overviews** To manage heterogeneous data sources and diverse types of queries, we introduce hierarchical planning, featuring the supervisor ReAct agent overseeing assistant ReAct agents (Figure 1). This design offers three major advantages over flat planning commonly implemented in previous works [3, 21]. (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 reducing the cognitive load on any individual agent; (3) by offloading specific functions to specialized agents, we minimize the context window consumption and schema parsing.

**Supervisor agent** The supervisor agent acts as a router and decision-maker, handling abstract logic between user requests and assistant agents. Here, we adopt ReAct on GPT-4 [42] to augment the agent's action space  $\mathcal{A}$  with a language space  $\mathcal{L}$  to create an expanded action space of  $\hat{\mathcal{A}} = \mathcal{A} \cup \mathcal{L}$ . This expanded action space empowers the agent to take action  $\hat{a}_t \in \mathcal{L}$  in language space that facilitate the collaboration with assistant agents to retrieve domain specific information and achieve complex downstream tasks such as molecular dynamics simulations.

Assistant agent The efficient function calling in LLMs is often hindered by the need to process complex API schemas, which can consume a significant portion of the context window. To address this, we assign a specialized ReAct 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 use of ReAct agents enables them to refine their API calls based on feedback, significantly improving task completion rates through ReAct's iterative self-correcting mechanism.

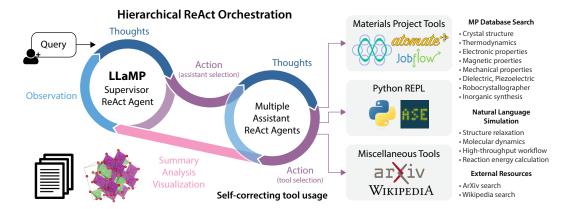


Figure 1: Hierarchical ReAct agent planning in LLaMP. Two levels of agents are deployed using a standardized LangChain interface. Supervisor ReAct agent oversees assistant ReAct agents at the bottom-level, each equipped with distinct toolkits and data/document stores to accomplish various tasks, including high-fidelity materials information retrieval, atomistic modeling and simulations, and literature search. For a detailed example, refer to Figure A.1.

The full list of agents are defined in 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 enable a seamless integration of new assistant agents, allowing for extensibility to various materials discovery methods and experimental techniques [20, 27, 36, 37, 44].

#### 4.2 Self-consistency of response (SCoR)

When LLMs are integrated in scientific workflows and deployed in high-stakes settings (*i.e.* self-driving labs), it is important for these models to have consistent and predictable behaviors [19]. For numeric knowledge retrieval tasks, we define the following metrics:

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

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

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

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

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.

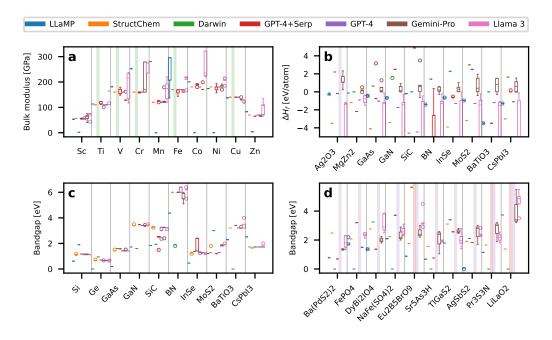


Figure 2: LLaMP RAG responses, baseline methods, and LLM intrinsic knowledge on material properties. (a) Bulk moduli, K, of 3d transition metals. (b) Formation energies,  $\Delta H_f$ , of common compounds. (c) Electronic bandgaps,  $E_g$ , of common intrinsic semiconductors. (d) Electronic bandgaps of multi-element (ternary or quaternary) materials. Missing predictions are marked by shaded areas. Fliers are marked in circles. All LLaMP results use GPT-4 as backend language provider.

#### 5 Experiments

#### 5.1 Multimodal ReAct Augmentation

Materials design often involves multi-objective property optimization. These properties span a Pareto front where optimizing one factor incurs deterioration in others. To succeed in such tasks, combining different modalities of materials properties is necessary. LLaMP achieves this through the hierarchical orchestration of multiple ReAct agents [42]. For the example question "What's the stiffest material with the lowest formation energy in Si-O system?" (Figure A.1), when a query requires multimodal information and compound logic, the supervisor agent decomposes the query into multiple subtasks, delegates them to assistant agents (MPThermoExpert and MPElasticityExpert) for information retrieval, and in the final stage of reasoning integrates information from both modalities, drawing on the context in episodic memory retrieved from the assistant agents (Figure 1). This enables LLaMP to achieve various tasks step-by-step by combining multiple data sources from the Materials Project (MP) (e.g. 3D crystal structures, thermodynamic, mechanical, magnetic properties, and more listed in and Appendix A.1) in a single query. In fact, we found that flat planning implementation struggles to accomplish tool usages and RAG because the single agent sees too much information at once and often fails to follow the API schema (see Section 4.1 for details).

#### 5.2 Performance Benchmarks

Response quality and consistency We evaluate the performance of LLaMP, StructChem [26], Darwin [38], and vanilla LLMs (gpt-4, llama3-8b, gemini-1.0-pro) on material properties such as bulk modulus, formation energy, and bandgap (Figure 2, Table 1). Performance is assessed through Precision, CoP, SCoR, and MAE metrics, as defined in Section 4.2. We argue that any useful LLM agents to be included in the scientific workflow should have high SCoR and low error on the materials properties. Notably, LLaMP consistently outperforms other models, achieving the highest SCoR and the lowest errors across material properties, making it highly suitable for scientific

Table 1: Performance metrics of LLaMP and LLM baselines on the prediction of material properties. 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), where Materials Project are taken as the ground truth. All the tabulated values are the average metrics over five runs and the sampled materials. All LLaMP and StructChem results use GPT-4 as backend language provider.

		Bulk	Modulus $K$ (0	GPa)			Formatio	on Energy $\Delta H$	f(eV)	
	Precision↓	CoP	Confidence	SCoR↑	MAE↓	Precision↓	CoP	Confidence	SCoR↑	MAE↓
LLaMP	2.698	0.900	1.000	0.900	14.574	0.007	0.993	0.960	0.953	0.009
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	2.222	0.352	1.000	0.352	41.742	5.947	0.745	1.000	0.745	8.214
GPT-4	0.186	0.910	1.000	0.910	41.225	0.000	1.000	0.180	0.180	1.680
Gemini-Pro	6.065	0.169	1.000	0.169	43.429	0.334	0.737	1.000	0.737	1.630
Llama 3	11.222	0.010	1.000	0.010	41.874	2.203	0.162	0.940	0.153	4.501
	Elec	tronic Ba	andgap $E_g$ - Co	ommon (e <sup>7</sup>	V)	Electronic Bandgap $E_g$ - Multi-element (eV)				
	Precision↓	CoP	Confidence	SCoR↑	MAE↓	Precision↓	CoP	Confidence	SCoR↑	MAE↓
LLaMP	0.000	1.000	1.000	1.000	0.000	0.013	0.988	0.950	0.938	0.182
StructChem	0.017	0.984	1.000	0.984	0.986	0.000	1.000	0.200	0.200	0.973
Darwin	0.002	0.998	1.000	0.998	1.224	0.000	1.000	1.000	1.000	1.951
GPT-4+Serp	0.040	0.963	1.000	0.963	1.012	0.000	1.000	0.660	0.660	0.576
GPT-4	0.027	0.975	1.000	0.975	0.978	-	-	0.000	0.000	-
Gemini-Pro	0.037	0.965	1.000	0.965	0.976	0.229	0.808	0.600	0.485	0.959
Llama 3	0.042	0.960	1.000	0.960	1.053	0.182	0.836	0.860	0.719	1.091

workflows. StructChem, despite extensive prompting strategies, often fails due to a lack of necessary domain knowledge, resulting in high refusal rates when it cannot validate outputs.

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 reference in the responses despite 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.

Our results demonstrate that vanilla LLMs fail to provide accurate formation energy predictions, with SCoR and 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).

In evaluating bandgaps, 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 (Figure 2c), with expected systematic deviation from MP values retrieved by LLaMP<sup>3</sup>. This is likely due to the extensive literature on experimental semi-conductor bandgaps, which have been studied and reported for decades. 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 (Figure 2d, Table B5.8), whereas LLaMP retrieves accurate data with a SCoR of 0.938 and correctly identifies the stable polymorph's bandgap when multiple forms are present.

**Ablation** Our frameworks relies on two principal components: first, factual material informatics on MP database; second, stable function calling mechanism that allows assistant agent to interact with tools. In Table 4, we examine three variants of LLaMP: (1) ReAct with MP tools; (2) 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 up-to-date, high-fidelity materials databases. In section Section 4.1, we mentioned the importance of hierarchical planning for robust function call. Evaluating several backbone models on bulk moduli and formation energy prediction, we found LLaMP's grounding performance correlates with the function-calling

<sup>&</sup>lt;sup>3</sup>Bandgaps calculated from generalized gradient approximation (GGA) functional are known to underestimate the experimental values by 40-50% [4]. Strategies to improve bandgap prediction at moderate or low computational cost will be included in MP in the future.

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 [40].

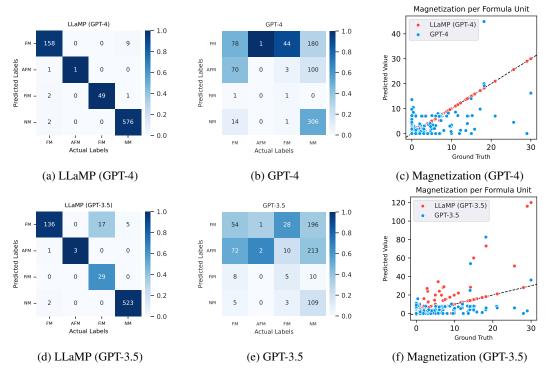


Figure 3: Prediction of LLaMP, GPT-3.5, and GPT-4 on (a,b,d,e) magnetic orderings and (c,f) total magnetization per formula unit of randomly selected materials. Confusion matrix presents the number of entries in each class. Colormap represents the percentage of correct classification.

**High-fidelity and higher-order data retrieval** The challenge for LLMs in excelling at knowledge- and data-intensive tasks is well-documented [5, 11, 19]. Figure 3 shows the prediction of LLaMP, GPT-3.5, and GPT-4 on the magnetic orderings and total magnetization of 800 materials randomly selected from all unary, binary, and ternary compounds in MP. Our result indicates that without RAG, vanilla LLMs suffer from hallucinations and misclassify the magnetic orderings of materials. LLaMP with GPT-4 as backend can counteract 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, while effective for classification and other information retrieval tasks, struggles to distinguish total\_magnetization from magnetization\_per\_formula\_unit in magnetism API schema and often requests the wrong field and 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 3a, d).

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 hightlights the limitation of

Table 2: Prediction 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	Magnetization		
	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	

intrsinic knowledge in LLMs to recall higherorder, more complex data for more comprehensive, holistic response.

#### 5.3 Real-world Applications

**Inorganic synthesis recipes** Empowered by the MP synthesis endpoint, 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 GPT-3.5 gives a seemingly correct and verbose synthesis procedure for YMnO<sub>3</sub> in Table B5.9, inferring possible reaction pathways from two common oxides as precursors (Y<sub>2</sub>O<sub>3</sub> and MnO<sub>2</sub>). However, it pulls irrelevant lithium compounds (Li<sub>2</sub>CO<sub>3</sub> and LiOH) into the recipe and overlooks the fact that metathesis reactions [18, 33] require less applied energy than high-temperature sintering, which relies on solid-state diffusion [23]. Vanilla GPT-3.5 also exhibits uncertainty about specific synthesis details, such as heating temperature, duration, cooling rate, *etc*.

Consider the example of LiFePO<sub>4</sub> presented in Table B5.10. Explicit instruction is provided: "*Please provide a detailed step-by-step procedure and reference*." While GPT-3.5 does offer both a procedure and reference as asked, and the reference is indeed associated with a real paper, the paper itself contains no information about the synthesis procedure of LiFePO<sub>4</sub>. The procedure listed in Table B5.10 is dissociated from the title and is hallucinated from the pre-training corpus.

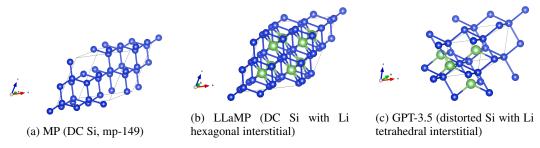


Figure 4: Generation and manipulation of crystal structures using LLMs to insert an additional lithium atom at the interstitial site in diamond cubic silicon structure. Blue: Si. Green: Li. Questionanswer pairs are listed in Table B5.11. Additional atoms extended through bonds are visualized.

Table 3: Structural parameters of the generated crystals compared with diamond cubic (DC) silicon. From left to right are fractional coordinates of inserted Li atom  $(x,y,z)_{\rm Li}$ , total cell volume V, average Si-Si bond lengths  $\ell_{\rm SiSi}$ , Si-Si-Si angles  $\theta_{\rm SiSiSi}$ , and Si-Li-Si angles  $\theta_{\rm SiLiSi}$ . GPT-4 refuses to respond due to their safeguard against the lack of atomic structure information.

	$(x,y,z)_{\mathrm{Li}}$	$\ell_{SiSi}  (\text{Å})$	Error (%)	$V(\text{\AA}^3)$	Error (%)	$\theta_{\rm SiSiSi} \ (^{\circ})$	Error (%)	$\theta_{\mathrm{SiLiSi}}$ (°)
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 generation and editing Fine-tuned LLMs for text-encoded atomistic information have shown the capability to generate stable crystals under the constraints of atomic positions and charges [9]. In this context, we delve into the examination and comparison of the crystal generation capabilities between LLaMP and GPT-3.5, without resorting to fine-tuning or tailored prompt messages in previous work. Figure 4 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 3), 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 pre-trained universal machine learning interatomic potential MACE-MP-0 [2] through language instruction. As demonstrated in Appendix C.1 and Appendix C.2, LLaMP is able to follow multi-step instruction to fetch stable crystal structure from MP, generate a supercell of atomic structure, and run annealing molecular dynamics simulation with varying temperature 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).

#### 6 Discussion

**Robustness** The hierarchical ReAct framework implemented here is essentially a graph of agents, or *language graph*, with one central node (supervisor) in connection with many satellite nodes (assistants). The implementation of ReAct for the assistant agents enables self-correcting tool usages and fortifies the robustness of data retrieval. As presented in Figure A.1c, MPThermoExpert initially misunderstood the schema at the first trial and filled in the formula field with Si-O, an invalid input but a valid one for chemical system (chemsys) field. The observation step (step 4) allows MPThermoExpert to handle exceptions and to refine the correct input fields after adaptation (step 6). Storing (Retrieving) question-answer and query-argument pairs to (from) vector databases could further reduce the number of trial-and-error steps, and the stored pairs can be used to refine foundation LLMs to improve function calling quality.

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. For example, sort\_fields argument allows sorting the returned documents in ascending order or descending order if the field is prefixed with —, but LLMs sometimes mistake the sign and sort in the opposite order. 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 magnetic configurations is also not complete. 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 available crystal structures on MP are not exhaustive but continuously expanding [24]. Furthermore, Kohn-Sham DFT theory is insufficient in some cases, and a higher level of theory is needed. Currently LLaMP only supports a few atomate2 workflows with machine learning force fields and VASP calculations. More diverse electronic calculation methods and workflows will be supported in the future work.

**Summary** We present a hierarchical agentic framework, LLaMP, based on ReAct to extract and manipulate material informatics through few-shot generalization. By grounding thoughts and actions with high-fidelity information, LLaMP showcases the ability to integrate various modalities of material properties and perform logical inferences to accomplish assigned tasks, all without the need for fine-tuning. In essence, the proposed LangChain framework holds the potential to expand its applicability to multiple data sources, encompassing both theoretical computations and experimental data, and real-world laboratories by incorporating additional assistant agents for data retrieval and robot control. LLaMP functions as a knowledge-aware agent, empowering users to navigate and manipulate complex materials informatics. In the context of self-driving labs [3, 32], LLM agents with multimodal data sources, sensors, and actors may improve their decision making and operation [25]. As new tools continue to emerge, there is an exciting avenue for further exploration to ascertain if this framework can effectively facilitate scientific hypothesis generation and guide data-driven experiments.

#### References

- [1] Y. Bang et al. A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity, Nov. 2023. URL http://arxiv.org/abs/2302.04023.arXiv:2302.04023 [cs].
- [2] I. Batatia et al. A foundation model for atomistic materials chemistry, Dec. 2023. URL http://arxiv.org/abs/2401.00096. arXiv:2401.00096 [cond-mat, physics:physics].
- [3] D. A. Boiko, R. MacKnight, B. Kline, and G. Gomes. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578, Dec. 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06792-0. URL https://www.nature.com/articles/s41586-023-06792-0. Number: 7992 Publisher: Nature Publishing Group.
- [4] P. Borlido et al. Exchange-correlation functionals for band gaps of solids: benchmark, reparametrization and machine learning. *npj Computational Materials*, 6(1):1–17, July 2020. ISSN 2057-3960. doi: 10.1038/s41524-020-00360-0. URL https://www.nature.com/articles/s41524-020-00360-0. Number: 1 Publisher: Nature Publishing Group.
- [5] K. Cobbe et al. Training Verifiers to Solve Math Word Problems, Nov. 2021. URL http://arxiv.org/abs/2110.14168. arXiv:2110.14168 [cs].
- [6] J. Dagdelen et al. Structured information extraction from scientific text with large language models. *Nature Communications*, 15(1):1418, Feb. 2024. ISSN 2041-1723. doi: 10.1038/s41467-024-45563-x. URL https://www.nature.com/articles/s41467-024-45563-x. Publisher: Nature Publishing Group.
- [7] N. Frey et al. Neural Scaling of Deep Chemical Models, May 2022. URL https://chemrxiv.org/engage/chemrxiv/article-details/627bddd544bdd532395fb4b5.
- [8] A. M. Ganose and A. Jain. Robocrystallographer: automated crystal structure text descriptions and analysis. *MRS Communications*, 9(3):874–881, Sept. 2019. ISSN 2159-6867. doi: 10. 1557/mrc.2019.94. URL https://doi.org/10.1557/mrc.2019.94.
- [9] N. Gruver et al. Fine-Tuned Language Models Generate Stable Inorganic Materials as Text. Nov. 2023. URL https://openreview.net/forum?id=0r5DE2ZSwJ.
- [10] T. Gupta, M. Zaki, N. M. A. Krishnan, and Mausam. MatSciBERT: A materials domain language model for text mining and information extraction. *npj Computational Materials*, 8(1):1–11, May 2022. ISSN 2057-3960. doi: 10.1038/s41524-022-00784-w. URL https://www.nature.com/articles/s41524-022-00784-w. Number: 1 Publisher: Nature Publishing Group.
- [11] D. Hendrycks et al. Measuring Massive Multitask Language Understanding, Jan. 2021. URL http://arxiv.org/abs/2009.03300. arXiv:2009.03300 [cs].
- [12] K. M. Jablonka, P. Schwaller, A. Ortega-Guerrero, and B. Smit. Leveraging large language models for predictive chemistry. *Nature Machine Intelligence*, 6(2):161–169, Feb. 2024. ISSN 2522-5839. doi: 10.1038/s42256-023-00788-1. URL https://www.nature.com/articles/s42256-023-00788-1. Publisher: Nature Publishing Group.
- [13] A. Jain et al. Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Materials*, 1(1):011002, July 2013. ISSN 2166-532X. doi: 10.1063/1.4812323. URL https://doi.org/10.1063/1.4812323.
- [14] O. Kononova et al. Text-mined dataset of inorganic materials synthesis recipes. *Scientific Data*, 6(1):203, Oct. 2019. ISSN 2052-4463. doi: 10.1038/s41597-019-0224-1. URL https://www.nature.com/articles/s41597-019-0224-1. Number: 1 Publisher: Nature Publishing Group.
- [15] G. Kresse and J. Furthmüller. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. *Physical Review B*, 54(16):11169–11186, Oct. 1996. doi: 10.1103/PhysRevB.54.11169. URL https://link.aps.org/doi/10.1103/PhysRevB.54. 11169. Publisher: American Physical Society.

- [16] P. Lewis et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html.
- [17] P. Lewis et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks, Apr. 2021. URL http://arxiv.org/abs/2005.11401. arXiv:2005.11401 [cs] version: 4.
- [18] J. Li et al. True Composition and Structure of Hexagonal "YAlO3", Actually Y3Al3O8CO3. *Inorganic Chemistry*, 54(3):837–844, Feb. 2015. ISSN 0020-1669. doi: 10.1021/ic502027k. URL https://doi.org/10.1021/ic502027k. Publisher: American Chemical Society.
- [19] P. Liang et al. Holistic Evaluation of Language Models, Oct. 2023. URL http://arxiv.org/abs/2211.09110. arXiv:2211.09110 [cs].
- [20] Y. Luo, C. Liu, and S. Ji. Towards Symmetry-Aware Generation of Periodic Materials, July 2023. URL https://arxiv.org/abs/2307.02707v2.
- [21] A. M. Bran et al. Augmenting large language models with chemistry tools. *Nature Machine Intelligence*, 6(5):525–535, May 2024. ISSN 2522-5839. doi: 10.1038/s42256-024-00832-8. URL https://www.nature.com/articles/s42256-024-00832-8. Publisher: Nature Publishing Group.
- [22] A. Mallen et al. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories, July 2023. URL http://arxiv.org/abs/2212.10511.arXiv:2212.10511 [cs].
- [23] A. L. Maximenko and E. A. Olevsky. Effective diffusion coefficients in solid-state sintering. Acta Materialia, 52(10):2953–2963, June 2004. ISSN 1359-6454. doi: 10.1016/j.actamat.2004.02.042. URL https://www.sciencedirect.com/science/article/pii/S1359645404001326.
- [24] A. Merchant et al. Scaling deep learning for materials discovery. *Nature*, pages 1–6, Nov. 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06735-9. URL https://www.nature.com/articles/s41586-023-06735-9. Publisher: Nature Publishing Group.
- [25] S. Miret and N. M. A. Krishnan. Are LLMs Ready for Real-World Materials Discovery?, Sept. 2024. URL http://arxiv.org/abs/2402.05200. arXiv:2402.05200 [cond-mat].
- [26] S. Ouyang et al. Structured Chemistry Reasoning with Large Language Models, Feb. 2024. URL http://arxiv.org/abs/2311.09656. arXiv:2311.09656 [cs].
- [27] G. Pilania, J. E. Gubernatis, and T. Lookman. Multi-fidelity machine learning models for accurate bandgap predictions of solids. *Computational Materials Science*, 129:156–163, Mar. 2017. ISSN 0927-0256. doi: 10.1016/j.commatsci.2016.12.004. URL https://www.sciencedirect.com/science/article/pii/S0927025616306188.
- [28] M. P. Polak and D. Morgan. Extracting Accurate Materials Data from Research Papers with Conversational Language Models and Prompt Engineering, June 2023. URL http://arxiv.org/abs/2303.05352. arXiv:2303.05352 [cond-mat].
- [29] A. Radford et al. Learning Transferable Visual Models From Natural Language Supervision, Feb. 2021. URL http://arxiv.org/abs/2103.00020. arXiv:2103.00020 [cs].
- [30] Y. Shao et al. Advances in molecular quantum chemistry contained in the Q-Chem 4 program package. *Molecular Physics*, 113(2):184–215, Jan. 2015. ISSN 0026-8976. doi: 10. 1080/00268976.2014.952696. URL https://www.tandfonline.com/doi/full/10.1080/00268976.2014.952696. Publisher: Taylor & Francis.
- [31] E. W. C. Spotte-Smith et al. A database of molecular properties integrated in the Materials Project, Oct. 2023. URL https://chemrxiv.org/engage/chemrxiv/article-details/651a3302ade1178b247ee6d1.

- [32] N. J. Szymanski et al. An autonomous laboratory for the accelerated synthesis of novel materials. *Nature*, 624(7990):86–91, Dec. 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06734-w. URL https://www.nature.com/articles/s41586-023-06734-w. Number: 7990 Publisher: Nature Publishing Group.
- [33] P. K. Todd et al. Selectivity in Yttrium Manganese Oxide Synthesis via Local Chemical Potentials in Hyperdimensional Phase Space. *Journal of the American Chemical Society*, 143(37):15185–15194, Sept. 2021. ISSN 0002-7863. doi: 10.1021/jacs.1c06229. URL https://doi.org/10.1021/jacs.1c06229. Publisher: American Chemical Society.
- [34] V. Tshitoyan et al. Unsupervised word embeddings capture latent knowledge from materials science literature. *Nature*, 571(7763):95–98, July 2019. ISSN 1476-4687. doi: 10.1038/s41586-019-1335-8. URL https://www.nature.com/articles/s41586-019-1335-8. Number: 7763 Publisher: Nature Publishing Group.
- [35] W. Wang et al. Augmenting Language Models with Long-Term Memory, June 2023. URL http://arxiv.org/abs/2306.07174. arXiv:2306.07174 [cs].
- [36] M. Wen et al. Chemical reaction networks and opportunities for machine learning. *Nature Computational Science*, 3(1):12–24, Jan. 2023. ISSN 2662-8457. doi: 10.1038/s43588-022-00369-z. URL https://www.nature.com/articles/s43588-022-00369-z. Publisher: Nature Publishing Group.
- [37] M. Wen et al. An equivariant graph neural network for the elasticity tensors of all seven crystal systems. *Digital Discovery*, 3(5):869–882, May 2024. ISSN 2635-098X. doi: 10. 1039/D3DD00233K. URL https://pubs.rsc.org/en/content/articlelanding/2024/dd/d3dd00233k. Publisher: RSC.
- [38] T. Xie et al. DARWIN Series: Domain Specific Large Language Models for Natural Science, Aug. 2023. URL http://arxiv.org/abs/2308.13565. arXiv:2308.13565 [cond-mat, physics:physics].
- [39] Z. Xu, S. Jain, and M. Kankanhalli. Hallucination is Inevitable: An Innate Limitation of Large Language Models, Jan. 2024. URL http://arxiv.org/abs/2401.11817. arXiv:2401.11817 [cs].
- [40] F. Yan et al. Berkeley Function Calling Leaderboard. 2024. URL https://gorilla.cs.berkeley.edu/blogs/8\_berkeley\_function\_calling\_leaderboard.html.
- [41] S. Yang et al. Accurate Prediction of Experimental Band Gaps from Large Language Model-Based Data Extraction. Nov. 2023. URL https://openreview.net/forum?id=oRKWhmtUG6.
- [42] S. Yao et al. ReAct: Synergizing Reasoning and Acting in Language Models, Mar. 2023. URL http://arxiv.org/abs/2210.03629. arXiv:2210.03629 [cs].
- [43] D. Zagorac et al. Recent developments in the Inorganic Crystal Structure Database: theoretical crystal structure data and related features. *Journal of Applied Crystallography*, 52 (5):918–925, Oct. 2019. ISSN 1600-5767. doi: 10.1107/S160057671900997X. URL https://journals.iucr.org/j/issues/2019/05/00/in5024/. Publisher: International Union of Crystallography.
- [44] C. Zeni et al. MatterGen: a generative model for inorganic materials design, Jan. 2024. URL http://arxiv.org/abs/2312.03687. arXiv:2312.03687 [cond-mat].
- [45] Z. Zheng et al. ChatGPT Chemistry Assistant for Text Mining and Prediction of MOF Synthesis. *Journal of the American Chemical Society*, 145(32):18048–18062, Aug. 2023. ISSN 0002-7863, 1520-5126. doi: 10.1021/jacs.3c05819. URL http://arxiv.org/abs/2306.11296. arXiv:2306.11296 [cond-mat, physics:physics].

#### **A Supplementary Information**

#### A.1 List of Implemented Assistant Agents and Tools

Here we provide the comprehensive list of implemented **assistant agents** and tools. Note that **MP Assistants** are highly modular so it is very trivial to support extra 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.

Table 4: Performance of LLaMP with different backbone LLMs and ablation of ReAct agent with MP API and SerpAPI tools.

		Bulk Moduli $K$ (GPa)				Formation Energies $\Delta H_f$ (eV)				
	Precision↓	CoP	Confidence	SCoR↑	MAE↓	Precision↓	CoP	Confidence	SCoR↑	MAE↓
LLaMP	2.698	0.900	1.000	0.900	14.574	0.007	0.993	0.960	0.953	0.009
LLaMP (Sonnet)	1.345	0.702	0.867	0.608	16.121	0.000	1.000	1.000	1.000	0.000
LLaMP (Gemini-1.5)	21.586	0.217	1.000	0.217	64.324	0.390	0.718	0.560	0.402	0.531
LLaMP (Llama 3)	7.314	0.327	0.800	0.261	47.386	0.413	0.721	1.000	0.721	3.062
GPT-4+Serp GPT-4	2.222 0.186	0.352 0.910	1.000 1.000	0.352 0.910	41.742 41.225	5.947 0.000	0.745 1.000	1.000 0.180	0.745 0.180	8.214 1.680

#### a. User Query c. Specialized ReAct (Assistant) What's the stiffest material with the lowest formation energy in Si-O system? **MPThermoExpert** b. LLaMP ReAct (Supervisor) "action": "search\_materiasls\_thermo\_\_get", "action\_input": { "formula": "Si-0", 1-Thought: To answer this question, I need to find materials in the Si-O sort\_fields": "formation\_energy\_per\_atom", system with the lowest formation energy and the highest stiffness. I will "limit": 5, "fields": "material\_id,formula\_pretty,formation\_energy\_per\_atom" use the MPThermoExpert tool to search for materials in the Si-O system and sort them by formation energy. After finding candidates, I will need to use the MPElasticityExpert tool to determine the stiffness of these 4-Observation: Error on search\_materials\_thermo\_\_get: REST query returned with error status code 400 on URL `json https://api.materialsproject.org/materials/thermo/?\_fields=material\_id... Problem processing one or more provided formulas. "action": "MPThermoExpert", 5-Thouaht: action\_input": { The query may be too broad or ambiguous. I will refine the search by specifying the chemical system "input": "What are the materials with the lowest formation energy in the Si-O system? 6-Action: } ison "action": "search\_materials\_thermo\_\_get", "action\_input": { "chemsys": "Si-O", 9-Thought: < "sort\_fields": "formation\_energy\_per\_atom", "limit": 5, The previous work has identified materials in the Si-O system with the lowest formation energy. The next step is to determine the stiffness of these materials. I will use the MPElasticityExpert tool to find the stiffness (typically represented by "fields": "material\_id,formula\_pretty,formation\_energy\_per\_atom" the bulk modulus or Young's modulus) for each of these materials 10-Action: 7-Observation: [('formula\_pretty': 'Si205', 'material\_id': 'mp-862998', 'formation\_energy\_per\_atom': -2.289431383571429}, {'formula\_pretty': 'Si02', 'material\_id': 'mp-733790', 'formation\_energy\_per\_atom': -3.039244523055555}, ...] `json action": "MPElasticityExpert", "action\_input": { "input": "What are the bulk and Young's moduli for The materials with the lowest formation energy in the Si-O system are: Si2O5 (mp-862998) with a Si2O5 (mp-862998), SiO2 (mp-733790), SiO2 (mp-6922), formation energy per atom of -2.289 eV/atom, SiO2 (mp-733790) with a formation energy per atom of -3.039 eV/atom, SiO2 (mp-6922) with a formation energy per atom of -3.037 eV/atom, SiO2 (mp-556985) $SiO2 \ (mp-556985), \ and \ SiO2 \ (mp-556994)?"$ with a formation energy per atom of -2.993 eV/atom, and SiO2 (mp-556994) with a formation energy per } atom of -2.958 eV/atom d. Specialized ReAct (Assistant) 14-Thought: The output from the MPElasticityExpert tool provides bulk modulus values for MPElasticityExpert the Si-O materials with the lowest formation energy. However, it seems there was ison a mismatch between the material IDs requested and the material IDs returned in the output. The Young's modulus values are not provided, but the bulk modulus "action": "search\_materials\_elasticity\_\_get", can be a good indicator of stiffness. The material with the highest bulk modulus "action\_input": { "formula": "Si205,Si02,Si02,Si02,Si02,Si02", would be considered the stiffest. "fields": "material\_id,formula\_pretty,bulk\_modulus,young\_modulus" From the provided output, the SiO2 with material ID mp-9258 has the highest bulk modulus of 310.27 GPa, which suggests it is the stiffest among the listed materials. Since this material also belongs to the Si-O system and has a high bulk modulus, it is likely to be one of the stiffest materials with low formation energy 12-Observation: [{'formula\_pretty': 'SiO2', 'material\_id': 'mp-1021503', 'elastic\_tensor': { 'raw': [[355.5072134787483, 236.22725456386155, 236.22725456386155, 2.081668171172168e-16, 7.771561172376094e-16, 7.771561172376094e-16], ...] in the Si-O system 13-Answer: The stiffest material in the Si-O system with low formation energy is likely SiO2 The bulk and Young's moduli for the requested materials are as follows: Si2O5 (mp-1021503) - Bulk ial ID mp-9258, which has a bulk modulus of 310.27 ailable; SiO2 (mp-11684) - Bulk modulus: 110.695 GPa, modulus values were not available, but the bulk modulus is a strong indicator of Young's modulus: Not available; SiO2 (mp-2374583) - Bulk modulus: 89.352 GPa, ... stiffness.

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 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.

	Question		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

			Answer
	Question	LLaMP	GPT-3.5
B5.2	What is the full elastic tensor of NaCl?	The elastic tensor for NaCl is available for two different entries on Materials Project. For material id 'mp-22851', the elastic tensor in IEEE format is [[76.0, 1.0, 1.0, -0.0, 0.0, -0.0], [1.0, 76.0, 1.0, -0.0, -0.0], [-0.0, -0.0, 0.0, -3.0, -0.0, -0.0], [-0.0, -0.0, -0.0, -0.0, -0.0, -0.0, -0.0], [-0.0, -0.0, -0.0, -0.0, -0.0, -0.0], [-0.0, 0.0, -0.0, -0.0], [-0.0, -0.0, -0.0], [12.0, 47.0, 12.0, 0.0, -0.0, -0.0], [12.0, 47.0, 12.0, 0.0, -0.0, -0.0], [0.0, 0.0, 0.0, 12.0, -0.0, -0.0], [-0.0, 0.0, -0.0, 12.0, -0.0, -0.0], [-0.0, -0.0, -0.0, 12.0, -0.0, -0.0], [-0.0, -0.0, -0.0, 12.0, -0.0, -0.0], [-0.0, -0.0, -0.0, 12.0], [-0.0, -0.0, -0.0, -0.0], [-0.0, -0.0, -0.0, 12.0]]	The elastic tensor describes the relationship between stress and strain in a material. It is a 4th-rank tensor that characterizes the elastic properties of a crystal. In the case of NaCl (sodium chloride), the elastic tensor is typically represented by a 6x6 matrix.  The elastic constants of NaCl can vary slightly depending on the specific conditions and crystallographic orientation. Here is an approximate form of the elastic tensor for NaCl:  C11 = 289.2 GPa C12 = 97.8 GPa C44 = 81.6 GPa C13 = C14 = C15 = C16 = C23 = C24 = C25 = C26 = 0 GPa  Note: The values are given in gigapascals (GPa). The elastic constants Cij represent the components of the elastic tensor. In this case, C11, C12, and C44 are the non-zero elastic constants. The off-diagonal terms C13, C14, C15, C16, C23, C24, C25, and C26 are typically zero due to cubic symmetry.  Keep in mind that these values are approximate, and there may be slight variations depending on the source and the specific conditions under which the measurements were made. Experimental techniques, such as Brillouin scattering or ultrasonic measurements, are often used to determine the elastic constants of materials.
B5.3	What are the bandgaps of the following materials: Si, Ge, and GaAs, GaN, SiC, BN, InSe, MoS2, BaTiO3, CsPbI3? Only look for ground-state bandgaps.	The ground state bandgaps for the requested materials are as follows: Si: 1.12 eV, Ge: 0.66 eV, GaAs: 1.42 eV, GaN: 3.28 eV, SiC: 2.36 eV, BN: 6.4 eV, InSe: 1.26 eV, MoS2: 1.29 eV, BaTiO3: 3.2 eV, CsPbI3: 1.73 eV. Please note that these values are typical for bulk materials and can vary depending on the specific structure and calculation method.	As of my last knowledge update in January 2022, I can provide you with approximate values for the ground state bandgaps of the materials you mentioned. Keep in mind that these values can vary depending on factors like temperature and crystal structure. Additionally, more recent research may have refined these values. Here are the approximate ground state bandgaps in electron volts (eV): Silicon (Si): 1.1 eV Germanium (Ge): 0.67 eV Gallium Arsenide (GaAs): 1.42 eV Gallium Nitride (GaN): 3.4 eV Silicon Carbide (SiC): 3.2 eV (for 4H-SiC) Boron Nitride (BN): 6 eV Indium Selenide (InSe): 1.25 eV Molybdenum Disulfide (MoS2): 1.3 eV (indirect bandgap) Barium Titanate (BaTiO3): 3.2 eV Cesium Lead Iodide (CsPbI3): 1.7 eV Please note that these values are general approximations and may not represent the specific conditions or variations in different forms of these materials. For the most accurate and up-to-date information, I recommend checking recent scientific literature or databases.

18

Table 5 – continued from previous page

		140.10 0 00.000	Answer
	Question	LLaMP	GPT-3.5
B5.4	Give me the decomposition enthalpies of the following materials: Ag2O3, MgZn2, GaAs, GaN, SiC, BN, InSe, MoS2, BaTiO3, CsPbI3. If there are multiple values, please give me the largest value.	The maximum decomposition enthalpies for the requested materials are as follows: - Ag2O3: 0.2565836414999989 eV/atom - MgZn2: 0.1356 eV/atom - GaAs: 0.075 eV/atom - GaN: 0.07162 eV/atom - SiC: -0.172 eV/atom - BN: -0.358 eV/atom - InSe: -0.01816 eV/atom - MoS2: 0.267	The decomposition enthalpies (in kJ/mol) of the following materials are: Ag2O3: -335 MgZn2: -82 GaAs: -317 GaN: -505 SiC: -40 BN: -248 InSe: -125 MoS2: -393
		eV/atom - BaTiO3: 3.055 eV/atom - CsPbI3: 0.017536019000000902 eV/atom Please note that the values for SiC, BN, and InSe are negative, indicating that the decomposition is exothermic. All other values are positive, indicating endothermic decompositions.	BaTiO3: -112 CsPbI3: -153
B5.5	Give me the formation energies of the following materials: Ag2O3, MgZn2, GaAs, GaN, SiC, BN, InSe, MoS2, BaTiO3, CsPbI3. If there are multiple values, please give me the smallest value.	The smallest formation energies of the materials you asked for are as follows: Ag2O3: -0.3601 eV/atom, MgZn2: -0.1306 eV/atom, GaAs: -0.446 eV/atom, GaN: -0.052 eV/atom, SiC: -0.204 eV/atom, BN: -1.41 eV/atom, InSe: -0.742582154921875 eV/atom, MoS2: -1.202 eV/atom, BaTiO3: -3.493 eV/atom, CsPbI3: -1.296 eV/atom.	Ag2O3: -2.63 eV/atom MgZn2: -0.27 eV/atom GaAs: -4.13 eV/atom GaN: -3.51 eV/atom SiC: -4.98 eV/atom BN: -3.8 eV/atom InSe: -0.89 eV/atom MoS2: -3.31 eV/atom BaTiO3: -4.45 eV/atom CsPbI3: -2.32 eV/atom
			Continued on next page

19

Table 5 – continued from previous page

	Table 3 – continued from previous page						
Question		Answer					
Question	LLaMP	GPT-3.5					
B5.6 Give me the decomposition enthalpies of the following materials: Ag2O3, MgZn2, GaAs, GaN, SiC, BN, InSe, MoS2, BaTiO3, CsPbI3. If there are multiple values, please give me the smallest value.	The smallest decomposition enthalpy values for the requested materials are as follows: Ag2O3: 0.019373648083333528 eV/atom, MgZn2: -0.0197 eV/atom, GaAs: -0.446 eV/atom, GaN: -0.05231 eV/atom, SiC: -0.204 eV/atom, BN: -1.41 eV/atom, InSe: -0.06336 eV/atom, MoS2: -0.669 eV/atom, BaTiO3: -0.063 eV/atom, CsPbI3: -0.018602 eV/atom.	The decomposition enthalpies of the materials are as follows: - Ag2O3: -1138.5 kJ/mol - MgZn2: -54.2 kJ/mol - GaAs: -452.2 kJ/mol - GaN: -510.9 kJ/mol - SiC: -2133.6 kJ/mol - BN: -2551.0 kJ/mol - InSe: -209.3 kJ/mol - MoS2: -739.3 kJ/mol - BaTiO3: -393.5 kJ/mol - CsPbI3: -65.3 kJ/mol					
B5.7 Give me the energy above hull values of the following materials: Ba2MnSe3, KAISiO4, LiCuO, H8C3N2O, Dy2ZnS4, Rb2BeF4, K9Fe2S7, Tb2FeC4, ScAlAu, CaCrO4. If there are multiple values, please give me the lowest one.	eV/atom H8C3N2O: 0.0735 eV/atom Dy2ZnS4: 0.00966 eV/atom Rb2BeF4: 0.0 eV/atom K9Fe2S7: 0.0 eV/atom Tb2FeC4: 0.0 eV/atom ScAlAu: 0.0 eV/atom CaCrO4: 0.0 eV/atom. Please note that these are the lowest values found in the Materials Project database and actual values may	The energy above hull values for the given materials are as follows: - Ba2MnSe3: 0.0000 eV/atom - KAlSiO4: 0.0000 eV/atom - LiCuO: 0.0000 eV/atom - H8C3N2O: 0.0000 eV/atom - Dy2ZnS4: 0.0000 eV/atom - Rb2BeF4: 0.0000 eV/atom - K9Fe2S7: 0.0000 eV/atom - Tb2FeC4: 0.0000 eV/atom - ScAlAu: 0.0000 eV/atom - CaCrO4: 0.0000 eV/atom					
	vary based on different entries.						

Table 5 – continued from previous page

		Tuble 5 continu	ued from previous page
	Question		Answer
	Question	LLaMP	GPT-3.5
B5.8	What are the bandgaps of the following materials: Ba(PdS2)2, FePO4, DyBi2IO4, NaFe(SO4)2, Eu2B5BrO9, Sr5As3H, Tl-GaS2, AgSbS2, Pr3S3N, LiLaO2. If there are polymorphs, give me only the bandgap of the stable material.	The bandgaps for the materials are: - Ba(PdS2)2: 0.779 eV - Most stable FePO4: 2.1932 eV - DyBi2IO4: 1.367 eV - Most stable NaFe(SO4)2: 2.094 eV - Eu2B5BrO9: 0.88 eV - Sr5As3H: 0.676 eV - TIGaS2: 1.7885 eV - Most stable AgSbS2: 1.372 eV - Pr3S3N: 1.15 eV - Most stable LiLaO2: 3.7336 eV Please note that these values are obtained from the Materials Project database and might be different in experimental results due to various factors.	I don't have real-time data access, and my training only includes information available up to January 2022. As a result, I may not have specific bandgap values for the materials you mentioned. However, I can provide you with general information about some of these materials and their known properties up to my last update:  Ba(PdS2)2: I don't have specific bandgap information for this compound. The bandgap may depend on factors such as crystal structure and synthesis conditions. FePO4: Iron(II) phosphate, known as ferrous phosphate, typically has a bandgap in the range of 2.5-3.0 eV.  DyBi2IO4: Dysprosium bismuth iodate. I don't have specific bandgap information for this compound.  NaFe(SO4)2: Sodium iron sulfate. I don't have specific bandgap information for this compound.  Sr5As3H: Strontium arsenide hydride. I don't have specific bandgap information for this compound.  TIGaS2: Thallium gallium sulfide. I don't have specific bandgap information for this compound.  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 compound.  LiLaO2: Lithium lanthanum oxide. I don't have specific bandgap information for this compound.  The most stable structure for each material can vary based on synthesis conditions and other factors. For accurate and up-to-date bandgap values and structural information, you may want to consult recent scientific literature or databases that specialize in material properties.

Table 5 – continued from previous page

			Answer
	Question	LLaMP	GPT-3.5
B5.9	How to synthesize YMnO <sub>3</sub> ?	To synthesize YMnO <sub>3</sub> , 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 <sub>3</sub> .	Synthesizing YMnO <sub>3</sub> (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 <sub>2</sub> O <sub>3</sub> ) Manganese oxide (MnO <sub>2</sub> or Mn <sub>2</sub> O <sub>3</sub> ) Highpurity manganese carbonate (MnCO <sub>3</sub> ) Lithium carbonate (Li <sub>2</sub> CO <sub>3</sub> ) 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 <sub>2</sub> O <sub>3</sub> ), manganese oxide (MnO <sub>2</sub> or Mn <sub>2</sub> O <sub>3</sub> ), and lithium carbonate (Li2CO3) or dithium hydroxide (LiOH) according to the stoichiometric ratio of YMnO3.  Step 2: Mixing the Powders  2. 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 <sub>3</sub> .  Step 5: Cooling 6. Slowly cool the sample to room temperature to avoid thermal stress.  Step 6: Characterization 7. Characterize the synthesized YMnO <sub>3</sub> 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 <sub>2</sub> or Mn <sub>2</sub> O <sub>3</sub> ) 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 adjusted based

Table 5 – continued from previous page

0 "		Answer
Question	LLaMP	GPT-3.5
B5.10 How to synthesize LiFePO <sub>4</sub> ? Please provide a detailed step-by-step procedure and the reference.	LiFePO <sub>4</sub> can be synthesized using various methods. One method involves a high-temperature	Synthesizing LiFePO <sub>4</sub> (lithium iron phosphate) typically involves a solid-state reaction method. Below is a step-by-step procedure for synthesizing LiFePO <sub>4</sub> :  Materials: - Lithium carbonate (Li2CO3) - Iron (II) acetate (Fe(CH <sub>3</sub> COO) <sub>2</sub> ) - Ammonium dihydrogen phosphate (NH4H2PO4) - Phosphoric acid (H3PO4) - Deionized water - Ethanol - Mortar and pestle - Crucible - Furnace  Procedure: 1. Weigh out the appropriate amounts of lithium carbonate (Li2CO3) and iron (II) acetate (Fe(CH <sub>3</sub> COO) <sub>2</sub> ) in a 1:1 molar ratio. For example, if you want to synthesize 1 mole of LiFePO <sub>4</sub> , you would weigh out 74.55 grams of Li <sub>2</sub> CO <sub>3</sub> and 173.93 grams of Fe(CH <sub>3</sub> COO) <sub>2</sub> .  2. Grind the lithium carbonate and iron (II) acetate together using a mortar and pestle to ensure a homogenous mixture.  3. Weigh out the appropriate amount of ammonium dihydrogen phosphate (NH4H2PO4) and add it to the lithium carbonate and iron (II) acetate mixture in a 1:1 molar ratio. For example, if you are synthesizing 1 mole of LiFePO <sub>4</sub> , you would weigh out 122.00 grams of NH4H2PO4.  4. Add a few drops of phosphoric acid (H3PO4) to the mixture to adjust the pH to around 7. This will help in the formation of LiFePO <sub>4</sub> .  5. Add a small amount of deionized water to the mixture and mix it thoroughly to form a slurry.  6. Slowly evaporate the water from the slurry by placing it on a hot plate or in a furnace at a low temperature.  7. Once the water has evaporated, heat the mixture in a crucible in a furnace at around 600-700°C for several hours to initiate the solid-state reaction and form LiFePO <sub>4</sub> .  8. After the reaction is complete, allow the crucible to cool and collect the LiFePO <sub>4</sub> powder.  9. Wash the LiFePO <sub>4</sub> powder with ethanol to remove any impurities and then dry it in a desiccator.  Reference: P. G. Bruce, B. Scrosati, JM. Tarascon, "Nanomaterials for Rechargeable Lithium Batteries," Angew. Chem. Int. Ed. 2008, 47, 2930-2946. DOI: 10.1002/anie.200702505.

Table 5 – continued from previous page

			Answer		
	Question	LLaMP	GPT-3.5		
B5.11	Could you insert one Li atom at the tetrahedral interstitial site in diamond cubic Si? Pull out the most stable DC Si structure from MP directly in JSON string, modify the retrieved structure directly, and give me the final pymatgen structure as a JSON string.	The final structure after inserting a Li atom at the tetrahedral interstitial site in the diamond cubic Si structure is as follows:  {'@module': 'pymatgen.core.structure', '     @class': 'Structure', 'charge': 0,     'lattice': ('matrix': [[3.333573,	Here's the JSON string for the most stable diamond cubic Si structure retrieved directly from MP:  "'json ('sites': [{'species': [{'element': 'Si', 'occu': 1]},		
			the state of the s		

Table 5 – continued from previous page

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

## 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 [2].

```
Query:
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
'''json
```

```
"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, 'O': 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
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:
```

```
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)})
```

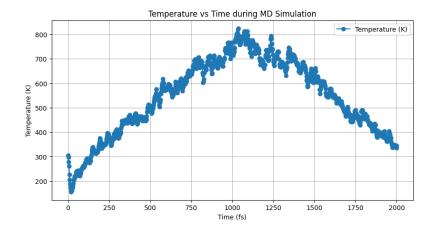
```
{'supercell_path': 'supercell.vasp', 'num_atoms': 270}
```

```
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-9044-0
    ed3cf9d0eef)
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-9044-0
    ed3cf9d0eef)
INFO: jobflow.core.job:Finished job - Forcefield MD (040866d1-fbcb-4e0b-9044-0
    ed3cf9d0eef)
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')
```

```
fmax
     Step Time
                       Energy
BFGS: 0 01:37:31
                      -0.048066
                                    8.967521
                      1.542253
       1 01:37:31
                                    33.456077
BFGS:
BFGS:
       2 01:37:31
                      -0.154186
                                    1.250352
BFGS:
       3 01:37:31
                      -0.156157
                                    0.259261
                                    0.000522
BFGS:
     4 01:37:31
                      -0.156246
Bulk modulus: 3.923492594026645e-23 GPa
```

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')
```

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