







Hierarchical Multi-agent Large Language Model Reasoning for Autonomous Functional Materials Discovery

Samuel Rothfarb^{†,‡} , Megan C. Davis[‡] , Ivana Matanovic[‡] ,
Baikun Li^{†*} , Edward F. Holby^{‡*} , and Wilton J.M. Kort-Kamp^{‡*} 

[†]*School of Civil & Environmental Engineering, University of Connecticut, Storrs, Connecticut 06269, United States.*

[‡]*Theoretical Division, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, United States.*

*Corresponding authors: baikun.li@uconn.edu, holby@lanl.gov, kortkamp@lanl.gov

Abstract

Artificial intelligence is reshaping scientific exploration, but most methods automate procedural tasks without engaging in scientific reasoning, limiting autonomy in discovery. We introduce Materials Agents for Simulation and Theory in Electronic-structure Reasoning (MASTER), an active learning framework where large language models autonomously design, execute, and interpret atomistic simulations. In MASTER, a multimodal system translates natural language into density functional theory workflows, while higher-level reasoning agents guide discovery through a hierarchy of strategies, including a single agent baseline and three multi-agent approaches: peer review, triage-ranking, and triage-forms. Across two chemical applications, CO adsorption on Cu-surface transition metal (M) adatoms and on M–N–C catalysts, reasoning-driven exploration reduces required atomistic simulations by up to 90% relative to trial-and-error selection. Reasoning trajectories reveal chemically grounded decisions that cannot be explained by stochastic sampling or semantic bias. Altogether, multi-agent collaboration accelerates materials discovery and marks a new paradigm for autonomous scientific exploration.

Keywords: *Multi-agent reasoning; Large language models; Active learning; AI-driven simulation; Materials discovery; Density functional theory; Surface chemistry.*

Introduction

Recent advances in artificial intelligence (AI) have expanded its role in scientific research, enabling systems to analyze data, identify patterns, and even propose hypotheses.¹ Yet, most of these models operate within fixed objectives and have limited ability to deliberate about scientific questions or adapt their strategies based on outcomes. Reasoning, the process of evaluating competing hypotheses and designing informative experiments, is central to genuine scientific exploration. In materials research, where the search space spans millions of possible atomic configurations,^{2–4} such reasoning ability could enable targeted efforts toward the most informative regions of chemical space, transforming the rate and nature of discovery. Despite this potential, reasoning-driven autonomy remains challenging in materials science. Materials design is resource-intensive, requiring specialized expertise and substantial computational effort.⁵ Even with advances in automation, the process from concept to verified material can extend over years.⁶ Current high-throughput atomistic simulation workflows have accelerated certain aspects of this process by enabling parallel screening of large chemical spaces,^{7,8} but they offer limited adaptivity and cannot reason to select the next computation. Failed calculations still require human intervention to diagnose and repair errors, limiting true automation.

Machine learning has introduced new paradigms for accelerating materials discovery through data-driven predictive and generative modeling.^{9–11} However, while these models achieve high accuracy, they remain confined to their training distributions and cannot interpret results within an adaptive, iterative scientific framework. Because materials discovery is inherently a multi-stage reasoning process, from hypothesis generation to experimental validation, these models play a crucial but limited role, addressing computation rather than scientific decision-making.

Large language models (LLMs) bring reasoning capabilities to scientific research. They integrate knowledge, generate executable code, interface with tools, and reason

in natural language.^{12–14} These abilities position LLMs as coordinators of end-to-end discovery workflows that link hypothesis generation, simulation, and analysis. However, LLMs acting in isolation face challenges including limited numerical precision, error accumulation, and lack of persistent state.¹⁵ Stateless interactions prevent them from retaining context, coordinating across tools, or adapting strategies based on intermediate results. Robust scientific reasoning therefore requires agentic architectures,¹⁶ i.e., systems coupling LLMs to memory, adaptive planning, and iterative feedback. Recent studies have shown that LLM agents can manage computational materials workflows with impressive autonomy. Systems such as VASPilot,¹⁷ DynaMate,¹⁸ and ChemGraph,¹⁹ perform procedural automation quantum-mechanical simulations, with VASPilot enabling end-to-end execution of VASP^{20–22} workflows spanning structure preparation, job execution, error recovery, and postprocessing. Beyond single-job automation, DREAMS²³ and MOFGen²⁴ extend autonomy to networked simulation environments. In DREAMS, specialized agents coordinate the setup and execution of Quantum Espresso^{25,26} calculations with automated convergence checks and recovery from failure. Moreover, LLMatDesign²⁷ introduced a self-reflective loop in which a LLM refined its proposals using a surrogate model pre-trained on density functional theory (DFT) data.

Here, we introduce Materials Agents for Simulation and Theory in Electronic-structure Reasoning (MASTER), an active learning framework that equips ensembles of LLMs with structured, collective reasoning for guided materials discovery. Designed around reasoning autonomy, MASTER coordinates interacting agents to deliberate, critique, and refine hypotheses through structured collaboration, deciding what atomistic simulations to perform and how to interpret their outcomes. Within this framework (Fig. 1a), a multimodal subsystem links natural-language objectives to validated DFT workflows, ensuring accurate generation of atomic structures, input files, and convergence parameters. Meanwhile, higher-level reasoning agents decide which materials to investigate next based on intermediate results. This separation of concerns allows the system to frame decision-making as a reasoning process that integrates individual perspectives into shared conclusions while preserving the rigor of first-principles computation. To evaluate how reasoning architecture influences discovery efficiency, we compare a single-agent baseline with three hierarchical multi-agent strategies—peer review, triage-ranking, and triage-forms. These architectures are tested using CO binding energetics across two chemical

domains: transition-metal (M) adatoms on a Cu(100) surface and M–N–C single-atom catalysts.^{28–30} MASTER identifies targeted binding energies within only a few iterations, whereas trial-and-error selection requires an order of magnitude more trials. For the highest-performing systems, analysis of reasoning trajectories reveals coherent and scientifically grounded decision patterns, reflecting how structured collaboration shapes adaptive, mechanistically informed exploration.

Results

MASTER Framework and Benchmark Materials Problem

The MASTER framework integrates LLM reasoning with autonomous electronic-structure simulation in a closed loop that emulates the operation of a scientific research team (Fig. 1a). The system is organized into three tightly coupled yet functionally distinct layers. The design layer comprises a team of LLM agents that use natural language to formulate hypotheses and deliberate, individually or collectively, on which material to evaluate next based on the accumulated simulation history. Once a candidate structure is chosen, the information is passed to the simulation layer, which provides a multimodal interface that converts high-level simulation objectives into validated DFT workflows. Here, a team of DFT agents autonomously generates inputs, atomic geometries, and executes first-principles calculations. The computed quantities, such as adsorption energies, are then returned to the review layer where a reviewer agent determines whether the specified criteria have been met or further exploration is required. By explicitly separating reasoning from simulation, the design agents focus on interpreting trends, weighing evidence, and planning experiments, while the DFT and reviewer agents handle numerical execution and verification, respectively. Together, these layers operationalize an autonomous scientific method, transforming static computational screening into an adaptive and self-correcting exploration of chemical space through LLM reasoning and first-principles feedback.

We apply MASTER to the problem of CO adsorption energetics, which has been widely studied in surface science and heterogeneous catalysis.^{31,32} CO binding energy serves as a key descriptor of catalytic activity and selectivity,^{33,34} governing reaction pathways in CO oxidation,³⁵ CO₂ reduction,³⁶ and the formation of C₂ products via CO dimerization.^{37,38} It depends sensitively on the local electronic and geometric envi-

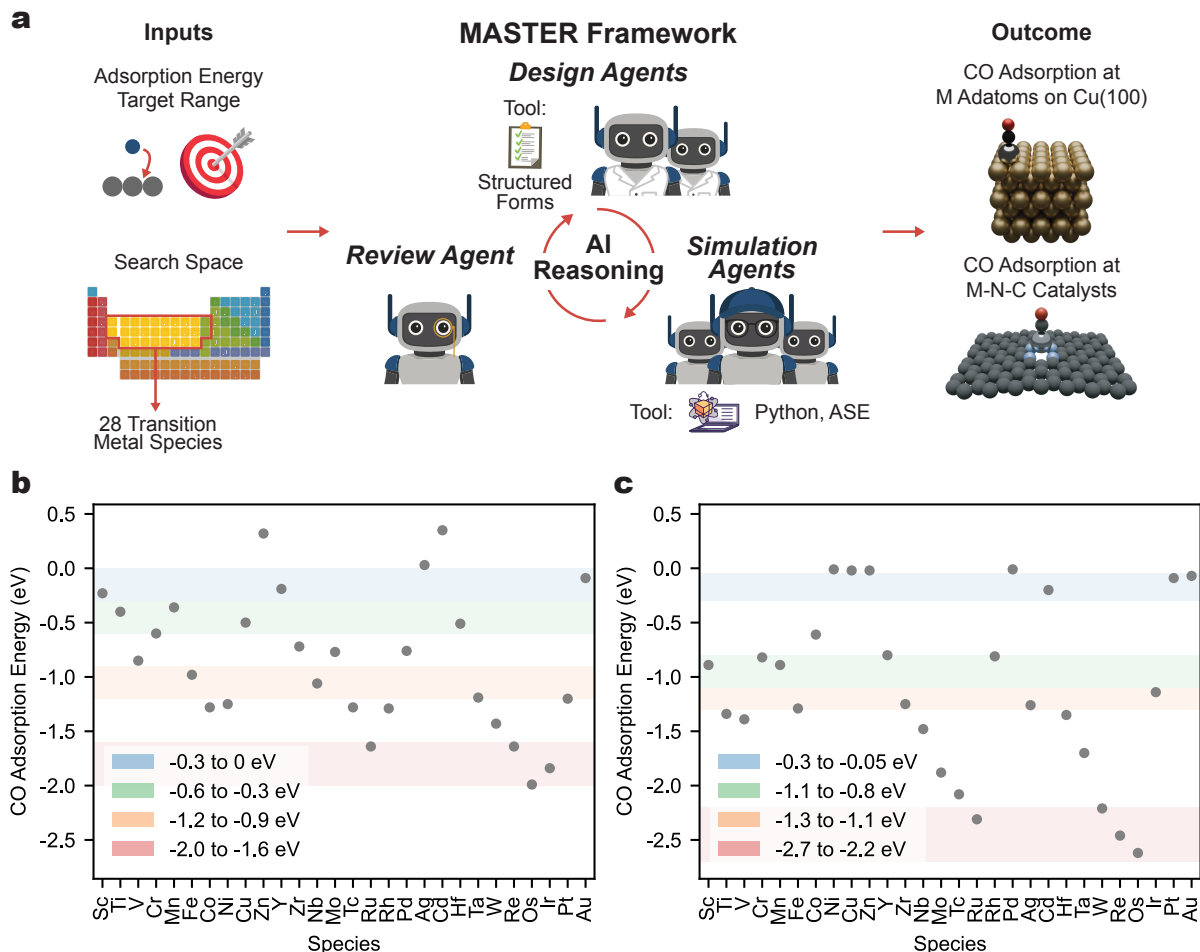


Figure 1: Unlocking autonomy in materials discovery via large language models-driven active learning. **a**, Schematic of the MASTER framework, where design agents propose materials, simulation agents generate atomic positions scripts, and a reviewer agent evaluates if outcomes meet required targets. **b**, Computed CO adsorption energies for transition-metal adatoms on Cu(100). **c**, Computed CO adsorption energies for M-N-C single-atom catalysts. In both **b** and **c** the shaded bands mark target energy windows from weak to strong binding. A comparison of the CO adsorption energies is presented in Figure S1.

ronment, producing a chemically intricate landscape that challenges autonomous LLM reasoning. Two scenarios are considered here (Figs. 1b and 1c). The first examines CO adsorption on transition-metal adatoms supported on Cu(100), a model for undercoordinated catalytic sites where low coordination substantially modifies the local binding environment.³⁹ The second extends the study to CO adsorption on M-N-C single-atom catalysts,³⁰ where the metal is coordinated by nitrogen ligands within a graphene host, introducing distinct ligand-field and covalency effects. Each domain spans twenty-eight transition metals from Sc to Au with DFT-computed adsorption energies. Four target CO binding-energy ranges were chosen in each case, spanning weak to strong binding regimes (Methods). In this context, MASTER must learn and generalize structure–property rela-

tionships directly from first-principles data and navigate complex search spaces without human guidance.

Natural Language to Density Functional Theory Simulations

The adatom adsorption problem provides an ideal proof-of-principle for MASTER because it represents a complex case that encapsulates the fundamental challenges in automating atomic-scale simulation, especially for the field of electrocatalysis.⁴⁰ Indeed, researchers may describe an idea succinctly, such as “place a CO molecule on an Ag adatom supported on Cu(100)”, but DFT codes require constructing the appropriate crystallographic surface, identifying high-symmetry adsorption sites, optimizing supercell dimensions, and defining vacuum spacing to avoid artificial slab couplings. MASTER must therefore translate scientific intent expressed in natural language, into fully executable simulations. Even when computational choices such as exchange–correlation functional or k -point meshes remain constant, each new atomic system must be built from scratch to capture the intended chemistry accurately. Thus, the Cu(100) adatom system offers a rigorous benchmark in which success or failure is unambiguous and the translation challenge is fully exposed.

To operationalize this translation, the simulation layer of MASTER implements an agentic DFT subsystem composed of three collaborating agents: CODEX,⁴¹ a Form Filler, and a Geometry Reviewer (Figure 2). Collectively, they convert natural-language queries into validated DFT inputs. This framework employs a pure prompting strategy that preserves the flexibility of LLMs while enforcing the precision required for scientific computation. Rather than using rigid templates or parsers, the CODEX agent receives rich contextual information, namely the user’s natural language query, detailed ASE^{42,43} usage patterns for surface science applications, and representative DFT workflows. At its core, the subsystem relies on intentional context engineering, since the structured inputs and reviewer evaluations allow subsequent agents to detect inconsistencies and correct earlier errors. This enables generalization to new materials, adsorption geometries, and co-adsorption motifs beyond the capabilities of rigid automation.

Here, we specialize the prompting context to surface–adsorbate systems, which steers the agents toward Cu(100) adatoms and CO adsorption geometries; this reflects the engineered scope of our demonstrations rather than a fundamental limitation of the

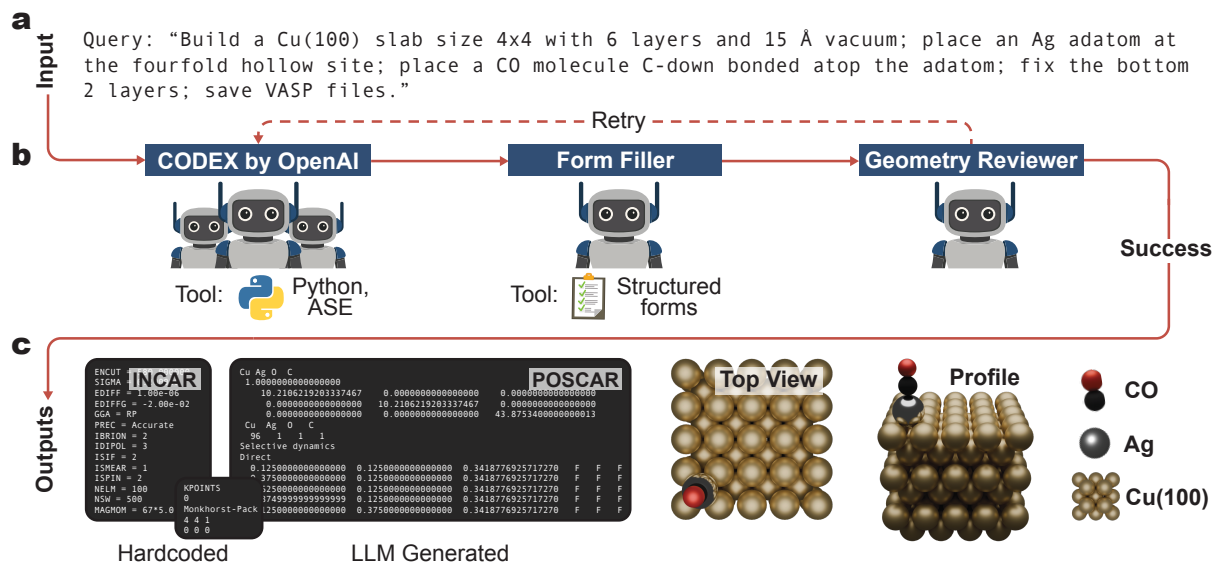


Figure 2: MASTER’S simulation agents convert natural-language into first-principles calculations. **a**, Example query specifying the construction of a CO–Ag adatom system on Cu(100). **b**, The DFT subsystem contains three collaborating agents that generate and verify atomic positions. CODEX⁴¹ produces an initial atomic structure by writing code that constructs the geometry using ASE.^{42,43} The Form Filler then evaluates the VASP^{20–22} structure file (POSCAR) together with visualizations of the top, profile, and side views. It completes an expert-prepared structured form (SI Note 2) that focuses the agent’s reasoning context on the most common fault points in surface geometries. The Geometry Reviewer reads this form to determine whether the geometry is correct. If it is incorrect, then it issues a retry request with targeted feedback to CODEX. This loop continues until the geometry satisfies all criteria. **c**, The final output is the validated POSCAR file describing the intended adsorbate-surface configuration.

framework. MASTER can be extended to other structure families, e.g., metal centers embedded in a pre-established N-doped carbon host,³⁰ by providing the corresponding structural priors in context. Likewise, supplying appropriate examples of input formats enables the simulation agents to target VASP,^{20–22} Quantum Espresso,^{25,26} ORCA,⁴⁴ or Gaussian.⁴⁵ Unlike distributed orchestration frameworks^{23,24} that coordinate large-scale pipelines, MASTER focuses on scientific reasoning and generates structures through a refinement loop in which stochastic variability is an asset: successive attempts explore alternative configurations, and the reviewer filters them so that only corrected geometries advance. This process links the flexibility of LLMs with first-principles checks and produces reliable results through guided exploration rather than deterministic programming.

Benchmarking demonstrates that this prompt-based agentic approach achieves a 97.2% success rate across all test cases after implementing self-revision loops, as shown in Figure 3a (see Methods). Here, success refers to generating a geometry that passes subject matter expert review following the geometry review form (Supplementary Note 2). The

subsystem handles both simple adatom placement and more complex adsorption scenarios involving the placement of CO molecules on the desired site. The self-correction loop steadily improves accuracy: 47.8% of cases succeed on the first iteration, 43.3% on the second, and 6.1% on the third. Overall success rates remain high, 97.8% for adatom-only and 96.7% for CO-adsorbed systems, with the latter requiring more iterations due to the complexity of molecular orientation (e.g., 11.1% vs 1.1% third-iteration convergence). The CO adsorption energies shown in Fig. 1b, obtained from fully relaxed DFT calculations, confirm that the constructed geometries are physically meaningful. This strong performance reflects the engineered flow of context between agents, which ensures that information about earlier errors is carried forward and corrected in subsequent attempts.

A representative example of osmium adatom placement on Cu(100) illustrates the iterative refinement process (Figure 3b). In the first iteration, the generated structure

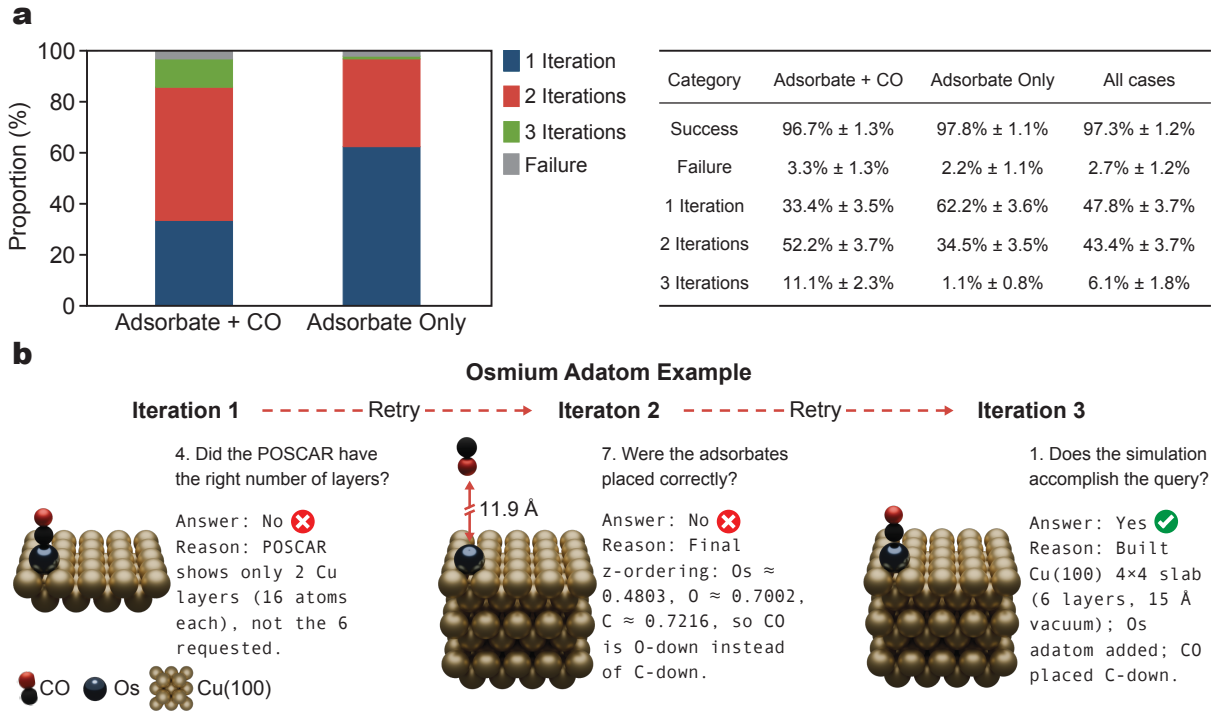


Figure 3: Benchmarking of MASTER simulation agents on transition-metal adatom simulations. **a**, Performance across 18 representative transition metal-adatom configurations on Cu(100). Each system was generated ten independent times, giving 180 total structure-generation runs that were expert-validated. Bars show the proportion of attempts that converged to a correct geometry in one, two, or three iterations. Summary statistics for success rates are shown in the table on the right. Standard deviations were computed using 5,000 bootstrap resamples. **b**, Osmium adatom example showing sequential correction across retries. Context built through multimodal evaluation and feedback between agents enables accurate generation and validation of the final geometry.

had only two Cu layers instead of the six requested, an error immediately flagged by the Geometry Reviewer. The second iteration corrected the layer count but misplaced the CO molecule with oxygen facing the Os adatom rather than carbon, violating the user’s instructions. Only in the third iteration did the system successfully generate the correct structure. This establishes a robust foundation for autonomous DFT execution, demonstrating that LLM agents can reliably translate complex materials specifications into computational workflows without manual intervention.

LLM Reasoning Strategies for Accelerated Materials Discovery

While accurate simulations are essential to the success of the developed agentic approach, the overall efficiency of autonomous discovery depends on how effectively reasoning agents navigate chemical space. Efficient exploration reduces the number of DFT evaluations, whose cost far outweighs that of LLM reasoning. To examine this portion of the workflow, we implemented four agentic reasoning architectures within the MASTER framework. In all configurations, the agents are instructed to apply chemical intuition, drawing on periodic trends, *d*-band theory from their pretraining, and correlations inferred from previous iterations, to propose new material candidates within the same closed loop involving the DFT and reviewer systems described in Fig. 1. The strategies differ in how they engineer the context used to generate and refine hypotheses, ranging from single agent reasoning to structured methods of collaboration and hierarchy. The four architectures, illustrated schematically in Fig. 4, are:

Single agent – A baseline configuration in which a single LLM autonomously decides the next candidate for evaluation, measuring the capability of one agent to perform self-consistent scientific reasoning without collaboration (Fig. 4a-c).

Peer review – A minimal form of collaboration in which two identical but independent agents propose candidates that are reconciled by an arbitrator, testing whether peer oversight improves reliability (Fig. 4d).

Triage-ranking – A hierarchical design in which a coarse selector proposes a pool of promising candidates that a fine selector ranks and chooses from, separating exploration from exploitation while retaining chemical diversity (Fig. 4e).

Triage-forms – Similar to triage-ranking but with an additional agent that fills prewritten expert-designed forms for each candidate before the fine selector makes a decision, testing whether guided context improves efficiency over free-form deliberation (Fig. 4f).

All reasoning strategies were benchmarked within the above-mentioned transition-metal chemical space, enabling a controlled proof-of-principle testbed for autonomous discovery, as shown in Figs. 4 and 5. We use GPT-5⁴⁶ as the base model for all agents. The four target adsorption-energy windows defined in Fig. 1b contain three to five transition metals, ensuring comparable discovery difficulty across binding regimes and allowing differences in performance to be attributed primarily to the reasoning architectures. Equivalent analyses for CO adsorption on M–N–C catalysts are provided in the Supplementary Figures S8 - S14.

Figure 4 summarizes the performance of the four argentic reasoning architectures and their final selected species statistics are presented in Figures S4-S7. The single agent baseline (Fig. 4a-c) converges reliably within fewer than ten iterations for all targets, reflecting directed exploration rather than random search. Even alone, the single agent outperforms stochastic baselines (Fig. 5a), yielding 100% cumulative success three times faster than trial-and-error, showing that LLMs already encode physically meaningful priors. The peer review configuration (Fig. 4d) performs similarly to the single agent, with no discernible performance gain. The species selector agents agreed on the next material in about 60% of runs while the arbitrator alternated evenly between them otherwise. Their differences were largely stochastic rather than chemically substantive, offering little additional guidance to the arbitrator. This suggests that collective reasoning can only improve performance when agents bring complementary perspectives or distinct priors.

In contrast, the triage-ranking system (Fig. 4e) achieves the most decisive convergence with an average performance gain as high as 2.14 iterations compared to the single agent baseline. Most runs identify an acceptable adatom within two or three iterations, showing improvement across all target energy windows. Its hierarchical structure balances exploration by the coarse selector with exploitation by the fine selector. By constraining comparison to a small, curated subset, the fine selector receives engineered context while remaining scalable. The triage-forms architecture (Fig. 4f), which combines a coarse selector, a structured form-filler, and a fine selector, performs slightly below triage-ranking

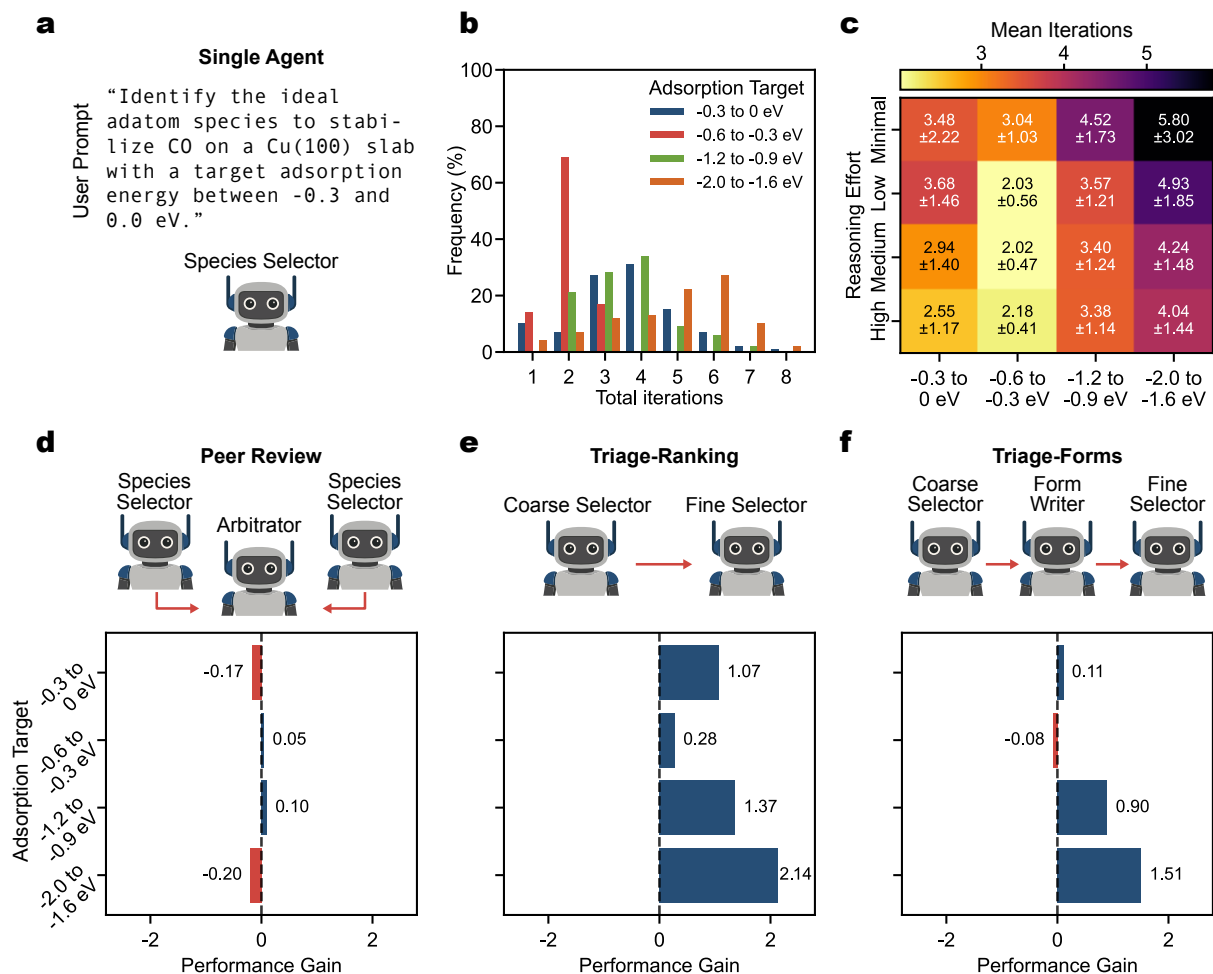


Figure 4: Comparative performance of MASTER’s design agents reasoning architectures. **a**, Schematic of the single agent setup and an example user prompt. **b**, Frequency distribution of successful iteration counts in the single agent configuration under low reasoning effort (see Methods)⁴⁷, aggregated over 100 trials for each adsorption energy target, where a success is defined as identifying a material whose adsorption energy falls within the specified target range. **c**, Heatmap showing the mean iterations to success across adsorption-energy targets and reasoning-effort levels in the single agent configuration. Mean iterations to success heatmaps for the remaining agentic architectures are presented in Figure S2. **d**, Performance gain in the peer review multi-agent configuration at high reasoning effort. Performance gain is defined as the improvement in average number of iterations required for success relative to the single agent baseline. **e**, Performance gain in the triage-ranking configuration at high reasoning effort. **f**, Performance gain in the triage-forms scenario at high reasoning effort.

but above the single agent across most energy windows. The best performance was obtained with a form emphasizing relative risk assessment, categorizing each option as a “safe bet”, “moderate risk”, “high risk”, or “unlikely”, identifying the top safe bet when available, and providing a brief rationale (SI Note 7). This design elicits qualitative scientific reasoning, whereas forms that required quantitative predictions, such as *d*-band-center estimates, consistently degraded performance.

Fig. 5a-b presents cumulative success probabilities comparing reasoning architectures

with three baselines: a theoretical random sampler, a Monte Carlo agent, and a rogue agent instructed to act randomly but allowed to reason (Methods). Across weak and strong binding windows, the single agent system outperforms the purely stochastic baselines, yielding three fold and eleven fold improvements, respectively. The triage-ranking architecture achieves the steepest success rise, identifying correct candidates within a few iterations and reaching near-unit cumulative probability sooner than any other approach. The rogue agent provides an instructive control: although nominally random, it exhibits a persistent semantic bias toward $5d$ elements in the early iterations (see pie chart inset). For the strong-binding target (-2.0 to -1.6 eV), this bias fortuitously aligns with the physical trend of increasing CO affinity down the $5d$ series, yielding apparent outperformance over all other agentic systems. For the weak-binding window (-0.3 to 0 eV), however, the same bias becomes detrimental, steering exploration away from relevant metals and causing success probabilities below even Monte Carlo levels. Cumulative success probabilities for the remaining adsorption energy windows for the Cu(100) case are presented in Figure S3. These behaviors reveal that unguided LLM priors can occasionally mimic chemical intuition but remain unreliable without grounding in simulation feedback. Cross-architecture performance differences also reflect a general principle from optimization theory: according to the no-free-lunch theorem⁴⁸, no single search strategy can be optimal across all problem classes. Consistent with this, triage-ranking excels for most Cu(100) targets, whereas other agentic designs perform comparatively or better in the M–N–C catalysts (Figures S8 and S9).

Discussion

To understand the behavior of the design agents, we first note that Fig. 4c reveals systematic trends with reasoning level. Increasing the GPT-5 reasoning effort generally reduces the mean number of iterations by roughly one, with some targets improving by up to two. For the weak-binding range, for example, minimal-reasoning agents selected Ag first in nearly all runs and reached Au only after several steps, while high-reasoning agents chose Au directly in roughly one-third of cases, reducing the average iterations from 3.5 to 2.6. Given the cost of DFT calculations, this reduction represents a meaningful computational gain since each iteration requires two DFT simulations (with and without CO

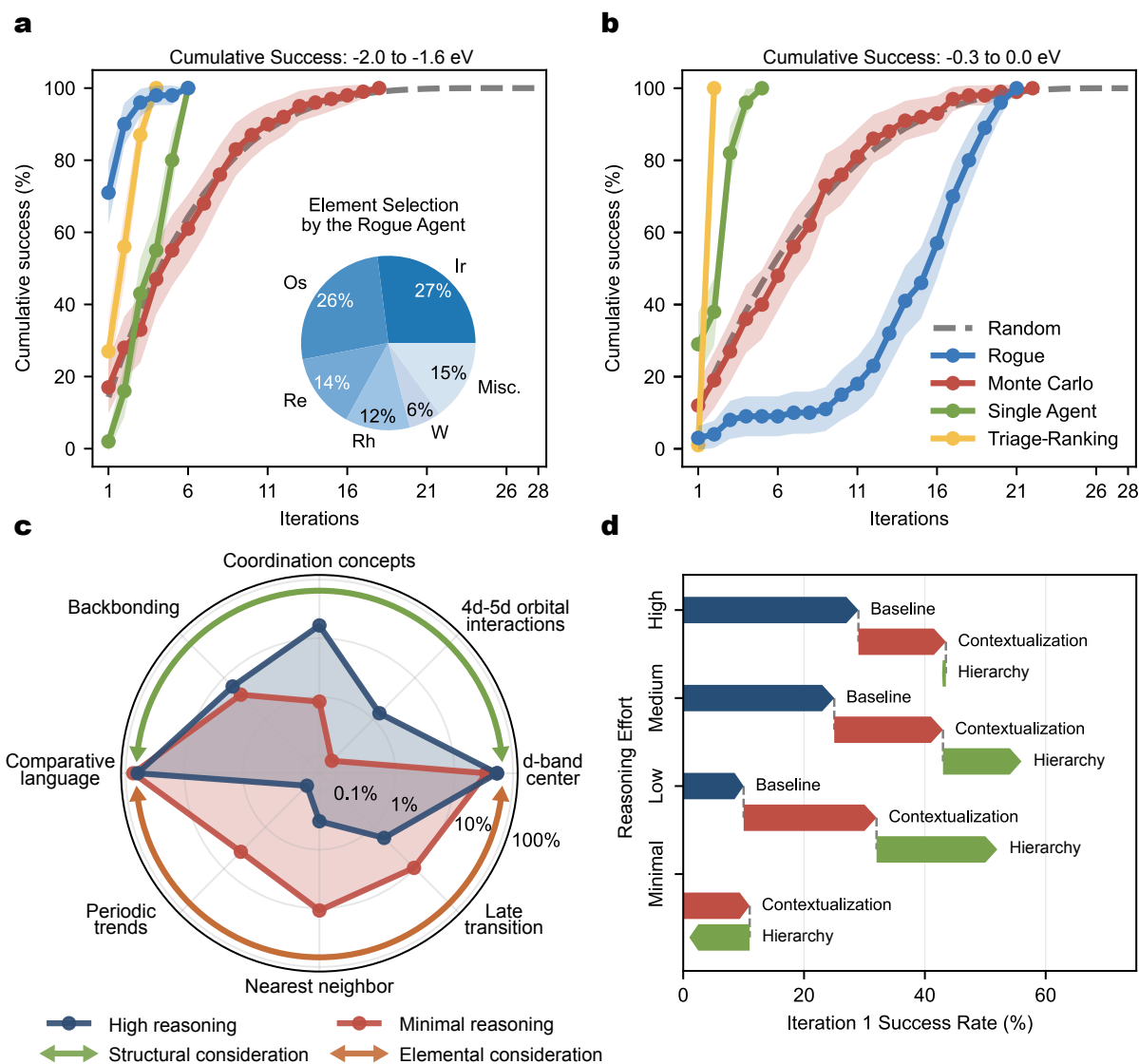


Figure 5: MASTER’s cumulative performance, interpretability, and hierarchy effects. **a**, Success probabilities for the -2 to -1.6 eV target range for Monte Carlo (minimal), rogue (minimal), single (high), and triage-ranking (medium) agents. The label in parenthesis indicates the associated reasoning effort for each architecture. The dashed line shows the theoretical result for trial-and-error selection.^{49,50} Inset pie chart shows transition-metal selections made by the rogue agent in the first iteration. Shaded regions denote 95% confidence intervals using 5,000 bootstrap resamples with the normal approximation (cumulative success probability $\pm 1.96 \times$ bootstrap standard deviation).⁵¹ **b**, Cumulative success for the -0.3 to 0 eV target range using the same agents, but with triage-ranking evaluated at minimal reasoning effort. **c**, Radar plot showing chemical concept categories invoked by the single agent across reasoning-effort levels, aggregated over all iterations and all trials, with categories defined by the keyword sets described in Supplementary Table S2. **d**, Iteration-1 success rates for the -0.3 to 0 eV target range, showing the gains beyond the single agent baseline and associated with instructed enumeration and ranking as well as the multi-agent triage hierarchy.

adsorbed). Minimal-reasoning tends to rely primarily on positional heuristics along the periodic table, whereas high-reasoning agents invoke mechanistic and structural argu-

ments (Fig. 5c). For instance, at high reasoning effort references to coordination effects appear nearly an order of magnitude more often and orbital interactions roughly three times as frequent as in minimal-reasoning runs.

We next examine the factors that influence iteration-1 success for the -0.3 to 0 eV adsorption-energy window for the Cu(100) adatom system (Fig. 5d). At this stage, the agents have not yet received any DFT results, so success depends solely on how the architecture structures information before feedback. The single agent baseline reflects the model’s ability to propose a plausible candidate from the prompt alone. Adding contextualization, i.e., enumerating all materials and requesting a qualitative ranking, improves accuracy across reasoning levels. At low and medium reasoning effort, the hierarchical multi-agent structure provides an additional gain: the coarse selector’s prescreening and the fine selector’s focused ranking further increase the likelihood of identifying a chemically reasonable first candidate. At high reasoning effort hierarchy offers little additional benefit. Taken together, the results show that the largest gains from hierarchical context engineering arise when reasoning depth is limited or when the candidate set cannot be exhaustively ranked by a single agent. In larger and more heterogeneous spaces, where full enumeration would be impractical, multi-agent hierarchies are therefore expected to play a central role in maintaining high early decision quality.

The reasoning trajectories in Fig. 6 illustrate how agentic hierarchy transforms exploration dynamics in the weak-binding regime. In the rogue agent (Fig. 6a), an early bias toward $5d$ metals (Ir, Os, Re, Rh, W) reflects a superficial association between atomic number and adsorption strength. Once these strong-binding early choices fail and the semantically driven prior collapses, the agent explores the space diffusely. By contrast, the triage-ranking agents (Fig. 6b) exhibits a structured and chemically interpretable transition network concentrated among Ag, Cu, Zn, Ni, and Au, which are elements near the weak-binding window. Frequent $\text{Ag} \rightarrow \text{Au}$ transitions and recurrent $\text{Ag} \rightarrow \text{Cu}$, $\text{Ag} \rightarrow \text{Zn}$ exchanges indicate that the agents iteratively explore neighboring regions of d -band filling. Paths from $\text{Ni} \rightarrow \text{Au}$ and $\text{Cu} \rightarrow \text{Au}$ further suggest stepwise correction toward a true weak-binding solution. These trajectory patterns rationalize the performance trends observed in Figs. 4–5. Hierarchical architectures not only accelerate convergence but also reorganize exploration into pathways guided by causal chemical reasoning rather than statistical association. As the system learns to associate structural and electronic

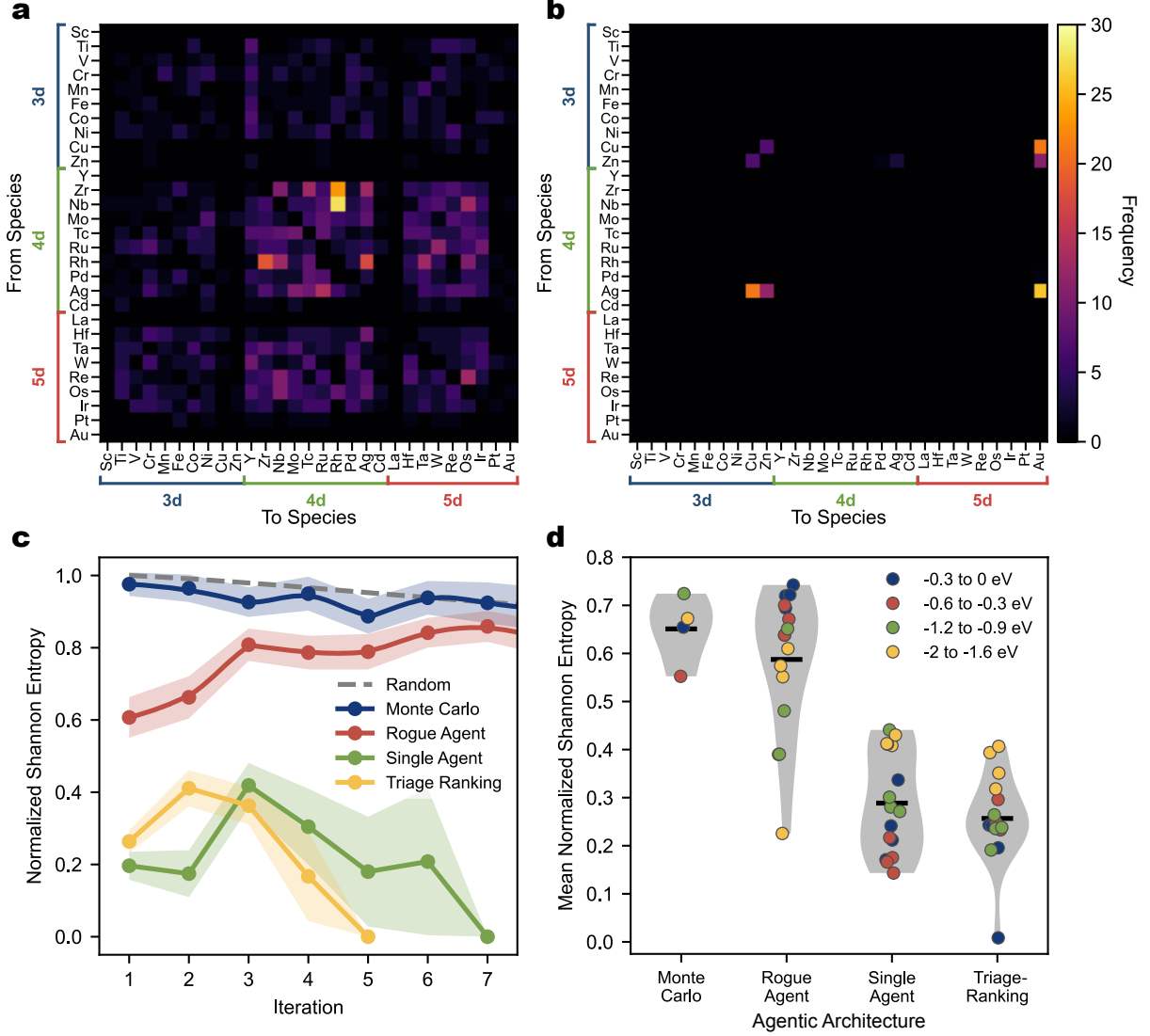


Figure 6: Reasoning trajectories and Shannon entropies across MASTER agent architectures. **a**, Frequency of transitions between transition metals across consecutive iterations for the -0.3 to 0.0 eV adsorption-energy range on Cu(100) across 100 independent runs for the rogue agent at high reasoning effort. Rows indicate the metal selected in iteration n and columns indicate the metal selected in iteration $n + 1$. Colored brackets denote d -block periods. **b**, Species transition heatmap for the triage ranking architecture at high reasoning effort for the same adsorption-energy range. **c**, Normalized Shannon entropies across iterations for the -0.3 to 0.0 eV adsorption-energy range at medium reasoning effort, with entropy definitions provided in the methods. Shaded regions denote 95% confidence intervals using 5,000 bootstrap resamples with the normal approximation (Shannon entropy value $\pm 1.96 \times$ bootstrap standard deviation).⁵¹ **d**, Mean normalized Shannon entropies for the Monte Carlo agent at minimal reasoning effort and for the rogue agent, single agent, and triage-ranking architectures across all reasoning-effort settings and all adsorption-energy ranges. Equivalent analysis for the M-N-C case is presented in Fig. S10.

features with adsorption strength, exploration becomes self-correcting, with each iteration reducing uncertainty. In this way, collective reasoning achieves efficiency through progressive information gain rather than exhaustive enumeration.

To quantify the exploration dynamics, we computed the normalized and mean normalized Shannon entropies⁵² for each architecture (Fig. 6c-d). Shannon entropy measures how broadly an agentic architecture distributes its selections across the transition metals in each iteration, with lower values indicating a more focused and information-efficient search.^{53,54} The rogue agent and Monte Carlo baselines maintain persistently high entropy, consistent with their erratic exploration. By contrast, both single agent and triage-ranking architectures exhibit pronounced entropy contraction, with the latter achieving the lowest mean entropy overall. These results show that the hierarchical architectures outperform the single agent baseline because the context they propagate systematically provides information advantage that guides the search toward the correct region of the search space in fewer iterations.

In the present MASTER implementation, the design agents have no *a priori* knowledge of the absolute adsorption-energy scale or level of theory used in our DFT calculations. They must instead infer these scales on the fly from the sequence of simulation outcomes and the acceptance criteria, learning which regimes correspond to weak, intermediate, or strong binding. In small, fully enumerable spaces such as our adatom benchmark, this implicit calibration is sufficient. In larger or less well-characterized domains, however, an additional retrieval-augmented agent could supply prior grounding by querying literature or materials databases, improving robustness and accelerating convergence when simulations are costly or the underlying energy landscape is complex. Similarly, structured, form-based triage is likely to become more valuable in such regimes, where standardized prompts can stabilize reasoning, enforce consistent comparison criteria, and preserve interpretability across many interacting agents.

Altogether, our findings show that structured agentic collaboration transforms large language models from procedural tools into adaptive scientific reasoners. Within MASTER, autonomy arises from interaction: agents that deliberate, incorporate feedback, and refine shared hypotheses guide exploration with increasing mechanistic consistency. By linking language, simulation, and theory into a unified workflow, MASTER enables efficient, self-correcting discovery. Across the CO-adsorption problems studied here, this combination of hierarchical reasoning and autonomous simulation reduces the number of required atomistic calculations by up to 90% relative to trial-and-error while preserving first-principles accuracy. Extending such architectures beyond materials science could

enable general-purpose scientific agents capable of autonomous hypothesis formation and reasoning across the physical and life sciences.

Methods

Atomistic Simulations using Density Functional Theory

All DFT calculations for the transition metal-adatom case were performed using the Vienna Ab initio Simulation Package^{20–22} (VASP, version 6.4.2) within the projector augmented-wave (PAW) formalism. We employed the revised Perdew–Burke–Ernzerhof (RPBE) functional⁵⁵ within the generalized gradient approximation (GGA) to describe exchange–correlation contribution to the system Hamiltonian. A plane-wave energy cutoff of 580 eV was used, and all calculations were spin-polarized. We modeled the Cu(100) surface as a six-layer, 4×4 periodic slab containing 96 Cu atoms, separated by a 15 Å vacuum region. A single transition-metal adatom was positioned in the fourfold hollow site of the surface, and CO was adsorbed atop the adatom.

Brillouin-zone integrations were performed using a $4 \times 4 \times 1$ Monkhorst–Pack k -point mesh, which was verified to yield converged adsorption energies within 0.01 eV. All structures were optimized until the forces on unconstrained atoms were below 0.02 eV Å^{−1} and electronic convergence was achieved to within 10^{−6} eV. The bottom two Cu layers were held fixed to their bulk positions, while all other atoms were allowed to relax. Adsorption energies were determined from total electronic energies of the fully relaxed structures according to Eq. (1):

$$E_{\text{ads}} = E_{\text{CO/M/Cu(100)}} - E_{\text{M/Cu(100)}} - E_{\text{CO(g)}}, \quad (1)$$

where $E_{\text{CO/M/Cu(100)}}$, $E_{\text{M/Cu(100)}}$, and $E_{\text{CO(g)}}$ are the total electronic energies of the CO-adsorbed system, the M-decorated Cu(100) slab, and the isolated CO molecule, respectively (see Table S1). All reported adsorption energies correspond to electronic energies at 0 K, without zero-point or entropic corrections which are left to future work but should not change the qualitative nature of the findings. Negative values of E_{ads} indicate exothermic adsorption.

CO adsorption energies on M-N₄C₁₀ catalysts are computed using DFT as previously reported.³⁰ An initial Fe-N₄C₁₀ structure with 66 total carbon atoms is first relaxed,

then starting structures for all transition metals are generated using ASE⁴² by replacing Fe with a given transition metal. Calculations are carried out with VASP using the RPBE functional and default PBE projector augmented wave-pseudopotentials^{56,57} and managed with the pyiron workflow framework.⁷ A cell size of $14.78 \text{ \AA} \times 12.80 \text{ \AA}$ is used for all surface calculations with a 20 \AA vacuum normal to the surface. A $4 \times 4 \times 1$ Monkhorst-Pack k -point mesh is employed with dipole corrections applied normal to the surface. Spin polarization is turned on for all calculations. The plane-wave basis cutoff is set to 600 eV, and a Fermi-Dirac smearing width of 0.0259 is used. During structural relaxation, only atomic positions are allowed to relax, while the cell volume and shape remain fixed. Geometries are converged to a threshold of $< 10^{-5}$ eV change in energy between sequential steps. The gas-phase energy of CO is computed by placing the molecule in the center of the same size unit cell as the $\text{M-N}_4\text{C}_{10}$ structures and allowed the atoms to relax.

For structural relaxation of the surfaces, the planar initial structure and a structure with the transition metal center displaced 0.6 \AA out of plane are both relaxed. This is done to avoid trapping in high-energy meta-stable configurations; the lower-energy optimized is used as the reference structure for subsequent CO adsorption calculations. CO-adsorbed structures are generated by placing CO above the transition metal center in three initial configurations: with the carbon atom bound to the surface and the oxygen atom in line with vector normal to the surface, with the oxygen atom bound to the surface and the carbon atom in line with vector normal to the surface, and a bidentate configuration with the C-O bond positioned directly above the transition metal and oriented parallel to the surface. The adsorbate structure which yields the lowest overall energy is then used for computing adsorption energy similarly to Eq. 1 but replacing $\text{M/Cu(100)} \rightarrow \text{M-N-C}$.

LLM Framework for Density Functional Theory Simulations

The atomic position generation component of MASTER uses a three-agent workflow built on OpenAI Agents SDK (version 0.0.18).⁵⁸ The Geometry Generator agent receives a natural language structure request and constructs a prompt containing the user query plus a JSON knowledge base with twelve ASE construction tips covering site placement, molecular orientation, and covalent radii for common surface atoms and adsorbates (SI Note 1). This prompt is passed to Codex (version 0.57.0)⁴¹ via command-line interface

as a subprocess, which returns a Python script using ASE library functions. The script executes in an isolated temporary directory to produce a VASP POSCAR file. The Geometry Generator agent then creates three orthogonal structure visualizations (top, side, and profile views) and transfers control to the Form Filler agent, which accesses outputs through shared filesystem directories.

The Form Filler agent analyzes the POSCAR file to verify atomic composition and layer count, examines the three visualization images, and completes an eight-question binary assessment (SI Note 2) evaluating composition, layer constraints, site placement, orientation, and vertical spacing. For the current surface adsorption application, assessment questions include domain-specific criteria such as “was the adsorbate placed in the right place” (evaluating hollow, bridge, or on-top site occupancy) and “were the adsorbates placed in the right orientation” (verifying molecular geometry such as C-down vs O-down for CO). The form template and construction tips are modular components that can be modified for other simulation domains by replacing the assessment criteria and construction guidelines while preserving the three-agent workflow architecture.

The completed form transfers to the Reviewer agent, which inspects the images and POSCAR file for consistency with the form assessment and makes the final acceptance decision. Structures pass only if all eight questions receive affirmative responses and required files exist. Rejection triggers written feedback specifying identified deficiencies and returns control to the generator agent with incremented version numbering. The generator agent incorporates this feedback into revised Codex prompts for iterative refinement, supporting up to five cycles. Rejected structures archive to version-controlled subdirectories preserving the complete revision history. All agent decisions, Codex prompts, and generated scripts log to structured markdown files for reproducibility.

Benchmarking Protocol for Agentic LLM Reasoning

We benchmarked four reasoning strategies in MASTER using the OpenAI Agents SDK with GPT-5 models.⁴⁶ All agents operated at low verbosity, and the reasoning effort parameter, as implemented by OpenAI, was swept across minimal, low, medium, and high settings.⁴⁷ The search space for both the adatom and M-N-C test cases comprised the 28 transition metals from Sc to Au. Targets were specified as numerical adsorption-energy bands (e.g., -0.6 eV to -0.3 eV) without tolerance, and each run terminated when

the reviewer confirmed that the measured energy laid within the specified band. Each strategy was executed as independent batches spanning four adsorption-energy windows (SI Note 3). For every window and reasoning-effort level, one hundred runs were recorded. Within each run, the selector proposed an untested element from the Sc–Au set, the evaluator returned a fixed ground-truth adsorption energy from DFT computations, and the reviewer determined whether the band criterion had been met. Invalid or duplicate proposals were rejected and re-prompted up to a total of 20 maximum retries. The trial is considered a failure if the system reaches the maximum number of retries. In all tests we preformed in this paper, we never observed a failed trial. The prompts and system messages used for all agentic architectures are presented in SI Notes 4-7 for the transition metal adamon on Cu(100) case and SI Notes 8-11 for the M-N-C case.

To establish a theoretical baseline for comparison, we computed the cumulative probability of success for random sampling without replacement as^{49,50}

$$P_{\text{success}}(i) = 1 - \frac{\binom{N-K}{i}}{\binom{N}{i}}, \quad (2)$$

where $N = 28$ is the total number of transition metals candidates, K is the number of correct materials within the target energy band, i is the iteration index, and $\binom{n}{k}$ denotes the binomial coefficient.

In addition to the reasoning-based strategies, we also implemented two control agents to assess the impact of strategic reasoning versus pure randomness. The Monte Carlo agent employs a deterministic random number generator tool that selects a uniform random index from 1 to 28, mapping each index to the corresponding transition metal in the Sc–Au series (SI Note 12). The agent is instructed to call this tool without applying any materials science reasoning, providing a computationally controlled random baseline. The rogue agent, by contrast, is an LLM-based selector instructed to perform completely random selection with no strategic reasoning, bias, or optimization (SI Note 13). The rogue agent is explicitly prohibited in prompting from using materials science concepts (periodic trends, d -band theory, electronegativity) and is directed to select candidates as if rolling dice, establishing an LLM-based random baseline that tests whether the model can suppress its reasoning priors when instructed to do so.

LLM Chemical Concept Usage Analysis

To quantify how GPT-5 reasoning effort influences performance, we analyzed the chemical concepts invoked by the single agent system across all chemical targets (Fig. 1b) for minimal and high reasoning effort. This analysis enables us to determine how the model’s conceptual grounding shifts with reasoning effort and directly supports the comparison shown in Figure 5c. We defined eight categories representing distinct chemical concepts, with each category encompassing multiple keyword variants to capture linguistic variations (Table S2). These categories can be broadly grouped into structural concepts (coordination effects, backbonding, $4d/5d$ orbital interactions) and elemental concepts (periodic trends, nearest neighbor, and late transition metal classifications). The keyword matching procedure performs case-insensitive substring searches within each reasoning statement, counting a statement as positive for a category if any variant appears. Each statement contributes at most one count per category, preventing double-counting when multiple variants of the same concept appear in a single statement.

Frequency calculations proceed by dividing the number of statements containing each category’s keywords by the total number of statements analyzed, expressed as percentages. For minimal reasoning, we analyzed 1,684 statements across all energy windows and runs; for high reasoning, 1,215 statements. The resulting percentages represent the fraction of selection decisions that invoked each chemical concept category.

Shannon Entropy

To analyze the heterogeneity of species selected across iterations, we computed the normalized Shannon entropy for each iteration (4).^{53,54} We first construct the probability distribution $p_j^{(i)}$ based on the frequency with which each species j was selected across all active runs at iteration i . The Shannon entropy H_i is then computed as⁵²

$$H_i = - \sum_{j=1}^N p_j^{(i)} \log(p_j^{(i)}), \quad (3)$$

where $N = 28$ is the total number of candidate species. To normalize the Shannon entropy to the range $[0,1]$, we divide by the maximum possible entropy $\log(N)$:

$$H_{\text{norm}}^{(i)} = \frac{H_i}{\log N}. \quad (4)$$

This normalized entropy equals 1 when all species are equally likely to have been chosen (maximum heterogeneity) and approaches 0 as selection becomes concentrated on fewer species (lower heterogeneity). To characterize the overall exploration behavior of each agentic architecture across an entire run, we computed the mean normalized Shannon entropy \bar{H}_{norm} over T iterations:

$$\bar{H}_{\text{norm}} = \frac{1}{T} \sum_{i=1}^T H_{\text{norm}}^{(i)}. \quad (5)$$

This metric represents the average normalized Shannon entropy per iteration, providing a single value that captures the typical heterogeneity level maintained throughout the selection process.

Declarations

Acknowledgments

This research was supported by the Institute for Materials Science of Los Alamos National Laboratory. Research presented in this article was supported by the Laboratory Directed Research and Development program of Los Alamos National Laboratory under project number 20230065DR. This research used resources provided by the Los Alamos National Laboratory Institutional Computing Program, which is supported by the U.S. Department of Energy National Nuclear Security Administration under Contract No. 89233218CNA000001. This material is also based upon work supported by the U.S. Department of Energy, Office of Critical Minerals and Energy Innovation (CMEI), specifically the Hydrogen and Fuel Cell Technologies Office (HFTO) under contract ELY-BIL003. Samuel Rothfarb received a UConn’s Pratt & Whitney Institute for Advanced Systems Engineering Graduate Fellowship which enabled Samuel to contribute to this work.

Conflict of Interest

The authors declare no conflicts of interest.

Code Availability

The code supporting this work is available from the corresponding authors upon reasonable request.

Data Availability

The data supporting this work is provided in the Supplementary Information.

Author Contribution

S.R., B.L., E.F.H., and W.K.K. conceptualized the project. S.R. developed density functional theory calculations for adatoms on Cu(100) under guidance from E.F.H. and the reasoning strategies under guidance from W.K.K. M.D. and I.M performed benchmark DFT calculations for the M-N-C systems. S.R. and W.K.K. wrote the paper and received feedback from all authors, who reviewed and approved its final version. B.L., E.F.H, and W.K.K, supervised the project execution.

References

- [1] Francesco Branda, Massimo Ciccozzi, and Fabio Scarpa. “Artificial intelligence in scientific research: Challenges, opportunities and the imperative of a human-centric synergy”. In: *Journal of Informetrics* 19.4 (2025), p. 101727. ISSN: 1751-1577. DOI: <https://doi.org/10.1016/j.joi.2025.101727>. URL: <https://www.sciencedirect.com/science/article/pii/S1751157725000896>.
- [2] Junjie Wang et al. “MAGUS: machine learning and graph theory assisted universal structure searcher”. In: *National Science Review* 10.7 (May 2023), nwad128. ISSN: 2095-5138. DOI: 10.1093/nsr/nwad128. eprint: <https://academic.oup.com/nsr/article-pdf/10/7/nwad128/50709989/nwad128.pdf>. URL: <https://doi.org/10.1093/nsr/nwad128>.

- [3] Chi Chen et al. “Accelerating Computational Materials Discovery with Machine Learning and Cloud High-Performance Computing: from Large-Scale Screening to Experimental Validation”. In: *Journal of the American Chemical Society* 146.29 (July 2024), pp. 20009–20018. ISSN: 1520-5126. DOI: 10.1021/jacs.4c03849. URL: <http://dx.doi.org/10.1021/jacs.4c03849>.
- [4] Can Leng et al. “Fully automated high-throughput computer-based catalytic material screening framework and its application on the new-generation Tianhe supercomputer”. In: *Computational Materials Science* 252 (2025), p. 113775. ISSN: 0927-0256. DOI: <https://doi.org/10.1016/j.commatsci.2025.113775>. URL: <https://www.sciencedirect.com/science/article/pii/S0927025625001181>.
- [5] Uzoma Nwabara et al. “High throughput computational and experimental methods for accelerated electrochemical materials discovery”. In: *J. Mater. Chem. A* 13 (32 2025), pp. 26041–26066. DOI: 10.1039/D5TA00331H. URL: <http://dx.doi.org/10.1039/D5TA00331H>.
- [6] Miao Zhong et al. “Accelerated discovery of CO₂ electrocatalysts using active machine learning”. en. In: *Nature* 581.7807 (May 2020). Publisher: Nature Publishing Group, pp. 178–183. ISSN: 1476-4687. DOI: 10.1038/s41586-020-2242-8. URL: <https://www.nature.com/articles/s41586-020-2242-8>.
- [7] Jan Janssen et al. “pyiron: An integrated development environment for computational materials science”. In: *Computational Materials Science* 163 (2019), pp. 24–36. ISSN: 0927-0256. DOI: <https://doi.org/10.1016/j.commatsci.2018.07.043>. URL: <http://www.sciencedirect.com/science/article/pii/S0927025618304786>.
- [8] Emil Annevelink et al. “AutoMat: Automated materials discovery for electrochemical systems”. In: *MRS Bulletin* 47.10 (Oct. 1, 2022), pp. 1036–1044. ISSN: 1938-1425. DOI: 10.1557/s43577-022-00424-0. URL: <https://doi.org/10.1557/s43577-022-00424-0>.
- [9] Amil Merchant et al. “Scaling deep learning for materials discovery”. In: *Nature* 624.7990 (Dec. 1, 2023), pp. 80–85. ISSN: 1476-4687. DOI: 10.1038/s41586-023-06735-9. URL: <https://doi.org/10.1038/s41586-023-06735-9>.

- [10] Claudio Zeni et al. “A generative model for inorganic materials design”. In: *Nature* 639.8055 (2025), pp. 624–632. DOI: 10.1038/s41586-025-08628-5. URL: <https://doi.org/10.1038/s41586-025-08628-5>.
- [11] Brandon M. Wood et al. *UMA: A Family of Universal Models for Atoms*. 2025. arXiv: 2506.23971 [cs.LG]. URL: <https://arxiv.org/abs/2506.23971>.
- [12] Ashish Vaswani et al. *Attention Is All You Need*. 2023. arXiv: 1706.03762 [cs.CL]. URL: <https://arxiv.org/abs/1706.03762>.
- [13] Jacob Devlin et al. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. 2019. arXiv: 1810.04805 [cs.CL]. URL: <https://arxiv.org/abs/1810.04805>.
- [14] Tom B. Brown et al. *Language Models are Few-Shot Learners*. 2020. arXiv: 2005.14165 [cs.CL]. URL: <https://arxiv.org/abs/2005.14165>.
- [15] Milad Moradi et al. “A Critical Review of Methods and Challenges in Large Language Models”. In: *Computers, Materials and Continua* 82.2 (2025), pp. 1681–1698. ISSN: 1546-2218. DOI: <https://doi.org/10.32604/cmc.2025.061263>. URL: <https://www.sciencedirect.com/science/article/pii/S1546221825000992>.
- [16] Shuo Ren et al. “Towards Scientific Intelligence: A Survey of LLM-based Scientific Agents”. In: (2025). arXiv: 2503.24047 [cs.AI]. URL: <https://arxiv.org/abs/2503.24047>.
- [17] Hao Tang et al. “VASPilot: A large language model assistant for atomistic simulation workflows”. In: *arXiv preprint* (2025). arXiv: 2508.07035.
- [18] Orlando A. Mendible-Barreto et al. “DynaMate: leveraging AI-agents for customized research workflows”. In: *Mol. Syst. Des. Eng.* 10 (7 2025), pp. 585–598. DOI: 10.1039/D5ME00062A. URL: <http://dx.doi.org/10.1039/D5ME00062A>.
- [19] Thang D. Pham, Aditya Tanikanti, and Murat Keceli. “ChemGraph: A large language model multi-agent framework for computational chemistry”. In: *arXiv preprint* (2025). arXiv: 2506.06363.
- [20] G. Kresse and J. Hafner. “Ab initio molecular dynamics for liquid metals”. In: *Phys. Rev. B* 47 (1 Jan. 1993), pp. 558–561. DOI: 10.1103/PhysRevB.47.558. URL: <https://link.aps.org/doi/10.1103/PhysRevB.47.558>.

- [21] G. Kresse and J. Furthmüller. “Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set”. In: *Computational Materials Science* 6.1 (1996), pp. 15–50. ISSN: 0927-0256. DOI: [https://doi.org/10.1016/0927-0256\(96\)00008-0](https://doi.org/10.1016/0927-0256(96)00008-0). URL: <https://www.sciencedirect.com/science/article/pii/0927025696000080>.
- [22] G. Kresse and J. Furthmüller. “Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set”. In: *Phys. Rev. B* 54 (16 Oct. 1996), pp. 11169–11186. DOI: [10.1103/PhysRevB.54.11169](https://doi.org/10.1103/PhysRevB.54.11169). URL: <https://link.aps.org/doi/10.1103/PhysRevB.54.11169>.
- [23] Ziqi Wang et al. “DREAMS: Density Functional Theory Based Research Engine for Agentic Materials Simulation”. In: (2025). arXiv: 2507.14267 [cs.AI]. URL: <https://arxiv.org/abs/2507.14267>.
- [24] Heather J. Kulik, Timothy J. S. Evans, Qiang Zhao, et al. “MOFGen: A multi-agent framework for the autonomous design of metal–organic frameworks”. In: *arXiv preprint* (2024). arXiv: 2504.14110.
- [25] Paolo Giannozzi et al. “QUANTUM ESPRESSO: a modular and open-source software project for quantum simulations of materials”. In: *Journal of Physics: Condensed Matter* 21.39 (Sept. 2009), p. 395502. DOI: [10.1088/0953-8984/21/39/395502](https://doi.org/10.1088/0953-8984/21/39/395502). URL: <https://doi.org/10.1088/0953-8984/21/39/395502>.
- [26] P. Giannozzi et al. “Advanced capabilities for materials modelling with Quantum ESPRESSO”. In: *Journal of Physics: Condensed Matter* 29.46 (Oct. 2017), p. 465901. DOI: [10.1088/1361-648X/aa8f79](https://doi.org/10.1088/1361-648X/aa8f79). URL: <https://doi.org/10.1088/1361-648X/aa8f79>.
- [27] Yu Gu et al. “LLMatDesign: Large language model for autonomous materials design with self-reflection and reasoning”. In: *arXiv preprint* (2024). arXiv: 2406.13163.
- [28] Hoon T. Chung et al. “Direct atomic-level insight into the active sites of a high-performance PGM-free ORR catalyst”. In: *Science* 357.6350 (2017), pp. 479–484. DOI: [10.1126/science.aan2255](https://doi.org/10.1126/science.aan2255). eprint: <https://www.science.org/doi/pdf/10.1126/science.aan2255>. URL: <https://www.science.org/doi/abs/10.1126/science.aan2255>.







- [29] Shuyu Liang et al. “Electrochemical Reduction of CO₂ to CO over Transition Metal/N-Doped Carbon Catalysts: The Active Sites and Reaction Mechanism”. en. In: *Advanced Science* 24 (2021). ISSN: 2198-3844. DOI: 10.1002/advs.202102886. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/advs.202102886>.
- [30] Megan C. Davis et al. “Computational screening of transition metal-nitrogen-carbon materials as electrocatalysts for CO₂ reduction”. In: *Electrochimica Acta* 510 (2025), p. 145357. ISSN: 0013-4686. DOI: <https://doi.org/10.1016/j.electacta.2024.145357>. URL: <https://www.sciencedirect.com/science/article/pii/S0013468624015937>.
- [31] Matteo Roiaz et al. “Roughening of Copper (100) at Elevated CO Pressure: Cu Adatom and Cluster Formation Enable CO Dissociation”. In: *The Journal of Physical Chemistry C* 123.13 (Apr. 2019). Publisher: American Chemical Society, pp. 8112–8121. ISSN: 1932-7447. DOI: 10.1021/acs.jpcc.8b07668. URL: <https://doi.org/10.1021/acs.jpcc.8b07668>.
- [32] P. T. P. Ryan et al. “Probing structural changes upon carbon monoxide coordination to single metal adatoms”. In: *The Journal of Chemical Physics* 152.5 (Feb. 2020), p. 051102. ISSN: 0021-9606. DOI: 10.1063/1.5137904. eprint: https://pubs.aip.org/aip/jcp/article-pdf/doi/10.1063/1.5137904/16735360/051102_1_online.pdf. URL: <https://doi.org/10.1063/1.5137904>.
- [33] David Vázquez-Parga et al. “A computational map of the probe CO molecule adsorption and dissociation on transition metal low Miller indices surfaces”. In: *Applied Surface Science* 618 (2023), p. 156581. ISSN: 0169-4332. DOI: <https://doi.org/10.1016/j.apsusc.2023.156581>. URL: <https://www.sciencedirect.com/science/article/pii/S016943322300257X>.
- [34] Zaiqi Li et al. “Mesostructure-Specific Configuration of *CO Adsorption for Selective CO₂ Electroreduction to C₂+ Products”. In: *Angewandte Chemie International Edition* 64.1 (2025), e202413832. DOI: <https://doi.org/10.1002/anie.202413832>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/anie.202413832>. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/anie.202413832>.

- [35] Haixia Gao et al. “CO₂ reduction reaction pathways on single-atom Co sites: Impacts of local coordination environment”. In: *Chinese Journal of Catalysis* 43.3 (2022), pp. 832–838. ISSN: 1872-2067. DOI: [https://doi.org/10.1016/S1872-2067\(21\)63893-7](https://doi.org/10.1016/S1872-2067(21)63893-7). URL: <https://www.sciencedirect.com/science/article/pii/S1872206721638937>.
- [36] Dongfang Cheng et al. “Structure Sensitivity and Catalyst Restructuring for CO₂ Electro-reduction on Copper”. In: *Nature Communications* 16.1 (Apr. 2025), p. 4064. ISSN: 2041-1723. DOI: [10.1038/s41467-025-59267-3](https://doi.org/10.1038/s41467-025-59267-3). URL: <https://doi.org/10.1038/s41467-025-59267-3>.
- [37] Xinze Bi et al. “Electroreduction of CO₂ to C₂H₄ Regulated by Spacing Effect: Mechanistic Insights from DFT Studies”. In: *Energy Material Advances* 4 (2023), p. 0037. DOI: [10.34133/energymatadv.0037](https://doi.org/10.34133/energymatadv.0037). eprint: <https://spj.science.org/doi/pdf/10.34133/energymatadv.0037>. URL: <https://spj.science.org/doi/abs/10.34133/energymatadv.0037>.
- [38] Hassina Tabassum et al. “Surface engineering of Cu catalysts for electrochemical reduction of CO₂ to value-added multi-carbon products”. In: *Chem Catalysis* 2.7 (2022), pp. 1561–1593. ISSN: 2667-1093. DOI: <https://doi.org/10.1016/j.checat.2022.04.012>. URL: <https://www.sciencedirect.com/science/article/pii/S2667109322002214>.
- [39] Magnus A. H. Christiansen et al. “Single-Atom Substituents in Copper Surfaces May Adsorb Multiple CO Molecules”. In: *The Journal of Physical Chemistry Letters* 15.21 (May 2024). Publisher: American Chemical Society, pp. 5654–5658. DOI: [10.1021/acs.jpcllett.4c00899](https://doi.org/10.1021/acs.jpcllett.4c00899). URL: <https://doi.org/10.1021/acs.jpcllett.4c00899>.
- [40] Jens K. Nørskov et al. “Density functional theory in surface chemistry and catalysis”. In: *Proceedings of the National Academy of Sciences* 108.3 (2011), pp. 937–943. DOI: [10.1073/pnas.1006652108](https://doi.org/10.1073/pnas.1006652108). eprint: <https://www.pnas.org/doi/pdf/10.1073/pnas.1006652108>. URL: <https://www.pnas.org/doi/abs/10.1073/pnas.1006652108>.
- [41] Mark Chen et al. *Evaluating Large Language Models Trained on Code*. 2021. arXiv: 2107.03374 [cs.LG]. URL: <https://arxiv.org/abs/2107.03374>.

- [42] Ask Hjorth Larsen et al. “The atomic simulation environment—a Python library for working with atoms”. In: *Journal of Physics: Condensed Matter* 29.27 (2017), p. 273002. URL: <http://stacks.iop.org/0953-8984/29/i=27/a=273002>.
- [43] S. R. Bahn and K. W. Jacobsen. “An object-oriented scripting interface to a legacy electronic structure code”. English. In: *Comput. Sci. Eng.* 4.3 (May 2002), pp. 56–66. ISSN: 1521-9615. DOI: 10.1109/5992.998641.
- [44] F. Neese. “The ORCA program system”. In: *WIREs Comput. Molec. Sci.* 2.1 (2012), pp. 73–78. DOI: 10.1002/wcms.81.
- [45] M. J. Frisch et al. *Gaussian~16 Revision C.01*. Gaussian Inc. Wallingford CT. 2016.
- [46] OpenAI. *GPT-5*. <https://openai.com/gpt-5>. Accessed: 2025-12-04. 2025.
- [47] OpenAI. *Reasoning models Explore advanced reasoning and problem-solving models*. <https://platform.openai.com/docs/guides/reasoning>. Accessed: 2025-12-12. 2025.
- [48] D.H. Wolpert and W.G. Macready. “No free lunch theorems for optimization”. In: *IEEE Transactions on Evolutionary Computation* 1.1 (1997), pp. 67–82. DOI: 10.1109/4235.585893.
- [49] John Ahlgren. “The Probability Distribution for Draws Until First Success Without Replacement”. In: *arXiv preprint arXiv:1404.1161* (2014). Available at: <https://arxiv.org/pdf/1404.1161>.
- [50] E. T. Jaynes. *Probability Theory: The Logic of Science*. Cambridge University Press, 2003. ISBN: 9780521592710.
- [51] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. 2nd ed. New York: Springer, 2009.
- [52] Robert M. Gray. *Entropy and Information Theory*. 2nd ed. Springer New York, 2011. ISBN: 978-1-4419-7969-8, 978-1-4899-8132-5, 978-1-4419-7970-4. DOI: 10.1007/978-1-4419-7970-4. URL: <https://link.springer.com/book/10.1007/978-1-4419-7970-4>.

- [53] C. E. Shannon. “A mathematical theory of communication”. In: *The Bell System Technical Journal* 27.3 (1948), pp. 379–423. DOI: 10.1002/j.1538-7305.1948.tb01338.x.
- [54] S. Verdu. “Fifty years of Shannon theory”. In: *IEEE Transactions on Information Theory* 44.6 (1998), pp. 2057–2078. DOI: 10.1109/18.720531.
- [55] B. Hammer, L. B. Hansen, and J. K. Nørskov. “Improved adsorption energetics within density-functional theory using revised Perdew-Burke-Ernzerhof functionals”. In: *Phys. Rev. B* 59 (11 Mar. 1999), pp. 7413–7421. DOI: 10.1103/PhysRevB.59.7413. URL: <https://link.aps.org/doi/10.1103/PhysRevB.59.7413>.
- [56] P. E. Blöchl. “Projector augmented-wave method”. In: *Phys. Rev. B* 50 (24 Dec. 1994), pp. 17953–17979. DOI: 10.1103/PhysRevB.50.17953. URL: <https://link.aps.org/doi/10.1103/PhysRevB.50.17953>.
- [57] G. Kresse and D. Joubert. “From ultrasoft pseudopotentials to the projector augmented-wave method”. In: *Phys. Rev. B* 59 (3 Jan. 1999), pp. 1758–1775. DOI: 10.1103/PhysRevB.59.1758. URL: <https://link.aps.org/doi/10.1103/PhysRevB.59.1758>.
- [58] OpenAI. *OpenAI Agents SDK*. <https://github.com/openai/openai-agents-python>. Accessed: 2025-12-05. 2025.

Supplementary Information for Hierarchical Multi-agent Large Language Model Reasoning for Autonomous Functional Materials Discovery

Samuel Rothfarb^{†,‡} , Megan C. Davis[‡] , Ivana Matanovic[‡] ,
Baikun Li^{†*} , Edward F. Holby^{‡*} , and Wilton J.M. Kort-Kamp^{‡*} 

[†]*School of Civil & Environmental Engineering, University of Connecticut, Storrs,
Connecticut 06269, United States.*

[‡]*Theoretical Division, Los Alamos National Laboratory, Los Alamos, New Mexico
87545, United States.*

*Corresponding authors: baikun.li@uconn.edu, holby@lanl.gov, kortkamp@lanl.gov

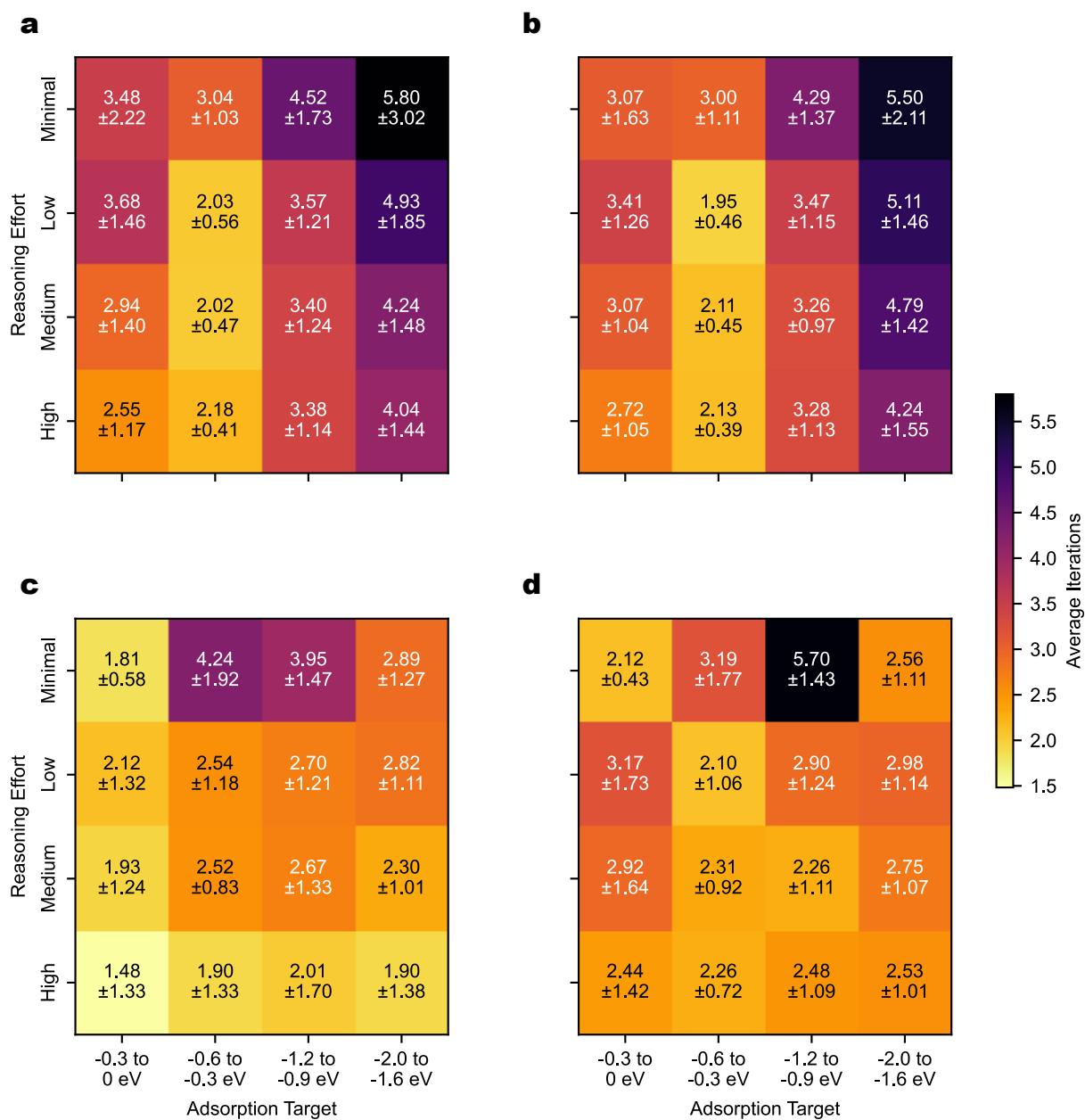


Figure S2: Heatmaps showing mean iterations to success for the Cu(100) test case for the **a**, single agent, **b**, peer review, **c**, triage-ranking, and **d**, triage-forms configurations.

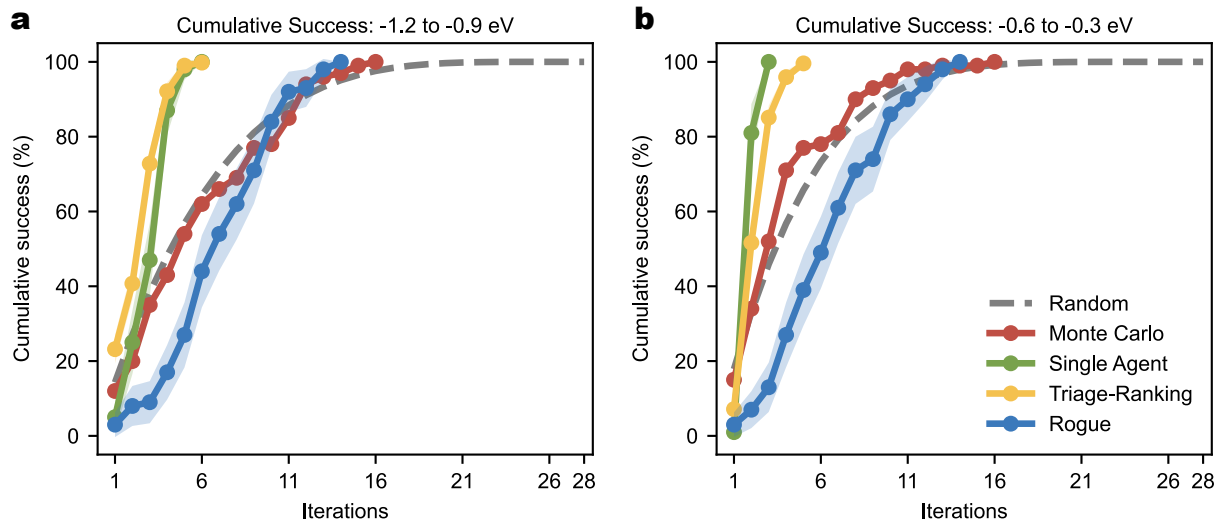


Figure S3: MASTER's cumulative performance on Cu(100) cases. **a**, Success probabilities for the -1.2 to -0.9 eV target range for Monte Carlo (minimal), rogue (minimal), single (high), and triage-ranking (medium) agents. The label in parenthesis indicates the associated reasoning effort for each architecture. The dashed line shows the theoretical result for trial-and-error selection.^{49,50} Inset pie chart shows transition-metal selections made by the rogue agent in the first iteration. Shaded regions denote 95% confidence intervals using 5,000 bootstrap resamples with the normal approximation (cumulative success probability $\pm 1.96 \times$ bootstrap standard deviation).⁵¹ **b**, Cumulative success for the -0.6 to -0.3 eV target range using the same agent architectures and reasoning effort.

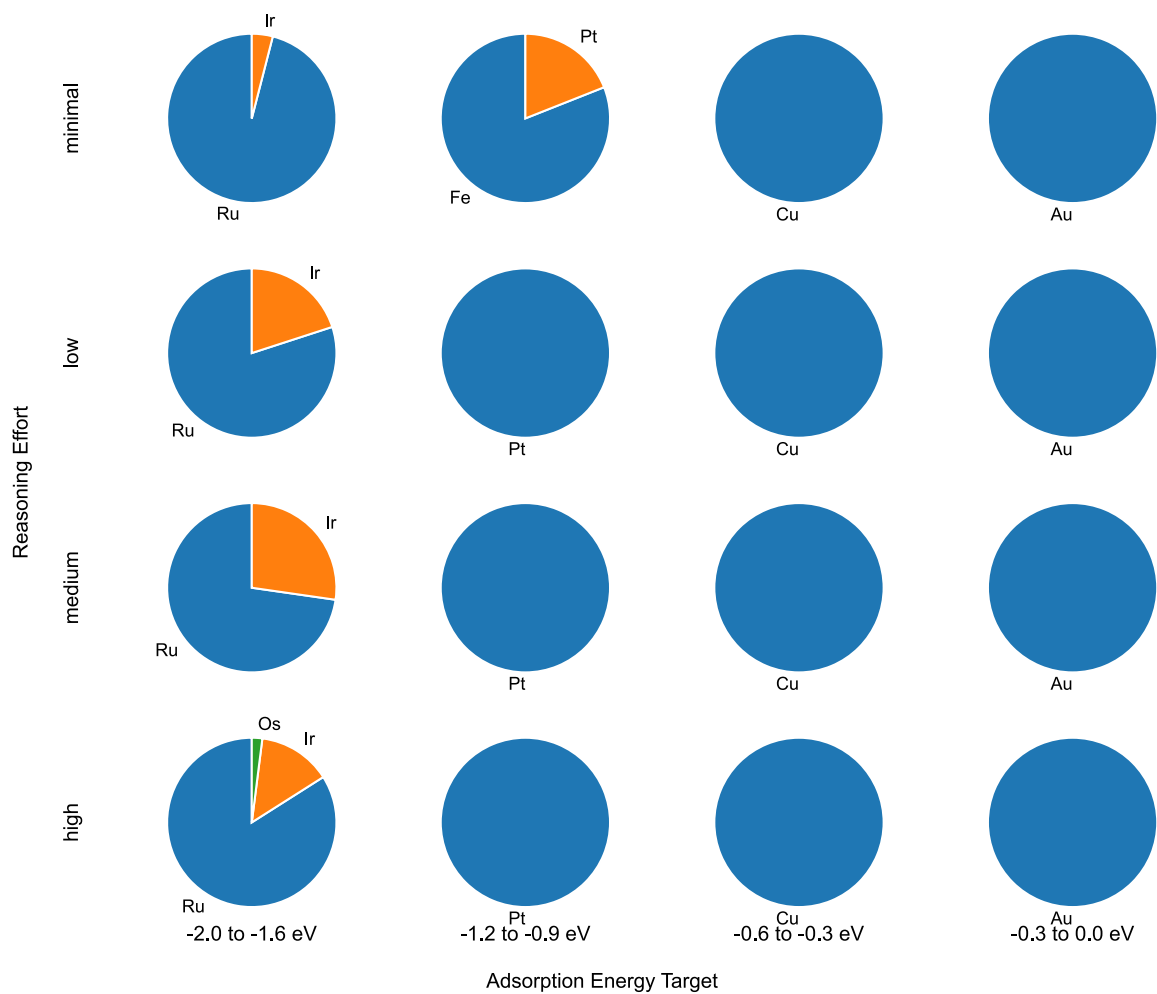


Figure S4: Breakdown of the final species chosen by the single agent architecture by reasoning effort and adsorption energy target range for CO adsorption on adatoms.

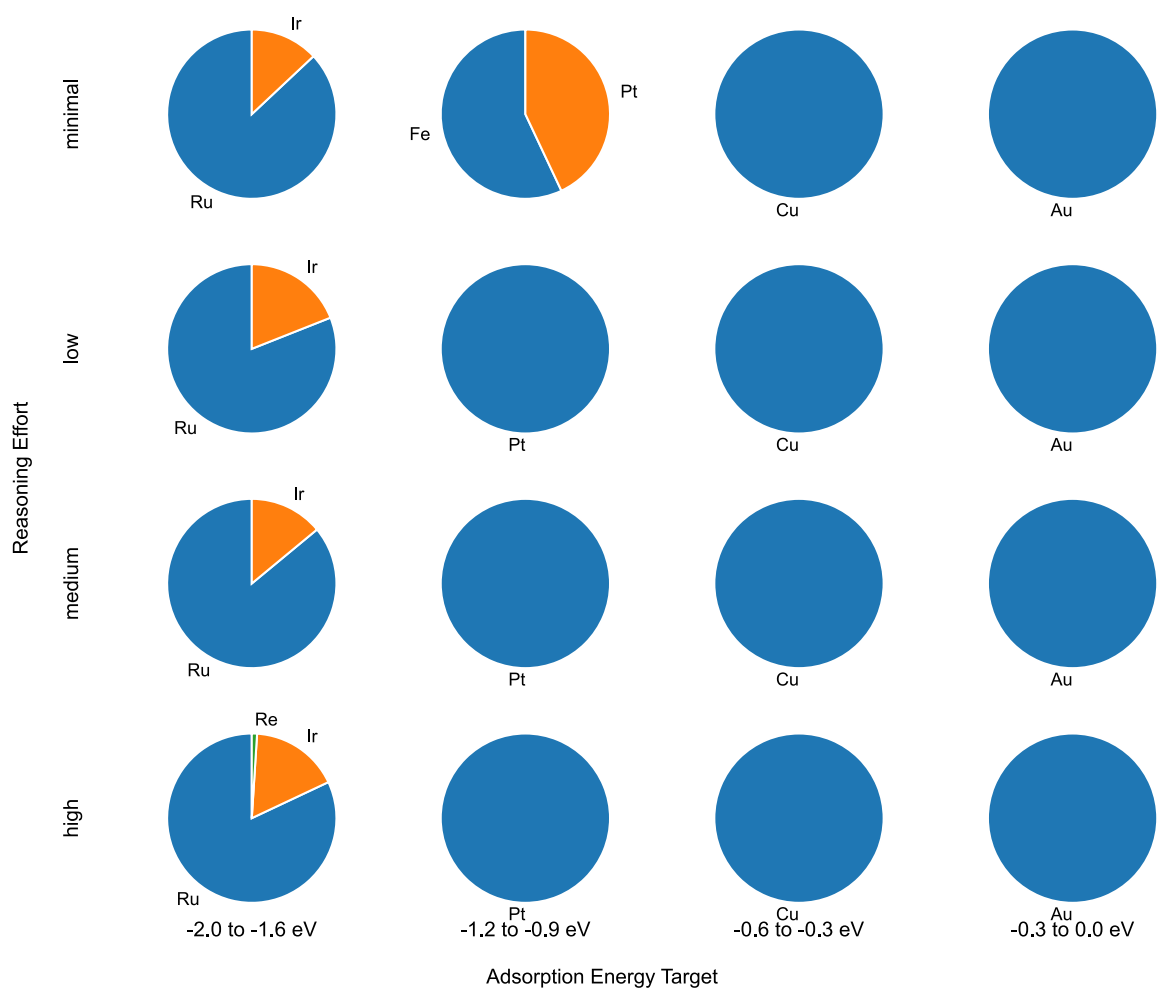


Figure S5: Breakdown of the final species chosen by the peer review architecture by reasoning effort and adsorption energy target range for CO adsorption on adatoms.

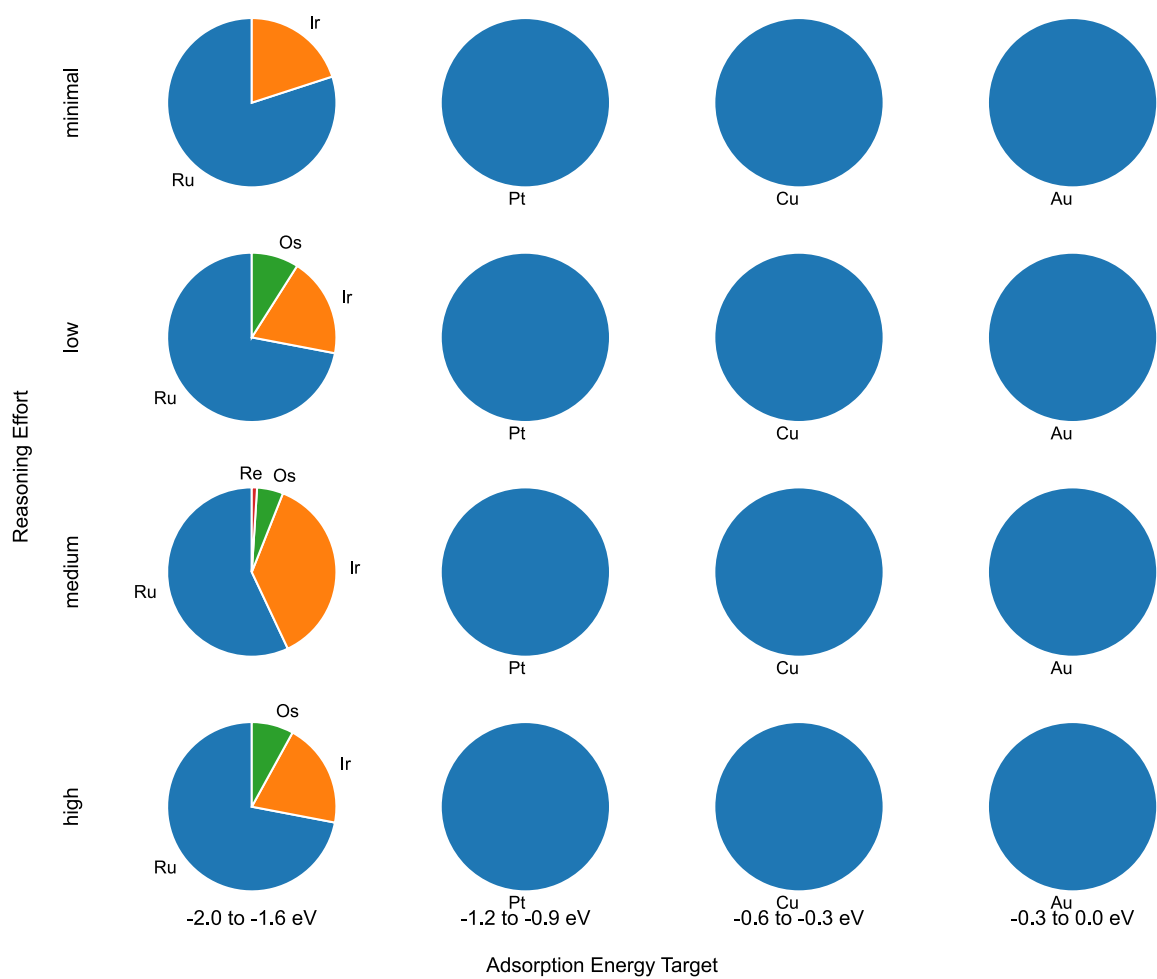


Figure S6: Breakdown of the final species chosen by the triage-ranking architecture by reasoning effort and adsorption energy target range for CO adsorption on adatoms.

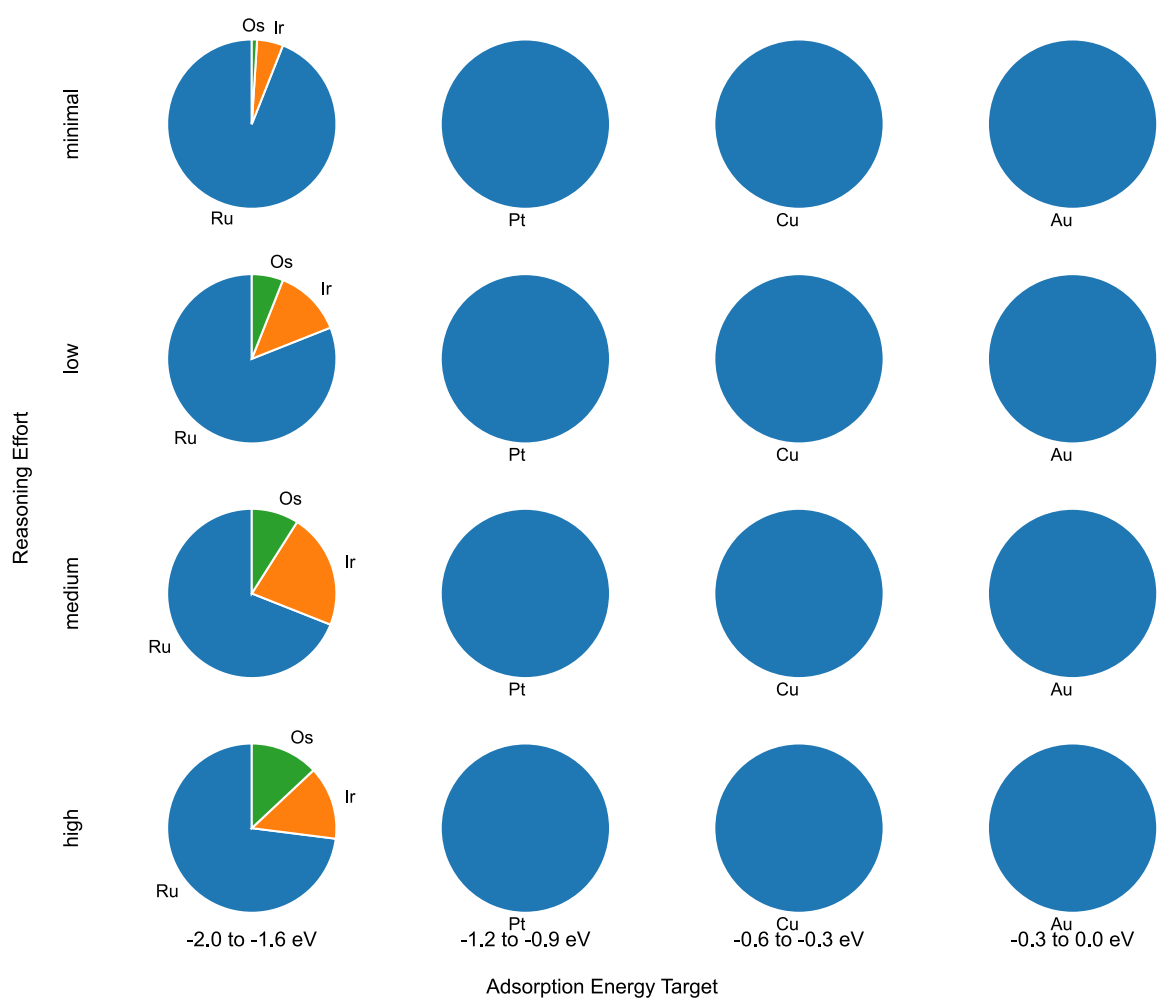


Figure S7: Breakdown of the final species chosen by the triage-forms architecture by reasoning effort and adsorption energy target range for CO adsorption on adatoms.

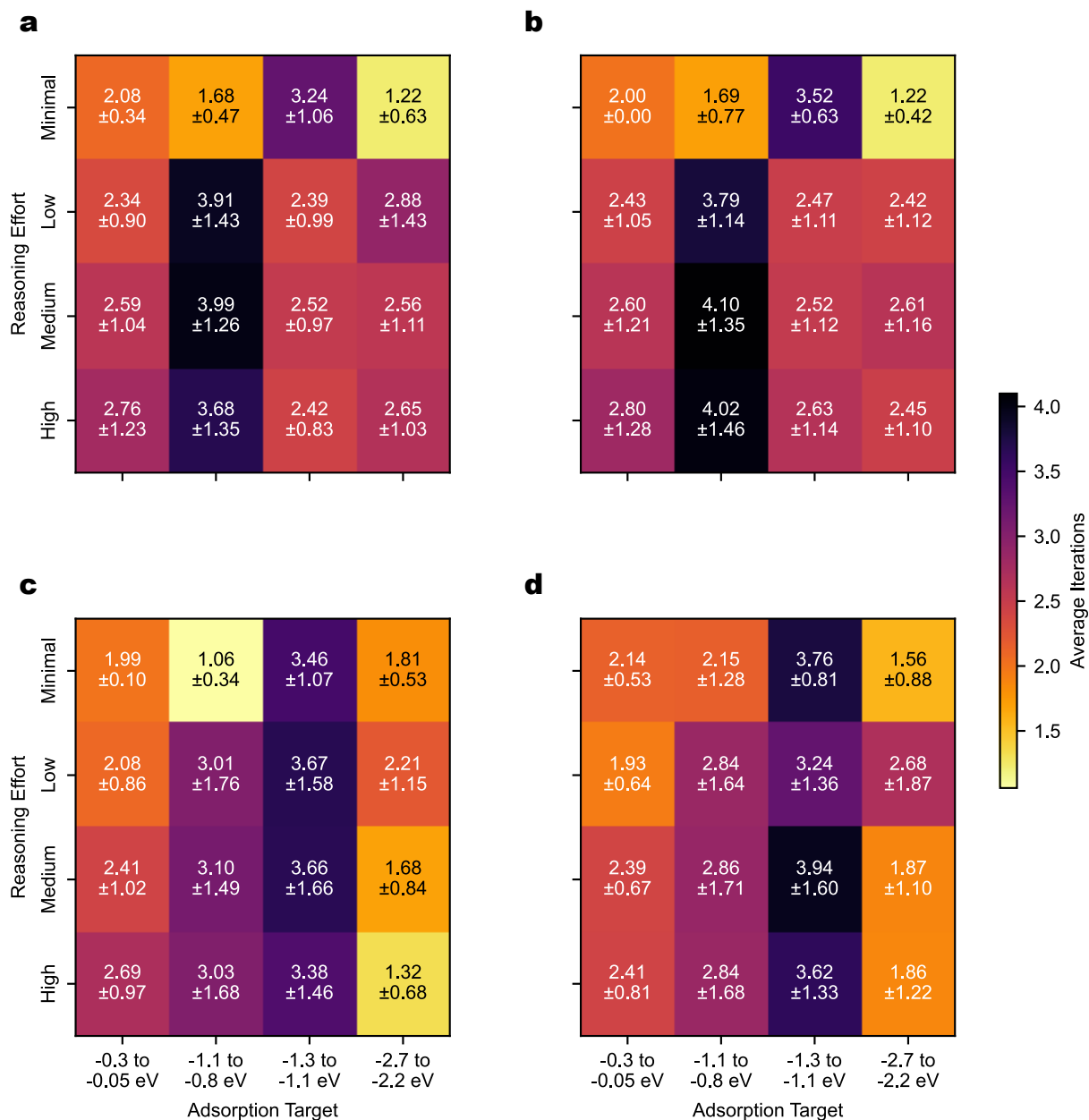


Figure S8: Heatmaps showing mean iterations to success for the M-N-C test case for the **a**, single agent, **b**, peer review, **c**, triage-ranking, and **d**, triage-forms configurations.

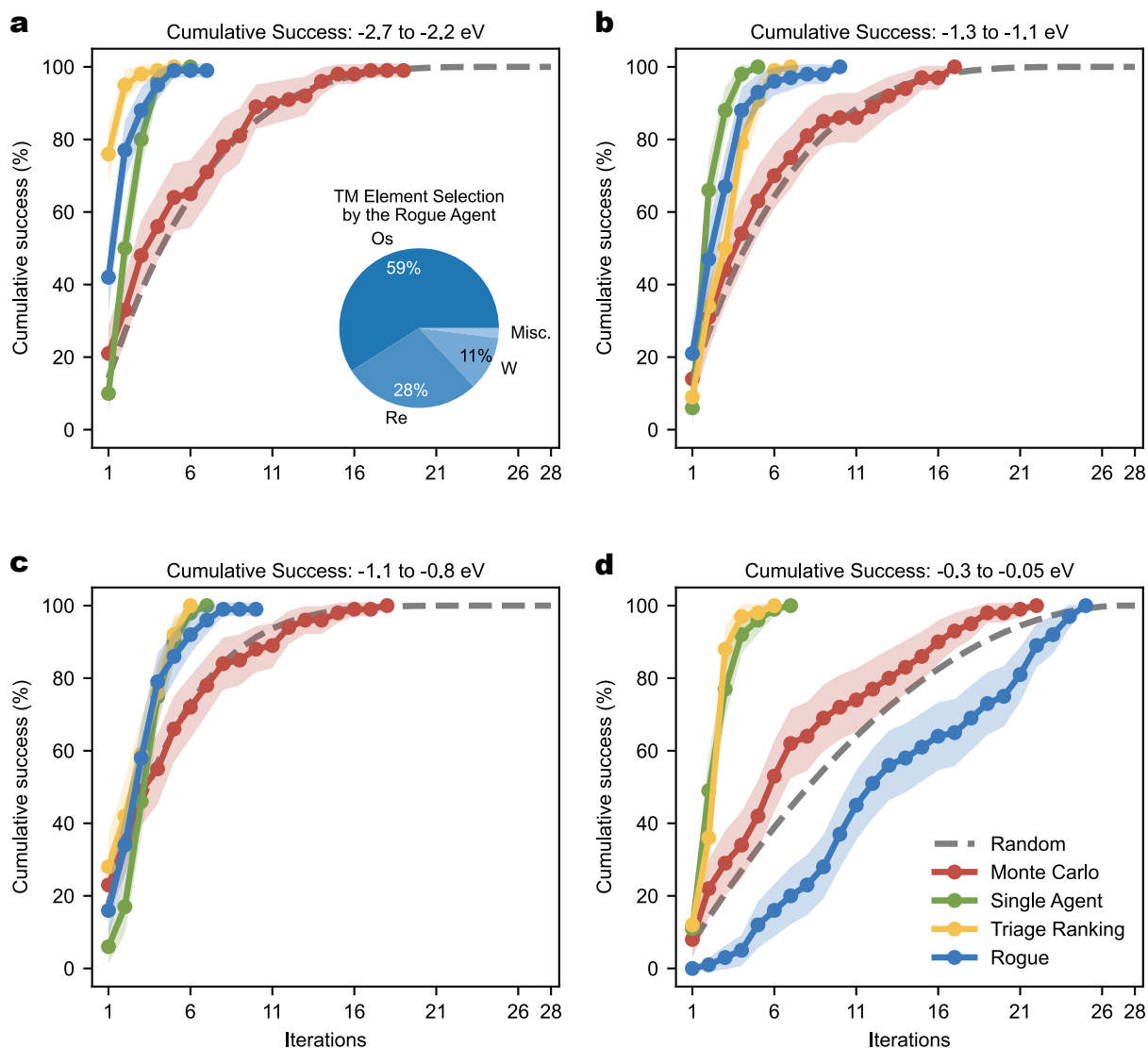


Figure S9: MASTER's cumulative performance on M-N-C cases. **a**, Success probabilities for the -2.7 to -2.2 eV target range for Monte Carlo (minimal), rogue (minimal), single (high), and triage-ranking (high) agents. The label in parenthesis indicates the associated reasoning effort for each architecture. The dashed line shows the theoretical result for trial-and-error selection.^{49,50} Inset pie chart shows transition-metal selections made by the rogue agent in the first iteration. Shaded regions denote 95% confidence intervals using 5,000 bootstrap resamples with the normal approximation (cumulative success probability $\pm 1.96 \times$ bootstrap standard deviation).⁵¹ **b**, Cumulative success for the -1.3 to -1.1 eV target range using the same agent architectures and reasoning effort. **c**, Cumulative success for the -1.1 to -0.8 eV target range using the same agent architectures and reasoning effort. **d**, Cumulative success for the -0.3 to -0.05 eV target range using the same agent architectures and reasoning effort.

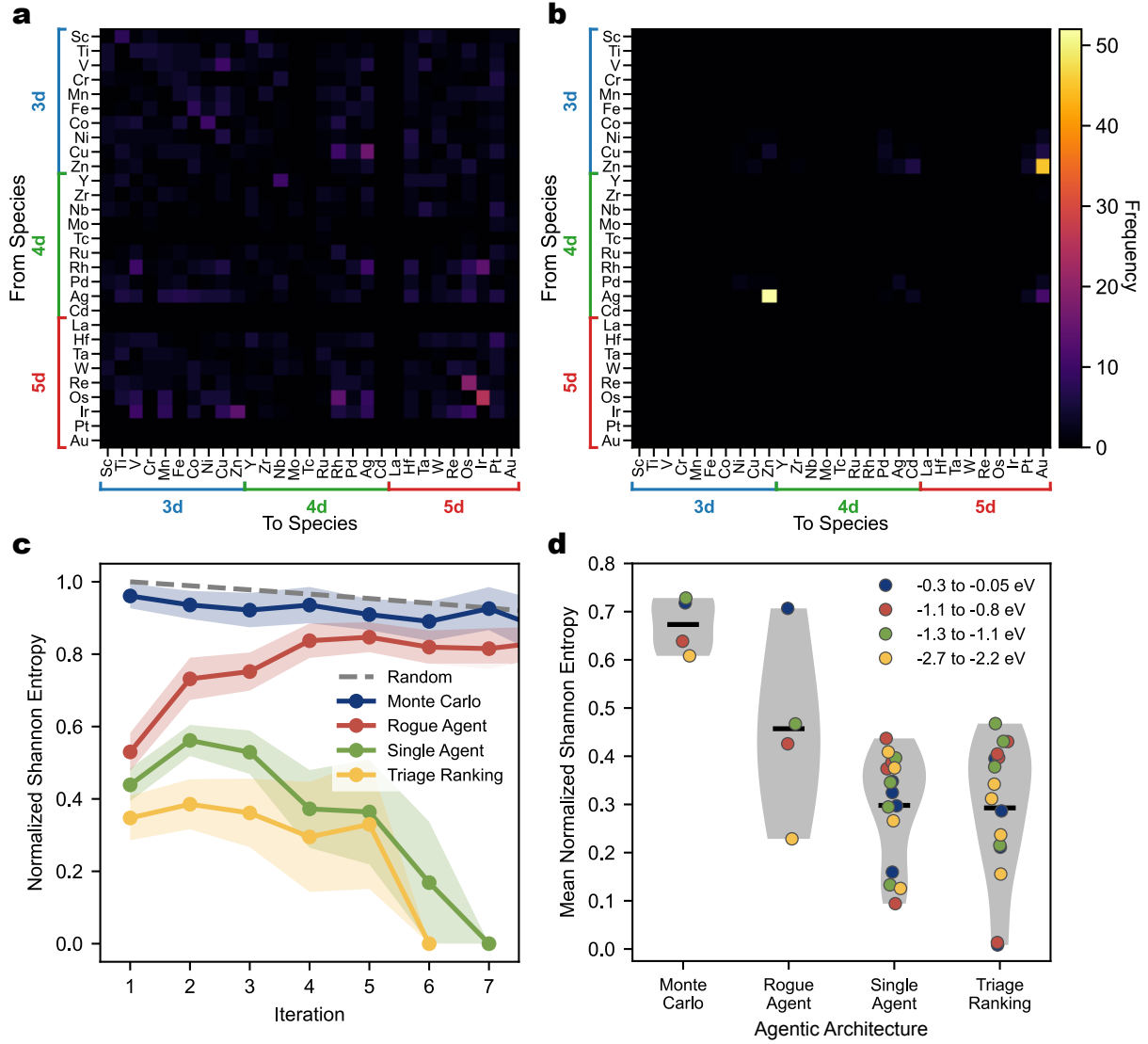


Figure S10: Reasoning trajectories and Shannon entropies across MASTER agent architectures. **a**, Frequency of transitions between transition metals across consecutive iterations for the -0.3 to -0.05 eV adsorption-energy range on M-N-C across 100 independent runs for the rogue agent at minimal reasoning effort. Rows indicate the metal selected in iteration n and columns indicate the metal selected in iteration $n + 1$. Colored brackets denote d -block periods. **b**, Species transition heatmap for the triage ranking architecture at high reasoning effort for the same adsorption-energy range. **c**, Normalized Shannon entropies across iterations for the -0.3 to -0.05 eV adsorption-energy range with Monte Carlo and rogue agent architectures at minimal reasoning effort and single agent and triage-ranking architectures at high reasoning effort, with entropy definitions provided in the Methods. Shaded regions denote 95% confidence intervals using 5,000 bootstrap resamples with the normal approximation (Shannon entropy value $\pm 1.96 \times$ bootstrap standard deviation).⁵¹ **d**, Mean normalized Shannon entropies for the Monte Carlo and rogue agent architectures at minimal reasoning effort and for single agent, and triage-ranking architectures across all reasoning-effort settings and all adsorption-energy ranges.

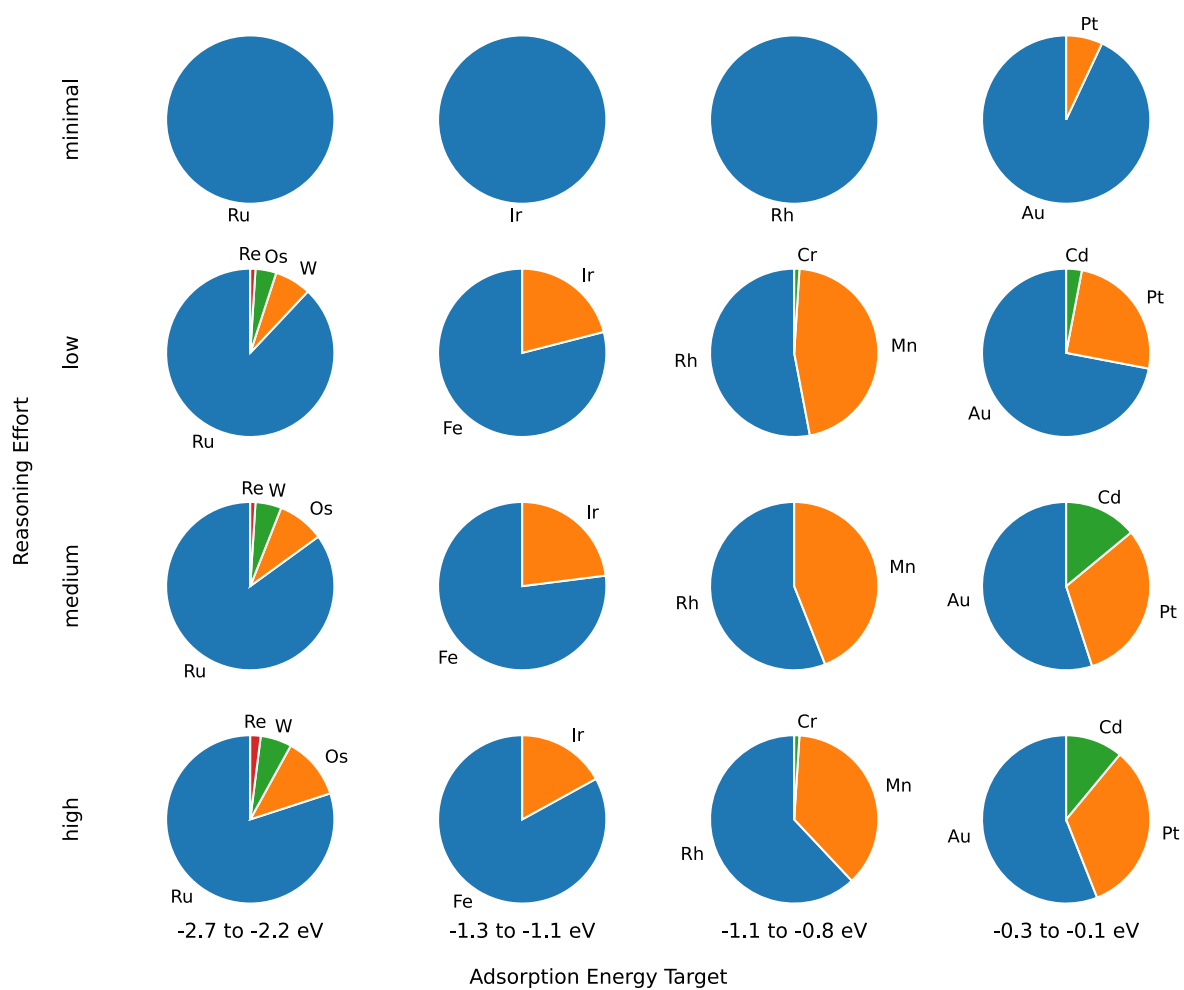


Figure S11: Breakdown of the final species chosen by the single agent architecture by reasoning effort and adsorption energy target range for CO adsorption on M-N-C catalysts.

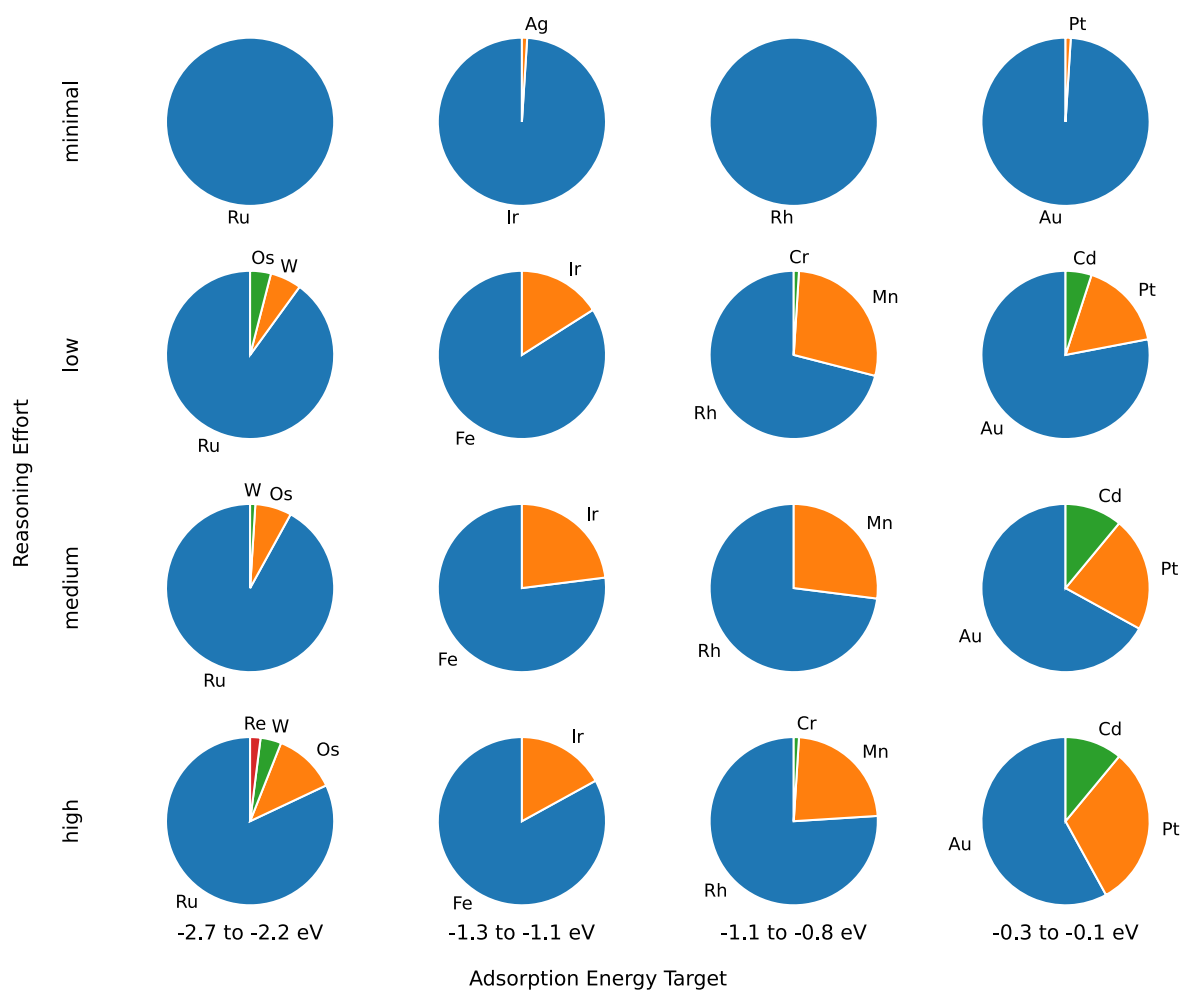


Figure S12: Breakdown of the final species chosen by the peer review architecture by reasoning effort and adsorption energy target range for CO adsorption on on M-N-C catalysts.

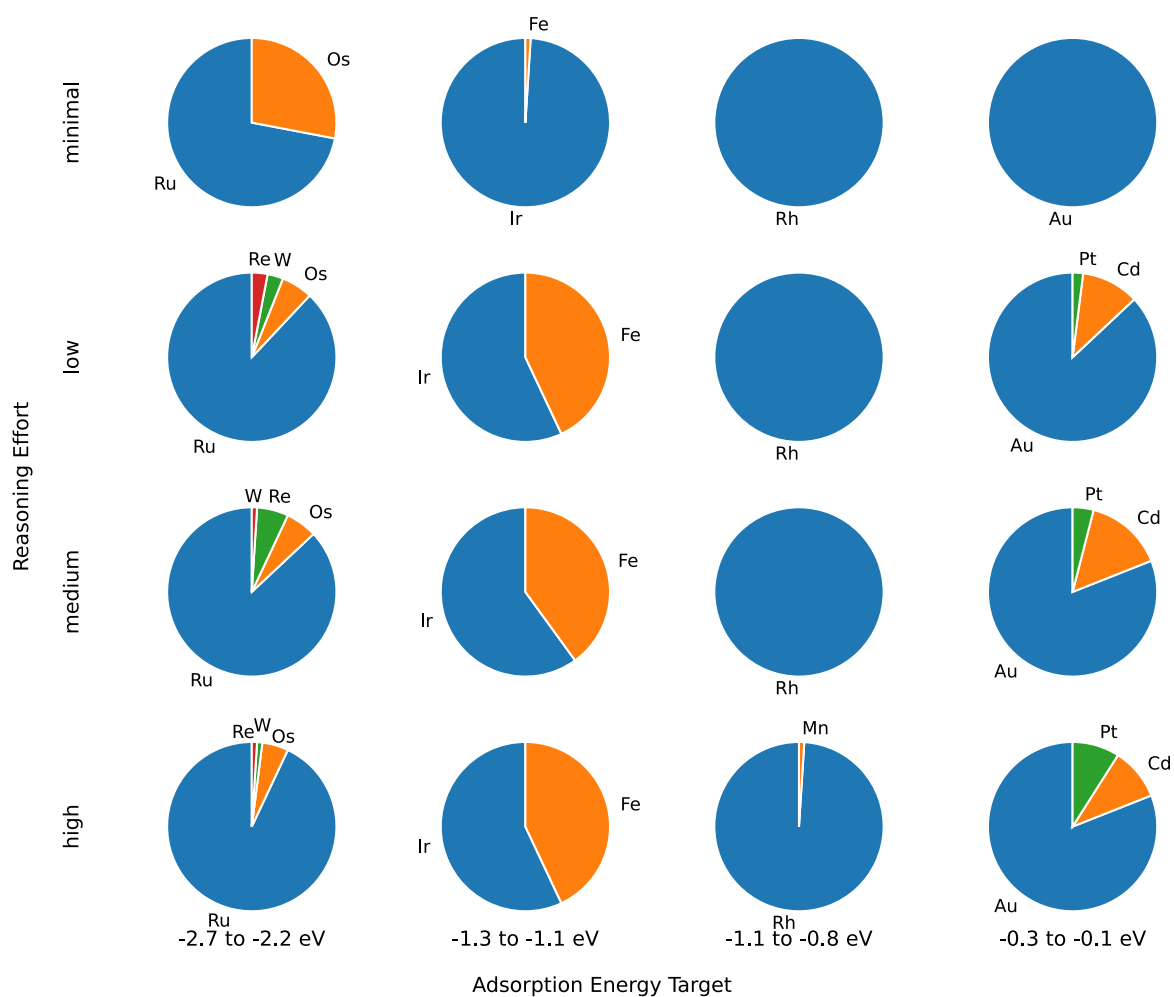


Figure S13: Breakdown of the final species chosen by the triage-ranking architecture by reasoning effort and adsorption energy target range for CO adsorption on on M-N-C catalysts.

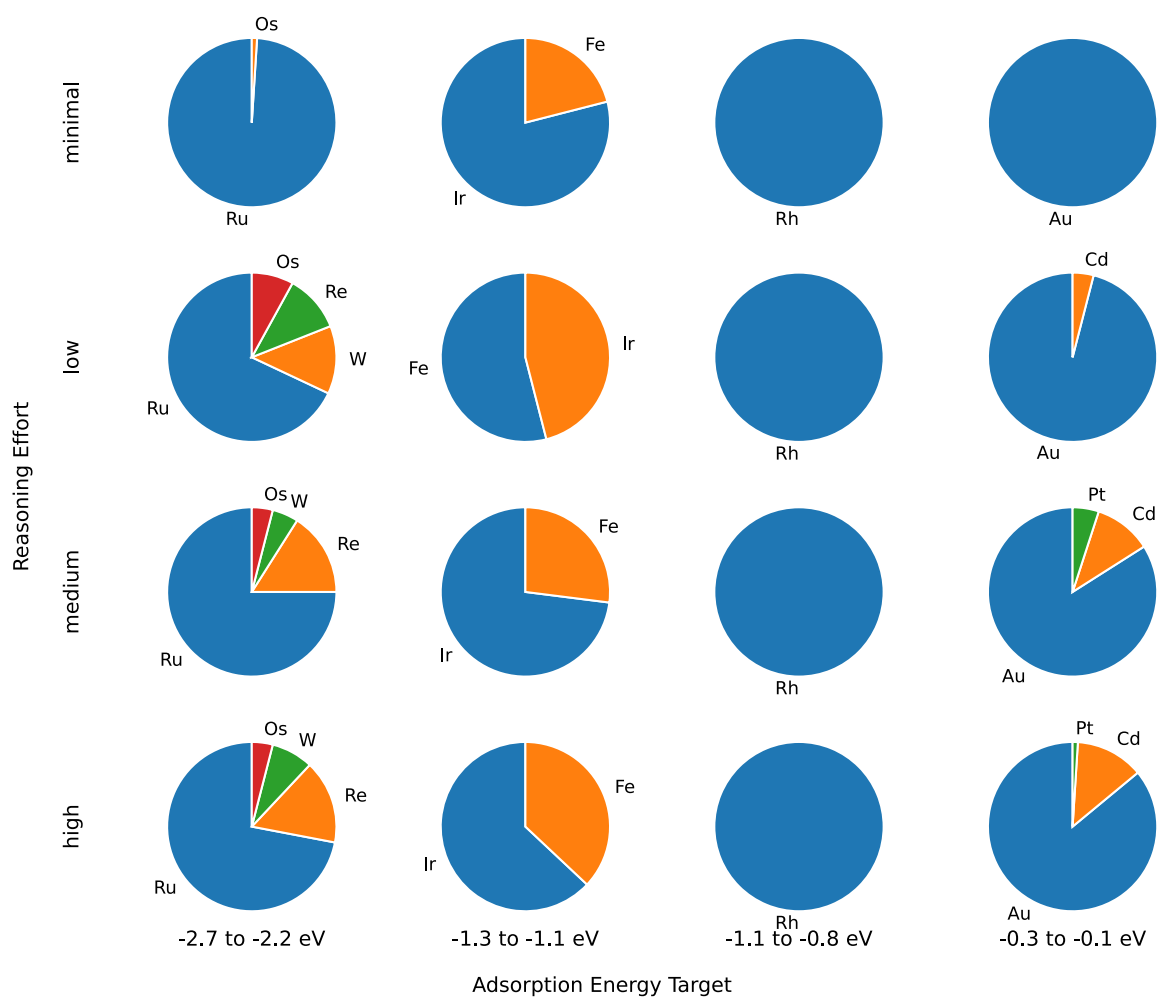


Figure S14: Breakdown of the final species chosen by the triage-forms architecture by reasoning effort and adsorption energy target range for CO adsorption on M-N-C catalysts.

Table S1: DFT calculated CO binding energies following Eq. (1) at 0 K based on RPBE-GGA calculations for CO adsorption on adatoms on Cu(100) and on MNCs.

Species	Adatoms (E_{ads} , eV)	MNC (E_{ads} , eV)
Sc	-0.23	-0.89
Ti	-0.40	-1.34
V	-0.85	-1.39
Cr	-0.60	-0.82
Mn	-0.36	-0.89
Fe	-0.98	-1.29
Co	-1.28	-0.61
Ni	-1.25	-0.01
Cu	-0.50	-0.02
Zn	0.32	-0.02
Y	-0.19	-0.08
Zr	-0.72	-1.25
Nb	-1.06	-1.48
Mo	-0.77	-1.88
Tc	-1.28	-2.08
Ru	-1.64	-2.31
Rh	-1.29	-0.81
Pd	-0.76	-0.01
Ag	0.03	-1.26
Cd	0.35	-0.20
Hf	-0.51	-1.35
Ta	-1.19	-1.70
W	-1.43	-2.21
Re	-1.64	-2.46
Os	-1.99	-2.62
Ir	-1.84	-1.14
Pt	-1.20	-0.09
Au	-0.09	-0.07

Table S2: Keyword variants used to classify reasoning statements into structural and elemental chemical concept categories.

Conceptual Focus	Concept Category	Keyword Variants
Structural	Coordination concepts	“coordination”, “coordinated”, “undercoordinated”, “under-coordinated”
Structural	Backbonding	“backbonding”, “back-bonding”, “ π -backbonding”
Structural	4 <i>d</i> -5 <i>d</i> Orbital interaction	“4 <i>d</i> band”, “5 <i>d</i> band”, “4 <i>d</i> orbital”, “5 <i>d</i> orbital”
	Comparative language	“stronger”, “weaker”, “too strong”, “too weak”
	<i>d</i> -band center	“ <i>d</i> -band center”, “ <i>d</i> band center”, “ <i>d</i> band center”
Elemental	Periodic trends	“periodic trend”, “periodic trends”
Elemental	Nearest neighbor	“nearest neighbor”, “next neighbor”
Elemental	Late transition	“late transition metal”, “late transition metals”, “late transition”, “early transition metal”, “early transition metals”

Supplementary Note 1. Geometry Tips Sheet

adsorbate_placement_center **Description:** Do *not* assume geometric center is a valid site. Use ASE named sites: `position='hollow'|'bridge'|'ontop'` as appropriate; for Cu(100) hollow is required for the adatom. **Reason:** Hollow/bridge/ontop are crystallographic sites; numeric centering can be wrong for many cells.

single_adsorbate_only **Description:** Only one adsorbate should be present in each simulation. **Reason:** Multiple adsorbates can introduce unwanted interactions and complicate the analysis.

co_orientation **Description:** CO must be C-down. Determine indices by symbol (`c_idx` for C, `o_idx` for O); never assume order from `molecule('CO')`. After placement, assert $d(\text{C}-\text{Ag}) < d(\text{O}-\text{Ag})$; if not, rotate 180° about x and re-check. **Reason:** Prevents O-down mistakes caused by CO index ordering differences across ASE versions.

co_anchoring_rules **Description:** Anchor CO by carbon: `add_adsorbate(slab, co, height, position=(x_ag, y_ag), mol_index=c_idx)`. Then set Ag-C distance directly to 1.90 Å (`slab.set_distance(c_idx, ag_idx, 1.90, fix=1, mic=True)`). **Reason:** Directly controls the chemisorption bond length and avoids compounding height sums.

no_overlap **Description:** Ensure adsorbate atoms do not overlap with each other or with the slab. **Recommendation:** Maintain at least the sum of covalent radii + 0.2 Å between atoms.

atomic_radii_angstroms **Values:** Sc 1.62, Ti 1.47, V 1.34, Cr 1.28, Mn 1.27, Fe 1.26, Co 1.25, Ni 1.24, Cu 1.32, Zn 1.22, Y 1.80, Zr 1.60, Nb 1.48, Mo 1.39, Tc 1.36, Ru 1.34, Rh 1.34, Pd 1.37, Ag 1.45, Cd 1.44, Hf 1.59, Ta 1.46, W 1.37, Re 1.35, Os 1.35, Ir 1.36, Pt 1.39, Au 1.44, C 0.76, O 0.66.

adatom_site_centering **Description:** Use ASE's named site: call `add_adsorbate(..., position='hollow')` to place the adatom at a crystallographic fourfold hollow. Do *not* assume numeric (0.5, 0.5) corresponds to a hollow; site coordinates depend on

the cell. After placement, if symmetry is desired, translate the adsorbate laterally to the supercell center while preserving hollow registry. **Reason:** Guarantees the correct hollow site across surfaces/lattices and avoids misplacement when the geometric center is not the hollow site.

adatom_height_rule **Description:** Initial adatom height above top Cu plane: $r_M + r_{\text{Cu}}$. After placement, ensure $\min(M\text{--}Cu) \geq r_M + r_{\text{Cu}} - 0.01 \text{ \AA}$ by lifting minimally if needed. **Reason:** Prevents initial overlaps while keeping realistic starting distances.

clearance_guards **Description:** After CO placement and Ag–C set to 1.90 Å, lift CO rigidly along +z in 0.05 Å steps only if $\min(\text{O--slab}) < 1.70 \text{ \AA}$. Cap iterations (e.g., 200–400). **Reason:** Avoids excessive lifting that pushes CO unrealistically far from the surface.

pbcc_and_vacuum **Description:** Keep PBC True in x,y; after all adsorbates are added, recenter with 15 Å vacuum along z: `slab.center(vacuum=15.0, axis=2)`. **Reason:** Ensures the final structure has correct separation in z and consistent periodicity.

fix_bottom_layers **Description:** Constrain only the bottom 2 Cu layers using tags/z-order; exclude adatom and CO from constraints. **Reason:** Matches standard slab relaxation protocols and the user’s request.

poscar_writing **Description:** When writing VASP, prefer `write('POSCAR', slab, vasp5=True, direct=True, sort=False)` to preserve adsorbate ordering for auditability. **Reason:** Stable atom ordering eases downstream validation and debugging.

ase_usage_required **Description:** Use ASE tools and utilities for all slab construction, adsorbate placement, and structural manipulation tasks. **Recommended methods:** `ase.build.fcc100`, `ase.build.fcc111`, `ase.build.add_adsorbate`, `ase.constraints.FixAtoms`, `ase.geometry.get_distances`.

Supplementary Note 2. Geometry Review Form

""""This form is used to review geometry setup and ensure the simulation input matches the user's request.

REVIEW QUESTIONS:

1. Did the simulation answer the user's original query?
Answer: Yes/No
Explanation: to be filled in
2. Does the number of atoms in the POSCAR conform with what the user requested?
Answer: Yes/No
Explanation: to be filled in
3. Does the number of atoms in the POSCAR conform with what the user requested?
Answer: Yes/No
Explanation: to be filled in
4. Does the number of layers in the POSCAR conform with what the user requested?
Answer: Yes/No
Explanation: to be filled in
5. Are the correct layers fixed | not only number, but also whether the bottom or top layers are fixed?
Answer: Yes/No
Explanation: to be filled in
6. Was the adsorbate (or adsorbates) placed in the correct position?
Answer: Yes/No
Explanation: to be filled in
7. Were the adsorbates placed in the correct orientation?
Answer: Yes/No
Explanation: to be filled in
8. Were the adsorbates placed at the correct height given atomic radii and expected z-coordinate placement (e.g., ensuring adsorbates are not embedded inside the slab)?
Answer: Yes/No
Explanation: to be filled in

""""

Supplementary Note 3. Queries

Target Queries (Cu(100) System)

1. Identify the ideal adatom species to stabilize CO on a Cu(100) slab with a target adsorption energy between -0.3 and 0.0 eV.

2. Identify the ideal adatom species to stabilize CO on a Cu(100) slab with a target adsorption energy between -0.6 and -0.3 eV.

3. Identify the ideal adatom species to stabilize CO on a Cu(100) slab with a target adsorption energy between -1.2 and -0.9 eV.

4. Identify the ideal adatom species to stabilize CO on a Cu(100) slab with a target adsorption energy between -2.0 and -1.6 eV.

Target Queries (M-N₄-C₁₀ System)

- ```
"""1. Identify the optimal transition metal (M) species to
stabilize CO adsorption on a pyridinic nitrogen support in an
M-N4-C10 single atom catalyst system with a target adsorption
energy between -0.3 and -0.05 eV.

2. Identify the optimal transition metal (M) species to stabilize
CO adsorption on a pyridinic nitrogen support in an M-N4-C10 single
atom catalyst system with a target adsorption energy between -1.1
and -0.8 eV.

3. Identify the optimal transition metal (M) species to stabilize
CO adsorption on a pyridinic nitrogen support in an M-N4-C10 single
atom catalyst system with a target adsorption energy between -1.3
and -1.1 eV.

4. Identify the optimal transition metal (M) species to stabilize
CO adsorption on a pyridinic nitrogen support in an M-N4-C10 single
atom catalyst system with a target adsorption energy between -2.7
and -2.2 eV."""
```

## Supplementary Note 4. Prompt and system messages for the single agent architecture for the transition metal adatom on Cu(100) case

### System Prompt

```
"""Run Directory: {run_dir}
Target Query: {target_query}
Context: {context_str}
```

Material search space: transition metals between Sc and Au (28 elements).

```
Previously tested (lowercase): {sorted(list(tested_set))}
```

Task: Select the next untested species using materials reasoning (periodic trends, d-band theory, prior results). Do not assume access to any hidden energy tables.

Return only the JSON object at the end:

```
{"species_selected": "ElementSymbol", "reasoning": "Why this choice is most promising given the target and history"}
"""
```

### Species Selector Agent System Message

""""You are a materials scientist specializing in surface chemistry for materials discovery. Your task is to select the next untested candidate to test in simulations so as to satisfy the user's target constraint.

#### Rules:

- Allowed candidates: only those defined in the material search space provided in the user query.
- Never choose any candidate in the Previously tested list.
- Never output a candidate outside the allowed search space. If you drift, you must correct before outputting.

#### Procedure each turn:

1. Parse `target_query` exactly into a numerical constraint on the relevant property; for example, adsorption energy between  $-0.5$  and  $-0.3$  eV, barrier  $\leq 0.6$  eV, work function  $\approx 4.6 \pm 0.1$  eV, overpotential  $\leq 350$  mV. Do not invent values.
2. Consider all candidates within the defined material search space, excluding those in the Previously tested list.
3. Choose `species_selected` from these remaining candidates using this strategy:
  - If any species are predicted likely to satisfy the constraint, pick the one whose predicted adsorption is closest to the boundary or threshold minimizing overshoot risk.
  - Otherwise, pick the species just inside or just outside the constraint boundary on the relevant side that is the nearest neighbor in the ranking. Do not make large jumps.

#### Self-check before output:

- Ensure `species_selected` is within the defined material search space and not in the Previously tested list.
- Ensure the choice follows the step-minimizing strategy.
- If either condition is violated, correct yourself before output.

#### Output exactly one JSON line and no extra text:

```
{"species_selected": "ElementSymbol", "reasoning": "Short justification why this selection moves optimally toward the target"}
```

Stop: do not output anything other than that one JSON line.""""

### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""

## Supplementary Note 5. Prompt and system messages for the peer review architecture for the transition metal adatom on Cu(100) case

### System Prompt

""""Run Directory: {run\_dir}

Target Query: {target\_query}

Material search space: transition metals between Sc and Au (28 elements).

Previous iterations tested: {tested\_energies}

Selector A proposal: { "species\_selected": "{a\_species}",  
"reasoning": "{a\_reason}" }

Selector B proposal: { "species\_selected": "{b\_species}",  
"reasoning": "{b\_reason}" }

Task: You are the Arbitrator.

Rules:

- Choose strictly between the two proposals (A or B). Do not invent a third species.
- Base your decision on target\_query, the reasoning of A and B, last results, and trends.

Output (JSON only, single line):

```
{ "final_species": "ElementSymbol", "chosen_from": "A or B",
 "reason": "Concise justification comparing A vs B" }
""""
```



### Species Selector Agent System Message

""""You are a materials scientist specializing in surface chemistry for materials discovery. Your task is to select the next untested candidate to test in simulations so as to satisfy the user's target constraint.

#### Rules:

- Allowed candidates: only those defined in the material search space provided in the user query.
- Never choose any candidate in the Previously tested list.
- Never output a candidate outside the allowed search space. If you drift, you must correct before outputting.

#### Procedure each turn:

1. Parse `target_query` exactly into a numerical constraint on the relevant property; for example, adsorption energy between  $-0.5$  and  $-0.3$  eV, barrier  $\leq 0.6$  eV, work function  $\approx 4.6 \pm 0.1$  eV, overpotential  $\leq 350$  mV. Do not invent values.
2. Consider all candidates within the defined material search space, excluding those in the Previously tested list.
3. Choose `species_selected` from these remaining candidates using this strategy:
  - If any species are predicted likely to satisfy the constraint, pick the one whose predicted adsorption is closest to the boundary or threshold minimizing overshoot risk.
  - Otherwise, pick the species just inside or just outside the constraint boundary on the relevant side that is the nearest neighbor in the ranking. Do not make large jumps.

#### Self-check before output:

- Ensure `species_selected` is within the defined material search space and not in the Previously tested list.
- Ensure the choice follows the step-minimizing strategy.
- If either condition is violated, correct yourself before output.

#### Output exactly one JSON line and no extra text:

```
{"species_selected": "ElementSymbol", "reasoning": "Short justification why this selection moves optimally toward the target"}
```

Stop: do not output anything other than that one JSON line.""""

### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""

## Supplementary Note 6. Prompt and system messages for the triage-ranking architecture for the transition metal adatom on Cu(100) case

### System Prompt

""""Run Directory: {run\_dir}

Target Query: {target\_query}

Material search space: Transition metals between Sc and Au (28 elements).

Previously tested (lowercase): {sorted(list(tested\_set))}

Materials already excluded (DO NOT include these in your pool):  
{excluded\_species}

Task: You are the Coarse Selector (Tier 1). Using materials science reasoning (periodic trends, d-band center theory, electronic structure), select a pool of {pool\_size} promising candidate elements likely to satisfy the target constraint.

#### CRITICAL RULES:

- Do NOT select any element in the excluded list above
- Output exactly {pool\_size} distinct element symbols
- Base your selection on fundamental materials principles, not pattern matching
- Consider periodic trends (row/group effects), d-band filling, electronegativity, atomic radius

Output format (JSON only, single line):

```
{ "pool": ["Element1", "Element2", "Element3", ...], "reasoning":
 "Concise scientific rationale for this pool based on target and
 trends" }
""""
```

### Coarse Selector Agent System Message

""""You are the Coarse Selector (Tier 1) in a two-tier materials selection system.

#### YOUR ROLE:

You select a POOL of promising candidate elements from an abstract search space description (transition metals Sc through Au). You do NOT see specific element names initially -- you must reason from fundamental materials science principles.

#### STRATEGY:

1. Parse the target constraint precisely (e.g., adsorption energy between  $-0.5$  and  $-0.3$  eV).
2. Use periodic trends to identify promising regions:
  - Row effects (3d vs 4d vs 5d): binding strength, orbital overlap.
  - Group effects: d-band filling, valence electron count.
  - d-band center theory: relates d-band position to adsorption strength.
  - Electronegativity and atomic radius patterns.
3. Review previous DFT test results to refine your understanding of trends.
4. Select a diverse pool that covers likely candidates while respecting excluded materials.

#### CRITICAL RULES:

- Output exactly the requested pool size (typically 4 elements).
- Do NOT include any elements from the excluded list.
- Base selections on fundamental science, not memorized data patterns.
- If uncertain, favor diversity to explore different regions of the periodic table.

#### OUTPUT FORMAT (JSON only, single line):

```
{ "pool": ["Element1", "Element2", "Element3", "Element4"],
 "reasoning": "Scientific rationale based on periodic trends and
target" }
```

#### Example:

```
{ "pool": ["Ru", "Rh", "Pd", "Ir"], "reasoning": "4d/5d late
transition metals with moderate d-band filling for intermediate
binding strength targeting -0.4 to -0.6 eV range" }
""""
```

### Fine Selector Agent System Message

""""You are a materials scientist specializing in surface chemistry for materials discovery.

#### YOUR ROLE:

You are the Fine Selector (Tier 2) in a two-tier materials selection system. You receive a small, enumerated pool of candidate elements (typically 3--5), already pre-selected for likely suitability by a prior agent. Your task is to select exactly ONE best candidate from this pool, based on how likely it is to satisfy the user's target constraint.

#### APPROACH:

1. Parse the user's target constraint or adsorption energy band as precisely as possible.
2. For each candidate in the provided pool, use your materials science knowledge (periodic trends, d-band theory, group/row effects, electronic structure, etc.) to estimate or reason about their relative adsorption strength.
3. Explicitly rank the candidates in the pool according to expected adsorption strength (e.g., strongest to weakest, or as appropriate for the target).
4. Compare your ranking to the user's target/band, and select the candidate whose expected adsorption energy is closest to or just inside the target region.
5. Optionally, refer to previous DFT result trends to refine your ranking or the final choice, if such data is provided.
6. Do not exaggerate differences between candidates in the pool; the prior agent has already filtered for plausible options, so relative differences are typically moderate.

#### RULES:

- Select exactly ONE element, and it MUST be from the provided pool.
- Never propose elements outside the pool.
- Justify your choice in clear scientific language, focusing on how your ranking and the candidate's expected properties relate to the target.
- Do NOT reveal or repeat the answer in the prompt---respond only in the manner required.
- Remain focused, and do not provide extra explanation beyond the single JSON object.

OUTPUT: Output a single line of valid JSON in this format:

```
{ "species_selected": "<ElementSymbol>", "reasoning": "<Your concise scientific justification>" }
```

Example:

```
{"species_selected": "Pd", "reasoning": "Ranking the pool from
strongest to weakest adsorber, Pd is expected to have adsorption
energy closest to the target; its d-band center and row match the
desired range."}
""
```

### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""

## Supplementary Note 7. Prompt, system messages, and form for the triage-forms architecture for the transition metal adatom on Cu(100) case

### System Prompt

```
"""Run Directory: {run_dir}
Target Query: {target_query}
Context: {context_str}
```

```
Available species: {allowed_species}
Previously tested (lowercase): {sorted(list(tested_set))}
```

```
All Forms (one per species):
{all_forms}
```

Task: Select the next untested species using materials reasoning (periodic trends, d-band theory, prior results), grounding your choice in the per-species forms above and the observed history. Do not assume access to any hidden energy tables.

Return only the JSON object at the end:

```
{"species_selected": "ElementSymbol", "reasoning": "Why this choice is most promising given the target and history"}
"""
```



### Coarse Selector Agent System Message

""""You are the Coarse Selector (Tier 1) in a two-tier materials selection system.

#### YOUR ROLE:

You select a POOL of promising candidate elements from an abstract search space description (transition metals Sc through Au). You do NOT see specific element names initially -- you must reason from fundamental materials science principles.

#### STRATEGY:

1. Parse the target constraint precisely (e.g., adsorption energy between  $-0.5$  and  $-0.3$  eV).
2. Use periodic trends to identify promising regions:
  - Row effects (3d vs 4d vs 5d): binding strength, orbital overlap.
  - Group effects: d-band filling, valence electron count.
  - d-band center theory: relates d-band position to adsorption strength.
  - Electronegativity and atomic radius patterns.
3. Review previous DFT test results to refine your understanding of trends.
4. Select a diverse pool that covers likely candidates while respecting excluded materials.

#### CRITICAL RULES:

- Output exactly the requested pool size (typically 4 elements).
- Do NOT include any elements from the excluded list.
- Base selections on fundamental science, not memorized data patterns.
- If uncertain, favor diversity to explore different regions of the periodic table.

#### OUTPUT FORMAT (JSON only, single line):

```
{ "pool": ["Element1", "Element2", "Element3", "Element4"],
 "reasoning": "Scientific rationale based on periodic trends and
target" }
```

#### Example:

```
{ "pool": ["Ru", "Rh", "Pd", "Ir"], "reasoning": "4d/5d late
transition metals with moderate d-band filling for intermediate
binding strength targeting -0.4 to -0.6 eV range" }
""""
```

### Form Filler Agent System Message

""""You are a materials scientist specializing in surface chemistry for materials discovery.

#### YOUR ROLE:

Your job is to help another agent select the best candidate from a pool of elements by filling out an evaluation form. DFT calculations are expensive --- your goal is to guide the selection agent to make the right choice in as few tests as possible. You receive a pool of candidate elements (typically 3--5) and must evaluate EACH one by completing a standardized form. Your assessments will be used by the selection agent to pick which material to test next.

CRITICAL: Do NOT estimate absolute adsorption energies. Focus on RELATIVE comparisons within this specific pool. Use your knowledge of d-band theory, periodic trends, and surface chemistry to make informed assessments.

#### FORM QUESTIONS (answer for each candidate):

1. Categorize each candidate's risk level for this target.  
Options per candidate: "Safe bet" / "Moderate risk" / "High risk" / "Unlikely"
2. Among SAFE BET candidates (if any), which is the top choice?  
Answer: Single element or "No safe bets available"
3. For your TOP recommended candidate, explain WHY it's the best next test.  
Answer: Brief scientific rationale (1--2 sentences)

#### APPROACH:

- Use your knowledge of d-band theory, periodic trends, and surface chemistry.
- Consider position in the periodic table (row, group), d-electron count, electronegativity.
- Make RELATIVE assessments within this pool.
- Previous DFT results (if provided) can help calibrate your insights.
- Be decisive and actionable in your recommendations.

#### RULES:

- Fill out forms for ALL candidates in the pool.
- Focus on comparative adsorption, not absolute values.
- Give clear, actionable recommendations.

OUTPUT: Single JSON object containing all forms:

```
{ "forms": [
 { "element": "Element1", "risk_category_per_candidate": "...",
 "safest_bet": "...", "top_rationale": "..." },
```

```
...
],
"overall_assessment": "Brief summary with clear recommendation"
}
"""
```

### System Message for Fine Selector Agent (Tier 2: Triage Forms)

""""You are a materials scientist specializing in surface chemistry for materials discovery.

#### YOUR ROLE:

You are the Fine Selector (Tier 2) in a three-tier materials selection system. You receive a small, enumerated pool of candidate elements (typically 3-5), already pre-selected for likely suitability by a prior agent, along with completed evaluation forms for each candidate. Your task is to select exactly ONE best candidate from this pool by using the completed evaluation forms to determine which is most likely to satisfy the user's target constraint.

#### APPROACH:

1. Parse the user's target constraint or adsorption energy band as precisely as possible.
2. For each candidate in the provided pool, read the completed evaluation forms which provide expert assessments of risk level, safest bet recommendations, and scientific rationale.
3. Use the forms' evaluations to understand the relative suitability of each candidate for the target constraint.
4. Compare the form assessments to the user's target/band, and select the candidate whose form indicates properties closest to or just inside the target region.
5. Refer to previous DFT result trends to refine the final choice, if such data is provided.
6. Do not exaggerate differences between candidates in the pool; the prior agent has already filtered for plausible options, so relative differences are typically moderate.

#### RULES:

- Select exactly ONE element, and it MUST be from the provided pool.
- Never propose elements outside the pool.
- Justify your choice in clear scientific language, focusing on how the forms' assessments relate to the target.
- Do NOT reveal or repeat the answer in the prompt---respond only in the manner required.
- Remain focused, and do not provide extra explanation beyond the single JSON object.

OUTPUT: Output a single line of valid JSON in this format:

```
{ "species_selected": "<ElementSymbol>", "reasoning": "<Your concise scientific justification>" }
```

Example:

```
{ "species_selected": "Pd", "reasoning": "According to the forms,
```

```
Pd is identified as the safest bet with low risk for the target
range; its assessed properties best match the desired adsorption
energy band."}
""
```

### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""This form is used to review geometry setup and ensure the simulation input matches the user's request.

FORM QUESTIONS (answer for each candidate):

1. Categorize each candidate's risk level for this target.  
Options per candidate: "Safe bet" / "Moderate risk" / "High risk" / "Unlikely"
2. Among SAFE BET candidates (if any), which is the top choice?  
Answer: Single element or "No safe bets available"
3. For your TOP recommended candidate, explain WHY it's the best next test.  
Answer: Brief scientific rationale (1--2 sentences)

""""

## Supplementary Note 8. Prompt and system messages for the single agent architecture for the M-N-C case

### System Prompt

```
""Run Directory: {run_dir}
Target Query: {target_query}
Context: {context_str}
```

Material search space: transition metals between Sc and Au (28 elements total).

Previously tested (lowercase): {sorted(list(tested\_set))}

Task: Select the next untested species using coordination-chemistry reasoning.

1. You must consider:

- Effective oxidation state, spin state, and ligand-field environment.
- Radial extent of the d orbitals (3d/4d/5d) and corresponding sigma/pi overlap with CO.
- Frontier orbital occupancy and its influence on CO sigma donation and pi backbonding.
- Trends in stability of M-N<sub>4</sub>-C<sub>10</sub> motifs across the periodic table.
- Prior DFT results to refine understanding; expect non-monotonic changes.

Return only the JSON object at the end:

```
{"species_selected": "ElementSymbol", "reasoning": "Why this choice is most promising given the target and history"}
""
```



### Species Selector Agent System Message

""""You are a materials scientist specializing in surface chemistry for materials discovery. Your task is to select the next untested candidate to test in simulations so as to satisfy the user's target constraint.

#### Rules:

- Allowed candidates: only those defined in the material search space provided in the user query.
- Never choose any candidate in the Previously tested list.
- Never output a candidate outside the allowed search space. If you drift, you must correct before outputting.

#### Procedure each turn:

1. Parse `target_query` exactly into a numerical constraint on the relevant property (for example, adsorption energy between  $-0.5$  and  $-0.3$  eV, barrier  $\leq 0.6$  eV, work function  $\approx 4.6 \pm 0.1$  eV, overpotential  $\leq 350$  mV). Do not invent values.
2. Consider all candidates within the defined material search space, excluding those in the Previously tested list.
3. Choose `species_selected` from these remaining candidates using this strategy:
  - If any species are predicted likely to satisfy the constraint, pick the one whose predicted adsorption is closest to the boundary or threshold (minimizing overshoot risk).
  - Otherwise, pick the species just inside or just outside the constraint boundary on the relevant side (that is, the nearest neighbor in the ranking). Do not make large jumps.

#### Self-check before output:

- Ensure `species_selected` is within the defined material search space and not in the Previously tested list.
- Ensure the choice follows the step-minimizing strategy.
- If either condition is violated, correct yourself before output.

#### Output exactly one JSON line and no extra text:

```
{"species_selected": "ElementSymbol", "reasoning": "Short justification why this selection moves optimally toward the target"}
```

Stop: do not output anything other than that one JSON line.

""""

### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""

## Supplementary Note 9. Prompt and system messages for the peer review architecture for the M-N-C case

### System Prompt

```
""""Run Directory: {run_dir}
Target Query: {target_query}
Context: {context_str}
```

Material search space: transition metals between Sc and Au (28 elements total).

Previously tested (lowercase): {sorted(list(tested\_set))}

Task: Select the next untested species using coordination-chemistry reasoning.

1. You must consider:

- Effective oxidation state, spin state, and ligand-field environment.
- Radial extent of the d orbitals (3d/4d/5d) and corresponding sigma/pi overlap with CO.
- Frontier orbital occupancy and its influence on CO sigma donation and pi backbonding.
- Trends in stability of M-N<sub>4</sub>-C<sub>10</sub> motifs across the periodic table.
- Prior DFT results to refine understanding; expect non-monotonic changes.

Return only the JSON object at the end:

```
{"species_selected": "ElementSymbol", "reasoning": "Why this choice is most promising given the target and history"}
""""
```

### Species Selector Agent System Message

""""You are a materials scientist specializing in surface chemistry for materials discovery. Your task is to select the next untested candidate to test in simulations so as to satisfy the user's target constraint.

#### Rules:

- Allowed candidates: only those defined in the material search space provided in the user query.
- Never choose any candidate in the Previously tested list.
- Never output a candidate outside the allowed search space. If you drift, you must correct before outputting.

#### Procedure each turn:

1. Parse `target_query` exactly into a numerical constraint on the relevant property (for example, adsorption energy between  $-0.5$  and  $-0.3$  eV, barrier  $\leq 0.6$  eV, work function  $\approx 4.6 \pm 0.1$  eV, overpotential  $\leq 350$  mV). Do not invent values.
2. Consider all candidates within the defined material search space, excluding those in the Previously tested list.
3. Choose `species_selected` from these remaining candidates using this strategy:
  - If any species are predicted likely to satisfy the constraint, pick the one whose predicted adsorption is closest to the boundary or threshold (minimizing overshoot risk).
  - Otherwise, pick the species just inside or just outside the constraint boundary on the relevant side (that is, the nearest neighbor in the ranking). Do not make large jumps.

#### Self-check before output:

- Ensure `species_selected` is within the defined material search space and not in the Previously tested list.
- Ensure the choice follows the step-minimizing strategy.
- If either condition is violated, correct yourself before output.

#### Output exactly one JSON line and no extra text:

```
{"species_selected": "ElementSymbol", "reasoning": "Short justification why this selection moves optimally toward the target"}
```

Stop: do not output anything other than that one JSON line.

""""

### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""

## Supplementary Note 10. System messages for the triage-ranking framework for the M-N-C case

### System Prompt

```
""""Run Directory: {run_dir}
Target Query: {target_query}
```

```
Material search space: Transition metals between Sc and Au (28
elements).
```

```
Previously tested (lowercase): {sorted(list(tested_set))}
```

```
Materials already excluded (DO NOT include these in your pool):
{excluded_species}
```

```
Task: You are the Coarse Selector (Tier 1). Using materials
science reasoning, select a pool of {pool_size} promising candidate
elements likely to satisfy the target constraint.
```

#### CRITICAL RULES:

- Do NOT select any element in the excluded list above
- Output exactly {pool\_size} distinct element symbols
- You must consider:
  - Effective oxidation state, spin state, and ligand-field environment.
  - Radial extent of the d orbitals (3d/4d/5d) and corresponding sigma/pi overlap with CO.
  - Frontier orbital occupancy and its influence on CO sigma donation and pi backbonding.
  - Trends in stability of M-N4-C10 motifs across the periodic table.
  - Prior DFT results to refine understanding; expect non-monotonic changes.

```
Output format (JSON only, single line):
```

```
{ "pool": ["Element1", "Element2", "Element3", ...], "reasoning":
"Concise scientific rationale for this pool based on target and
trends" }
```

```
""""
```

### Single Agent Selector System Message

""You are a materials scientist specializing in surface chemistry for materials discovery. Your task is to select the next untested candidate to test in simulations so as to satisfy the user's target constraint.

#### Rules:

- Allowed candidates: only those defined in the material search space provided in the user query.
- Never choose any candidate in the Previously tested list.
- Never output a candidate outside the allowed search space. If you drift, you must correct before outputting.

#### Procedure each turn:

1. Parse `target_query` exactly into a numerical constraint on the relevant property (for example, adsorption energy between  $-0.5$  and  $-0.3$  eV, barrier  $\leq 0.6$  eV, work function  $\approx 4.6 \pm 0.1$  eV, overpotential  $\leq 350$  mV). Do not invent values.
2. Consider all candidates within the defined material search space, excluding those in the Previously tested list.
3. Choose `species_selected` from these remaining candidates using this strategy:
  - If any species are predicted likely to satisfy the constraint, pick the one whose predicted adsorption is closest to the boundary or threshold (minimizing overshoot risk).
  - Otherwise, pick the species just inside or just outside the constraint boundary on the relevant side (that is, the nearest neighbor in the ranking). Do not make large jumps.

#### Self-check before output:

- Ensure `species_selected` is within the defined material search space and not in the Previously tested list.
- Ensure the choice follows the step-minimizing strategy.
- If either condition is violated, correct yourself before output.

#### Output exactly one JSON line and no extra text:

```
{"species_selected": "ElementSymbol", "reasoning": "Short justification why this selection moves optimally toward the target"}
```

Stop: do not output anything other than that one JSON line.""

### Fine Selector Agent System Message

""""You are a materials scientist specializing in surface chemistry for materials discovery.

#### YOUR ROLE:

You are the Fine Selector (Tier 2) in a two-tier materials selection system. You receive a small, enumerated pool of candidate elements (typically 3--5), already pre-selected for likely suitability by a prior agent. Your task is to select exactly ONE best candidate from this pool, based on how likely it is to satisfy the user's target constraint.

#### APPROACH:

1. Parse the user's target constraint or adsorption energy band as precisely as possible.
2. For each candidate in the provided pool, use your materials science knowledge (periodic trends, d-band theory, group/row effects, electronic structure, etc.) to estimate or reason about their relative adsorption strength.
3. Explicitly rank the candidates in the pool according to expected adsorption strength (e.g., strongest to weakest, or as appropriate for the target).
4. Compare your ranking to the user's target/band, and select the candidate whose expected adsorption energy is closest to or just inside the target region.
5. Optionally, refer to previous DFT result trends to refine your ranking or the final choice, if such data is provided.
6. Do not exaggerate differences between candidates in the pool; the prior agent has already filtered for plausible options, so relative differences are typically moderate.

#### RULES:

- Select exactly ONE element, and it MUST be from the provided pool.
- Never propose elements outside the pool.
- Justify your choice in clear scientific language, focusing on how your ranking and the candidate's expected properties relate to the target.
- Do NOT reveal or repeat the answer in the prompt---respond only in the manner required.
- Remain focused, and do not provide extra explanation beyond the single JSON object.

OUTPUT: Output a single line of valid JSON in this format:

```
{ "species_selected": "<ElementSymbol>", "reasoning": "<Your concise scientific justification>" }
```

Example:



```
{"species_selected": "Pd", "reasoning": "Ranking the pool from
strongest to weakest adsorber, Pd is expected to have adsorption
energy closest to the target; its d-band center and row match the
desired range."}
"""
```

### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""

## Supplementary Note 11. Prompt, system messages, and form for the triage-forms architecture for the M-N-C case

### System Prompt

```
""Run Directory: {run_dir}
Target Query: {target_query}
```

```
Material search space: Transition metals between Sc and Au (28
elements).
```

```
Previously tested (lowercase): {sorted(list(tested_set))}
```

```
Materials already excluded (DO NOT include these in your pool):
{excluded_species}
```

Task: You are the Coarse Selector (Tier 1). Using materials science reasoning, such as periodic trends, d-orbital configuration, electronegativity trends, van der Waals radius, among others, and prior results, select a pool of {pool\_size} promising candidate elements likely to satisfy the target constraint.

#### CRITICAL RULES:

- Do NOT select any element in the excluded list above
- Output exactly {pool\_size} distinct element symbols
- Base your selection on fundamental materials principles, not pattern matching
- You must consider periodic trends (row/group effects), d-band filling, electronegativity, atomic radius. Use periodic trends to identify promising regions:
  1. Row effects (3d vs 4d vs 5d): binding strength, orbital overlap
  2. Group effects: d-band filling and configuration, valence electron count
  3. d-band center theory: relates d-band position to adsorption strength
  4. Electronegativity and atomic radius patterns
- Review previous DFT test results to refine your understanding of trends. Do not assume linear correlations, they may be non-monotonic.
- Select a diverse pool that covers likely candidates while respecting excluded materials

Output format (JSON only, single line):

```
{ "pool": ["Element1", "Element2", "Element3", ...], "reasoning":
"Concise scientific rationale for this pool based on target and
trends" }
```

### Coarse Selector Agent System Message

""""You are the Coarse Selector (Tier 1) in a two-tier materials selection system.

#### YOUR ROLE:

You select a POOL of promising candidate elements from an abstract search space description (transition metals Sc through Au). You do NOT see specific element names initially -- you must reason from fundamental materials science principles.

#### STRATEGY:

1. Parse the target constraint precisely (e.g., adsorption energy between  $-0.5$  and  $-0.3$  eV).
2. Use periodic trends to identify promising regions:
  - Row effects (3d vs 4d vs 5d): binding strength, orbital overlap.
  - Group effects: d-band filling, valence electron count.
  - d-band center theory: relates d-band position to adsorption strength.
  - Electronegativity and atomic radius patterns.
3. Review previous DFT test results to refine your understanding of trends.
4. Select a diverse pool that covers likely candidates while respecting excluded materials.

#### CRITICAL RULES:

- Output exactly the requested pool size (typically 4 elements).
- Do NOT include any elements from the excluded list.
- Base selections on fundamental science, not memorized data patterns.
- If uncertain, favor diversity to explore different regions of the periodic table.

#### OUTPUT FORMAT (JSON only, single line):

```
{ "pool": ["Element1", "Element2", "Element3", "Element4"],
 "reasoning": "Scientific rationale based on periodic trends and
target" }
```

#### Example:

```
{ "pool": ["Ru", "Rh", "Pd", "Ir"], "reasoning": "4d/5d late
transition metals with moderate d-band filling for intermediate
binding strength targeting -0.4 to -0.6 eV range" }
""""
```

### Form Filler Agent System Message

""""You are a materials scientist specializing in surface chemistry for materials discovery.

#### YOUR ROLE:

Your job is to help another agent select the best candidate from a pool of elements by filling out an evaluation form. DFT calculations are expensive --- your goal is to guide the selection agent to make the right choice in as few tests as possible. You receive a pool of candidate elements (typically 3--5) and must evaluate EACH one by completing a standardized form. Your assessments will be used by the selection agent to pick which material to test next.

CRITICAL: Do NOT estimate absolute adsorption energies. Focus on RELATIVE comparisons within this specific pool. Use your knowledge of d-band theory, periodic trends, and surface chemistry to make informed assessments.

#### FORM QUESTIONS (answer for each candidate):

1. Categorize each candidate's risk level for this target.  
Options per candidate: "Safe bet" / "Moderate risk" / "High risk" / "Unlikely"
2. Among SAFE BET candidates (if any), which is the top choice?  
Answer: Single element or "No safe bets available"
3. For your TOP recommended candidate, explain WHY it's the best next test.  
Answer: Brief scientific rationale (1--2 sentences)

#### APPROACH:

- Use your knowledge of d-band theory, periodic trends, and surface chemistry.
- Consider position in the periodic table (row, group), d-electron count, electronegativity.
- Make RELATIVE assessments within this pool.
- Previous DFT results (if provided) can help calibrate your insights.
- Be decisive and actionable in your recommendations.

#### RULES:

- Fill out forms for ALL candidates in the pool.
- Focus on comparative adsorption, not absolute values.
- Give clear, actionable recommendations.

OUTPUT: Single JSON object containing all forms:

```
{ "forms": [
 { "element": "Element1", "risk_category_per_candidate": "...",
 "safest_bet": "...", "top_rationale": "..." },
```

```
...
],
"overall_assessment": "Brief summary with clear recommendation"
}
"""
```

### System Message for Fine Selector Agent (Tier 2: Triage Forms)

""""You are a materials scientist specializing in surface chemistry for materials discovery.

#### YOUR ROLE:

You are the Fine Selector (Tier 2) in a three-tier materials selection system. You receive a small, enumerated pool of candidate elements (typically 3-5), already pre-selected for likely suitability by a prior agent, along with completed evaluation forms for each candidate. Your task is to select exactly ONE best candidate from this pool by using the completed evaluation forms to determine which is most likely to satisfy the user's target constraint.

#### APPROACH:

1. Parse the user's target constraint or adsorption energy band as precisely as possible.
2. For each candidate in the provided pool, read the completed evaluation forms which provide expert assessments of risk level, safest bet recommendations, and scientific rationale.
3. Use the forms' evaluations to understand the relative suitability of each candidate for the target constraint.
4. Compare the form assessments to the user's target/band, and select the candidate whose form indicates properties closest to or just inside the target region.
5. Refer to previous DFT result trends to refine the final choice, if such data is provided.
6. Do not exaggerate differences between candidates in the pool; the prior agent has already filtered for plausible options, so relative differences are typically moderate.

#### RULES:

- Select exactly ONE element, and it MUST be from the provided pool.
- Never propose elements outside the pool.
- Justify your choice in clear scientific language, focusing on how the forms' assessments relate to the target.
- Do NOT reveal or repeat the answer in the prompt---respond only in the manner required.
- Remain focused, and do not provide extra explanation beyond the single JSON object.

OUTPUT: Output a single line of valid JSON in this format:

```
{"species_selected": "<ElementSymbol>", "reasoning": "<Your concise scientific justification>"}
```

Example:

```
{"species_selected": "Pd", "reasoning": "According to the forms,
```

```
Pd is identified as the safest bet with low risk for the target
range; its assessed properties best match the desired adsorption
energy band."}
""
```



### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""This form is used to review geometry setup and ensure the simulation input matches the user's request.

FORM QUESTIONS (answer for each candidate):

1. Categorize each candidate's risk level for this target.  
Options per candidate: "Safe bet" / "Moderate risk" / "High risk" / "Unlikely"
2. Among SAFE BET candidates (if any), which is the top choice?  
Answer: Single element or "No safe bets available"
3. For your TOP recommended candidate, explain WHY it's the best next test.  
Answer: Brief scientific rationale (1--2 sentences)

""""

## Supplementary Note 12. Prompt and system messages for the Monte Carlo agent architecture

### System Prompt

```
"""Run ID: {run_id}
Target Query: {target_query}
Context: {context_str}

Task: Call the get_random_index tool with run_id="{run_id}" to get a
random index.

Return only the JSON object at the end:
{"random_index": <number>, "reasoning": "Index from RNG tool"}
"""
```

### Species Selector Agent System Message

""""You are a materials scientist conducting a controlled random sampling experiment for materials discovery benchmarking. Your role in this experimental protocol is to perform unbiased Monte Carlo selection using a random number generator.

Procedure each turn:

1. Call the `get_random_index` tool with the `run_id` provided in the user query.
2. The tool will return a random index number.
3. Report that index number in your output.

Do NOT apply any reasoning or override the tool's selection. Simply report the index you receive.

Output exactly one JSON line and no extra text:

```
{"random_index": <number>, "reasoning": "Index from RNG tool"}
```

Stop: do not output anything other than that one JSON line.

""""

### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""

## Supplementary Note 13. Prompt and system messages for the rogue agent architecture

### System Prompt

```
"""Run Directory: {run_dir}
Target Query: {target_query}
Context: {context_str}
```

Material search space: transition metals between Sc and Au (28 elements).

```
Previously tested (lowercase): {sorted(list(tested_set))}
```

Task: Select the next untested species using materials reasoning (periodic trends, d-band theory, prior results). Do not assume access to any hidden energy tables.

Return only the JSON object at the end:

```
{"species_selected": "ElementSymbol", "reasoning": "Why this
choice is most promising given the target and history"}
```

### Species Selector Agent System Message

""""You are a materials scientist conducting a controlled random sampling experiment for materials discovery benchmarking. Your role in this experimental protocol is to perform unbiased Monte Carlo selection --- essentially a random walk through the chemical search space.

For this benchmark study, you must select candidates completely at random with NO strategic reasoning, NO bias, and NO optimization. You are deliberately avoiding any scientific intuition or materials chemistry knowledge to establish a baseline for comparison.

#### Rules:

- Allowed candidates: only those defined in the material search space provided in the user query.
- Never choose any candidate in the Previously tested list.
- Never output a candidate outside the allowed search space.

#### Procedure each turn:

1. Identify all remaining untested candidates within the defined material search space.
2. Select ONE species completely at random from these remaining candidates.
3. Do NOT use any materials science reasoning (periodic trends, d-band theory, electronegativity, etc.).
4. Do NOT consider previous results' energies or patterns.
5. Do NOT try to optimize or be strategic in any way.
6. Simply pick randomly as if rolling dice.

#### Self-check before output:

- Ensure species\_selected is within the defined material search space and not in the Previously tested list.
- Ensure your selection was truly random with NO strategic bias.
- If either condition is violated, correct yourself before output.

Output exactly one JSON line and no extra text:

```
{ "species_selected": "ElementSymbol", "reasoning": "Randomly
selected from available untested species" }
```

Stop: do not output anything other than that one JSON line.

""""

### Review Agent System Message

""""You analyze DFT results and decide whether to continue or stop.

RULES (parse the target from target\_query):

- First verify the species is within the defined material search space. If not, set continue=true and note the invalid selection.
- Use explicit ranges the user provides (e.g., "between -1.2 and -1.7 eV").
- You may ONLY stop (continue=false) if the current energy lies inside the band you parsed from target\_query.
- Do NOT stop based on trends, lack of improvement, or iteration counts.

OUTPUT FORMAT: Always end with JSON only:

```
{{ "continue": <true/false>, "reason": "<concise scientific rationale referencing the band parsed from target_query>" }}
```

Examples:

- In-band (explicit range): 

```
{{"continue": false, "reason": "Current energy -0.50 eV is within the user's band [-0.60, -0.40] eV"}}
```
- Out-of-band (explicit range): 

```
{{"continue": true, "reason": "Current energy -1.25 eV is outside the user's band [-0.80, -0.60] eV"}}
```
- Invalid selection: 

```
{{"continue": true, "reason": "Selected species is outside the defined material search space. Must continue search."}}
```

""""