

---

# CatAgent: Multi-Agent Orchestration for Electrocatalyst Discovery

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

1 The discovery of efficient hydrogen evolution reaction catalysts is hindered by  
2 the combinatorial scale of candidate material space, which limits conventional  
3 screening approaches. Here we present CatAgent, an autonomous multi-agent  
4 workflow driven by large language models that achieves a 2.27-fold increase in  
5 discovery rates over the random baseline. The framework integrates hypothesis  
6 generation, surrogate-model evaluation, and iterative critic feedback to navigate  
7 chemical space toward optimal hydrogen adsorption strength. We benchmark 13  
8 state-of-the-art language models in 'single-shot mode' and 'iterative CatAgent  
9 mode' across alloy compositions spanning  $L1_0$  and  $L1_2$  crystallographic phases.  
10 Critic-enabled iterations improve performance for most model architectures, though  
11 self-refinement effectiveness varies across model families. Top-performing models  
12 such as Gemini 3 Flash and GPT-5.2 concentrate proposals near zero theoretical  
13 overpotential, identifying successful candidates twice as efficiently as random  
14 sampling. In contrast, outdated models such as GPT-3.5 Turbo fail to formulate  
15 consistent hypotheses and exhibit unreliable search behavior. This work provides  
16 the first demonstration that state-of-the-art language models can autonomously  
17 reshape high-throughput catalyst screening into reasoning-guided exploration.

18 

## 1 Introduction

19 Scalable renewable energy storage demands efficient electrocatalysts for hydrogen production, yet  
20 state-of-the-art systems rely on scarce noble metals such as platinum and iridium [1]. This limitation  
21 drives an intensive search for earth-abundant alternatives, with bimetallic alloys emerging as  
22 promising candidates due to their tunable electronic properties. For the hydrogen evolution reaction  
23 (HER), the hydrogen adsorption energy ( $\Delta E_H$ ) serves as a foundational descriptor for catalytic  
24 activity, following a well-established volcano relationship where optimal performance correlates  
25 with near-thermoneutral binding [2, 3]. Specifically, targeting a free energy of adsorption near zero  
26 requires an electronic adsorption energy of approximately  $-0.27$  eV, once zero-point energy and  
27 entropy corrections are accounted for [4, 5].

28 However, navigating the vast chemical space of electrocatalysts presents a significant combinatorial  
29 challenge [6]. Conventional high-throughput screening is computationally intensive, often requiring  
30 exhaustive sampling of myriad adsorption sites, surface terminations, and adsorption configurations  
31 [7]. Such enumeration typically relies on density functional theory (DFT), which remains prohibitively  
32 expensive for large-scale exploration. To overcome these bottlenecks, researchers have increasingly  
33 turned to active learning [8] and machine learning interatomic potentials [9], which can accelerate  
34 discovery by an order of magnitude or more compared to manual exploration [10]. These frameworks  
35 allow for rapid navigation through thousands of intermetallic configurations by prioritizing promising  
36 candidates for high-fidelity evaluation.

37 Recently, the integration of large language models (LLMs) [11, 12] and transformer architectures  
38 [13] into materials science has introduced new paradigms for catalyst design [14–16]. Models such  
39 as CatBERTa and CatGPT have demonstrated that transformer-based representations can substitute  
40 expensive 3D descriptors in early-stage screening and enable inverse design of catalyst structures  
41 [17, 18]. Beyond these surrogate models, LLM-driven agentic workflows have emerged for fully  
42 autonomous experimentation. The A-Lab achieved a 63% success rate in synthesizing computationally  
43 predicted inorganic materials through closed-loop integration of ML-based synthesis planning, robotic  
44 experimentation, and active learning [19]. Similarly, multi-agent frameworks such as DREAMS have  
45 shown that LLM-based planners can orchestrate complex DFT simulation workflows with expert-level  
46 accuracy [20], while agentic systems for metal-organic framework discovery have demonstrated the  
47 capacity to autonomously navigate vast combinatorial spaces through iterative hypothesis generation  
48 and validation [21]. Frameworks such as Adsorb-Agent have further illustrated how LLM-driven  
49 reasoning can strategically identify stable adsorption configurations without exhaustive geometric  
50 enumeration by combining proposal engines with surrogate models such as graph neural networks  
51 [22]. These advancements demonstrate that LLM-based reasoning can drive autonomous navigation  
52 of chemical space, motivating the development of agentic workflows for catalyst discovery.

53 In this study, we investigate LLM-driven candidate selection for the screening of bimetallic HER  
54 catalysts within the  $L_{1_0}$  and  $L_{1_2}$  composition space. We compare the efficacy of a single-shot LLM  
55 approach against an iterative agentic workflow, CatAgent, with and without an integrated critic module  
56 for refined reasoning. Using a pretrained graph neural network potential (uma-s-1p1)[23] as surrogate  
57 model for rapid property evaluation, we assess how different LLM-orchestration strategies navigate  
58 bimetallic surface chemistry toward optimal hydrogen adsorption strength. The main contributions of  
59 this work are:

- 60 • A comprehensive benchmark of 13 state-of-the-art LLMs across three major providers  
61 (Google Gemini, OpenAI GPT, and Anthropic Claude), establishing quantitative metrics for  
62 chemical reasoning over 1,998 bimetallic compositions.
- 63 • CatAgent, a multi-agent framework integrating hypothesis generation, surrogate-model  
64 simulation, and critic feedback, achieving discovery rates up to 2.27-fold above random  
65 baseline.
- 66 • Systematic analysis revealing architecture-dependent search behaviors: top-performing  
67 models (Gemini 3 Flash, GPT-5.2) concentrate proposals near optimal adsorption energy,  
68 while outdated models (GPT-3.5 Turbo) exhibit dispersed exploration.

## 69 2 Methodology

### 70 2.1 Dataset

71 The search space for catalyst discovery is defined by bimetallic alloys composed of 37 metals: Al,  
72 Sc, Ti, V, Cr, Mn, Fe, Co, Ni, Cu, Zn, Ga, Y, Zr, Nb, Mo, Tc, Ru, Rh, Pd, Ag, Cd, In, Sn, La, Hf,  
73 Ta, W, Re, Os, Ir, Pt, Au, Hg, Tl, Pb, and Bi [24]. These metals are combined into two distinct  
74 composition families: the  $L_{1_2}$  structure with a stoichiometry of metal1:metal2 = 0.75:0.25, and the  
75  $L_{1_0}$  structure with a stoichiometry of 0.50:0.50. The total candidate space consists of 1,998 unique  
76 entries, comprising 1,332  $L_{1_2}$  and 666  $L_{1_0}$  candidates. For each candidate, adsorption sites are  
77 enumerated according to symmetry, resulting in 9 unique sites for  $L_{1_2}$  and 10 for  $L_{1_0}$ .

78 Adsorption energies are calculated for every enumerated site, and the representative value for a given  
79 composition is defined as the minimum adsorption energy across all sites. The adsorption energy  
80  $\Delta E_H$  is calculated relative to gas-phase H<sub>2</sub> as follows:

$$\Delta E_H = E_{\text{slab}+\text{H}} - E_{\text{slab}} - \frac{1}{2}E_{\text{H}_2}.$$

81 The target adsorption energy is  $-0.27$  eV, corresponding to a near-thermoneutral free energy of  
82 adsorption ( $\Delta G_H = \Delta E_H + 0.24$  eV  $\approx 0$ ) that maximizes HER activity according to Sabatier’s  
83 principle [4, 8]. A successful candidate is defined as falling within the window of  $-0.57$  eV  $\leq$   
84  $\Delta E_H \leq 0.03$  eV. The baseline dataset distribution corresponds to the minimum H adsorption energies  
85 across the 1,998 candidates, utilizing Universal Models for Atoms (UMA, uma-s-1p1 checkpoint) as a  
86 surrogate evaluator [23]. Among machine learning interatomic potentials benchmarked for hydrogen

87 adsorption energy prediction, UMA-s achieves the lowest total MAE (0.168 eV) combined with  
 88 high structural reliability (84.58% normal relaxation rate), making it well-suited for rapid screening  
 89 [25, 26].

90 **2.2 CatAgent workflow**

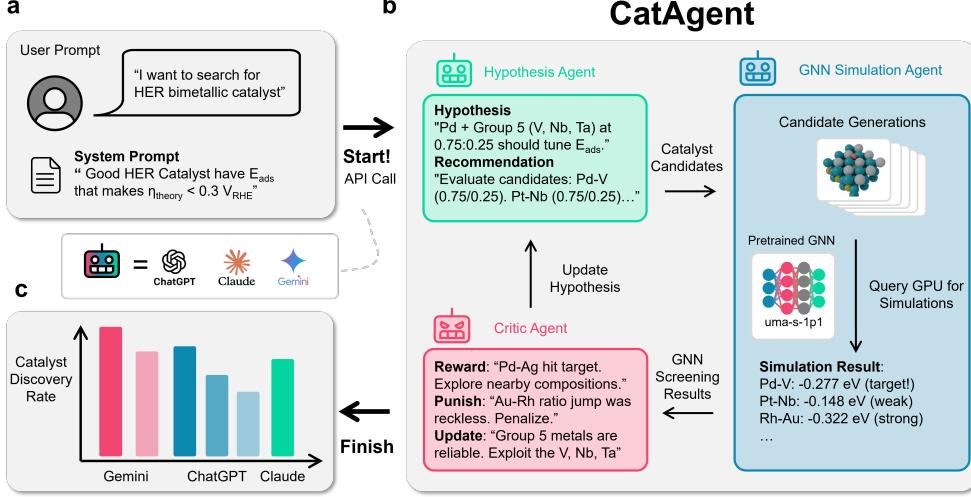


Figure 1: Schematic representation of the CatAgent multi-agent workflow architecture. (a) User and system prompts specify the screening objective and acceptance criterion. (b) Iterative loop of LLM-driven hypothesis formulation, candidate suggestion, UMA-based energy evaluation, and critic feedback for successive refinement. (c) Performance evaluation of screening outcomes.

91 Two primary screening approaches are employed: single-shot screening and the iterative CatAgent  
 92 workflow. In the single-shot mode, the large language model (LLM) performs a combined role of  
 93 hypothesis formulation and candidate generation to propose 200 candidates in a single non-iterative  
 94 pass (Fig. 2a). This mode was evaluated using 13 models: Gemini 3 Flash, Gemini 3 Pro, Gemini 2.5  
 95 Flash, Gemini 2.5 Pro, GPT-5.2, GPT-5, GPT-5 mini, GPT-5 nano, GPT-4o mini, GPT-3.5 Turbo,  
 96 Claude Opus 4.5, Claude Haiku 4.5, and Claude Sonnet 4.5.

97 The CatAgent workflow utilizes a multi-agent architecture comprising LLM-driven hypothesis,  
 98 simulation, and critic agents. After the user specifies the screening objective and acceptance criterion  
 99 (Fig. 1a), each run explores 200 candidates over 20 steps, proposing 10 candidates per step. We  
 100 investigate two variants of this workflow: a critic-enabled and a critic-disabled configuration. In  
 101 the critic-enabled variant, the critic agent evaluates stepwise outcomes to identify recurring success  
 102 and failure patterns, and its feedback is used by the hypothesis agent to refine the search strategy in  
 103 subsequent steps (Fig. 1b). After termination of the iterative loop, overall screening performance  
 104 is evaluated using the metrics defined in Section 2.3 (Fig. 1c), enabling systematic comparison of  
 105 screening outcomes across models and workflow configurations. The workflow was benchmarked  
 106 using Gemini 3 Flash, GPT-5.2, GPT-4o mini, and GPT-3.5 Turbo; for all models and providers,  
 107 default sampling hyperparameters were used.

108 The simulation agent executes an evaluation pipeline using UMA: (i) construction of the bulk phase,  
 109 (ii) bulk relaxation, (iii) construction of a  $(2 \times 2)$  three-layer slab with 20 Å of vacuum, (iv) slab relax-  
 110 ation, (v) adsorption site enumeration, (vi) relaxation of the H+slab adsorbed structure with the bottom  
 111 two layers fixed, and (vii) calculation of  $\Delta E_H$  using the gas-phase H<sub>2</sub> energy. Geometry optimizations  
 112 were performed using the Limited-memory Broyden-Fletcher-Goldfarb-Shanno (LBFGS) algorithm  
 113 with a force criterion of  $f_{max} = 0.05$  eV/Å and a maximum of 100 steps. The codes, exact model  
 114 snapshots, full prompts, and all raw outputs for the single-shot and CatAgent workflows are available  
 115 at [https://osf.io/t8he3/overview?view\\_only=1255c0a1cb1646bd9616d634314cd5e9](https://osf.io/t8he3/overview?view_only=1255c0a1cb1646bd9616d634314cd5e9).

116 **2.3 Metrics and statistics**

117 The performance of the screening process is evaluated using the success probability, defined as:

$$P_{\text{success}} = \frac{|\{i \mid \Delta E_{H,i} \in [-0.57, 0.03] \text{ eV}\}|}{200}$$

118 where the success window corresponds to the target adsorption energy of  $-0.27 \text{ eV} \pm 0.30 \text{ eV}$ . The  
 119 random-choice success probability is defined as the probability of selecting a successful candidate  
 120 from the baseline distribution. Among the 1,998 candidates in the search space, 699 fall within the  
 121 success window, yielding a random baseline of  $699/1998 = 35.0\%$ . The improvement factor is  
 122 calculated as the ratio of the observed success probability to this random baseline; a value of 1.0  
 123 indicates performance equivalent to random selection, while values exceeding unity suggest an ability  
 124 to identify compositions closer to the target adsorption energy.

125 Token usage is reported as the sum of input and output tokens. Cached tokens are excluded to ensure  
 126 consistent comparison across providers where caching metadata may not be available. For  $L1_0$   
 127 compositions, symmetric pairs are considered equivalent and deduplicated within and across steps.  
 128 All statistical values derived from triplicate runs are reported as mean  $\pm$  standard deviation.

129 **3 Results and discussion**

130 **3.1 Intrinsic performance comparison of language models**

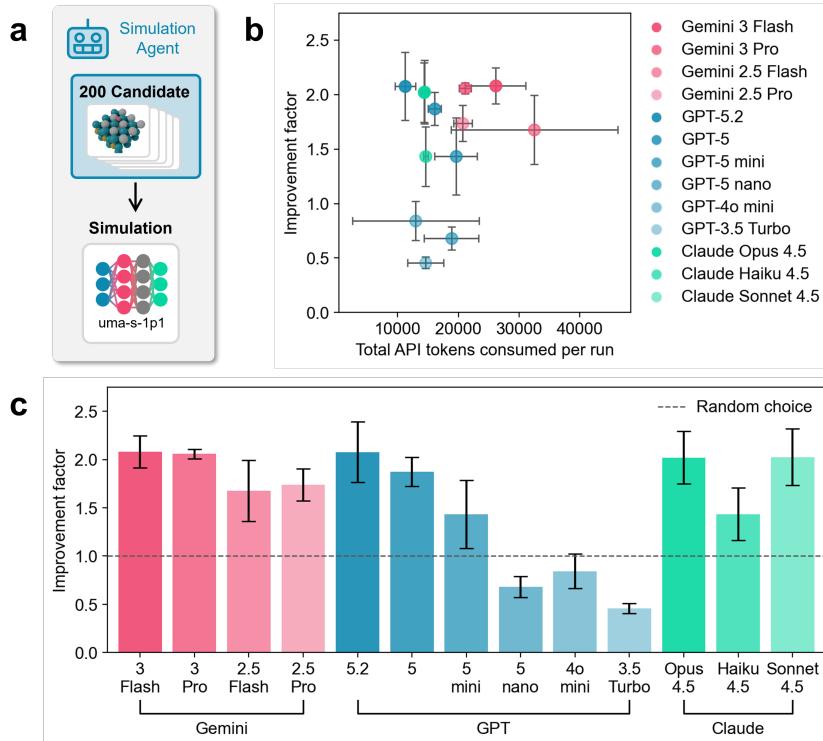


Figure 2: Single-shot screening performance analysis. (a) Workflow schematic. (b) Relationship between discovery efficiency and computational resource consumption across the 13 language models. (c) Improvement factors of the language models. Mean  $\pm$  standard deviation from triplicate runs.

131 Single-shot benchmarking reveals substantial variation in chemical reasoning capabilities across the  
 132 13 language models evaluated (Fig. 2a). Resource consumption differs markedly among models,  
 133 ranging from approximately 26,120 tokens per trial for Gemini 3 Flash to considerably lower values  
 134 for more efficient architectures (Fig. 2b). Despite this variation in computational cost, performance  
 135 differences are even more pronounced: Gemini 3 Flash achieves the highest mean improvement

136 factor of 2.08, demonstrating robust navigation of the chemical composition space in a single attempt,  
137 whereas GPT-3.5 Turbo exhibits an improvement factor of only 0.46, performing below the random  
138 baseline (Fig. 2c). This disparity across model architectures indicates that LLM training and design  
139 fundamentally determine single-shot predictive accuracy for catalyst discovery tasks.

140 Further analysis reveals that the number of unique candidates explored correlates strongly with the  
141 resulting improvement factor. While the experimental target is 200 unique catalyst candidates per  
142 trial, several models fail to meet this goal; GPT-4o mini averages only 116 candidates. Notably,  
143 resource efficiency does not track directly with overall performance: GPT-5.2 requires an average of  
144 only 78 tokens per successful candidate, representing a favorable balance between computational  
145 cost and discovery effectiveness. These metrics provide a comprehensive baseline for intrinsic model  
146 capabilities before introducing iterative agentic workflows.

### 147 3.2 CatAgent performance comparison across language models

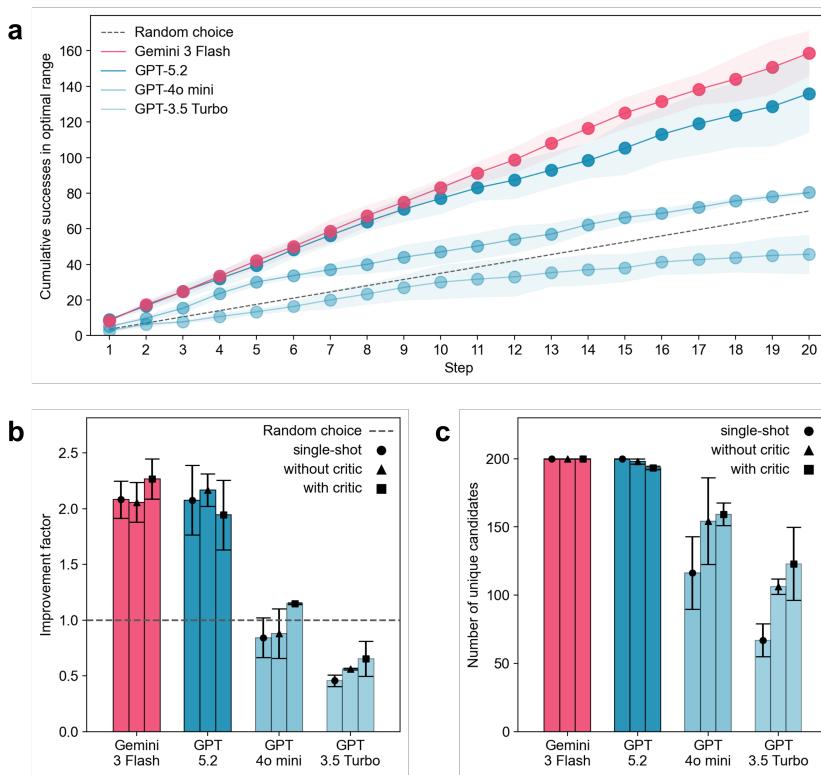


Figure 3: CatAgent workflow performance across model architectures. (a) Cumulative discovery of successful candidates during critic-guided iterative search. (b) Improvement factors and (c) unique candidate counts for single-shot, non-critic iterative, and critic-enabled iterative workflows across Gemini 3 Flash, GPT-5.2, GPT-4o mini, and GPT-3.5 Turbo. Mean  $\pm$  standard deviation from triplicate runs.

148 The introduction of the CatAgent workflow allows for a direct comparison of discovery performance  
149 across different architectural designs: the single-shot baseline, the CatAgent workflow without critic,  
150 and the CatAgent workflow with critic. This analysis is focused on a subset of four models: Gemini 3  
151 Flash, GPT-5.2, GPT-4o mini, and GPT-3.5 Turbo. As demonstrated in Fig. 3b, for Gemini 3  
152 Flash, the addition of a critic module is consistent with an increase in the improvement factor of  
153 approximately 0.21 compared to the iterative workflow lacking a critic. Similarly, GPT-4o mini  
154 demonstrates a notable improvement in its improvement factor, rising from 0.88 to 1.15 when  
155 the critic module is integrated. These results suggest that iterative self-correction and the inclusion of  
156 an evaluation step generally enhance the model's ability to narrow down the chemical search space  
157 toward the desired adsorption energy target.

158 However, the efficacy of the critic module is not universal across all model architectures. Contrary  
 159 to the general trends observed in other models, GPT-5.2 exhibits a lower improvement factor when  
 160 a critic is employed compared to iterative trials conducted without one. In some instances, success  
 161 probability can fall below the level of a random choice, particularly for models like GPT-3.5 Turbo.  
 162 This performance deficit can be attributed to unfocused search behaviors or a failure to generate  
 163 the required number of unique candidates. Additionally, the integration of a critic can influence the  
 164 volume of output (Fig. 3c); with a critic, GPT-4o mini explores approximately 159 unique candidates  
 165 per trial, which represents an increase over its single-shot baseline of 116.

166 **3.3 Analysis of catalyst compositions discovered via CatAgent**

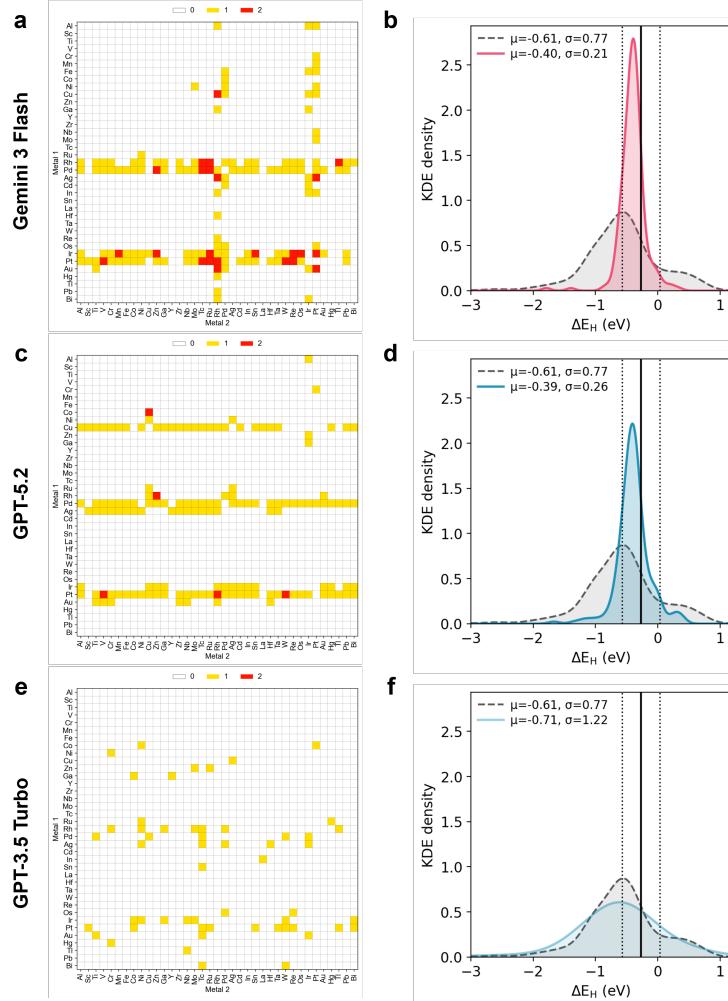


Figure 4: (a, c, e) Heatmaps showing the composition distributions of candidates with adsorption energies within the energy threshold and (b, d, f) KDE plots showing the adsorption-energy distributions of all catalysts identified in one repeat of the critic-enabled CatAgent. (a,b) Gemini 3 Flash; (c,d) GPT-5.2; (e,f) GPT-3.5 Turbo.

167 A qualitative analysis of the catalyst compositions discovered through the CatAgent workflows  
 168 reveals different search behaviors, categorized as either focused or dispersed. The performance of the  
 169 models is reflected in the adsorption-energy distribution mean and standard deviation relative to the  
 170 broader dataset distribution. Models that exhibit focused search behaviors generate candidates with an  
 171 energy distribution that has a smaller standard deviation and a mean shift toward the target, whereas  
 172 dispersed search behaviors result in candidates that more closely resemble the original dataset's broad  
 173 distribution.

174 Gemini 3 Flash exhibited a clear shift toward the target energy when using the CatAgent workflow  
175 with critic (Figs. 4a,b, and Fig. 5). The mean absolute deviation from the target decreased from 0.34  
176 eV at baseline to 0.12 eV across three trials. This focusing effect is also evident in the narrowing of  
177 the energy distribution, with the standard deviation of adsorption energies reduced to an average of  
178 36.7% of the baseline dataset value.

179 GPT-5.2 demonstrated the most precise alignment with the target energy among the models analyzed  
180 (Figs. 4c,d, and Fig. 6). In the CatAgent workflow with critic, GPT-5.2 achieved a mean absolute  
181 deviation of 0.08 eV from the target. The model shifted the mean candidate adsorption energy by  
182 0.26 eV relative to the dataset baseline and reduced the standard deviation to 47.4% of the dataset  
183 value. This combination of a significant mean shift and a reduced standard deviation indicates that  
184 GPT-5.2 employs a focused search strategy that identifies compositions near the desired adsorption  
185 energy.

186 In contrast, GPT-3.5 Turbo exhibited dispersed search behavior with limited success in approaching  
187 the HER target (Figs. 4e,f, and Fig. 8). In the CatAgent workflow with critic, GPT-3.5 Turbo shifted  
188 the mean candidate adsorption energy by -0.05 eV relative to the dataset baseline across three trials,  
189 representing a move away from the target direction. Coupled with its reduced candidate generation  
190 count of 123, the resulting energy distribution remained largely unfocused, consistent with its lower  
191 normalized success rate. The contrasting energy distribution characteristics across these models  
192 highlight the varying capacity of language models to interpret and act upon chemical design targets  
193 within an agentic framework.

## 194 4 Conclusion

195 CatAgent, a multi-agent workflow driven by large language models, demonstrates substantial effectiveness  
196 in navigating complex chemical spaces for catalyst discovery. Through systematic evaluation  
197 of 13 language models in single-shot screening and four models within iterative multi-agent configura-  
198 tions, we establish that LLM-guided exploration can strategically identify promising bimetallic  
199 hydrogen evolution reaction catalysts from a search space of 1,998 compositions spanning  $L1_0$  and  
200  $L1_2$  crystallographic phases.

201 The CatAgent framework integrates hypothesis generation, computational evaluation, and critic  
202 feedback mechanisms to achieve discovery rates up to 2.27-fold above random baseline selection.  
203 The integration of critic modules demonstrates architecture-dependent efficacy, with models such as  
204 GPT-4o mini exhibiting substantial performance improvements while GPT-5.2 showing decreased  
205 effectiveness, highlighting the heterogeneous nature of self-refinement capabilities across model fam-  
206 ilies. Notably, performance remains robust across three major API providers—Google, OpenAI, and  
207 Anthropic—validating the generalizability of the approach beyond specific model implementations.  
208 Leading models like Gemini 3 Flash demonstrate focused search behaviors that concentrate candidate  
209 proposals near the target adsorption energy while substantially reducing distribution variance relative  
210 to the baseline dataset. Resource efficiency varies significantly across LLM architectures, with  
211 optimal configurations requiring as few as 78 tokens per successful candidate.

212 However, several opportunities for advancement remain. First, the current search space is relatively  
213 constrained. Scaling to broader compositional ranges, alternative crystal structures, and multi-  
214 component systems would better test framework capabilities. Second, whether domain-specific  
215 fine-tuning on materials science data could enhance chemical reasoning remains unexplored. Third,  
216 extending CatAgent to multi-objective optimization would enhance practical utility, as real-world cat-  
217 alyst design requires balancing activity, selectivity, stability, cost, and synthesizability simultaneously.  
218 This work establishes LLM-powered autonomous reasoning as a promising paradigm for accelerating  
219 electrocatalyst discovery in vast chemical spaces. Continued development in LLM reasoning capacity,  
220 prompting strategies, and experimental integration will be essential for realizing the full potential of  
221 AI-assisted materials discovery.

222    **References**

- 223 [1] C Lawrence Zitnick, Lowik Chanussot, Abhishek Das, Siddharth Goyal, Javier Heras-Domingo, Caleb  
224 Ho, Weihua Hu, Thibaut Lavril, Aini Palizhati, Morgane Riviere, et al. An introduction to electrocatalyst  
225 design using machine learning for renewable energy storage. *arXiv preprint arXiv:2010.09435*, 2020.
- 226 [2] Hyeonjung Jung, Seokhyun Choung, and Jeong Woo Han. Design principles of noble metal-free electrocata-  
227 lysts for hydrogen production in alkaline media: combining theory and experiment. *Nanoscale Advances*,  
228 3(24):6797–6826, 2021.
- 229 [3] Seokhyun Choung, Heejae Yang, Jinuk Moon, Wongyu Park, Hyekjoon June, Chaesung Lim, and  
230 Jeong Woo Han. Theoretical tuning of local coordination environment of metal-nitrogen doped carbon  
231 catalysts for selective chlorine-evolution reaction. *Catalysis Today*, 425:114358, 2024.
- 232 [4] Jens Kehlet Nørskov, Thomas Bligaard, Ashildur Logadóttir, John R. Kitchin, Jingguang G. Chen, Svetlozar  
233 Pandelov, and Ulrich Stimming. Trends in the exchange current for hydrogen evolution. *Journal of The  
234 Electrochemical Society*, 152(2):J23–J26, 2005. doi: 10.1149/1.1856988.
- 235 [5] Xiangyun Xiao, Sungsu Kang, Seokhyun Choung, Jeong Woo Han, Jungwon Park, and Taekyung Yu.  
236 Synthesis of metal cation doped nanoparticles for single atom alloy catalysts using spontaneous cation  
237 exchange. *Journal of Materials Chemistry A*, 11(6):2857–2867, 2023.
- 238 [6] Seokhyun Choung, Wongyu Park, Jinuk Moon, and Jeong Woo Han. Rise of machine learning potentials  
239 in heterogeneous catalysis: Developments, applications, and prospects. *Chemical Engineering Journal*,  
240 494:152757, 2024.
- 241 [7] Zhi Wei Seh, Jakob Kibsgaard, Colin F Dickens, IB Chorkendorff, Jens K Nørskov, and Thomas F Jaramillo.  
242 Combining theory and experiment in electrocatalysis: Insights into materials design. *Science*, 355(6321):  
243 eaad4998, 2017.
- 244 [8] Kevin Tran and Zachary W. Ulissi. Active learning across intermetallics to guide discovery of  
245 electrocatalysts for co2 reduction and h2 evolution. *Nature Catalysis*, 1(9):696–703, 2018. doi:  
246 10.1038/s41929-018-0142-1.
- 247 [9] Ryan Jacobs, Dane Morgan, Siamak Attarian, Jun Meng, Chen Shen, Zhenghao Wu, Clare Yijia Xie,  
248 Julia H Yang, Nongnuch Artrith, Ben Blaiszik, et al. A practical guide to machine learning interatomic  
249 potentials—status and future. *Current Opinion in Solid State and Materials Science*, 35:101214, 2025.
- 250 [10] L. Kavalsky, V. I. Hegde, E. Muckley, M. S. Johnson, B. Meredig, and V. Viswanathan. By how much can  
251 closed-loop frameworks accelerate computational materials discovery? *Digital Discovery*, 2:1112–1125,  
252 2023. doi: 10.1039/D2DD00133K.
- 253 [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
254 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.  
255 *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 256 [12] Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed  
257 Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language models. *ACM  
258 Transactions on Intelligent Systems and Technology*, 16(5):1–72, 2025.
- 259 [13] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
260 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*,  
261 30, 2017.
- 262 [14] Kevin Maik Jablonka, Qianxiang Ai, Alexander Al-Feghali, Shruti Badhwar, Joshua D Bocarsly, Andres M  
263 Bran, Stefan Bringuer, L Catherine Brinson, Kamal Choudhary, Defne Circi, et al. 14 examples of how  
264 llms can transform materials science and chemistry: a reflection on a large language model hackathon.  
265 *Digital discovery*, 2(5):1233–1250, 2023.
- 266 [15] Negin Orouji, Jeffrey A. Bennett, Richard B. Canty, Long Qi, Shijing Sun, Paulami Majumdar, Chong Liu,  
267 Núria López, Neil M. Schweitzer, John R. Kitchin, Hongliang Xin, and Milad Abolhasani. Autonomous  
268 catalysis research with human–AI–robot collaboration. *Nature Catalysis*, 8:1135–1145, 2025. doi:  
269 10.1038/s41929-025-01430-6.
- 270 [16] Hongliang Xin, John R. Kitchin, and Heather J. Kulik. Towards agentic science for advancing scientific  
271 discovery. *Nature Machine Intelligence*, 7:1373–1375, 2025. doi: 10.1038/s42256-025-01110-x.

- 272 [17] Janghoon Ock, Chakradhar Guntuboina, and Amir Barati Farimani. Catalyst energy prediction with  
273 catberta: Unveiling feature exploration strategies through large language models. *ACS Catalysis*, 13(24):  
274 16032–16044, 2023. doi: 10.1021/acscatal.3c04956.
- 275 [18] Dong Hyeon Mok and Seoin Back. Generative pretrained transformer for heterogeneous catalysts. *Journal  
276 of the American Chemical Society*, 146(49):33712–33722, 2024. doi: 10.1021/jacs.4c11504.
- 277 [19] Nathan J. Szymanski, Bernardus Rendy, Yuxing Fei, Rishi E. Kumar, Tanjin He, David Milsted, Matthew J.  
278 McDermott, Max Gallant, Ekin Dogus Cubuk, Amil Merchant, Haegyeom Kim, Anubhav Jain, Christopher J. Bartel,  
279 Kristin Persson, Yan Zeng, and Gerbrand Ceder. An autonomous laboratory for the  
280 accelerated synthesis of inorganic materials. *Nature*, 624:86–91, 2023. doi: 10.1038/s41586-023-06734-w.
- 281 [20] Ziqi Wang, Hongshuo Huang, Hancheng Zhao, Changwen Xu, Shang Zhu, Jan Janssen, and Venkatasubra-  
282 manian Viswanathan. DREAMS: Density functional theory based research engine for agentic materials  
283 simulation. *arXiv preprint arXiv*, 2025.
- 284 [21] Theo Jaffrelot Inizan, Sherry Yang, Aaron Kaplan, Yen-hsu Lin, Jian Yin, Saber Mirzaei, Mona Abdelgaid,  
285 Ali H Alawadhi, KwangHwan Cho, Zhiling Zheng, et al. System of agentic ai for the discovery of  
286 metal-organic frameworks. *arXiv preprint arXiv:2504.14110*, 2025.
- 287 [22] Janghoon Ock, Tirtha Vinchurkar, Yayati Jadhav, and Amir Barati Farimani. Adsorb-agent: Autonomous  
288 identification of stable adsorption configurations via large language model agent. 2024. *arXiv:2410.16658*.
- 289 [23] Brandon M Wood, Misko Dzamba, Xiang Fu, Meng Gao, Mohammed Shuaibi, Luis Barroso-Luque,  
290 Kareem Abdelmaqsoud, Vahe Gharakhanyan, John R Kitchin, Daniel S Levine, et al. Uma: A family of  
291 universal models for atoms. *arXiv preprint arXiv:2506.23971*, 2025.
- 292 [24] Osman Mamun, Kirsten T. Winther, Jacob R. Boes, and Thomas Bligaard. High-throughput cal-  
293 culations of catalytic properties of bimetallic alloy surfaces. *Scientific Data*, 6:76, 2019. doi:  
294 10.1038/s41597-019-0080-z.
- 295 [25] Jinuk Moon, Uchan Jeon, Seokhyun Choung, and Jeong Woo Han. Catbench framework for benchmarking  
296 machine learning interatomic potentials in adsorption energy predictions for heterogeneous catalysis. *Cell  
297 Reports Physical Science*, 6(12), 2025.
- 298 [26] Seokhyun Choung, Miyeon Kim, Jinuk Moon, and Jeong Woo Han. From atomic motif to realistic single  
299 atom catalysts through machine learning interatomic potentials. *ACS Energy Letters*, 10(12):6288–6296,  
300 2025.
- 301 [27] LangChain AI. Langchain: License (mit license). GitHub repository, . URL <https://github.com/langchain-ai/langchain/blob/master/LICENSE>. License: MIT. Accessed 2026-01-31.
- 303 [28] LangChain AI. Langgraph: License (mit license). GitHub repository, . URL <https://github.com/langchain-ai/langgraph/blob/main/LICENSE>. License: MIT. Accessed 2026-01-31.
- 305 [29] FAIR-Chem. Fair-chem: License (mit license). GitHub repository. URL <https://github.com/FAIR-Chem/fairchem/blob/main/LICENSE.md>. License: MIT (software). Accessed 2026-01-31.
- 307 [30] Meta AI and FAIR Chemistry. Uma model weights: License (fair chemistry license v1). Hugging Face  
308 model repository (facebook/UMA). URL <https://huggingface.co/facebook/UMA/blob/main/LICENSE>. License: FAIR Chemistry License v1 (model weights). Accessed 2026-01-31.
- 310 [31] OpenAI. Openai services agreement. OpenAI Policies. URL <https://openai.com/policies/services-agreement/>. Terms governing use of OpenAI services for businesses/developers (incl. API).  
311 Accessed 2026-01-31.
- 313 [32] Google. Gemini api additional terms of service. Google AI for Developers. URL <https://ai.google.dev/gemini-api/terms>. Additional terms for Gemini API usage. Accessed 2026-01-31.
- 315 [33] Anthropic. Anthropic commercial terms of service. Anthropic Legal. URL <https://www.anthropic.com/legal/commercial-terms>. Terms governing use of Anthropic API keys and related commercial  
316 offerings. Accessed 2026-01-31.
- 317

318 **A Appendix**

319 **A.1 List of Assets**

- 320 • **LangChain** [27] [MIT License]
- 321 • **LangGraph** [28] [MIT License]
- 322 • **FAIR-Chem / fairchem-core** [29] [MIT License]
- 323 • **UMA pretrained potential (uma-s-1p1)** [30] [FAIR Chemistry License v1]
- 324 • **OpenAI API + GPT-family model endpoints** [31] [Proprietary (API Terms of Use)]
- 325 • **Google Gemini API endpoints** [32] [Proprietary (API Terms of Service)]
- 326 • **Anthropic Claude API endpoints** [33] [Proprietary (API Terms of Service)]

327 **A.2 Compositional distribution and energetic profiles of catalysts identified via CatAgent  
328 with critic**

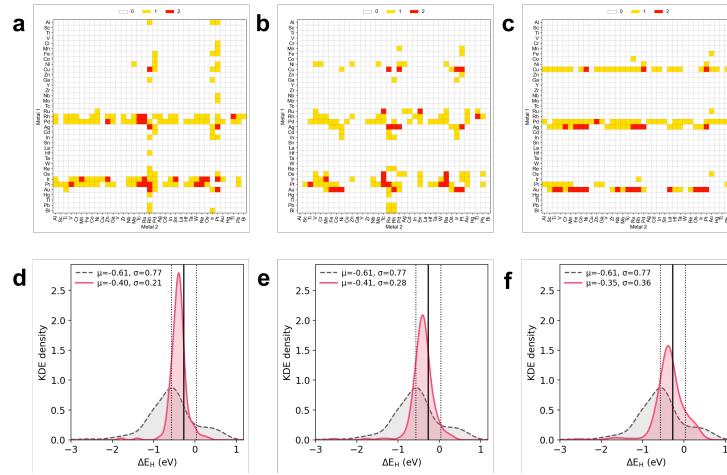


Figure 5: Catalyst discovery patterns for Gemini 3 Flash using critic-enabled CatAgent across three independent trials. (a-c) composition heatmaps of successful candidates; (d-f) kernel density estimates of adsorption energy distributions for all explored candidates. (a,d) Trial 1, (b,e) Trial 2, (c,f) Trial 3.

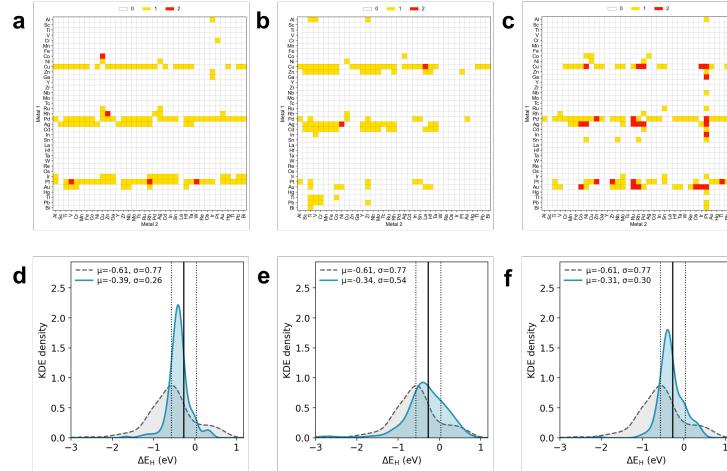


Figure 6: Catalyst discovery patterns for GPT-5.2 using critic-enabled CatAgent across three independent trials. (a-c) composition heatmaps of successful candidates; (d-f) kernel density estimates of adsorption energy distributions for all explored candidates. (a,d) Trial 1, (b,e) Trial 2, (c,f) Trial 3.

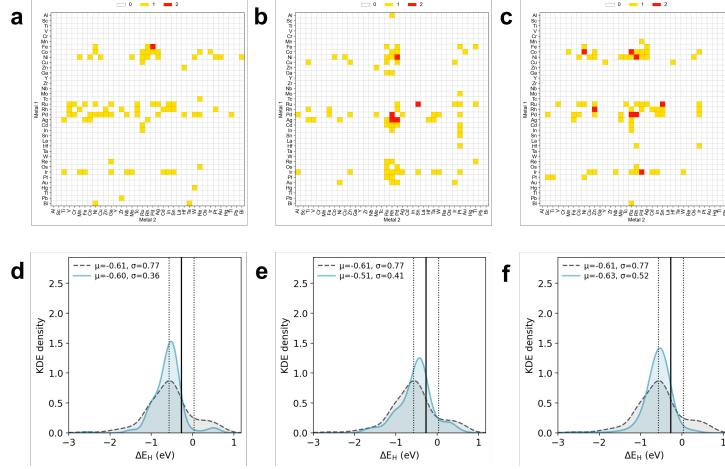


Figure 7: Catalyst discovery patterns for GPT-4o mini using critic-enabled CatAgent across three independent trials. (a-c) composition heatmaps of successful candidates; (d-f) kernel density estimates of adsorption energy distributions for all explored candidates. (a,d) Trial 1, (b,e) Trial 2, (c,f) Trial 3.

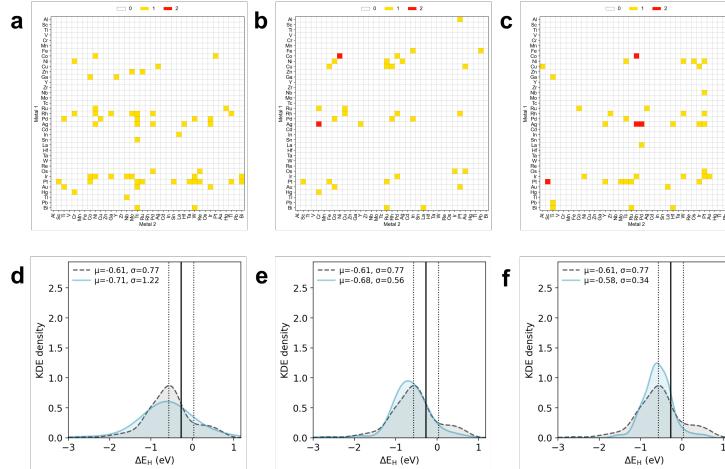


Figure 8: Catalyst discovery patterns for GPT-3.5 Turbo using critic-enabled CatAgent across three independent trials. (a-c) composition heatmaps of successful candidates; (d-f) kernel density estimates of adsorption energy distributions for all explored candidates. (a,d) Trial 1, (b,e) Trial 2, (c,f) Trial 3.

329 **AI Co-Scientist Challenge Korea Paper Checklist**

330 **1. Claims**

331 Question: Do the main claims made in the abstract and introduction accurately reflect the  
332 paper's contributions and scope?

333 Answer: [Yes]

334 Justification: The abstract and introduction correctly frame the core contribution — CatAgent  
335 (multi-agent iterative workflow) and a 13-model single-shot benchmark over bimetallic  
336 space — consistent with the Introduction's bullet-point contributions.

337 Guidelines:

- 338 • The answer NA means that the abstract and introduction do not include the claims  
339 made in the paper.
- 340 • The abstract and/or introduction should clearly state the claims made, including the  
341 contributions made in the paper and important assumptions and limitations. A No or  
342 NA answer to this question will not be perceived well by the reviewers.
- 343 • The claims made should match theoretical and experimental results, and reflect how  
344 much the results can be expected to generalize to other settings.
- 345 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
346 are not attained by the paper.

347 **2. Limitations**

348 Question: Does the paper discuss the limitations of the work performed by the authors?

349 Answer: [Yes]

350 Justification: In the final paragraph of the Conclusion, the paper explicitly acknowledges  
351 limitations and open gaps—most notably that the search space is relatively constrained  
352 (limited compositions and phases), that domain-specific fine-tuning is not explored, and that  
353 the current framework does not yet address multi-objective optimization, framing these as  
354 clear directions for advancement.

355 Guidelines:

- 356 • The answer NA means that the paper has no limitation while the answer No means that  
357 the paper has limitations, but those are not discussed in the paper.
- 358 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 359 • The paper should point out any strong assumptions and how robust the results are to  
360 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
361 model well-specification, asymptotic approximations only holding locally). The authors  
362 should reflect on how these assumptions might be violated in practice and what the  
363 implications would be.
- 364 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
365 only tested on a few datasets or with a few runs. In general, empirical results often  
366 depend on implicit assumptions, which should be articulated.
- 367 • The authors should reflect on the factors that influence the performance of the approach.  
368 For example, a facial recognition algorithm may perform poorly when image resolution  
369 is low or images are taken in low lighting. Or a speech-to-text system might not be  
370 used reliably to provide closed captions for online lectures because it fails to handle  
371 technical jargon.
- 372 • The authors should discuss the computational efficiency of the proposed algorithms  
373 and how they scale with dataset size.
- 374 • If applicable, the authors should discuss possible limitations of their approach to  
375 address problems of privacy and fairness.
- 376 • While the authors might fear that complete honesty about limitations might be used by  
377 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
378 limitations that aren't acknowledged in the paper. The authors should use their best  
379 judgment and recognize that individual actions in favor of transparency play an impor-  
380 tant role in developing norms that preserve the integrity of the community. Reviewers  
381 will be specifically instructed to not penalize honesty concerning limitations.

382     **3. Theory Assumptions and Proofs**

383     Question: For each theoretical result, does the paper provide the full set of assumptions and  
384     a complete (and correct) proof?

385     Answer: [N/A]

386     Justification: The paper is an empirical workflow/benchmark study.

387     Guidelines:

- 388       • The answer NA means that the paper does not include theoretical results.
- 389       • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
390        referenced.
- 391       • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 392       • The proofs can either appear in the main paper or the supplemental material, but if  
393        they appear in the supplemental material, the authors are encouraged to provide a short  
394        proof sketch to provide intuition.
- 395       • Inversely, any informal proof provided in the core of the paper should be complemented  
396        by formal proofs provided in appendix or supplemental material.
- 397       • Theorems and Lemmas that the proof relies upon should be properly referenced.

398     **4. Experimental Result Reproducibility**

399     Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
400     perimental results of the paper to the extent that it affects the main claims and/or conclusions  
401     of the paper (regardless of whether the code and data are provided or not)?

402     Answer: [Yes]

403     Justification: The Methodology section points to an anonymous artifact that contains the  
404     exact model snapshot identifiers, prompts, codes, and run commands needed to reproduce  
405     the reported figures. Detailed information can be found at [https://osf.io/t8he3/overview?view\\_only=1255c0a1cb1646bd9616d634314cd5e9](https://osf.io/t8he3/overview?view_only=1255c0a1cb1646bd9616d634314cd5e9).

407     Guidelines:

- 408       • The answer NA means that the paper does not include experiments.
- 409       • If the paper includes experiments, a No answer to this question will not be perceived  
410        well by the reviewers: Making the paper reproducible is important, regardless of  
411        whether the code and data are provided or not.
- 412       • If the contribution is a dataset and/or model, the authors should describe the steps taken  
413        to make their results reproducible or verifiable.
- 414       • Depending on the contribution, reproducibility can be accomplished in various ways.  
415        For example, if the contribution is a novel architecture, describing the architecture fully  
416        might suffice, or if the contribution is a specific model and empirical evaluation, it may  
417        be necessary to either make it possible for others to replicate the model with the same  
418        dataset, or provide access to the model. In general, releasing code and data is often  
419        one good way to accomplish this, but reproducibility can also be provided via detailed  
420        instructions for how to replicate the results, access to a hosted model (e.g., in the case  
421        of a large language model), releasing of a model checkpoint, or other means that are  
422        appropriate to the research performed.
- 423       • While AI Co-Scientist Challenge Korea does not require releasing code, the conference  
424        does require all submissions to provide some reasonable avenue for reproducibility,  
425        which may depend on the nature of the contribution. For example
  - 426           (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
427            to reproduce that algorithm.
  - 428           (b) If the contribution is primarily a new model architecture, the paper should describe  
429            the architecture clearly and fully.
  - 430           (c) If the contribution is a new model (e.g., a large language model), then there should  
431            either be a way to access this model for reproducing the results or a way to reproduce  
432            the model (e.g., with an open-source dataset or instructions for how to construct  
433            the dataset).

434 (d) We recognize that reproducibility may be tricky in some cases, in which case  
435 authors are welcome to describe the particular way they provide for reproducibility.  
436 In the case of closed-source models, it may be that access to the model is limited in  
437 some way (e.g., to registered users), but it should be possible for other researchers  
438 to have some path to reproducing or verifying the results.

439 **5. Open access to data and code**

440 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
441 tions to faithfully reproduce the main experimental results, as described in supplemental  
442 material?

443 Answer: [Yes]

444 Justification: The codes, full prompts, model snapshots, and raw data are available at [https://osf.io/t8he3/overview?view\\_only=1255c0a1cb1646bd9616d634314cd5e9](https://osf.io/t8he3/overview?view_only=1255c0a1cb1646bd9616d634314cd5e9).

446 Guidelines:

- 447 • The answer NA means that paper does not include experiments requiring code.
- 448 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 449 • While we encourage the release of code and data, we understand that this might not be  
450 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
451 including code, unless this is central to the contribution (e.g., for a new open-source  
452 benchmark).
- 453 • The instructions should contain the exact command and environment needed to run to  
454 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 455 • The authors should provide instructions on data access and preparation, including how  
456 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 457 • The authors should provide scripts to reproduce all experimental results for the new  
458 proposed method and baselines. If only a subset of experiments are reproducible, they  
459 should state which ones are omitted from the script and why.
- 460 • At submission time, to preserve anonymity, the authors should release anonymized  
461 versions (if applicable).
- 462 • Providing as much information as possible in supplemental material (appended to the  
463 paper) is recommended, but including URLs to data and code is permitted.

466 **6. Experimental Setting/Details**

467 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
468 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
469 results?

470 Answer: [Yes]

471 Justification: The paper reports the key experimental settings in Methodology, in-  
472 cluding the candidate budget (200), iterative schedule (20 steps  $\times$  10), UMA check-  
473 point (uma-s-1p1). The codes, experimental settings, full prompts, model snap-  
474 shots, and raw data are available at [https://osf.io/t8he3/overview?view\\_only=1255c0a1cb1646bd9616d634314cd5e9](https://osf.io/t8he3/overview?view_only=1255c0a1cb1646bd9616d634314cd5e9).

476 Guidelines:

- 477 • The answer NA means that the paper does not include experiments.
- 478 • The experimental setting should be presented in the core of the paper to a level of detail  
479 that is necessary to appreciate the results and make sense of them.
- 480 • The full details can be provided either with the code, in appendix, or as supplemental  
481 material.

482 **7. Experiment Statistical Significance**

483 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
484 information about the statistical significance of the experiments?

485 Answer: [Yes]

486 Justification: Uncertainty reporting is defined explicitly: all statistics from triplicate runs are  
487 reported as mean  $\pm$  standard deviation, and the main quantitative figures (Fig. 2 and Fig. 3)  
488 state this in their captions.

489 Guidelines:

- 490 • The answer NA means that the paper does not include experiments.  
491 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
492 dence intervals, or statistical significance tests, at least for the experiments that support  
493 the main claims of the paper.  
494 • The factors of variability that the error bars are capturing should be clearly stated (for  
495 example, train/test split, initialization, random drawing of some parameter, or overall  
496 run with given experimental conditions).  
497 • The method for calculating the error bars should be explained (closed form formula,  
498 call to a library function, bootstrap, etc.)  
499 • The assumptions made should be given (e.g., Normally distributed errors).  
500 • It should be clear whether the error bar is the standard deviation or the standard error  
501 of the mean.  
502 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
503 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
504 of Normality of errors is not verified.  
505 • For asymmetric distributions, the authors should be careful not to show in tables or  
506 figures symmetric error bars that would yield results that are out of range (e.g. negative  
507 error rates).  
508 • If error bars are reported in tables or plots, The authors should explain in the text how  
509 they were calculated and reference the corresponding figures or tables in the text.

## 510 8. Experiments Compute Resources

511 Question: For each experiment, does the paper provide sufficient information on the com-  
512 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
513 the experiments?

514 Answer: [Yes]

515 Justification: Detailed information can be found at [https://osf.io/t8he3/overview?  
516 view\\_only=1255c0a1cb1646bd9616d634314cd5e9](https://osf.io/t8he3/overview?view_only=1255c0a1cb1646bd9616d634314cd5e9).

517 Guidelines:

- 518 • The answer NA means that the paper does not include experiments.  
519 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
520 or cloud provider, including relevant memory and storage.  
521 • The paper should provide the amount of compute required for each of the individual  
522 experimental runs as well as estimate the total compute.  
523 • The paper should disclose whether the full research project required more compute  
524 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
525 didn't make it into the paper).

## 526 9. Code Of Ethics

527 Question: Does the research conducted in the paper conform, in every respect, with the  
528 NeurIPS Code of Ethics <https://nips.cc/public/EthicsGuidelines>?

529 Answer: [Yes]

530 Justification: The work is a computational study of LLM-guided catalyst screening and does  
531 not involve human subjects, crowdsourcing, or personally identifiable data, and it uses stan-  
532 dard scientific simulations without deceptive, surveillance, or discriminatory applications.  
533 In line with the NeurIPS Code of Ethics, the project's foreseeable impacts are predominantly  
534 beneficial (accelerating materials discovery).

535 Guidelines:

- 536 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
537 • If the authors answer No, they should explain the special circumstances that require a  
538 deviation from the Code of Ethics.

- 539           • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
540           eration due to laws or regulations in their jurisdiction).

541           **10. Broader Impacts**

542           Question: Does the paper discuss both potential positive societal impacts and negative  
543           societal impacts of the work performed?

544           Answer: [N/A]

545           Justification: The contribution of this work is methodological and empirical—benchmarking  
546           LLM-assisted workflows for electrocatalyst discovery.

547           Guidelines:

- 548           • The answer NA means that there is no societal impact of the work performed.
- 549           • If the authors answer NA or No, they should explain why their work has no societal  
550           impact or why the paper does not address societal impact.
- 551           • Examples of negative societal impacts include potential malicious or unintended uses  
552           (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
553           (e.g., deployment of technologies that could make decisions that unfairly impact specific  
554           groups), privacy considerations, and security considerations.
- 555           • The conference expects that many papers will be foundational research and not tied  
556           to particular applications, let alone deployments. However, if there is a direct path to  
557           any negative applications, the authors should point it out. For example, it is legitimate  
558           to point out that an improvement in the quality of generative models could be used to  
559           generate deepfakes for disinformation. On the other hand, it is not needed to point out  
560           that a generic algorithm for optimizing neural networks could enable people to train  
561           models that generate Deepfakes faster.
- 562           • The authors should consider possible harms that could arise when the technology is  
563           being used as intended and functioning correctly, harms that could arise when the  
564           technology is being used as intended but gives incorrect results, and harms following  
565           from (intentional or unintentional) misuse of the technology.
- 566           • If there are negative societal impacts, the authors could also discuss possible mitigation  
567           strategies (e.g., gated release of models, providing defenses in addition to attacks,  
568           mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
569           feedback over time, improving the efficiency and accessibility of ML).

570           **11. Safeguards**

571           Question: Does the paper describe safeguards that have been put in place for responsible  
572           release of data or models that have a high risk for misuse (e.g., pretrained language models,  
573           image generators, or scraped datasets)?

574           Answer: [N/A]

575           Justification: This work does not release any high-risk generative model weights or scraped  
576           user data; it releases only workflow code, configuration files, and precomputed catalyst-  
577           screening outputs. Access to external LLMs and UMA is mediated through their providers'  
578           APIs or gated checkpoints under existing terms.

579           Guidelines:

- 580           • The answer NA means that the paper poses no such risks.
- 581           • Released models that have a high risk for misuse or dual-use should be released with  
582           necessary safeguards to allow for controlled use of the model, for example by requiring  
583           that users adhere to usage guidelines or restrictions to access the model or implementing  
584           safety filters.
- 585           • Datasets that have been scraped from the Internet could pose safety risks. The authors  
586           should describe how they avoided releasing unsafe images.
- 587           • We recognize that providing effective safeguards is challenging, and many papers do  
588           not require this, but we encourage authors to take this into account and make a best  
589           faith effort.

590           **12. Licenses for existing assets**

591 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
592 the paper, properly credited and are the license and terms of use explicitly mentioned and  
593 properly respected?

594 Answer: [Yes]

595 Justification: The paper (A.1 List of Assets) credits all non-original assets that are essential  
596 to the study.

597 Guidelines:

- 598 • The answer NA means that the paper does not use existing assets.
- 599 • The authors should cite the original paper that produced the code package or dataset.
- 600 • The authors should state which version of the asset is used and, if possible, include a  
601 URL.
- 602 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 603 • For scraped data from a particular source (e.g., website), the copyright and terms of  
604 service of that source should be provided.
- 605 • If assets are released, the license, copyright information, and terms of use in the  
606 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
607 has curated licenses for some datasets. Their licensing guide can help determine the  
608 license of a dataset.
- 609 • For existing datasets that are re-packaged, both the original license and the license of  
610 the derived asset (if it has changed) should be provided.
- 611 • If this information is not available online, the authors are encouraged to reach out to  
612 the asset's creators.

### 613 13. New Assets

614 Question: Are new assets introduced in the paper well documented and is the documentation  
615 provided alongside the assets?

616 Answer: [Yes]

617 Justification: The new assets introduced by the paper (the CatAgent workflow code, the  
618 released raw experimental outputs, and the baseline dataset) are documented in <https://osf.io/t8he3/>?view\_only=1255c0a1cb1646bd9616d634314cd5e9.

620 Guidelines:

- 621 • The answer NA means that the paper does not release new assets.
- 622 • Researchers should communicate the details of the dataset/code/model as part of their  
623 submissions via structured templates. This includes details about training, license,  
624 limitations, etc.
- 625 • The paper should discuss whether and how consent was obtained from people whose  
626 asset is used.
- 627 • At submission time, remember to anonymize your assets (if applicable). You can either  
628 create an anonymized URL or include an anonymized zip file.

### 629 14. Crowdsourcing and Research with Human Subjects

630 Question: For crowdsourcing experiments and research with human subjects, does the paper  
631 include the full text of instructions given to participants and screenshots, if applicable, as  
632 well as details about compensation (if any)?

633 Answer: [N/A]

634 Justification: The study does not involve crowdsourcing or human-subject experiments.

635 Guidelines:

- 636 • The answer NA means that the paper does not involve crowdsourcing nor research with  
637 human subjects.
- 638 • Including this information in the supplemental material is fine, but if the main contribu-  
639 tion of the paper involves human subjects, then as much detail as possible should be  
640 included in the main paper.

- 641           • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
642           or other labor should be paid at least the minimum wage in the country of the data  
643           collector.

644       **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
645           Subjects**

646       Question: Does the paper describe potential risks incurred by study participants, whether  
647           such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
648           approvals (or an equivalent approval/review based on the requirements of your country or  
649           institution) were obtained?

650       Answer: [N/A]

651       Justification: This work does not involve human subjects or study participants.

652       Guidelines:

- 653           • The answer NA means that the paper does not involve crowdsourcing nor research with  
654           human subjects.
- 655           • Depending on the country in which research is conducted, IRB approval (or equivalent)  
656           may be required for any human subjects research. If you obtained IRB approval, you  
657           should clearly state this in the paper.
- 658           • We recognize that the procedures for this may vary significantly between institutions  
659           and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
660           guidelines for their institution.
- 661           • For initial submissions, do not include any information that would break anonymity (if  
662           applicable), such as the institution conducting the review.