Hierarchical Deep Research with Local-Web RAG: Toward Automated System-Level Materials Discovery

Rui Ding*

Pritzker School of Molecular Engineering University of Chicago Chicago, IL 60637 ruiding@uchicago.edu

Yuxin Chen

Department of Computer Science University of Chicago Chicago, IL 60637 yuxinchen@uchicago.edu

Junhong Chen†

Pritzker School of Molecular Engineering University of Chicago Chicago, IL 60637 junhongchen@uchicago.edu

Abstract

We present a long-horizon, hierarchical deep research (DR) agent designed for complex materials and device discovery problems that exceed the scope of existing ML surrogates or closed-source commercial agents. Our framework instantiates a locally deployable DR instance that integrates local retrieval-augmented generation (RAG) with large language model (LLM) reasoners, enhanced by a Deep Tree of Research (DToR) mechanism that adaptively expands and prunes research branches to maximize coverage, depth, and coherence. We evaluate across 21 nanomaterials/device topics using an *LLM-as-judge* rubric with 5 web-enabled SOTA models as jurors. In addition, we conduct dry-lab validations on five representative tasks, where human experts use domain simulations (e.g., DFT) to verify whether DR-agent proposals are actionable. Results show that our DR agent produces reports with quality comparable to commercial systems (ChatGPT-o3/o4-mini-high, Gemini Deep Research) at substantially lower cost, while enabling on-prem integration with local data and tools.

1 Introduction

Machine learning for nanomaterials and device discovery has advanced from tractable single-entity predictions to increasingly complex, system-level challenges. At the *molecular or crystal level* (S1), modern surrogates routinely predict intrinsic properties [1]. *Small assemblies* (S2)—interfaces, heterostructures, or

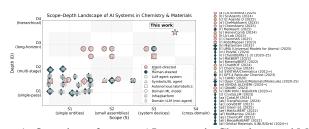


Figure 1: Overview of current AI systems in Chemistry and Materials in S (Scope)/D (Depth) landscape

simple 2–3-component devices—introduce coupling effects yet remain manageable. Real-world nano-devices, however, operate at *system scope* (S3): multi-layer perovskite stacks with interdependent layers, core–shell–doped catalysts with engineered interfaces, or nano-architected battery

^{*}Also Chemical Sciences and Engineering Division, Argonne National Laboratory, Lemont, IL 60439

[†] Also Chemical Sciences and Engineering Division, Argonne National Laboratory, Lemont, IL 60439

electrodes where transport across length scales governs performance [2]. The ultimate challenge lies in *cross-domain integration* (S4): platforms that span material families, recognition chemistries, and transduction modalities—for example, engineering 2D-material field-effect transistor (FET) sensors with aptamer functionalization for environmental detection [3, 4].

Meeting this escalation in scope requires a matching escalation in depth. We frame depth from *single-pass* prediction/forward screening (D1), to *multi-stage pipelines* with structured outputs (D2), to *long-horizon synthesis* with iterative planning (D3), and ultimately *hierarchical, breadth-depth* exploration (D4). As shown in Fig. 1, while classical high-throughput density functional theory (DFT)/molecular dynamics (MD) and data-driven screening deliver important wins at S1–S2 / D1–D2, they routinely struggle at S3–S4 where interfacial chemistry, kinetic pathways, and manufacturing constraints are dominated by cascading multiscale interactions [5–8] (see Appendix B).

To address this gap, we introduce a hierarchical deep-research (DR) agent purpose-built for long-horizon scientific inquiry. Given a natural-language inquiry, the system instantiates a locally deployable DR instance that iterates between local and web retrieval-augmented generation (RAG) and large language model (LLM) reasoners (summarizers, reflectors). We further expand the DR instance into a Deep Tree of Research (DToR) where each instance is treated as a node in a tree-structured workflow. DToR adaptively expands/prunes branches to improve coverage, depth, and coherence. The agent is self-directed/navigated (plans its own research path), runs locally for controllability and privacy, and targets S3–S4 hypotheses. Our contributions are

- 1. Methodology: DToR—an open, on-prem framework that casts long-horizon inquiry as a resource-bounded tree over research nodes (single-path DR instances), combining local-first RAG, diversity-aware complementary queries, gap-driven expansion/pruning, and provenance-centric synthesis. Runs with Ollama/LMStudio for controllability and privacy, democratizing DR in facing scientific problems at scopes S3–S4.
- 2. Evaluation: LLM-as-jury across 21 topics, a locally served open-source 32B LLM on consumer-grade devices driving DToR could achieve report quality comparable to leading commercial systems (DR conducted by ChatGPT-o3/o4-mini-high) while avoiding closed APIs, with consistent rankings from five web-enabled judges (rubric scoring and dueling). As for dry-lab validation in 5 representative applications, physics-based simulations on agent-proposed candidates match or surpass those from the best commercial baselines across multiple task metrics, demonstrating actionable, physically grounded hypotheses.

Related work Recent work on structured reasoning for LLMs has moved beyond linear prompting toward explicit search structures. Tree-of-Thoughts (ToT) casts inference as a branching search over intermediate 'thoughts', enabling exploration and backtracking [9]. Graph-of-Thoughts (GoT) generalizes the search space to arbitrary graphs to capture non-linear dependencies and coordination across subproblems [10, 11]. In parallel, reason—act and self-improvement paradigms—ReAct, Reflexion, and Self-Refine—provide strong baselines for multi-step decision-making and revision ([12–14]). Foundational web-agent environments/benchmarks—BrowseComp, WebLINX, WebArena, and BrowserGym—formalize browsing/navigation competencies and support reproducible evaluation [15–18]. Work on retrieval is also moving toward more adaptive strategies: Self-RAG [19] learns when to retrieve and how to critique generations; Corrective RAG (CRAG) [20] assesses retrieval quality and escalates to web search when local evidence is weak; and GraphRAG [21] leverages graph-structured knowledge to guide retrieval beyond flat vector search. Despite this progress, there is no framework that unifies structured branching and adaptive retrieval in a way that scales under practical compute and data limits for domain-specific research settings.

2 Method

D4 requires structured and budget-aware planning beyond short-term queries. We introduce DToR, a breadth-then-depth controller for DR that starts from multiple parallel perspectives, selectively expands or prunes based on explicit gap checks, runs evidence-first research steps (local RAG with adapted web search), and produces source-grounded reports.

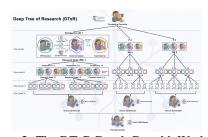


Figure 2: The DToR Depth-Breadth Workflow

Single DR Instance As shown in Fig. S1 (with comparison against common standard DR agent [22]), each DR instance is an evidence-first loop that (1) generates a search query; (2) performs local RAG; (3) summarizes local evidence; (4) generates a complementary query to induce topical diversity; (5) conducts web research for both queries; (6) integrates/updates previous retrieved segments in a running summary; (7) reflects to propose follow-up query. The process would loop until a set maximum round before (8) finalizing. Three design choices distinguish this instance from common local DR agents [22] which only conduct (1)-(5)-(6)-(7): local-first retrieval (before web) to reduce drift, hallucination, and surface domain priors; diversity-aware query generation (local summary derived complementary query) to guarantee comprehensiveness; and robust I/O discipline to prevent controller stalls on local-hosted LLMs. Compared with classic ReAct/Reflexion-style linear loops, this instance adds an explicit breadth mechanism and escalation policy akin to adaptive/corrective-RAG, while remaining compatible with the drop-in alternative RAG methods.

DToR: from perspectives to gap-driven expansion Based on the single DR instance, we treat each DR instance as a Research Node (RN) inside a branch-and-bound controller (Fig. 2, see Fig. S2 for full size schematic, Appendix C for pseudocode). Starting from a user query, a diversifier proposes several orthogonal Perspectives (P1, P2...)—each seeds one branch with a budget: maximum depth, nodes per branch, and total branches. For the active branch, the router (i) runs the next pending RN (the single DR instance), producing a RN summary; (ii) invokes analyst based on remaining depth/budget/current RN report quality to give a decision EXPAND, PRUNE; and, if EXPAND with budget remaining, (iii) invokes knowledge gap explorer to materializes new RNs from this RN targeted to those gaps (depth+1) and updates remaining budget. Branches that stagnate are pruned; (iv) branches that reach depth or converge are synthesized into a perspective report. When all branches reach synthesis, a final synthesizer reconciles cross-branch evidence, resolves conflicts, and outputs a provenance-rich, consolidated report. Conceptually, this is a tree-of-search over (query, evidence, summary) states rather than over purely symbolic "thoughts," inheriting the benefits of structured exploration/backtracking from tree/graph-based reasoning while remaining tailored to information-seeking and citation fidelity under explicit resource budgets.

3 Experiments

Experimental setup We evaluate across 21 human-expert—crafted, independent topics spanning multiple domains (e.g., sensors, batteries, catalysts, as detailed in Appendix D.1). For each topic, we consider 27 agents drawn from 9 commercial solutions, and 18 local settings deployed locally on consumer-level hardware (9 single DR-instance settings, and 9 DToR-mode settings, details in Appendix D.2). Using an LLM-as-judge protocol (details in Appendix D.3), we generate 27 distinct report variants per topic and have five web-enabled SOTA judges—Claude 4 Opus (thinking), Gemini 2.5 Pro, Grok-3 (thinking), ChatGPT-o3, and ChatGPT-o4-mini-high—score each report under a fixed, double-blind rubric with five equally weighted dimensions: (1) relevance, (2) depth, (3) clarity, (4) applicability, and (5) novelty. Judges assign 0–10 per dimension and provide a one-sentence justification; numeric ratings are extracted programmatically. We repeat the evaluation three times per topic-report-judge combination, yielding $21 \times 27 \times 5 \times 3 = 8,505$ rubric-based judgments. In a second stage, we run A/B dueling preferences using the same setup. For each topic, we rank 27 reports within commercial and local pools by rubric average, select the top-3 from each pool, and conduct all pairwise matchups (C(6,2)=15 pairs per topic). Duels are double-blind and orderrandomized (A vs. B) and use a concise "experienced materials scientist" prompt emphasizing clarity, gap prioritization, and actionable next steps. Each duel is repeated three times, yielding $21 \times 15 \times 3 = 945A/B$ trials and 4,725 individual judge decisions. We also conduct dry-lab validations to test whether DR-agent proposals are actionable. From the 21 topics, we select five representative tasks—per- and polyfluoroalkyl substances (PFAS) sensor probe, PFAS degradation catalysts, battery binder selection, oxygen evolution catalyst stability, and CO₂ sensor probe. For each, we take the best local and best commercial agent and deliver reports anonymously to a human expert, who extracts candidates and builds atomic-level simulation environments (details in Appendix E).

Experimental results

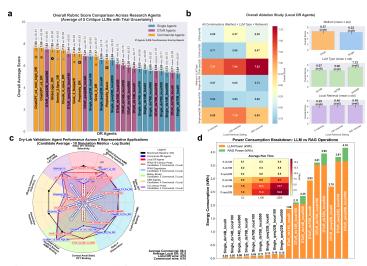


Figure 3: (a) Overall five-dimension rubric scores across 27 agents (21 topics; 5 judges; 3 trials). (b) Full-factorial ablation (Method × LLM × Local retrieval); heatmap shows overall means; side bars show each factor impact. (c) Dry-lab validation on five tasks: radar of 10 simulation metrics comparing best local vs. best commercial DR Agents. (d) Local deployment cost: stacked energy (kWh) split into LLM vs. RAG; inset = average runtime (hour).

Across the five LLM-as-judge scoring dimensions (Fig. 3a). the locally deployed agent DToR_qwq32B_local500³ achieved an average score 7.83/10, ranking 3rd of 27 agents—just below ChatGPT-o3 (7.98/10) and ChatGPT-o4-mini-high (7.96/10)—and outperforming other widely used commercial DR systems, including Gemini 2.5 Pro (7.78/10), Claude Opus 4 (7.73/10), Grok 3 (7.43/10), and Perplexity (7.40/10). resource-reduced settings, DToR_qwq32B_local100 (7.44/10) ranked 7th and single DR instance, DTor_qwq32B_local0 (7.21/10), ranked 10th. Upon detailed look on each of the five dimensions (Fig. S3), we

could observe that for the

most challenging dimensions of *novelty* and *depth*, though our local agent could not make to top-3, the difference between DToR and single DR instance have shown our DToR strategy could manage to narrow the gap. For instance, DToR enables qwq32B_local500 to increase from an average of 5.04/10 to 6.84/10 in *depth*, and from 8.11/10 to 8.61/10 in *novelty*.

Inter-rater consistency across the five LLM critics was high: average Pearson correlation p=0.96 on raw scores and Spearman $\rho=0.95$ on rank orders (Fig. S4). Reliability is further supported by agent-specific and inter-trial consistency (Fig. S5, Fig. S6). From ablation study perspective (Fig. 3b), enabling local RAG (6.49/10) and especially DToR mode (6.57/10) would increase the average score. Pairwise dueling (Fig. S7) confirms this: DToR_qwq32B_local500 achieves a 74% mean win rate across head-to-head comparisons on all 21 topics with great consistency (Fig. S8).

In dry-lab validation spanning 10 metrics over five application scenarios (Fig. 3c), the commercial baseline leads on 6 out of 10 metrics when judged by per-metric averages, but aggregating across metrics the local agent achieves a higher overall score (92.9 vs. 88.4). Although both approaches still trail domain-expert benchmarks (e.g., beta-cyclodextrin [23] for selective PFAS binding), the DToR workflow yields recommendations competitive with—and sometimes exceeding—commercial systems. Fig. 3d shows an average of only 4.16 kWh per completed report. Without high subscription fee (Fig. S10) for commercial solutions, this makes DToR practical for nanomaterials/device researchers without large compute or API budgets, as small open-source models on consumer hardware keep energy modest and runtime tunable.

4 Conclusion

We introduced DToR to enable a locally deployable DR agent with a resource-bounded tree controller that unifies local-first RAG, diversity-aware query generation, and gap-driven branching into provenance-rich, long-horizon information seeking. Across 21 topics and 27 agents, under 8,505 rubric judgments from 5 web-enabled judges, our best local model (DToR_qwq32B_local500) ranked 3rd (7.83/10) with a 74% mean duel win rate (4,725 dueling judgement) while using only 4.16 kWh/report. In dry-lab validation over five applications, its candidates also achieved close or better quality compared with the best commercial solutions. This work underscores actionable competitive guidance and positioning open, democratized, on-prem DR agents as a secure, controllable, and scalable path to automated accelerated scientific discovery at the S3-S4 level.

³Here, DToR_qwq32B_local500 refers to using qwq32B model to drive DToR framework and enable 500 volumes (all) of local corpus for RAG when running DR instances. See details in Appendix B.

References

- [1] Tian Xie and Jeffrey C. Grossman. Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties. *Physical Review Letters*, 120(14):145301, 2018.
- [2] Kristen A. Severson, Peter M. Attia, Norman Jin, Nicholas Perkins, Benben Jiang, Zi Yang, Michael H. Chen, Muratahan Aykol, Patrick K. Herring, Dimitrios Fraggedakis, Martin Z. Bazant, Stephen J. Harris, William C. Chueh, and Richard D. Braatz. Data-driven prediction of battery cycle life before capacity degradation. *Nature Energy*, 4(5):383–391, 2019.
- [3] Nathan J. Szymanski, Bernardus Rendy, Yuxing Fei, Rishi E. Kumar, Tanjin He, David Milsted, Matthew J. McDermott, Max Gallant, Ekin D. Cubuk, Amil Merchant, Haegyeom Kim, Anubhav Jain, Christopher J. Bartel, Kristin Persson, Yan Zeng, and Gerbrand Ceder. An autonomous laboratory for the accelerated synthesis of novel materials. *Nature*, 624:86–91, 2023.
- [4] Shun Wang, Mingyuan Sun, Yunhong Zhang, Hao Ji, Jianwei Gao, Shuai Song, Jun Sun, Hong Liu, Yu Zhang, and Lin Han. Ultrasensitive antibiotic perceiving based on aptamer-functionalized ultraclean graphene field-effect transistor biosensor. *Analytical Chemistry*, 94(42):14785–14793, 2022.
- [5] Joseph H. Montoya, Kirsten T. Winther, Raul A. Flores, Thomas Bligaard, Jens S. Hummelshøj, and Muratahan Aykol. Autonomous intelligent agents for accelerated materials discovery. *Chemical Science*, 11(32):8517–8532, 2020.
- [6] Amil Merchant, Simon Batzner, Samuel S. Schoenholz, Muratahan Aykol, Gowoon Cheon, Ekin Dogus Cubuk, et al. Scaling deep learning for materials discovery. *Nature*, 624:80–85, 2023.
- [7] Lowik Chanussot, Abhishek Das, Siddharth Goyal, Thibaut Lavril, Muhammed Shuaibi, Morgane Riviere, Kevin Tran, Javier Heras-Domingo, Caleb Ho, Weihua Hu, Aini Palizhati, Anuroop Sriram, Brandon M. Wood, Junwoong Yoon, Devi Parikh, C. Lawrence Zitnick, and Zachary W. Ulissi. Open catalyst 2020 (oc20) dataset and community challenges. *ACS Catalysis*, 11:6059–6072, 2021.
- [8] Daniel S. Levine, Muhammed Shuaibi, Evan Walter Clark Spotte-Smith, Michael G. Taylor, Muhammad R. Hasyim, Kyle Michel, Ilyes Batatia, Gábor Csányi, Misko Dzamba, Peter Eastman, Nathan C. Frey, Xiang Fu, Vahe Gharakhanyan, Aditi S. Krishnapriyan, Joshua A. Rackers, Sanjeev Raja, Ammar Rizvi, Andrew S. Rosen, Zachary Ulissi, Santiago Vargas, C. Lawrence Zitnick, Samuel M. Blau, and Brandon M. Wood. The open molecules 2025 (omol25) dataset, evaluations, and models. *arXiv preprint arXiv:2505.08762*, 2025.
- [9] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- [10] Michał Besta et al. Graph of thoughts: Solving elaborate problems with large language models. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pages 17682–17690, 2024.
- [11] Yao Yao, Zuchao Li, and Hai Zhao. Beyond chain-of-thought: Effective graph-of-thought reasoning in language models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2901–2921, Mexico City, Mexico, 2024. Association for Computational Linguistics.
- [12] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, and Karthik Narasimhan. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- [13] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *arXiv* preprint *arXiv*:2303.11366, 2023.
- [14] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651*, 2023.
- [15] J. Wei et al. Browsecomp: A simple yet challenging benchmark for browsing agents. *arXiv* preprint arXiv:2504.12516, 2025.

- [16] Y. Pan et al. Weblinx: Real-world website navigation with multi-turn dialogue. arXiv preprint arXiv:2402.05930, 2024.
- [17] S. Zhou et al. Webarena: A realistic web environment for building autonomous agents. *arXiv* preprint arXiv:2307.13854, 2023.
- [18] Thibault Le Sellier De Chezelles et al. The browsergym ecosystem for web agent research. *arXiv preprint arXiv:2412.05467*, 2024.
- [19] Akari Asai et al. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In Proceedings of the Twelfth International Conference on Learning Representations (ICLR), 2024.
- [20] S.-Q. Yan et al. Corrective retrieval augmented generation. arXiv preprint arXiv:2401.15884, 2024.
- [21] B. Peng et al. Graph retrieval-augmented generation: A survey. arXiv preprint arXiv:2408.08921, 2024.
- [22] LangChain. Local deep researcher. https://github.com/langchain-ai/local-deep-researcher, 2025. GitHub repository. Version pinned at commit 8a5e220 (2025-08-08).
- [23] Yuqin Wang, Hyun-June Jang, Max Topel, Siva Dasetty, Yining Liu, Mohamed Ateia, Aaron Tam, Vepa Rozyyev, Ellie Ouyang, Wen Zhuang, Haihui Pu, Sang Soo Lee, Jeffrey Elam, Andrew Ferguson, Seth Darling, and Junhong Chen. Reversible ppt-level detection of perfluorooctane sulfonic acid in tap water using field-effect transistor sensors. *Nature Water*, 3:1187–1197, 2025.
- [24] Richard Tran, Janice Lan, Muhammed Shuaibi, Brandon M. Wood, Siddharth Goyal, Abhishek Das, Javier Heras-Domingo, Adeesh Kolluru, Ammar Rizvi, Nima Shoghi, Anuroop Sriram, Felix Therrien, Jehad Abed, Oleksandr Voznyy, Edward H. Sargent, Zachary Ulissi, and C. Lawrence Zitnick. The open catalyst 2022 (oc22) dataset and challenges for oxide electrocatalysts. *ACS Catalysis*, 13:3066–3084, 2023.
- [25] Luis Barroso-Luque, Muhammed Shuaibi, Xiang Fu, Brandon M. Wood, Misko Dzamba, Meng Gao, Ammar Rizvi, C. Lawrence Zitnick, and Zachary W. Ulissi. Open materials 2024 (omat24) inorganic materials dataset and models. *arXiv preprint arXiv:2410.12771*, 2024.
- [26] Brandon M. Wood, Misko Dzamba, Xiang Fu, Meng Gao, Muhammed Shuaibi, Luis Barroso-Luque, Kareem Abdelmaqsoud, Vahe Gharakhanyan, John R. Kitchin, Daniel S. Levine, Kyle Michel, Anuroop Sriram, Taco Cohen, Abhishek Das, Ammar Rizvi, Sushree Jagriti Sahoo, Zachary W. Ulissi, and C. Lawrence Zitnick. Uma: A family of universal models for atoms. *arXiv preprint arXiv:2506.23971*, 2025.
- [27] Tanishq Gupta, Mohd Zaki, N. M. Anoop Krishnan, and Mausam. MatSciBERT: A materials domain language model for text mining and information extraction. *npj Computational Materials*, 8:102, 2022.
- [28] Amalie Trewartha, Nicholas Walker, Haoyan Huo, Sanghoon Lee, Kevin Cruse, John Dagdelen, Alexander Dunn, Kristin A. Persson, Gerbrand Ceder, and Anubhav Jain. Quantifying the advantage of domain-specific pre-training on named entity recognition tasks in materials science. *Patterns*, 3:100488, 2022.
- [29] Michiko Yoshitake, Fumitaka Sato, Hiroyuki Kawano, and Hiroshi Teraoka. Materialbert for natural language processing of materials science texts. Science and Technology of Advanced Materials: Methods, 2:372–380, 2022.
- [30] Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. Chemberta: Large-scale self-supervised pretraining for molecular property prediction. arXiv preprint arXiv:2010.09885, 2020.
- [31] Walid Ahmad, Elana Simon, Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. Chemberta-2: Towards chemical foundation models. *arXiv preprint arXiv:2209.01712*, 2022.
- [32] Riya Singh, Aryan Amit Barsainyan, Rida Irfan, Connor Joseph Amorin, Stewart He, Tony Davis, Arun Thiagarajan, Shiva Sankaran, Seyone Chithrananda, Walid Ahmad, Derek Jones, Kevin McLoughlin, Hyojin Kim, Anoushka Bhutani, Shreyas Vinaya Sathyanarayana, Venkat Viswanathan, Jonathan E. Allen, and Bharath Ramsundar. Chemberta-3: An open source training

- framework for chemical foundation models. *ChemRxiv preprint ChemRxiv: 10.26434/chemrxiv-2025-4glrl*, 2025.
- [33] NVIDIA. Megamolbart, 2022. GitHub repository. Release: Aug 30, 2022. Accessed: 2025-08-19.
- [34] Nathan C. Frey, Ryan Soklaski, Simon Axelrod, Siddharth Samsi, Rafael Gómez-Bombarelli, Connor W. Coley, and Vijay Gadepally. Neural scaling of deep chemical models. *nature machine intelligence*, 5:1297–1305, 2023.
- [35] Changwen Xu, Yuyang Wang, and Amir Barati Farimani. Transpolymer: a transformer-based language model for polymer property predictions. *npj Computational Materials*, 9:64, 2023.
- [36] Christopher Kuenneth and Rampi Ramprasad. polybert: a chemical language model to enable fully machine-driven ultrafast polymer informatics. *Nature Communications*, 14:4099, 2023.
- [37] Haoke Qiu, Lunyang Liu, Xuepeng Qiu, Xuemin Dai, Xiangling Ji, and Zhao-Yan Sun. Polync: a natural and chemical language model for the prediction of unified polymer properties. *Chemical Science*, 15:534–544, 2024.
- [38] Yeonghun Kang and Jihan Kim. Chatmof: an artificial intelligence system for predicting and generating metal—organic frameworks using large language models. *Nature Communications*, 15:4705, 2024.
- [39] Luis M. Antunes, Keith T. Butler, and Ricardo Grau-Crespo. Crystal structure generation with autoregressive large language modeling. *Nature Communications*, 15:10570, 2024.
- [40] Claudio Zeni, Robert Pinsler, Daniel Zügner, Andrew Fowler, Matthew Horton, Xiang Fu, Zilong Wang, Aliaksandra Shysheya, Jonathan Crabbé, Shoko Ueda, Roberto Sordillo, Lixin Sun, Jake Smith, Bichlien Nguyen, Hannes Schulz, Sarah Lewis, Chin-Wei Huang, Ziheng Lu, Yichi Zhou, Han Yang, Hongxia Hao, Jielan Li, Chunlei Yang, Wenjie Li, Ryota Tomioka, and Tian Xie. A generative model for inorganic materials design. *Nature*, 639:624–632, 2025.
- [41] Merck KGaA. Synthia® retrosynthesis software. https://www.synthiaonline.com/, 2025. Accessed 2025-08-19.
- [42] Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D. White, and Philippe Schwaller. Augmenting large language models with chemistry tools. *Nature Machine Intelli*gence, 6:525–535, 2024.
- [43] Huan Zhang, Yu Song, Ziyu Hou, Santiago Miret, and Bang Liu. Honeycomb: A flexible llm-based agent system for materials science. *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3369–3382, 2024.
- [44] Alireza Ghafarollahi and Markus J. Buehler. Sciagents: Automating scientific discovery through bioinspired multi-agent intelligent graph reasoning. *Advanced Materials*, 37:2413523, 2025.
- [45] Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici, Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat, Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan, Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R D Costa, José R Penadés, Gary Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan. Towards an ai co-scientist. arXiv preprint arXiv:2502.18864, 2025.
- [46] Nathan J. Szymanski, Bernardus Rendy, Yuxing Fei, Rishi E. Kumar, Tanjin He, David Milsted, Matthew J. McDermott, Max Gallant, E. Dogus Cubuk, Amil Merchant, Haegyeom Kim, Anubhav Jain, Christopher J. Bartel, Kristin Persson, Yan Zeng, and Gerbrand Ceder. An autonomous laboratory for the accelerated synthesis of novel materials. *Nature*, 624:86–91, 2023.
- [47] IBM Research. Ibm roborxn for chemistry. https://rxn.res.ibm.com/rxn/robo-rxn/, 2020.
- [48] Wen Jie Ong, Piero Altoè, Justin S. Smith, Dallas Foster, Melisa Alkan, and Harry Petty. Revolutionizing ai-driven material discovery using nvidia ALCHEMI. *AI4AM* 2025, 2025.
- [49] Hui Lin, Junfeng Niu, Shangtao Liang, Chong Wang, Yujuan Wang, Fangyuan Jin, Qi Luo, Szu-Yin Chiang, and Qingguo Huang. Development of macroporous magnéli phase ti4o7 ceramic materials: As an efficient anode for mineralization of poly- and perfluoroalkyl substances. *Chemical Engineering Journal*, 354:1058–1067, 2018.

- [50] Xuehu Zhong, Junwei Han, Lingling Chen, Wei Liu, Fei Jiao, Hong Zhu, and Wenqing Qin. Binding mechanisms of pvdf in lithium ion batteries. *Applied Surface Science*, 553:149564, 2021.
- [51] Rui Ding, Jianguo Liu, Kang Hua, Xuebin Wang, Xiaoben Zhang, Minhua Shao, Yuxin Chen, and Junhong Chen. Leveraging data mining, active learning, and domain adaptation for efficient discovery of advanced oxygen evolution electrocatalysts. *Science Advances*, 11(14):eadr9038, 2025.
- [52] T. V. K. Karthik, H. Martínez-García, F. Ortiz-Chi, C. G. Espinosa-González, J. G. Torres-Torres, A. G. Hernandez, S. Godavarthi, and M. K. Kesarla. Co2 gas sensing properties of graphitic carbon nitride (g-c3n4) thin films. *Diamond and Related Materials*, 133:109736, 2023.

Hierarchical Deep Research with Local-Web RAG: Toward Automated System-Level Materials Discovery

Technical Appendices and Supplementary Material

A Supplementary Figures

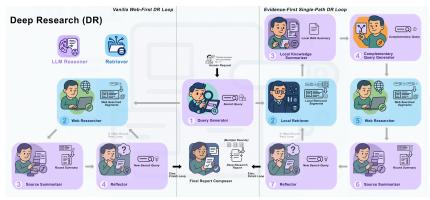


Figure S1: Schematic of the local evidence-first single DR instance in comparison with common web-first DR (Adapted from [22])

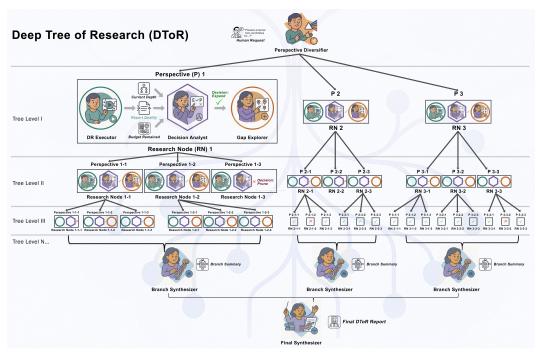


Figure S2: Schematic of the DToR workflow

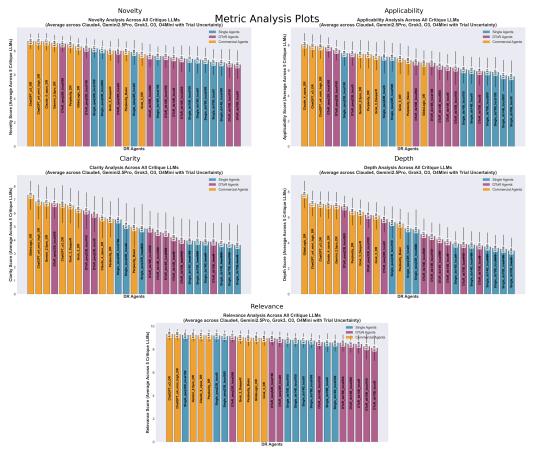


Figure S3: Overall rubric scores across 27 agents on different evaluation dimensions.

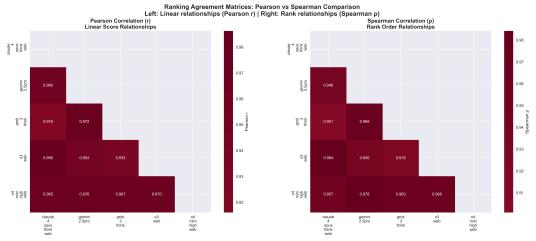
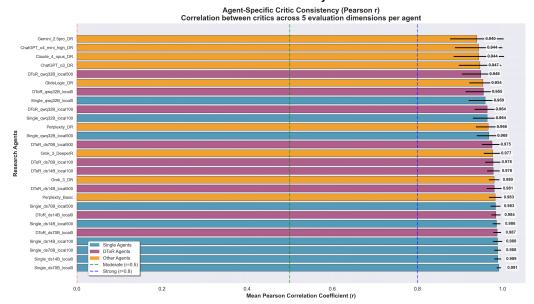


Figure S4: Rubric score agent agreement level over the five critic LLMs, computed by correlation factors.

Pearson Consistency



Spearman Consistency

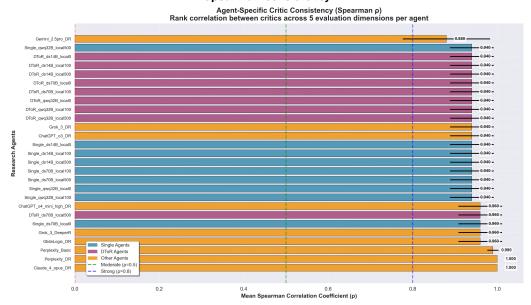


Figure S5: Agent-Specific Critic Consistency Analysis. Pearson correlation coefficients measuring linear agreement between critique LLMs across 5 evaluation dimensions (Relevance, Depth, Clarity, Applicability, Novelty) for each research agent. Spearman rank correlation coefficients measuring ordinal agreement between critics. Higher values indicate greater consensus among the 5 critique LLMs (Claude, Gemini, Grok, O3, O4) about each agent's relative performance strengths and weaknesses. Data aggregated from 21 research topics \times 3 trials per agent. Perfect correlation ($\rho = 1.00$) indicates all critics completely agree on the ranking of an agent's performance across the 5 dimensions.

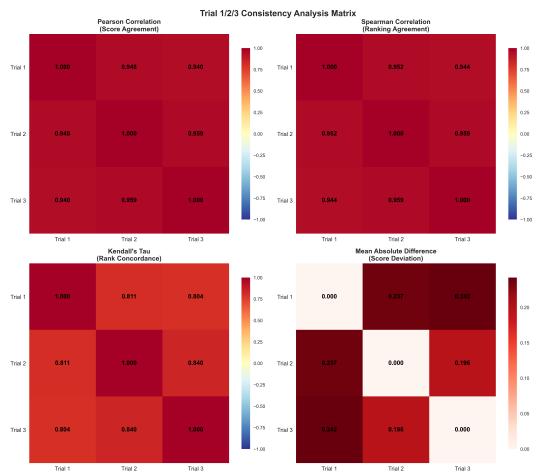


Figure S6: Trial-to-Trial Consistency Analysis Matrix. Four metrics assess reproducibility across three independent trials: (A) Pearson correlation measures score agreement, (B) Spearman correlation evaluates ranking consistency, (C) Kendall's tau provides alternative rank concordance, and (D) Mean absolute difference quantifies score deviations. Data aggregated across 21 topics, 27 agents, and 5 LLM critics. Higher correlations (>0.8) and lower deviations (<0.5) indicate excellent experimental reproducibility.

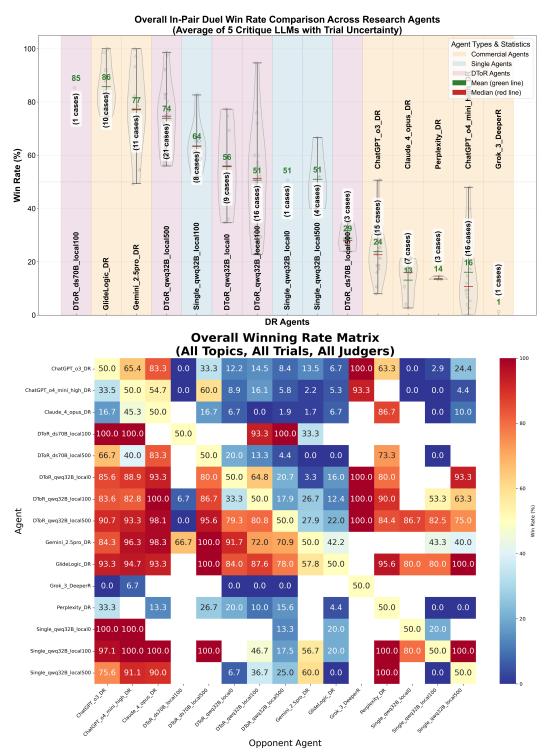


Figure S7: Upper: Win Rate Violin Plot - Shows the distribution of win rates across different agents and topics with detailed statistical information and performance variations. Bottom: Overall Winning Matrix All Topics - Displays the comprehensive winning matrix with head-to-head performance comparisons between all agents across all research topics.

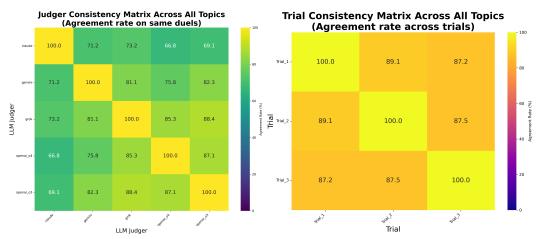


Figure S8: Consistency of in pair dueling. Left: Judger Consistency All Topics - Shows the consistency of judgment patterns across different research topics and agents, indicating how reliable the evaluation criteria are. Right: Trial Consistency All Topics - Displays the consistency of performance across multiple trials for each agent-topic combination, revealing the stability and reproducibility of results.

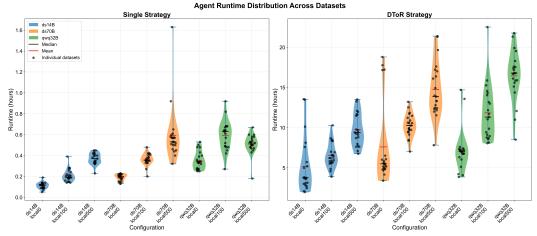


Figure S9: Violin plot visualizing the distribution of local DR agents runtime under different settings.

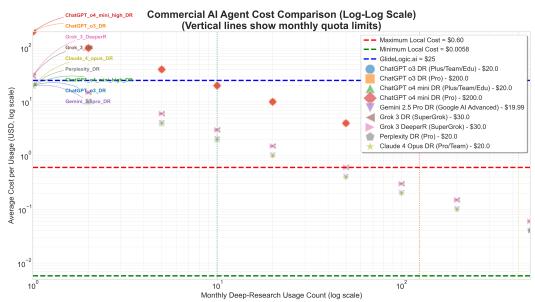


Figure S10: Comparison of cost over commercial DR agent solution with our local solution based on Chicago area electricity price.

B Related Work in Physics ML Surrogates, Domain LLMs, and Agentic Systems

Physics-aligned surrogates and classical ML. Physics-guided predictive surrogates now operate at web-scale, with graph networks and equivariant models trained on massive curated corpora. For crystalline stability, GNoME trains large graph models and active-learning pipelines to propose millions of candidate structures and flag hundreds of thousands as likely stable—establishing that data+scale can systematically widen the known inorganic universe [6]. In catalysis, the Open Catalyst efforts formalized realistic workloads: OC20 introduced 1.28M DFT relaxations and 265M single-point evaluations for adsorbate-surface energetics and created standardized tasks (IS2RE/IS2RS/S2EF) with stringent generalization splits; OC22 extended to oxide electrocatalysts relevant to OER with more realistic catalysts and targets [7, 24]. New foundation-style datasets push further: OMat24 aggregates >110M DFT entries and releases strong baselines (e.g., Equiformer-V2) that achieve near-chemical-accuracy on multiple targets [25]; OMol25 analogously scales molecular DFT data to >100M calculations spanning tens of millions of unique molecules, enabling broad pretraining for small-molecule physics [8]. Complementing datasets, universal model of atoms (UMA) explore single architectures that transfer across chemistry and tasks, aiming for cross-domain generalization with shared tokenization/featurization of atoms and local environments [26]. These surrogates excel at S1-S2 targets (bulk stability, adsorption energies, local structural relaxations) but are largely D1-D2 in operation: fixed-form predictors trained on equilibrium data, with limited capacity to reason about kinetics, synthesis windows, or device-level constraints that emerge at S3–S4. The net effect is that they remain superb "atoms-to-properties" engines, yet they do not autonomously help solve the long-horizon, multi-paper evidence assembly and process-aware planning needed, from atomic scale to device scale rationales [6–8, 24–26].

Domain LLMs (non-agent). Domain LLMs span three useful roles. (i) Materials-text models (MatSciBERT/MatBERT/MaterialBERT) are trained on materials corpora to power entity/relation extraction, document classification, and graph building; MatSciBERT outperforms SciBERT on named entity recognition (NER) and relation tasks across curated benchmarks, while MatBERT shows that domain-specific pretraining significantly lifts NER F1 and sample efficiency in materials texts; MaterialBERT introduces a domain tokenizer/vocabulary tuned to materials literature [27–29]. (ii) Chemical/polymer LLMs target sequence-level prediction/generation: ChemBERTa established transformer baselines and scaling trends for SMILES; ChemBERTa-2 pushed to 77M SMILES and showed consistent gains (and that multi-task regression pretraining can beat (Masked Language Modeling) MLM in transfer); ChemBERTa-3 contributes open, end-to-end training infrastructure for chemical foundation models; MegaMolBART provides a large-scale generative sequence model optimized for molecular design [30–33]. ChemGPT demonstrates generative scaling to >1B parameters with clear neural-scaling laws for chemical LLMs, clarifying how model/data scale trade off in practice [34]. In polymers, TransPolymer introduces a chemically-aware tokenizer and MLM pretraining that improves regression performances across ten benchmarks; polyBERT builds an endto-end, ultrafast pipeline and screens 100M hypothetical polymers; PolyNC fuses natural-language and chemical-sequence inputs to improve unified polymer property prediction [35–37]. (iii) Generative/design LLMs for materials include ChatMOF, which conditions MOF proposals/properties from text and structural descriptors; Crystal-LLM (autoregressive) generates crystal structures from composition; and MatterGen scales inorganic structure generation with validated discoveries and robust property control [38–40]. Collectively these models deliver flexible text/code/structure interfaces and strong S1–S2 performance—ideal for curation, candidate suggestion, or surrogate-assisted screening—but they are typically short-horizon (D1–D2): they answer prompts, not conduct hours-long, tool-grounded, hierarchical inquiry across conflicting literature [27–40].

Agentic systems for scientific inquiry. Agent frameworks aim to orchestrate end-to-end inquiry: decompose goals, retrieve and reconcile evidence, use domain tools/simulators, and iteratively critique/plan. Early exemplars used rule/graph planning and Bayesian loops within constrained scopes—Chematica/SYNTHIA operationalized expert retrosynthesis rules and heuristic search; CAMD framed autonomous design/make/test/learn for materials screening [5, 41]. LLM-based agents extend this with language-native decomposition and tool-use: ChemCrow wires models to chemistry tools (e.g., RDKit, databases) to complete realistic tasks; HoneyComb provides a flexible materials-agent stack with retrieval/tooling and HPC hooks; SciAgents formalizes multi-agent graph reasoning for scientific planning; "AI co-scientist" frameworks couple planning with web/code tools

and emphasize evaluation and provenance [42–45]. Execution back-ends close the loop: A-Lab demonstrates a self-driving lab discovering and characterizing new materials; IBM RXN/RoboRXN links digital retrosynthesis with automated synthesis; NVIDIA ALCHEMI targets scaled infrastructure (data, models, and accelerated simulation) for materials discovery [46–48]. Despite progress, most agents today operate at D2–D3 depth and S1–S3 scope; they rarely sustain thousand-query retrieval runs, explicit topic-tree governance, or tight coupling to cross-domain process constraints that S3–S4 device problems demand. This is precisely the niche for deep-research agents: to deliver hierarchical planning, multi-stage local+web retrieval at scale, targeted tool/simulator calls, and systematic self-critique with dense provenance until a decision-ready synthesis emerges [5,41–48].

Gap to deep research at S3–S4/D4. Despite existing progress, the transition to S3–S4 nano-device complexity poses a fundamental barrier. The combinatorial explosion across materials, interfaces, processing parameters, and operational constraints creates hypothesis spaces where decisive insights lie at the intersections of disparate literature, conflicting reports, and cross-domain phenomena. Crucially, long-horizon, hierarchical evidence accumulation (D4) remains under-documented: many agents cap at D2–D3, with limited iterative loops, shallow coordination, and no explicit topic-tree control for breadth-then-depth synthesis. Commercial deep-research offerings introduced in 2025 enable multi-step inquiry but remain closed-source with limited controllability and local integration, reinforcing the need for open, on-prem frameworks tailored to proprietary materials workflows.

C Pseudocode

```
Algorithm 1: Vanilla Web-First DR Agent(Adapted from [22])
Require: Inputs: research_topic \in \mathbb{S}; configuration C = \{ \text{llm_provider} \in \{ \text{lmstudio}, \text{ollama} \},
                             {tavily, perplexity, duckduckgo, searxng}, fetch_full_page∈
           search_api∈
           strip_thinking_tokens \in \{0, 1\}, use_tool_calling \in \{0, 1\}, max_web_research_loops \in \mathbb{N} \}.
Ensure: final running_summary (with deduplicated sources appended).
 1: State: running_summary \( \lefta \) "; search_query \( \lefta \)"; web_research_results \( \lefta \) [];
    sources_gathered \leftarrow []; research_loop_count \leftarrow 1.
 2: function GENERATEQUERY(state, C)
         Build system prompt from research_topic and current_date.
 3:
 4:
        if C.use_tool_calling = 1 then
 5:
             Ask LLM to return tool payload {query, rationale}.
 6:
        else
 7:
             Ask LLM in JSON mode to return {query, rationale}.
        Parse JSON/tool output; let q be the extracted query.
 8:
 9:
        if q is empty or parsing fails then
10:
             Strip \langle \text{think} \rangle blocks if C.strip_thinking_tokens = 1.
             Fallback to first nonempty line; if still empty, set q \leftarrow "Tell me more about
11:
             research_topic".
        search_query \leftarrow q; return q.
12:
13: function WEBRESEARCH(state, C)
        Select backend b \leftarrow C.search_api.
        Issue web search with q = \text{search} query via b; obtain raw results \mathcal{R}.
15:
        Format bullets S \leftarrow \texttt{format\_sources}(\mathcal{R}).
16:
        Build long string W \leftarrow \text{deduplicate\_and\_format\_sources}(\mathcal{R}, m, C.\text{fetch full page}),
17:
        where m is the max tokens per source (vanilla uses \approx 1000).
18:
        Append S to sources_gathered; append W to web_research_results.
19:
        research\_loop\_count \leftarrow research\_loop\_count + 1.
20:
        return W.
21: function SUMMARIZESOURCES(state, C)
        Let E \leftarrow running summary (possibly empty).
        Let W \leftarrow last element of web research results.
23:
        if E \neq "" then
24:
                                                                                          > Update existing
            Build human message with \langle \text{Existing Summary} \rangle E \langle / \rangle and \langle \text{New Context} \rangle W \langle / \rangle.
25:
            Instruct: "Update the Existing Summary with the New Context on: research_topic."
26:
27:
                                                                                              ▷ Create fresh
28:
             Build human message with \langle \text{Context} \rangle W \langle / \rangle.
             Instruct: "Create a Summary using the Context on: research topic."
29:
30:
        Call LLM with temperature 0 to get S^*.
31:
        if C.strip thinking tokens = 1 then strip \langle \text{think} \rangle blocks in S^*.
        running_summary \leftarrow S^*; return S^*.
32:
33: function REFLECTONSUMMARY(state, C)
        Provide running_summary and research_topic to LLM with a reflection prompt.
34:
        if C.use_tool_calling = 1 then
35:
             Request tool payload {follow up query, knowledge gap}.
36:
37:
38:
             Request JSON {follow up query, knowledge gap}.
        Extract q' \leftarrow follow_up_query; on failure, strip (think) then fallback to "Tell me more about
39:
        research topic".
40:
        search_query \leftarrow q'; return q'.
41: function FINALIZESUMMARY(state)
        Initialize seen\leftarrow \emptyset, unique_lines\leftarrow [].
42:
        for all bullet block ∈ sources gathered do
43:
44:
             for all line \in lines(bullet block) do
                 if line is nonempty \wedge line \notin seen then
45:
                     add line to unique lines; seen\leftarrowseen\cup{line}.
46:
        all_sources← join(unique_lines, "\n").
47:
        running summary ← "## Summary" | running summary || "\n\n ### Sources:\n" || all sources.
48:
```

```
49: return running summary.
    Main flow
50: search query \leftarrow GENERATEQUERY(state, C).
52:
        \_\leftarrow WebResearch(state, C).
        \leftarrow SUMMARIZESOURCES(state, C).
53:
        search_query \leftarrow REFLECTONSUMMARY(state, C).
54:
        if research_loop_count > C.max_web_research_loops then break
56: \overline{\text{final\_summary}} \leftarrow \text{FINALIZESUMMARY}(state).
57: return final summary.
Algorithm 2: Local Evidence-First DR Agent
Require: Inputs: research_topic \in S; configuration C = \{ lm_provider \in \{ lmstudio, ollama \}, \}
           search\_api \in \quad \{tavily, perplexity, duckduckgo, searxng\}, \quad fetch\_full\_page \in \quad \{0,1\},
           strip_thinking_tokens\in \{0,1\}, use_local_rag\in \{0,1\}, vector_store_paths: List,
           embedding_model \in \mathbb{S}, local_results_count \in \mathbb{N}, research_mode \in \{\text{single}, \text{DToR}\},
           max web research loops \in \mathbb{N}.
Ensure: final running summary; in single mode it includes deduplicated sources (tagged by chan-
          nel), in DToR mode it omits sources.
 1: State: running_summary \( \lefta \) ""; summary_history \( \lefta \) []; search_query \( \lefta \) ""; complementary_search_query \( \lefta \) ""; local_research_results \( \lefta \) []; local_rag_summary \( \lefta \) ""; web_research_results \( \lefta \) []; complementary_web_research_results \( \lefta \) [];
    local_sources_gathered \leftarrow []; sources_gathered \leftarrow []; complementary_sources_gathered \leftarrow [];
    research loop count \leftarrow 1.
 2: function GENERATEQUERY(state, C)
        Build system prompt from research topic and current date.
 4:
        Gate JSON mode by provider/model; attempt JSON parse of {query,rationale}.
 5:
        if JSON parse fails then
            if C.strip_thinking_tokens = 1 then strip \langle \text{think} \rangle
 6:
 7:
            Set q \leftarrow first nonempty line of content, else q \leftarrow research topic.
        search_query \leftarrow q; return q.
 9: function LOCALRAGRESEARCH(state, C)
10:
        if C.use_local_rag = 0 then return {local_research_results : [], local_sources_gathered :
        Expand C.vector_store_paths into discovered volumes; query Chroma stores with device-
11:
        aware embeddings.
12:
        Extract raw contents list R; format citations S_{loc}.
        return {local_research_results : R, local_sources_gathered : [S_{loc}] }.
13:
14: \bar{\mathbf{function}} SUMMARIZELOCALRAGRESULTS(state, C)
        if local_research_results = \emptyset then return {local_rag_summary : ""}
15:
16:
        Join local results into one string; prompt concise (< 500 chars) summary with fixed lead-in.
        return {local_rag_summary : \ell }.
17:
18: function GENERATECOMPLEMENTARYQUERY(state, C)
19:
        if C.use_local_rag = 0 or local_research_results = \emptyset then
20:
             return {complementary_search_query : ""}
                                                                              ⊳ skip when no local signal
        Ask LLM (JSON) using local_rag_summary context; recursively search any JSON key
21:
        containing "query".
22:
        if extracted query empty ∨ identical to search query then
             apply backup prompt(s) and safe template fallbacks.
        return {complementary_search_query : q_c }.
25: function WEBRESEARCH(state, C)
        Call configured backend with q = \text{search\_query}.
26:
        Build W \leftarrow \texttt{deduplicate\_and\_format\_sources}(\cdot, max\_tokens\_per\_source =
27:
        2500, C.fetch_full_page).
28:
        Append formatted bullets to sources_gathered and W to web_research_results.
        research loop count \leftarrow research loop count +1; return W.
30: function COMPLEMENTARY WEBRESEARCH(state, C)
31:
        if complementary_search_query = "" then
32:
            return {complementary_sources_gathered : [], complementary_web_research_results :
```

```
Query backend with q_c; build W_c as above (budget 2500 tokens per source).
       return {complementary_sources_gathered : [S_c], complementary_web_research_results :
34:
        [W_c].
35: function SUMMARIZESOURCES(state, C)
36:
       Let E \leftarrow \text{running\_summary}; assemble context string X by concatenating
            (i) last of local_research_results if any as \langle Local \rangle \dots \langle / \rangle,
            (ii) last of web research results as \langle Main \rangle \dots \langle / \rangle,
            (iii) last of complementary_web_research_results as (Complementary) . . . (/).
       if E \neq "" then prompt to update E with X;
37:
       elseprompt to create summary from X.
38:
39:
       Get S^* (temperature 0.3); if C.strip_thinking_tokens = 1 then strip tokens.
       Append {iteration = research_loop_count, summary = S^*, query =
40:
       search_query, complementary_query = complementary_search_query} to
       summary_history.
       running summary \leftarrow S^*; return S^*.
41:
42: function REFLECTONSUMMARY(state, C)
       Request JSON with keys {follow_up_query, knowledge_gap}; use model-specific JSON
       gating.
44:
       if parse fails then
45:
           if C.strip thinking tokens = 1 then strip \langle \text{think} \rangle
           fallback to a safe heuristic query.
46:
       search query \leftarrow q'; return q'.
47:
48: function ROUTERESEARCH(state, C)
       if research loop count < C.max web research loops then return local_rag_research
       else return finalize_summary
50:
51: function FinalizeSummary(state, C)
       Build history text from summary_history if available, else use running_summary.
53:
       Ask LLM (temperature 0.3) for final mini-report.
54:
       if C.research mode = DToR then return final running summary \triangleright sources handled at
       DToR synthesis
55:
       else
           Deduplicate lines from sources gathered (main), complementary sources gathered (tag
56:
           as complementary), and local sources gathered (tag as local); join as ### Sources:
           block.
57:
           Append sources block to final summary and return.
    Main flow
58: INIT_SESSION; search_query \leftarrow GENERATEQUERY(state, C).
59: LOCALRAGRESEARCH(state,C) \rightarrow SUMMARIZELOCALRAGRESULTS(state,C) \rightarrow GENER-
    ATECOMPLEMENTARY QUERY (state, C).
60: WebResearch(state,C) \rightarrow ComplementaryWebResearch(state,C) \rightarrow Summarize-
    SOURCES(state,C).
61: search_query \leftarrow REFLECTONSUMMARY(state, C).
62: if ROUTERESEARCH(state,C) = local_rag_research then goto LocalRAGResearch else
    FINALIZESUMMARY(state,C).
Algorithm 3: DToR Orchestrator Over Research Nodes (Breadth–Depth Controller)
Require: Inputs: research_topic \in S; configuration C = \{ llm_provider, search_api, \}
          fetch_full_page, strip_thinking_tokens, use_local_rag, vector_store_paths, embed-
          ding_model, max_branches \in \mathbb{N}, max_branch_depth \in \mathbb{N}, nodes_per_branch \in \mathbb{N} \}.
Ensure: final_summary and (optionally) all_sources.
 1: State (ToTState):
      branches: Dict[branch id \rightarrow ResearchBranch], where each
         ResearchBranch = {perspective, depth, remaining_budget, research_nodes:
    List<SummaryState>, branch summary, is complete};
      active branch id; knowledge gaps; all sources; is complete \in \{0, 1\}.
 2: function DIVERSIFYINITIALQUERY(state, C)
                                                                                ▷ Breadth seeding
       Generate up to C.max_branches orthogonal perspectives (titles, descriptions, seed queries).
 3:
       For each perspective i: create a branch with one RN (SummaryState) seeded by its query;
 4:
```

set depth $\leftarrow 0$, remaining budget $\leftarrow C$.nodes per branch.

5:

```
Choose initial active branch id.
7: function SELECTNEXTBRANCH(state)
       if active branch has actionable work then return active branch id
       elsereturn first incomplete branch id
9:
10: function RUNRESEARCHNODE(state, C)
                                                                 \triangleright Depth step = RN from Alg. 2
       In the active branch, pick the next RN with processing status = "pending".
       Invoke Algorithm 2 on that RN (local evidence \rightarrow web \rightarrow summarize \rightarrow reflect).
12:
       Update the RN with its running summary and mark completed; clear per-node sources in
13:
       DToR mode.
14: function DECISIONANALYST(state, C)
                                                                    ⊳ EXPAND or PRUNE only
       Read the latest completed RN in the active branch and its context (depth, remaining_budget).
15:
       Return decision ∈ {EXPAND, PRUNE} and (if EXPAND) a small list of knowledge_gaps
16:
       with web queries.
       if depth \geq C.\text{max\_branch\_depth} then
17:
           force decision ← PRUNE; knowledge_gaps ← []
18:
       Mark the analyzed RN as processing_status = "analyzed".
19:
20: function GENERATEFOLLOWUPQUERIES(state, C)
                                                                             ⊳ Depth expansion
       Let k \leftarrow \min(|\text{knowledge\_gaps}|, \text{remaining\_budget}).
       For the first k gaps: create new RN nodes with refined topics/queries (avoid compounded
22:
       prefixes).
23:
       remaining_budget \leftarrow remaining_budget -k; depth \leftarrow depth +1.
24: function SYNTHESIZEBRANCH(state)
25:
       Compress all RN running_summary in the branch into branch_summary;
26:
       deduplicate and append sources into all sources; persist branch report.
27: function SYNTHESIZEFINAL(state)
       Integrate all branch summary across branches into final summary; mark is complete \leftarrow 1.
28:
29: function ROUTENEXTACTION(state, C)
                                                                       ⊳ Breadth–depth control
       if not in DToR mode or state.is_complete then return END
30:
31:
       if no branches then return DIVERSIFYINITIALQUERY
       Let b \leftarrow active branch.
32:
33:
       if knowledge_gaps \neq \emptyset and b.remaining_budget > 0 then return GENERATEFOLLOWUP-
       OUERIES
34:
       if exists pending RN in b then return RUNRESEARCHNODE
35:
       if exists completed but not analyzed RN in b then return DECISIONANALYST
36:
       if b.is_complete = 1 and b.branch_summary empty then return SYNTHESIZEBRANCH
37:
       if all branches complete and summarized then return SYNTHESIZEFINAL
       return SELECTNEXTBRANCH
38:
   Main flow (controller loop)
39: DIVERSIFYINITIALQUERY(state, C); active_branch_id \leftarrow SELECTNEXTBRANCH(state).
40: while not state is complete do
       a \leftarrow \text{ROUTENEXTACTION}(state, C).
41:
       if a = RUNRESEARCHNODE then RUNRESEARCHNODE(state,C)
42:
       if a = DECISIONANALYST then
43:
44:
           (decision, knowledge gaps) \leftarrow DECISIONANALYST(state, C).
           if decision = EXPAND and b.depth < C.max_branch_depth then GENERATEFOL-
45:
           LOWUPQUERIES(state,C)
                                                        \triangleright PRUNE \equiv stop expanding this branch
46:
           else
47:
              b.is complete \leftarrow 1
       if a = GENERATEFOLLOWUPQUERIES then GENERATEFOLLOWUPQUERIES(state,C)
48:
       if a = SYNTHESIZEBRANCH then SYNTHESIZEBRANCH(state)
49:
       if a = SYNTHESIZEFINAL then SYNTHESIZEFINAL(state); break
50:
       if a = SELECTNEXTBRANCH then active_branch_id \leftarrow SELECTNEXTBRANCH(state)
52: return final summary (and all sources).
```

D Supplemental Results: Experimental Setup

D.1 Nanomaterials/Devices Topic Suite

PFAS FET Sensor Probe

Which two-dimensional (2D) nanomaterials or molecular probes-including novel, \rightarrow underexplored candidates-offer the greatest potential to achieve highly \rightarrow sensitive and selective detection of per- and polyfluoroalkyl substances \rightarrow (PFAS), or more broadly chemical and biological analytes, when integrated \rightarrow into a field-effect transistor (FET)-based sensor platform?

Battery Binder

Which advanced binder technologies are being developed to improve the \hookrightarrow performance and longevity of lithium-ion batteries?

Staphylococcus aureus FET Biosensor

How can emerging probe chemistries that are commercially available, $\begin{array}{ll} \hookrightarrow & \text{cost-effective, and exhibit minimal batch-to-batch variability be used to} \\ \hookrightarrow & \text{develop a novel FET-based biosensor for detecting Staphylococcus aureus?} \end{array}$

Resource Recovery Electrocatalysts

Which nanostructured electrocatalyst materials demonstrate the highest selectivity and efficiency for electrochemical detection and recovery of critical resources (e.g., Li^+ , $\operatorname{PO_4}^{3-}$, $\operatorname{NH_4}^+$) from complex wastewater matrices, and what key performance metrics-such as selectivity, sensitivity, recovery rate, energy consumption, and operational stability-distinguish them?

OER Stability

Identify the top nanostructured electrocatalyst materials for driving the oxygen evolution reaction (OER) in complex wastewater matrices—such as high-chloride, high-organic—load, or multi—ion streams—, for each the material class/composition, comprehensively consider key performance metrics (overpotential at 10 mA cm $^{-2}$, Faradaic efficiency, stability), and the surface—engineering strategies that confer corrosion resistance and sustained activity in real effluent conditions.

Photothermal Wastewater Evaporation

What photothermal materials and system designs are most effective for solar-driven water evaporation in complex wastewater matrices-such as high-organic-load or multi-ion streams-and how do they compare in terms of material composition, solar-to-vapor conversion efficiency under one-sun illumination, evaporation rate, fouling resistance, and integrated resource-recovery functionalities?

Nutrient Sensor Interference

Which sensor probe materials/chemicals and designs offer the best performance

→ in minimizing interference factors-such as competing ions, dissolved

→ organic matter, pH fluctuations, and temperature variations-for accurate

→ and selective detection of nutrients (e.g., nitrate, phosphate, ammonium)

→ in complex water matrices, and what mitigation strategies do they employ?

CO2 Sensor Probe

Which two-dimensional materials-such as graphene derivatives, transition \rightarrow metal dichalcogenides, or MXenes-offer the highest CO_2 sensing \rightarrow performance in complex gas or aqueous environments, and how do they \rightarrow compare in terms of detection limit (ppm), selectivity against common \rightarrow interferents (e.g., O_2 , O_2 , O_2 , O_3), response/recovery time, and long-term stability, including any functionalization or structural modifications O_3 that enhance these metrics? Think of novel candidates.

Printed FET Variability

Which printed-electronics fabrication parameters and post-processing

⇒ strategies-such as ink viscosity, printing speed and resolution,

⇒ substrate surface energy, annealing temperature profiles, and in-line

⇒ calibration methods-have been shown to minimize device-to-device

⇒ variability in FET sensor arrays, and what specific process windows

⇒ achieve low variation in threshold voltage and field-effect mobility?

2D Synaptic Transistors

Which two-dimensional material platforms (e.g., MoS_2 , WSe_2 , black phosphorus, \rightarrow h-BN), device architectures (e.g., floating-gate, ionic-gated, dual-gate), and fabrication protocols (e.g., channel thickness control, dielectric engineering, contact metallurgy) have been shown to optimize synaptic transistor performance-specifically in terms of energy per event, weight-update linearity, retention time, and cycling endurance-for neuromorphic sensing applications? Think of most practical and promising candidates.

Microplastics 2D Sensors

Which two-dimensional material platforms (e.g., graphene derivatives,

→ transition metal dichalcogenides, MXenes), molecular recognition elements

→ (e.g., molecularly imprinted polymers, aptamers, peptide receptors), and

→ device integration strategies (e.g., FET, electrochemical impedance,

→ photonic transduction) have demonstrated the highest sensitivity,

→ selectivity against organic matter and ionic interferents, and rapid

→ response times for detecting micro- and nanoplastic particles in complex

→ water matrices? Think of novel candidates.

Antibiotics 2D Sensors

Which commercially available, cost-effective probe chemistries (e.g., \hookrightarrow thiolated DNA aptamers, antibody mimetics, conductive MIPs) combined with two-dimensional nanomaterial transducer platforms (e.g., WS $_2$ FETs, \hookrightarrow graphene field-effect sensors, nanotube-extended gate FETs) deliver the \hookrightarrow lowest detection limits, minimal batch variability, and robust \hookrightarrow performance for sensing trace levels of pharmaceutical antibiotics in \hookrightarrow diverse aqueous waterbody.

Li–Na Selective Membranes

What material or membrane exhibits the highest selectivity for Li⁺ over Na⁺

→ in aqueous systems, given their nearly identical hydrated ionic radii and

→ solvation environments? Beyond crown-ether-functionalized polymers, what

→ novel or unexpected materials-such as bioinspired ultrahigh-selectivity

→ membranes or covalent organic framework nanochannels-might provide

→ breakthrough Li⁺ discrimination over Na⁺? Find novel and promising

→ candidates

PFAS Degradation

Which novel electrode materials can achieve efficient PFAS degradation under ambient aqueous electrochemical conditions, delivering both high mineralization and defluorination rates? What intrinsic properties-such as PFAS adsorption affinity, reactive oxygen species generation capacity, and C-F bond activation energy-should be optimized to guide their discovery?

Li–Co–Ni Separation

Which membrane materials can effectively separate Li^+ , Co^{2+} , and Ni^{2+} ions \to from aqueous solutions by leveraging selective transport properties—such as tailored pore sizes, specific surface functionalizations, and charge \to affinities—and what membrane design principles optimize both selectivity and permeability? Think of potential novel, effective, and practical \to candidates.

LC-TEM MoS2 FET Sensing

How can in situ liquid cell TEM be employed to directly visualize the \hookrightarrow real-time adsorption and structural changes of 2D MoS $_2$ nanosheets used in aqueous FET sensors during analyte binding? What fluid cell \hookrightarrow configurations and electron dose parameters are necessary to preserve \hookrightarrow native water-material interfaces while capturing high-resolution sensing \hookrightarrow events without beam-induced artifacts?

CO2RR Catalyst Leaders

Given the current landscape of CO_2 electroreduction, which state-of-the-art catalyst platforms-such as oxide-derived copper, single-atom catalysts on nitrogen-doped carbon supports, or metal-organic framework-derived attended materials-demonstrate the highest activity and selectivity? Moreover, what emerging catalytic systems or novel heterostructures beyond these examples could feasibly outperform today's leading electrocatalysts in terms of faradaic efficiency and stability?

Unconventional PV Platforms

Perovskite-silicon tandem cells, organic photovoltaics, and quantum-dot solar

→ cells currently represent the forefront of next-generation solar

→ technologies. Beyond these established platforms, which unconventional

→ material classes or innovative device architectures-such as chalcogenide

→ perovskites, 2D semiconductor heterostructures, or ferroelectric

→ photovoltaic systems-offer the most unexpected promise for achieving

→ breakthroughs in efficiency, stability, and scalability?

Ambient-Pressure Diamond Growth

How can ambient-pressure diamond synthesis using a Ga-Ni-Fe-Si liquid-metal

- \rightarrow alloy at 1 atm and ~1025 °C be adapted to produce larger-area diamond
- \hookrightarrow films or oriented single crystals, and what key mechanistic steps govern
- \rightarrow nucleation and growth kinetics under these mild-pressure conditions?

Diamane Property Tuning

What chemical functionalizations or reaction pathways are most promising for

- $\,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\,$ tuning the electronic bandgap and mechanical stability of diamane-like 2D
- $\,\,\hookrightarrow\,\,$ diamond films formed via $\text{sp}^2\text{-to-sp}^3$ conversion, and how might those
- → modifications influence their integration into nanoelectronic devices?

Anion-Receptor Ether Electrolytes

Search for novel fluorinated ether-based electrolyte candidates that

- \rightarrow incorporate covalently bound anion-receptor motifs (e.g. boron- or
- → phosphorus-centered groups) to deliver oxidative stability beyond 5.6 V,
- \rightarrow Li⁺ transference numbers above 0.8, and minimal interfacial impedance on
- \hookrightarrow lithium-metal anodes.

D.2 Details of DR Agent Experimental Settings

Open-source LLMs and local RAG corpus We implement the framework in LangGraph and deploy three low-deployment-cost, open-source LLMs—DeepSeek-R1-14B, DeepSeek-R1-70B, and QwQ-32B—served locally via Ollama. All agent prompts and tool I/O use a JSON/tool-calling discipline for determinism and easy parsing, and the same research node (RN) described in Algorithm 2 is the unit of work in the controller (Algorithm 3 DToR). To provision representative, applicationdriven knowledge for nanomaterials and devices, we build a standardized Chroma vector store spanning 1,136,297 full-text papers organized around four pillars: sensors, batteries, catalysis, and semiconductors. The corpus is queried first under a local-first RAG policy; complementary web retrieval is invoked only when Algorithm 2 (local evidence-first)'s reflection detects unresolved knowledge gaps. Consistent with our emphasis on report-level usefulness rather than information retrieval (IR) benchmarking, the local-first RAG layer is treated as a routine, swappable component; DToR only triggers complementary web evidence upon explicit knowledge gaps and resolves crosssource conflicts by evidence majority \rightarrow recency \rightarrow domain authority. To make domain coverage portable and auditable without exposing any local texts or indices, we provide exact domain-scoping query templates corresponding to these four pillars (sensors, energy/environment covering batteries, semiconductors, and catalysis) so others can target the same scope with their own holdings while keeping the orchestration identical.

Corpus construction Raw full texts are normalized with a cleaning pipeline designed for scientific documents at scale. The cleaner (i) preserves section structure and headers, (ii) retains LaTeX/MathJax equations and common math environments via stable markers, (iii) standardizes tables into a compact textual schema, (iv) normalizes URLs/DOIs and citation markers instead of stripping them, and (v) removes XML/HTML artifacts and noisy back matter conservatively. For long inputs, we apply sliding-window tokenization (default window 8192 tokens, stride 4096) or a section-aware variant that respects detected boundaries, yielding JSONL suitable for downstream training or indexing. As our focus is on the DToR orchestration, we adhere to widely used, off-the-shelf preparation and vector-indexing practices; the corpus layer functions as a routine, swappable component rather than a methodological contribution. To enable scope-equivalent replication across different holdings, we release the exact domain-scoping query templates for sensors, energy/environment (including batteries), semiconductors, and catalysis; users may instantiate a local store with their customary institutional pipelines and preferred indexers.

Indexing and retrieval Cleaned texts are split with a recursive character splitter (default chunk size 2500, overlap 500) and embedded using BAAI/bge-m3. Each chunk carries rich metadata (paper id, DOI, chunk id/index, source path, shard id) to support provenance and de-duplication. Embeddings are persisted in Chroma, enabling fast k-NN retrieval and incremental persistence during large ingestions. For robustness across machines, the store is partitioned into 500 volumes (shards). This organization serves two purposes: (1) controllable breadth—DToR can expand to additional volumes when the analyst decides to broaden coverage—and (2) bounded depth/latency—branches can be pruned once depth or evidence sufficiency is reached, without touching unrelated volumes. To mitigate licensing and privacy risk, our local corpus is restricted to institutionally licensed or open sources, and we release only the domain-scoping query templates—not texts or indices—so others can replicate scope without redistributing copyrighted content or PII.

Policy integration In each research node, local retrieval is attempted first; summaries are synthesized from local evidence. Only when reflection surfaces concrete gaps does the controller expand (selecting additional volumes and spawning new RNs) or prune (stopping further growth of a branch). This strategy couples breadth (multi-volume exploration across pillars) with depth (iterative drilling within a volume) while keeping costs predictable on commodity hardware and avoiding external dependencies during the main research loop.

Agent Configuration By default, we run all agents in a local-first configuration with cost-aware web fallback. Concretely, both single-agent DR and DToR branches use a research depth of 3 iterations, each RN queries the local Chroma first with top-k = 5, and only then invokes web search—by default DuckDuckGo—to avoid API fees. In DToR mode, breadth–depth is bounded with max branches = 3, max branch depth = 3, nodes per branch = 100, and a global cap of 500 nodes per experiment, ensuring predictable latency and compute. Models are served locally via Ollama; RAG is enabled by default with BAAI/bge-m3 embeddings, and we preserve document fidelity via full-page fetch and

safe output hygiene. All parameters remain user-tunable through the configuration/RunnableConfig to trade off coverage, latency, and cost—for example, increasing depth or top-k for harder queries, or switching web backends when budgets permit.

D.3 Rubric Template for LLM-as-Judge Evaluation Experiment

LLM-as-Judge Rubric Prompt (applied to DR Report)

{Anonymous DR Report content}

Now please analyze and evaluate the Report above on five equally weighted dimensions:

- 1. Relevance to the stated objective,
- 2. Depth and rigor of validation and benchmarking protocols,
- 3. Clarity and reproducibility of methodology,
- 4. Practical applicability and implementation readiness,
- 5. Novelty and breadth of candidate coverage.

For each dimension, assign a score from 0–10 and supply a one-sentence justification. This rubric places equal emphasis on procedural rigor and clear implementation steps alongside overall alignment, ensuring that reports with well-defined benchmarks and actionable guidance receive full credit.

Example raw response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe) — Claude 4 Opus (thinking)

```
Based on my analysis of the research report on 2D nanomaterials and molecular
\,\,\,\,\,\,\,\,\,\,\,\, probes for FET-based PFAS detection, here are the evaluation scores across
   the five equally weighted dimensions:
## **Dimension 1: Relevance to the Stated Objective**
**Score: 9/10**
The report directly addresses the core question about identifying 2D
\hookrightarrow nanomaterials and molecular probes with the greatest potential for highly
\hookrightarrow sensitive and selective PFAS detection in FET-based sensors. The
\hookrightarrow comprehensive analysis covers both established and novel candidates

→ including MXenes, graphene, TMDCs, aptamers, and molecularly imprinted

    polymers, with specific focus on their integration into FET platforms for
   PFAS detection.
## **Dimension 2: Depth and Rigor of Validation and Benchmarking Protocols**
**Score: 6/10**
While the report provides specific performance metrics such as sensitivity
\rightarrow ranges (100 fM to 10 pM) and specificity percentages (99% for PFOS), it
\rightarrow lacks detailed descriptions of the experimental protocols, standardized
   testing conditions, and validation methodologies used to obtain these

ightharpoonup benchmarks. The report acknowledges this gap by calling for "standardized
    protocols for probe-material interface characterization" in its future
   research recommendations.
## **Dimension 3: Clarity and Reproducibility of Methodology**
**Score: 7/10**
The report presents a well-structured synthesis approach with clear

→ categorization of materials, probes, and integration strategies, but

    provides limited detail on the specific research methodologies used to
   generate the findings from the five research branches (A-E). The

→ contradiction analysis and trade-off evaluations demonstrate systematic

→ thinking, though the underlying experimental procedures for material

   characterization and performance testing are not explicitly described.
## **Dimension 4: Practical Applicability and Implementation Readiness**
**Score: 8/10**
```

The report excels in addressing real-world implementation challenges, \hookrightarrow specifically highlighting manufacturing scalability through roll-to-roll → processing, environmental stability concerns for MXenes, and the need for → protective coatings. The comprehensive candidate inventory with performance → metrics and the identification of hybrid systems like MXene/graphene \hookrightarrow heterostructures provide actionable pathways for practical deployment. ## **Dimension 5: Novelty and Breadth of Candidate Coverage** **Score: 8/10** The report demonstrates excellent breadth by covering both conventional ightharpoonup materials (graphene, TMDCs) and emerging candidates (MXenes, fluorinated \hookrightarrow phospholipid bilayers), while exploring novel integration approaches and \hookrightarrow creative solutions like machine learning optimization and bioinspired → coatings, though some candidates like fluorinated phospholipid bilayers are \rightarrow noted as requiring further validation. ## **Overall Assessment** **Total Score: 38/50 (76%) ** The report provides a comprehensive and well-structured analysis that \hookrightarrow effectively addresses the stated objective with strong practical focus and \hookrightarrow protocols and methodological transparency to enhance reproducibility and scientific rigor.

Example raw response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe) — Gemini 2.5 Pro

Based on the provided research report, here is an analysis and evaluation across \hookrightarrow the five requested dimensions:

- * **Relevance to the stated objective**
 - * **Score**: 10/10
 - * **Justification**: The report directly and comprehensively answers the
 - $\,\hookrightarrow\,$ core research question by identifying and evaluating specific 2D
 - \rightarrow materials and molecular probes for FET-based PFAS detection.
- * **Depth and rigor of validation and benchmarking protocols**
 - * **Score**: 7/10
 - * **Justification**: While the report synthesizes key quantitative
 - → performance metrics, it correctly identifies the lack of standardized
 - $\,\hookrightarrow\,$ testing protocols as a critical gap in the field.
- * **Clarity and reproducibility of methodology**
 - * **Score**: 6/10
 - * **Justification**: The report clearly outlines high-level integration
 - $\,\hookrightarrow\,$ strategies but lacks the detailed experimental parameters required for
 - $\,\,\hookrightarrow\,\,\, \text{direct methodological reproduction.}$
- * **Practical applicability and implementation readiness**
 - * **Score**: 9/10
 - * **Justification**: The analysis heavily focuses on practical barriers to
 - $\,\hookrightarrow\,$ commercialization, such as stability and scalability, and proposes a
 - \rightarrow clear roadmap for real-world implementation.
- * **Novelty and Breadth of Candidate Coverage**

```
* **Score**: 10/10
* **Justification**: The report covers a wide breadth of candidates,

→ including novel materials like MXenes and underexplored probes like

→ fluorinated phospholipid bilayers, while also considering innovative

→ hybrid systems.
```

Example raw response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe) — Grok 3 (thinking)

```
Below is the evaluation of the provided research report titled *"Integrated
\hookrightarrow Research Report: 2D Nanomaterials and Molecular Probes for FET-Based PFAS
\hookrightarrow Detection"* based on the five equally weighted dimensions specified in the
\hookrightarrow rubric. Each dimension is scored from 0 to 10, with a one-sentence
  justification for each score.
### **Evaluation of the Research Report**
#### **1. Relevance to the Stated Objective**
- **Score**: 9
- **Justification**: The report directly addresses the detection of PFAS and
\hookrightarrow other analytes using FET-based sensors with 2D nanomaterials and molecular
\rightarrow probes, covering material selection, probe design, and integration
→ strategies; however, it could have explored the specific challenges of

→ detecting diverse PFAS compounds in greater depth.

#### **2. Depth and Rigor of Validation and Benchmarking Protocols**
- **Score**: 7
- **Justification**: The report provides specific performance metrics such as
⇒ sensitivity (100 fM-10 pM) and specificity (99%), but it lacks detailed
\hookrightarrow descriptions of validation methods or comparisons with existing

→ technologies, limiting the rigor of its claims.

#### **3. Clarity and Reproducibility of Methodology**
- **Score**: 8
- **Justification**: The methodology is clearly outlined for integration
→ strategies like layer-by-layer assembly and roll-to-roll manufacturing, but
\,\,\,\,\,\,\,\,\,\,\,\,\,\,\, it falls short on detailed, reproducible steps for material
\hookrightarrow functionalization or probe design.
#### **4. Practical Applicability and Implementation Readiness**
- **Score**: 7
- **Justification**: The report discusses scalability through techniques like
actionable details such as cost estimates or specific implementation steps,
\hookrightarrow keeping it somewhat theoretical.
#### **5. Novelty and Breadth of Candidate Coverage**
- **Score**: 10
- **Justification**: The report stands out for its comprehensive coverage of 2D
→ materials (e.g., MXenes, graphene, TMDCs, BN) and molecular probes (e.g.,
→ aptamers, MIPs, fluorinated phospholipid bilayers), including innovative
  concepts like hybrid systems and machine learning optimization.
### **Summary of Scores**
- **Relevance to the Stated Objective**: 9
- **Depth and Rigor of Validation and Benchmarking Protocols**: 7
- **Clarity and Reproducibility of Methodology**: 8
- **Practical Applicability and Implementation Readiness**: 7
```

```
- **Novelty and Breadth of Candidate Coverage**: 10

---

### **Overall Assessment**
The report is highly relevant to the objective of identifying 2D nanomaterials

and molecular probes for FET-based PFAS detection and excels in its

innovative and broad coverage of candidates. However, it would benefit from

stronger validation protocols, including detailed benchmarking, and more

specific guidance on practical implementation to enhance its real-world

applicability. These improvements would bolster its procedural rigor and

readiness while preserving its strengths in relevance and novelty.
```

Example raw response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe) — ChatGPT o3

```
Overall, the report captures the current consensus that MXene- and
   aptamer-functionalised FETs can drive PFAS detection into the femto- to
   picomolar regime, mirroring recent experimental demonstrations of ultra-low
  limits of detection and hybrid material strategies.
  ([pmc.ncbi.nlm.nih.gov][1], [mdpi.com][2]) However, it gives only a
  qualitative treatment of validation and field-deployment protocols that
  agencies such as the U.S. EPA and DoD now require, leaving practical
   benchmarking details thin. ([epa.gov][3], [denix.osd.mil][4]) The table
   below assigns balanced scores across the five requested dimensions.
## Evaluation Table
                                                       | Score (0-10) |
| Dimension

→ One-sentence justification

          ______
   .
|------|
   -
       _____
\hookrightarrow
| **1. Relevance to Objective**
\hookrightarrow discussion centres on MXenes, graphene, TMDCs and probe chemistries
\hookrightarrow (aptamers, MIPs) that literature identifies as the most sensitive FET
\,\hookrightarrow\, transducers for PFAS and analogous trace contaminants, directly matching the

    stated goal. ([pmc.ncbi.nlm.nih.gov][1], [mdpi.com][5],

  [pmc.ncbi.nlm.nih.gov][6]) |
| **2. Depth & Rigor of Validation/Benchmarking**

→ While impressive detection limits are quoted, the report lacks specific

→ inter-laboratory comparisons, LOD calculation methods, or reference to
  emerging EPA/DoD validation workflows, so its rigor is only moderate.
  ([epa.gov][3], [denix.osd.mil][4])
| **3. Clarity & Reproducibility of Methodology**
                                                      l **7**
                                                                    l A
  clear LbL/R2R/CVD integration scheme is laid out, but omission of key
   experimental parameters (e.g., gate bias, buffer composition, chip layout)
   hampers full reproducibility expected in contemporary FET biosensor reports.
  ([exploration-journals.com][7], [pubs.rsc.org][8])
| **4. Practical Applicability & Implementation Readiness** | **6**
→ narrative recognises roll-to-roll graphene films and protective coatings to
\hookrightarrow curb MXene oxidation, yet stops short of cost modelling, encapsulation, or
\rightarrow shelf-life testing that determine commercial readiness.
  ([link.springer.com][9], [pubs.rsc.org][10], [sciencedirect.com][11])
```

```
| **5. Novelty & Breadth of Candidate Coverage**
                                                               | **8**
\hookrightarrow Inclusion of under-explored fluorinated bilayers, machine-learning-optimised
→ aptamers, and MXene/graphene heterostructures shows commendable breadth

→ beyond mainstream materials and highlights emerging cross-disciplinary

→ avenues. ([nature.com][12], [sciencedirect.com][13], [mdpi.com][2])
### Key Take-aways
* **Strengths: ** High topical relevance, comprehensive material/probe coverage,

→ and acknowledgment of scalability challenges.

* **Gaps to Address: ** Explicit validation protocols, detailed experimental
\hookrightarrow conditions, and economic/packaging analyses would improve rigor and
\hookrightarrow real-world impact.
By closing these gaps, the report could rise to a top-tier reference for
\hookrightarrow PFAS-oriented FET sensor development.
[1]: https://pmc.ncbi.nlm.nih.gov/articles/PMC12109899/?utm_source=chatgpt.com
→ "Aptamer-Based Planar Electric Double-Layer Field-Effect Transistor"
[2]: https://www.mdpi.com/2079-6374/14/10/497?utm_source=chatgpt.com "Recent
\hookrightarrow Advancements in MXene-Based Biosensors for Health and ..."
[3]: https://www.epa.gov/water-research/pfas-analytical-methods-development-and
\rightarrow -sampling-research?utm_source=chatgpt.com "PFAS Analytical Methods

→ Development and Sampling Research"

[4]: https://www.denix.osd.mil/edqw/denix-files/sites/43/2023/02/Module-6-Data-
→ Validation-Guidelines-1633-PFAS-Final-1.pdf?utm_source=chatgpt.com "[PDF]
→ Module 6 Data Validation Guidelines - 1633 PFAS Final - DOD DENIX"
[5]: https://www.mdpi.com/2227-9040/12/7/116?utm_source=chatgpt.com "First
→ Direct Gravimetric Detection of Perfluorooctane Sulfonic Acid ..."
[6]: https://pmc.ncbi.nlm.nih.gov/articles/PMC7944395/?utm_source=chatgpt.com
→ "MXene-Graphene Field-Effect Transistor Sensing of Influenza Virus ..."
[7]: https://www.exploration-journals.com/article/aptamer-functionalized-field-
\hookrightarrow effect-transistor-biosensors-for-disease-diagnosis-and-environmental-monito
\hookrightarrow ring/?utm_source=chatgpt.com "Aptamer-functionalized field-effect transistor \hookrightarrow biosensors for disease ..."
[8]: https://pubs.rsc.org/en/content/articlehtml/2022/ma/d1ma00929j?utm_source=|
\hookrightarrow chatgpt.com "Microfluidic sensors based on two-dimensional materials for
[9]: https://link.springer.com/article/10.1007/s40820-024-01534-x?utm_source=ch |
→ atgpt.com "Flexible Graphene Field-Effect Transistors and Their Application
\hookrightarrow \quad \text{in} \ \dots \text{"}
[10]: https://pubs.rsc.org/en/content/articlehtml/2024/tc/d4tc00884g?utm_source |
\hookrightarrow =chatgpt.com "Shielding 2D MXenes against oxidative degradation - RSC
[11]: https://www.sciencedirect.com/science/article/abs/pii/S1385894724055864?u

→ ScienceDirect"

[12]: https://www.nature.com/articles/s41467-021-22555-9?utm_source=chatgpt.com
→ "Machine learning guided aptamer refinement and discovery - Nature"
[13]: https://www.sciencedirect.com/science/article/abs/pii/S0360319921047455?u
\hookrightarrow tm_source=chatgpt.com "Heterostructured two dimensional materials of MXene
\hookrightarrow and graphene ..."
```

Example raw response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe) — ChatGPT o4 mini high

```
Here are the scores and one-sentence justifications for each dimension:

## 1. Relevance to the Stated Objective

**Score: 9/10**
```

The report aligns strongly with the goal of identifying 2D materials and \hookrightarrow molecular probes for PFAS detection in FET sensors, covering both

 \hookrightarrow sensitivity and selectivity considerations, with only minor omissions around

 \hookrightarrow alternative non-FET platforms.

2. Depth and Rigor of Validation and Benchmarking Protocols

Score: 6/10

While key performance metrics (e.g., detection limits, specificity) are cited, \hookrightarrow the report lacks detailed descriptions of experimental controls, statistical \hookrightarrow analyses, and comparative benchmarks across different sensing platforms.

3. Clarity and Reproducibility of Methodology

Score: 7/10

Integration strategies and trade-offs are clearly delineated, but the absence of \hookrightarrow step-by-step protocols (e.g., exact functionalization procedures, probe \hookrightarrow immobilization conditions) limits precise reproduction.

4. Practical Applicability and Implementation Readiness

Score: 8/10

The discussion of scalability (R2R, LbL, CVD) and durability solutions offers \hookrightarrow actionable guidance for commercialization, though cost-analysis and \hookrightarrow regulatory considerations could be more fully addressed.

5. Novelty and Breadth of Candidate Coverage

Score: 8/10

The inventory spans established (graphene, TMDCs) and emerging (fluorinated \hookrightarrow bilayers, MXENE hybrids) candidates, but could further explore under-studied \hookrightarrow 2D systems (e.g., phosphorene, borophene) and alternative probe chemistries.

Table 1: Example raw response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe; Trial #1) — cross-model score matrix (0–10).

Model	Relevance	Depth	Clarity	Applicability	Novelty	Overall (mean)
Claude 4 Opus (thinking)	9	6	7	8	8	7.6
Gemini 2.5 Pro	10	7	6	9	10	8.4
Grok 3 (thinking)	9	7	8	7	10	8.2
ChatGPT o3	9	6	7	6	8	7.2
ChatGPT o4 mini high	9	6	7	8	8	7.6
Average across models	9.2	6.4	7.0	7.6	8.8	7.8

A/B Dueling Prompt (DR Report A vs. DR Report B)

Report A: {Anonymous DR Report content} Report B: {Anonymous DR Report content}

As an experienced materials scientist who bridges the gap between research and pilot lines, please review these two reports. I need your professional judgment: which one offers the clearer high-level strategic roadmap for our team's R&D over the next six months? We are especially interested in the report that presents a well-structured narrative, systematically identifies and prioritizes knowledge gaps, and synthesizes actionable next steps—beyond performance metrics. Please give me your recommendation in a brief, plain-language note—straight to the point—ending with a clear verdict on which report to back and the single most decisive reason why.

A/B Duel — Metadata

agent_a: DToR_qwq32B_loca1500_RD3_RAG5

agent_b: ChatGPT_o3_DR

trial: trial_1

Raw dueling response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe)— Claude 4 Opus (thinking)

Professional Assessment: Report A vs Report B

After reviewing both reports as a materials scientist with pilot-line
→ experience, **Report A is the clear choice** for your R&D roadmap.

Why Report A wins:

Report A provides a structured, decision-ready framework that directly addresses

your team's needs. It synthesizes findings across five research branches,

explicitly identifies contradictions (like MXene sensitivity vs. stability),

 \hookrightarrow and most importantly, provides **resolution strategies** for each challenge. \hookrightarrow The report prioritizes MXene-aptamer systems while acknowledging their

 $\hookrightarrow\,$ limitations and offering concrete solutions like hybrid coatings and R2R

 \hookrightarrow manufacturing approaches.

Report A also includes a practical "Candidate Inventory" table that ranks

→ materials by performance metrics and a clear action plan focusing on three

→ critical barriers: stability testing, standardized protocols, and

→ cross-disciplinary collaboration.

Why Report B falls short:

Report B is comprehensive but reads more like a literature review than a

strategic roadmap. While it contains valuable technical depth, it lacks
prioritization and actionable synthesis. The extensive material descriptions
don't clearly guide decision-making, and the conclusions are more
boservational than directive.

The decisive factor: Report A explicitly addresses the "how to move forward" \hookrightarrow question with ranked recommendations and specific next steps, while Report B \hookrightarrow provides extensive "what exists" information without clear strategic \hookrightarrow direction.

Verdict: Back Report A - it's the roadmap your R&D team can immediately act \hookrightarrow upon.

Raw dueling response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe)— Gemini 2.5 Pro

To the R&D Lead,

After reviewing both documents, my recommendation is to use **Report A** as the \hookrightarrow strategic roadmap for your team's R&D efforts over the next six months.

While Report B is an excellent and exhaustive technical review, Report A is structured as an executive-level strategic plan. It moves beyond a simple survey of materials and performance metrics to actively synthesize the information into a coherent strategy. It identifies the most critical trade-offs (e.g., Sensitivity vs. Stability), frames them as solvable problems in its "Contradiction Analysis," and provides a clear, prioritized conclusion with actionable next steps for your team.

My verdict is to back **Report A**.

The single most decisive reason is: **Report A is a decision-making tool, not \hookrightarrow just a literature review; its structure forces a strategic focus on \hookrightarrow resolving key conflicts, which is exactly what a short-term R&D plan \hookrightarrow requires.**

Raw dueling response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe)— Grok 3 (thinking)

After reviewing both reports as a materials scientist with experience in scaling

R&D from lab to pilot, Report A stands out for its concise, integrated

structure that weaves together multi-branch insights into a cohesive

narrative. It systematically flags key gaps (e.g., stability in humid

environments, scalability via manufacturing methods) and prioritizes them

through contradiction resolutions and a focused conclusion, while distilling

actionable steps like stability testing, protocol standardization, and

collaborative innovations to guide near-term R&D.

Report B, while thorough and data-rich, reads more like a broad literature → review with less emphasis on gap prioritization and explicit roadmapping, → making it better for background reference than strategic planning.

Verdict: Back Report A. The single most decisive reason is its synthesis of \hookrightarrow actionable next steps into a clear, prioritized R&D framework that directly \hookrightarrow addresses real-world deployment barriers beyond just metrics.

Raw dueling response — ChatGPT o4 mini high

After reviewing both documents, Report A delivers a concise, goal-oriented or roadmap that explicitly prioritizes knowledge gaps (e.g. MXene stability, scalable integration, standardized testing) and pairs each with actionable next steps (hybrid coatings development, R2R pilot runs, cross-disciplinary protocol design). In contrast, Report B, while rich in technical depth, reads more like an encyclopedic survey than a six-month R&D playbook.

Verdict: Back Report A-because it alone translates material insights into a

→ structured, gap-driven action plan suitable for immediate project planning.

Raw dueling response (DToR_qwq32B_local500; Topic: PFAS FET sensor probe)— Chat-GPT o3

Both reports demonstrate deep technical insight, but they serve different $\ensuremath{\hookrightarrow}$ purposes.

- Report A reads like an executive playbook. It starts with a concise problem

 → statement, then walks us through material/probe choices, trade-offs,

 → contradictions, and cross-disciplinary innovations in a clearly labeled,

 → logical sequence. Crucially, it flags three "blocking" knowledge gaps (MXene

 → stability, manufacturing scale-up, and standardized testing) and ties each

 → gap to concrete, near-term actions (e.g., hybrid coatings, roll-to-roll

 → trials, protocol drafting). The Candidate Inventory table distills options

 → at a glance and can be plugged directly into a six-month work-plan.
- Report B is an encyclopedic review. It is rich in context and citations, but

 → the narrative meanders, mixes PFAS with many other analytes, and buries

 → action items inside long technical passages. The sheer length makes it

 → harder to extract priorities or sequence tasks.

 $\label{thm:prop} \mbox{Verdict: Back Report A-because it delivers a structured, priority-driven roadmap} \hookrightarrow \mbox{we can act on immediately.}$

E First-Principles Simulation

E.1 DFT simulation settings

Static Simulation

To simulate the binding interactions of diverse combinations under periodic boundary conditions, we employed ab initio calculations employing parameters optimized for precision and computational efficiency. The plane-wave energy cutoff was set to 400 eV, and spin-polarized density functional theory DFT with the Perdew-Burke-Ernzerhof (PBE) functional was utilized to describe exchangecorrelation interactions. The Brillouin zone integration employed a Gamma centered $1 \times 1 \times 1$ k-point grid and van der Waals corrections were incorporated using the DFT-D3 method to capture dispersion effects. Ionic relaxations were conducted with a convergence criterion of 0.02 eV/Å for the forces, and the electronic structure iterations were set to converge to an energy difference of 10^{-4} eV. Implicit solvent field is added for static simulations by setting the dielectric constant at 78.4 if applicable (namely excluding the ab initio molecular dynamics (AIMD) with explicit solvent molecules, and the CO₂ sensor probe binding scenarios). The modeling of the slabs is directly conducted by human expert based on anonymously reading the DR report. To ensure fair comparison, we asked the expert to create same number of candidates from each report (e.g. 4 from DR report A, and 4 from DR report B), and also pick one domain commonly chosen benchmark material/chemical for each of the 5 applications. To prevent interactions between periodic images, we also include >20 Å vacuum gap along the z-, ensuring a reliable comparison between adsorbed and separated configurations.

Binding Energy Convention.

Following a standard DFT convention (Phys. Rev. B, 2008, 77, 12) for adsorption on typical studied 2D material systems like graphene, we report the binding energy ΔE_{bind} of an adsorbate at a given coverage as

$$\Delta E_{\text{bind}} = (E_{\text{clean slab}} + E_{\text{adsorbate}}) - E_{\text{adsorbed}} \tag{1}$$

where $E_{\rm adsorbed}$ is the total energy of the fully relaxed adsorbate—slab complex (for certain applications, bottom atoms of bulk crystals are frozen and surface area are free to mimic realistic surface interactions), $E_{\rm clean_slab}$ is the energy of the corresponding pristine clean slab supercell, and $E_{\rm adsorbed}$ is the energy of the isolated adsorbate in the same supercell. All terms are computed with identical exchange—correlation functional, plane-wave cutoff, and k-point mesh. Unless stated specifically (like for battery binder normalized to be contact areal) We report $\Delta E_{\rm bind}$ for only one adsorbate in the slab super cell for qualitative comparison. This definition mirrors the adsorption-energy convention widely used for molecules on typical substrates like graphene or MXene (we adopt the term "binding energy" throughout). By this convention, a more positive $\Delta E_{\rm bind}$ indicates an exothermic, favorable adsorption (energy released upon binding).

Ab Initio Molecular Dynamics (AIMD) Simulation

To improve the physical realism of our calculations and especially considering solvent effect for surface events, we embedded explicit water molecules in the supercell according to conformations identified in preceding static simulations (Mol. Syst. Des. Eng. 2025, 10, 345–356; Sci. Adv. 2025, 11, eadr9038) using an end-state protocol. Specifically, the optimized adsorbed configuration corresponding to $E_{\rm adsorbed}$ was used as the initial snapshot for the solvated adsorbed state, with explicit water molecules added thereafter. For the pulled state, the adsorbate was translated by at least 15 Å away from the relaxed slab (e.g., graphene)—without additional geometry optimization—and then solvated under identical conditions. These two equilibrated states form the basis for directly comparing adsorption energetics in explicit solvent, consistent with end-state comparisons commonly applied in MD umbrella-sampling studies (J. Phys. Chem. C 2023, 127, 5920–5930; ChemPhysChem 2018, 19, 690–702), and in line with benchmark explicit-solvent AIMD treatments of metal/water interfacial solvation (J. Chem. Phys. 2020, 152, 144703) and first-principles determinations of interfacial hydration free energies for adsorbates on Pt (J. Chem. Phys. 2021, 154, 094107).

The water density was maintained at 1 g cm⁻³ (for battery binder, solvent composition is different as describe in corresponding section). Our goal was to compare the evolution of the system's potential (electronic) energy upon equilibration (Electrocatalysis 2017, 8, 577–586; Commun. Chem. 2022, 5, 76; J. Chem. Phys. 2020, 152, 144703; J. Chem. Phys. 2021, 155, 134703).

To ensure that the solvent adopted physically meaningful configurations, each system was first equilibrated for 2 ps at 300 K with a 0.5 fs time step. Temperature was regulated via a Langevin thermostat to suppress extraneous inertial effects while preserving canonical sampling. After equilibrium, we further performed an 3 ps production run, from which we extracted time-averaged electronic potential energies.

We therefore report an explicit-solvent end-state potential-energy gap, denoted $\Delta E_{\rm pot}^{\rm ES}$ (AIMD), defined as the difference between the time-averaged electronic potential energies of the solvated adsorbed and pulled basins in identical cells:

$$\Delta E_{\text{pot}}^{\text{ES}} \equiv \langle E_{\text{pot}} \rangle_{\text{pulled}} - \langle E_{\text{pot}} \rangle_{\text{ads}}.$$
 (2)

Unless otherwise noted, all comparisons use the same supercell, water count, ion composition, thermostatting, and slab-electrostatics settings. Hence, the more positive $\Delta E_{\rm pot}^{\rm ES}$ (AIMD), the more energetically favorable the adsorbed state—i.e., a potential-energy analogue of the binding energy—whereas a negative $\Delta E_{\rm pot}^{\rm ES}$ indicates that explicit solvation stabilizes the pulled state and adsorption is disfavored.

Scope and limitations of $\Delta E_{\mathrm{pot}}^{\mathrm{ES}}$:

 $\Delta E_{
m pot}^{
m ES}$ is not a rigorously converged adsorption free energy. By construction it (i) neglects configurational and translational entropy; (ii) folds solvent reorganization and water (or other solvent)–surface exchange energetics into the potential-energy baseline; and (iii) depends on finite-time sampling of each basin rather than a reversible pathway or an alchemical thermodynamic cycle. Consequently, its absolute magnitudes can reflect not only adsorbate–surface interactions but also large-scale differences in hydrogen-bond networks and long-range electrostatics associated with hundreds of explicit waters. Nonetheless, because all systems are treated consistently within the same AIMD framework, $\Delta E_{
m pot}^{\rm ES}$ serves as a heuristic yet informative comparative descriptor for explicit-solvent binding trends across molecules and candidate surfaces, revealing solvent-mediated, surface-dependent effects that are not captured in vacuum or implicit-solvent calculations. Formal free-energy methods (umbrella sampling/PMF, metadynamics, and double-decoupling/TI with standard-state and slab-electrostatics corrections) are acknowledged as the rigorous routes to ΔG , as exemplified for interfacial adsorbates at water–Pt interfaces (J. Chem. Phys. 2021, 154, 094107).

Computational feasibility rationale:

Rigorous free-energy calculations become computationally prohibitive for our systems, which routinely exceed 500 atoms and include numerous transition-metal centers in addition to light elements; fully converged interfacial free-energy landscapes are beyond feasible AIMD timescales under such conditions. Extensive explicit-solvent AIMD benchmarks of metal/water interfaces have likewise highlighted the costliness of such simulations (J. Chem. Phys. 2020, 152, 144703). In this context, the explicit-solvent end-state potential-energy gap provides a practical dynamic compromise: it captures key solvent-mediated effects absent from static DFT binding energies while remaining tractable for large, chemically realistic slabs. Thus, although approximate, $\Delta E_{\rm pot}^{\rm ES}$ offers higher fidelity and greater interpretive value than purely static vacuum or implicit-solvent descriptions for the screening-level comparisons emphasized in this work.

E.2 PFAS FET Sensor Probe

For the first application: PFAS FET Sensor Probe. The query as shown in D.1 is:

PFAS FET Sensor Probe

Which two-dimensional (2D) nanomaterials or molecular probes-including novel, \hookrightarrow underexplored candidates-offer the greatest potential to achieve highly \hookrightarrow sensitive and selective detection of per- and polyfluoroalkyl substances \hookrightarrow (PFAS), or more broadly chemical and biological analytes, when integrated \hookrightarrow into a field-effect transistor (FET)-based sensor platform?

The best commercial solution is: ChatGPT_o3_DR. The best local solution is: DToR_qwq32B_local500.

The expert based on the two reports, provided 3 commercial DR report derived candidates (A1-A3), and 3 local DR report derived candidates (B1-B3), visualized in DL 11:

A1 — 2D rGO + β -Cyclodextrin (β -CD)

Hexagonal P1 slab with a = b = 36.90 Å, $\gamma = 120^\circ$, c = 45.00 Å (in-plane area ≈ 1179.19 Ų). Atom counts (C/H/O = 483/102/94) indicate a 2D graphene-derived sheet with oxygenated defects (epoxy/hydroxyl) on one face and a carbohydrate torus consistent with β-CD mounted above it. The populated thickness is ≈ 15.02 Å (vacuum ≈ 29.98 Å); the 2D carbon plane sits near $z \approx 7.5$ Å. The β-CD presents a toroidal pocket with a projected mouth \sim 15.3 Å diameter, oriented normal to the basal plane for PFOS capture. Oxygen sites on rGO act as plausible anchoring points while the underlying 2D lattice remains largely sp² and planar.

A2 — 2D MoS₂/2D hBN heterojunction + methacrylic-acid (MIP-like) cavity, pyrene-linked

Hexagonal P1 with a = b = 28.44 Å, $\gamma = 120^\circ$, c = 46.13 Å (area ≈ 700.47 Ų). Composition shows a **2D MoS**₂/**2D hBN heterojunction** stack capped by an aromatic (pyrene-type) tether and a **methacrylic-acid cluster forming a** β **-CD-like cavity** on one side (organics C/H/O sit above the dielectric). The populated thickness is ≈ 18.68 Å (vacuum ≈ 27.45 Å). The cavity's projected mouth is \sim **13.6** Å **in diameter**, positioned above hBN so PFOS can approach without perturbing the MoS₂ channel. The pyrene linker bridges π - π to the 2D cap while keeping the cavity axis roughly perpendicular to the surface, preserving electronic decoupling from the channel.

A3 — 2D Fluorographene (CF) slab

Hexagonal P1 with a = b = 36.90 Å, γ = 120°, c = 45.00 Å (area \approx 1179.19 Ų). Stoichiometry $C_{450}F_{450}$ yields a single-layer **2D** CF sheet with chair-like alternating C–F buckling; the "solid" thickness is \approx 3.26 Å and vacuum \approx 41.74 Å. The surface is uniform and free of additional adsorbates, providing a continuous fluorinated interface. Fluorine is arranged on both sides of the 2D carbon lattice, giving a symmetric fluorophilic environment. The wide lateral extent minimizes edge effects, making it suitable for baseline PFOS adsorption comparisons.

B1 — 2D Ti₃C₂O₂ MXene + methacrylic-acid (MIP-like) cavity

Hexagonal P1 with a=b=30.67 Å, $\gamma=120^\circ$, c=46.70 Å (area ≈ 814.63 Ų). The slab is an O-terminated **2D Ti**₃C₂ MXene bearing, on one side, a **methacrylic-acid cluster organized** as a β -CD-like cavity (additional C/H/O located at higher z). Populated thickness ≈ 16.14 Å (vacuum ≈ 30.56 Å). The cavity mouth projects \sim 13.5 Å in diameter, perched above the O-terminated surface while the bottom MXene layer serves as a rigid support. The one-sided functionalization creates an asymmetric interface—electronegative terminations and the cavity's polar rim cooperate to stabilize PFOS headgroup interactions.

B2 — 2D Graphene + methacrylic-acid (MIP-like) cavity

Hexagonal P1 with a = b = 36.90 Å, $\gamma = 120^{\circ}$, c = 45.00 Å (area $\approx 1179.19 \text{ Å}^2$). A $15 \times 15 \text{ 2D}$ graphene sheet supports, on one face, a **methacrylic-acid cluster arranged into a** β **-CD-like cavity** (C/H/O entirely above the carbon plane). The populated thickness is $\approx 9.58 \text{ Å}$ with $\approx 35.42 \text{ Å}$ vacuum. The cavity presents a projected $\sim 13.6 \text{ Å}$ **diameter** single-sided pocket aligned normal to the basal plane for PFOS engagement. The 2D sp² lattice remains flat and conductive, making this a clean platform for decoupling recognition from transduction.

B3 — 2D Graphene + fluorinated phospholipid monolayer

Hexagonal P1 with a = b = 36.90 Å, γ = 120°, c = 45.00 Å (area \approx 1179.19 Ų). Elemental makeup (C/F/H/O with trace N and P) is consistent with phosphatidylcholine headgroups near the **2D** graphene plane and perfluorinated tails extending outward to form a dense brush. The populated thickness is \approx 14.84 Å (vacuum \approx 30.16 Å), with F spanning roughly z \approx 9.21–22.34 Å above the basal plane. The assembly yields a hydrophilic-to-fluorophilic gradient, presenting a fluorinated outer surface for PFOS tail affinity. Lipid tails show a collective tilt relative to the surface normal, producing a compact, low-defect fluorinated interface on one side while the opposite side remains bare graphene.

Selectivity Scoring for PFAS Sensor Performance The practical deployment of PFAS sensors demands not only strong binding to target per- and polyfluoroalkyl substances but also selective discrimination against common environmental interferents. To quantitatively evaluate this selectivity, we developed a unified scoring framework applicable to both static DFT calculations and explicit-solvent AIMD simulations, enabling direct comparison between computational methodologies while maintaining thermodynamic rigor.

Selectivity Metric Derivation: The fundamental selectivity metric $\Delta \Delta E$ quantifies the differential binding preference between PFAS compounds and interferent molecules:

$$\Delta \Delta E = \overline{\Delta E}_{PFAS} - \overline{\Delta E}_{Interf}$$
 (3)

where $\overline{\Delta E}_{PFAS}$ represents the average binding energy across four representative PFAS compounds (PFOA, PFOS, PFNA, and PFHxS), and $\overline{\Delta E}_{Interf}$ denotes the average binding energy for interferent species (humic acid and SDS), as illustrated in Figure DL_1 2. A larger positive $\Delta \Delta E$ value indicates superior selectivity, signifying stronger preferential binding to PFAS over interferents.

For static DFT calculations, ΔE values correspond to the conventional binding energies $\Delta E_{\rm bind}$ computed under implicit solvent conditions (Figure DL_1 3). In explicit-solvent AIMD simulations, these are replaced by the time-averaged end-state potential energy gaps $\Delta E_{\rm pot}^{\rm ES}$, which capture solvent reorganization effects and hydrogen-bonding network rearrangements absent from static calculations (Figures DL_1 4 and DL_1 5).

Linear Selectivity Scoring Framework: To enable meaningful comparison across candidate systems and computational methods, we normalize the selectivity metrics to a common scoring scale:

Score =
$$100 \times (1 + \beta(\Delta \Delta E - \Delta \Delta E_{\beta\text{-CD}}))$$
 (4)

where $\Delta\Delta E_{\beta\text{-CD}}$ represents the selectivity of β -cyclodextrin, our domain benchmark reference with established PFAS binding capability [23]. The sensitivity parameter $\beta=0.2$ was selected to provide appropriate score discrimination while maintaining numerical stability across both DFT and AIMD datasets. This linear transformation preserves the physical meaning of energy differences while establishing β -CD as the baseline reference (Score = 100).

Comparative Analysis of Selectivity Under Different Solvation Models: The application of this unified scoring framework reveals striking differences between implicit and explicit solvation treatments. In static DFT calculations with implicit solvent, system A1 achieves the highest selectivity score of 102.4, outstandingly exceeding the β -CD baseline (Figure DL_1 6). The B-series candidates also prove competitive as B2 and B3 also showed slightly passing 100 score.

However, the inclusion of explicit water molecules and dynamic sampling in AIMD simulations dramatically alters the selectivity landscape (Figure DL_17). Most notably, all candidate systems

score below the β -CD baseline in explicit solvent, with B1 achieving the highest relative performance at 95.7. This systematic reduction in scores reflects the critical role of solvent-mediated effects: competitive water adsorption at binding sites, strong hydration of anionic PFAS headgroups, and the entropic penalties associated with constraining flexible molecules in adsorbed configurations.

The reversal of relative rankings between static DFT and AIMD calculations—particularly the superior performance of B-series candidates under explicit solvation—underscores the importance of including solvent dynamics in sensor design evaluation. While static calculations provide valuable initial screening, the explicit-solvent AIMD results more accurately capture the complex interplay between adsorbate—surface interactions and solvent reorganization that governs selectivity under realistic aqueous conditions.

Implications for Sensor Design: The divergence between implicit and explicit solvent results carries significant implications for PFAS sensor development. Systems optimized solely based on static DFT calculations may exhibit diminished selectivity in aqueous environments due to unforeseen solvent competition effects. Conversely, candidates that appear suboptimal in vacuum or implicit solvent calculations may benefit from favorable solvent-mediated interactions that enhance their practical selectivity. This highlights the necessity of multi-scale computational approaches that progress from efficient static screening to rigorous explicit-solvent validation for promising candidates.

The linear selectivity scoring framework presented here provides a quantitative, thermodynamically grounded metric for evaluating sensor performance across computational methodologies. By maintaining consistent notation and normalization across static and dynamic simulations, this approach enables systematic comparison of candidate materials while preserving the physical interpretability of the underlying energy differences.

PFAS Sensing Candidate Structures

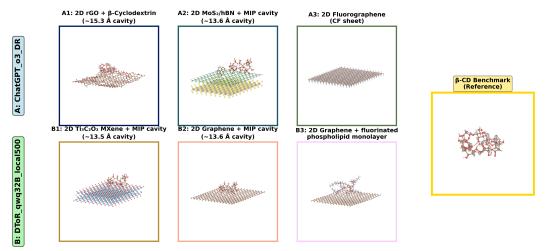


Figure DL_1 1: Structural configurations of candidates proposed by DR agents for selective PFAS binding, with β -cyclodextrin [23] included as the domain benchmark reference.

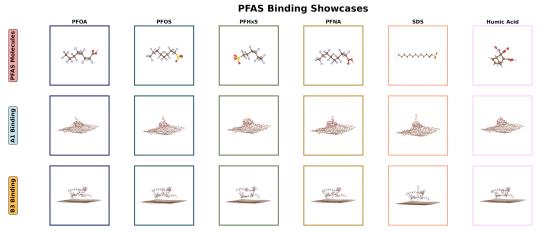


Figure DL_1 2: Representative PFAS compounds (PFOA, PFOS, PFNA, PFHxS) and interferent molecules (SDS and humic acid) employed in selectivity assessment, with corresponding optimized binding configurations. Two candidates' relaxed structures are shown.

PFAS Binding Energy (eV) Radar Chart

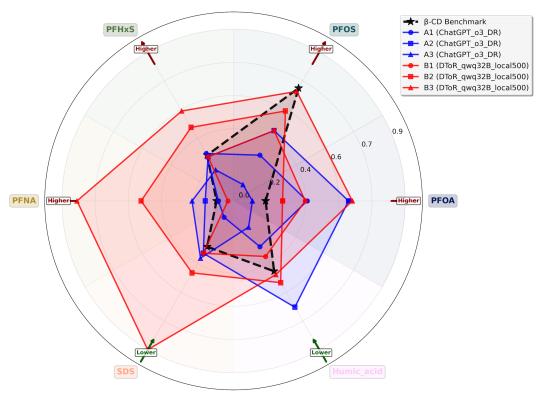


Figure DL_1 3: Radar plot comparing binding energies ($\Delta E_{\rm bind}$) across all candidate systems and benchmark β -CD obtained from static DFT calculations with implicit solvent.

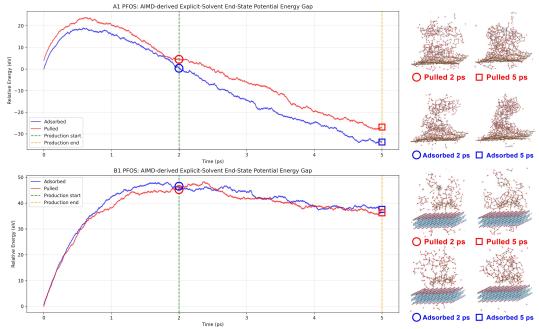


Figure DL_1 4: Representative AIMD trajectory showing relative energy evolution under explicit solvent conditions. The green dashed line denotes the 2 ps equilibration period. Right panels display system configurations at highlighted time points along the production trajectory.

AIMD-derived Explicit-Solvent End-State Potential Energy Gap (eV) Radar Chart (All Compounds)

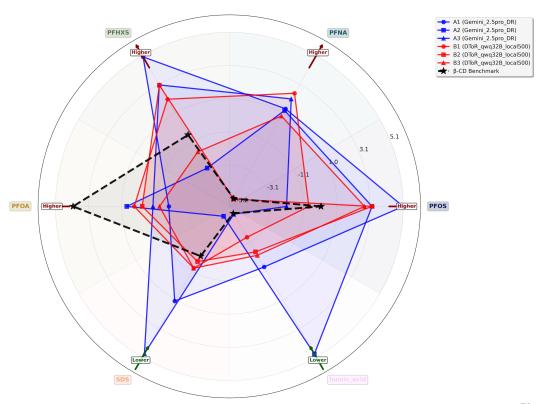


Figure DL_1 5: Radar plot comparing explicit-solvent end-state potential energy gaps ($\Delta E_{\rm pot}^{\rm ES}$) averaged over the 3 ps production period for all candidate systems and benchmark β -CD from AIMD simulations.

DFT Static: PFAS Sensor Selectivity Scores

Implicit solvent, linear selectivity metric

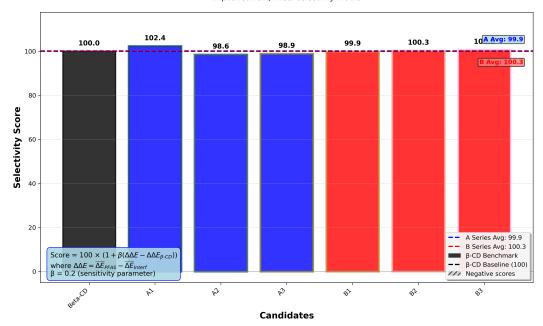


Figure DL_1 6: Selectivity scores for all candidate systems from static DFT calculations, normalized to β -CD baseline (100). Scores are calculated using the linear selectivity metric with sensitivity parameter $\beta=0.2$, where $\Delta\Delta E=\overline{\Delta E}_{\text{PFAS}}-\overline{\Delta E}_{\text{Interf}}$ represents the difference between average PFAS and interferent binding energies.

AIMD Explicit Solvent: PFAS Sensor Selectivity Scores

Explicit water molecules, captures solvent reorganization effects

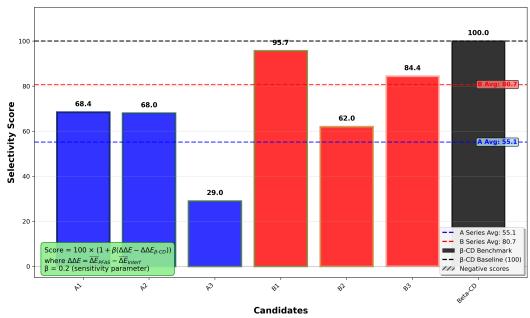


Figure DL_1 7: Selectivity scores for all candidate systems from explicit-solvent AIMD simulations, normalized to β -CD baseline (100). The scoring methodology remains identical to static DFT, except ΔE values are replaced by $\Delta E_{\rm pot}^{\rm ES}$, the time-averaged potential energy gaps from the 3 ps production period.

E.3 PFAS Degradation

For the second application: PFAS Degradation. The query as shown in D.1 is:

PFAS Degradation

Which novel electrode materials can achieve efficient PFAS degradation under

- \rightarrow ambient aqueous electrochemical conditions, delivering both high
- \hookrightarrow mineralization and defluorination rates? What intrinsic properties-such as
- \hookrightarrow PFAS adsorption affinity, reactive oxygen species generation capacity, and
- → C-F bond activation energy-should be optimized to guide their discovery?

The best commercial solution is: **Gemini_2.5_Pro_DR**.

The best local solution is: DToR_qwq32B_local500.

The expert based on the two reports, provided 5 commercial DR report derived candidates (A1-A5), and 5 local DR report derived candidates (B1-B5), visualized in DL_2 1:

Ti₄O₇ — Pristine Magnéli (001) bench slab

Triclinic P1 slab with a = 22.376 Å, b = 27.651 Å, γ = 75.306°, c = 42.208 Å (in-plane area \approx 598.50 Ų). Atom counts (Ti/O = 256/448; 704 total) match a Ti₄O₇ (001) termination. The populated thickness spans \approx 12.13 Å (z \approx 3.41–15.53 Å), leaving \approx 30.08 Å vacuum; the slab mid-plane lies near z \approx 9.47 Å. Surfaces are stoichiometric and flat, preserving Magnéli shearplane connectivity and metallic conduction pathways—our baseline conductive-oxide reference for PFAS tests. Build/relax notes: 4 atomic layers; bottom two constrained. (Adsorbates for subsequent runs were seeded \approx 3 Å above the most protruding Ti site, tail parallel to the surface.)

A1 — Ti_4O_7 (001) with through-pore + rim O-vacancies

Triclinic P1 cell with a = 22.376 Å, b = 27.651 Å, γ = 75.306°, c = 42.208 Å (in-plane area \approx 598.50 Ų). Atom counts (Ti/O = 249/437; 686 total) reflect removal of a top-side cluster and additional rim O-vacancies relative to the bench (Δ = -7 Ti, -11 O). The populated thickness is \approx 13.44 Å ($z \approx 3.00$ –16.45 Å), leaving \approx 28.77 Å vacuum; the slab mid-plane sits near $z \approx$ 9.73 Å. A circular through-pore was introduced by deleting atoms within a \sim 5 Å-radius cylinder from the upper layers, yielding an effective mouth \sim 1.0 nm diameter; the rim exposes under-coordinated Ti with a locally reduced O environment, while the surrounding lattice retains Magnéli shear-plane connectivity. Relaxation protocol used: bottom half constrained; top region (and any adsorbates) relaxed so the pore edge dominates binding/activation. Connection to the Ti₄O₇ bench: same lattice, cell vectors, and basal (001) termination as the pristine Ti₄O₇ bench slab; bottom region unchanged. All edits are localized to the top surface (pore machining + rim O-vacancies). Away from the pore, Ti/O stoichiometry per layer and shear-plane topology match the bench reference.

A2 — Pd-overlayer/ ${ m Ti_4O_7}$ (001) composite (amorphous Pd on the bench slab)

Triclinic P1 cell with a = 22.376 Å, b = 27.651 Å, γ = 75.306°, c = 42.21 Å (in-plane area \approx 598.50 Ų). Atom counts (Ti/O/Pd = 256/448/94; 798 total) indicate a stoichiometric Ti₄O₇ substrate capped by a Pd overlayer. The populated thickness spans \approx 16.45 Å (z \approx 3.14–19.59 Å), leaving \approx 25.76 Å vacuum; mid-plane at z \approx 11.37 Å. Pd atoms populate the top terrace broadly (no single island), forming a disordered sheet with Pd heights \sim 16.12–19.59 Å (mean \approx 17.71 Å), i.e., \sim 3.1 Å above the top oxide layer (z_{ox,max} \approx 16.49 Å). Connection to the Ti₄O₇ bench: identical lattice and Ti/O stoichiometry to the pristine Magnéli (001) slab; modifications are limited to a Pd overlayer on the top surface (bottom half retained as in the bench), with no pores or O-vacancies introduced.

$A3 - Pd_{55}$ cluster on Ti_4O_7 (001) (icosahedral-like nanoparticle on the bench slab)

Triclinic P1 cell with a = 22.376 Å, b = 27.651 Å, γ = 75.306°, c = 42.21 Å (in-plane area \approx 598.50 Ų). Atom counts (Ti/O/Pd = 256/448/55; 759 total) show a Ti₄O₇ substrate supporting a compact Pd nanoparticle. The populated thickness spans \approx 22.65 Å (z \approx 3.23–25.87 Å), leaving \approx 19.56 Å vacuum; mid-plane at z \approx 14.55 Å. The Pd cluster (initially icosahedral; allowed to relax) (n = 55) is localized on the top terrace with a projected footprint \sim 9.2 \times 9.6 Å (max in-plane radius \approx 5.33 Å; projected diameter \approx 10.65 Å) and vertical extent placing the apex \sim 9.7 Å above the top oxide layer (z_{ox,max} \approx 16.15 Å). Connection to the Ti₄O₇ bench: same lattice and stoichiometric oxide as the pristine (001) slab; the only alteration is a single Pd₅₅ nanoparticle anchored on the top surface (no pore machining or O-vacancies), with the bottom region kept as in the bench reference.

A4 — F-doped rutile SnO₂ (110) slab (FTO, O \rightarrow F substitution near surface)

Orthogonal P1 cell with a = 26.8022 Å, b = 25.4880 Å, γ = 90.0°, c = 42.1012 Å; in-plane area \approx 683.13 Ų. The metrics correspond to a 4×8 SnO₂(110) supercell (a \approx 4×a_(rutile) $\sqrt{2}$, b \approx 8×c_(rutile)). Atom counts (Sn/O/F = 256/497/15; 768 total) indicate \sim 2.9% O \rightarrow F substitution on the anion sublattice (15/512), i.e., a lightly F-doped "FTO" slab. The populated thickness spans \approx 14.28 Å (z \approx 4.32–18.60 Å), leaving \approx 27.83 Å vacuum; Sn layers occupy z \approx 4.38–15.29 Å, while F sites are confined to a narrow band near the top (z \approx 16.33–17.17 Å; mean \approx 16.75 Å), just beneath the outermost O ridges (O_{max} \approx 18.60 Å). This places F primarily in near-surface O positions on (110), preserving the characteristic bridging-O/in-plane-O relief of rutile and providing a slightly more electron-rich, less hydroxylated top termination expected for FTO.

A5 — B-doped four-layer graphene

Hexagonal P1; a = b = 19.68 Å, $\gamma = 120^\circ$, $c \approx 36.0$ Å. Four graphene sheets stacked with the standard interlayer spacing $d(002) \approx 3.35 \pm 0.05$ Å (three gaps $\rightarrow \sim 10.05$ Å slab thickness). Substitutional B is confined to the top layer (≈ 10 at% within that plane; $\approx 2-3$ at% overall), preserving an sp^2 framework. The remaining ~ 26 Å is vacuum along c to suppress image interactions. During relaxation, we fixed the bottom two layers; relax the top two (and any adsorbates).

B1 — F-grafted four-layer graphene (consensus spacing; A5-consistent)

Hexagonal P1; keep A5's in-plane metrics (a = b, γ = 120°) unchanged and c \approx 36.0 Å. Stack four graphene sheets with the standard interlayer spacing d(002) = 3.35 \pm 0.05 Å (three gaps $\rightarrow \sim 10.05$ Å slab thickness), with the remaining \sim 26 Å as vacuum along c to decouple periodic images. Chemistry/dopant pattern exactly as in A5 (no changes to which layer(s) are doped/functionalized); the carbon network remains predominantly sp². Same Relaxation protocol as A5: fix the bottom two layers; relax the top two (and any adsorbates). Connection to A5: identical spacing convention and vacuum budget as A5 to ensure cross-comparable adsorption energies and electrostatics across the graphene series, while leaving only the dopant/functional differences to drive behavior.

B2 — **2D** Ti₃C₂ MXene (O/OH-terminated)

Hexagonal P1 slab with a = b = 24.536 Å, $\gamma = 120^\circ$, c = 36.700 Å (in-plane area ≈ 521.36 Ų). Atom counts (Ti/C/O/H = 192/128/128/114; 562 total) correspond to **64 f.u. of Ti**₃C₂ with full two-sided terminations (128/site total): $\approx 89\%$ –**OH** (114 H) and $\approx 11\%$ =**O** (14 O), distributed symmetrically per face (57 –**OH** + 7 =**O** each side). The vacuum layer is around **27.21** Å. The sheet preserves the characteristic **2D MXene** core (Ti–C–Ti trilayers) with mixed O/OH terminations pointing outward, yielding a hydrophilic, electronically conductive surface appropriate for adsorption tests (e.g., PFAS headgroup coordination via –OH domains).

B3 — 2D MnO₂-like (birnessite-style) nanosheet with sparse one-sided –OH

Orthogonal P1 cell with a = b = 29.3628 Å, γ = 90.0°, c = 33.595 Å (in-plane area \approx 862.17 Ų). Atom counts (Mn/O/H = 216/411/18; 645 total) are consistent with a near-stoichiometric MnO₂ layer carrying light hydroxylation (H:Mn \approx 0.083). The populated thickness is \approx 9.32 Å (z \approx 0.82–10.14 Å), leaving \approx 24.28 Å vacuum. Element-resolved spans: Mn \approx 1.52–9.59 Å; O \approx 0.82–9.61 Å; H confined to \approx 7.12–10.14 Å on the outer face, indicating **one-sided –OH/water terminations** on a **2D** octahedral sheet. Relaxation note: we kept the lower half constrained and relaxed the top half/adsorbates; head-group binding is expected at surface –OH/O sites with the basal Mn–O framework remaining robust.

B4 — 2D MnO₂-like (birnessite-style) nanosheet, smaller supercell; one-sided –OH

Orthogonal P1 cell with a = b = 22.000 Å, γ = 90.0°, c = 33.610 Å (in-plane area = 484.00 Ų). Composition (Mn/O/H = 150/285/12; 447 total) again reflects a MnO₂ sheet with light hydroxylation (H:Mn \approx 0.08). The populated thickness is \approx 9.40 Å (z \approx 1.77–11.17 Å), leaving \approx 24.21 Å vacuum. Element-resolved spans: Mn \approx 2.12–10.48 Å; O \approx 1.77–10.99 Å; H \approx 8.43–11.17 Å, i.e., **–OH/water groups localized on the top face** of the **2D** sheet. Relaxation matches B3 (bottom region constrained, top free) to keep comparisons clean across the MnO₂ series.

B5 — Graphene-capped 2D MXene (derived from B2)

Derived directly from **B2** (same cell vectors and terminations), a **single-layer graphene** sheet is added parallel to the top face of the MXene and centered to span the full in-plane cell. The initial MXene \leftrightarrow graphene separation is set to \sim 3.35 Å (graphite d(002)), with **no covalent links**—pure van der Waals contact. Graphene is pristine sp² and allowed to relax laterally/vertically; the MXene beneath follows B2's relaxation protocol (bottom constrained, top free). The **c-axis vacuum is kept as in B2** to avoid inter-image coupling. **Connection to B2:** identical substrate; the only modification is the graphene capping layer.

Multi-Metric Evaluation of PFAS Degradation Catalysts

Effective catalytic degradation of PFAS compounds requires both strong surface adsorption to concentrate the pollutants and sufficient activation of the exceptionally stable C–F bonds. To comprehensively evaluate candidate catalysts for PFAS degradation, we employed a three-pronged computational assessment framework encompassing static binding energies, explicit-solvent dynamics, and C–F bond dissociation energetics, with Ti_4O_7 serving as the domain benchmark reference [49].

Integrated Performance Metrics

The degradation performance evaluation integrates three complementary metrics that capture distinct aspects of the catalytic process. First, the binding energy $\Delta E_{\rm bind}$ from static DFT calculations quantifies the thermodynamic driving force for PFAS adsorption onto catalyst surfaces (Figure DL_2 3). Strong binding ensures sufficient surface coverage and residence time for subsequent degradation reactions. Second, the explicit-solvent potential energy gap $\Delta E_{\rm pot}^{\rm ES}$ from AIMD simulations reveals how aqueous solvation affects the adsorption thermodynamics (Figures DL_2 7 and DL_2 8), accounting for competitive water adsorption and hydrogen-bonding network reorganization that static calculations cannot capture. Third, the C–F bond dissociation energy quantifies the activation barrier

for the rate-limiting defluorination step, computed by displacing fluorine atoms from adsorbed PFAS configurations to isolated positions within the supercell (Figure DL_2 4).

Scoring Framework for Degradation Performance

To enable quantitative comparison across these diverse metrics, we applied a unified linear scoring framework:

Score =
$$100 \times \left(1 + \beta \left(\Delta \bar{E} - \Delta \bar{E}_{\text{Ti}_4\text{O}_7}\right)\right)$$
 (5)

where $\Delta \bar{E}$ represents the mean energy value averaged over four representative PFAS compounds (PFOA, PFOS, PFNA, PFHxS), and $\beta=0.2$ serves as the sensitivity parameter. For binding metrics (both static DFT and AIMD), more negative $\Delta \bar{E}$ values indicate stronger adsorption and thus higher scores. Conversely, for C–F dissociation, lower energy barriers (less positive values) yield higher scores, reflecting more facile bond activation.

The distribution of C–F dissociation energies across different PFAS compounds and candidate surfaces exhibits substantial variability (Figure DL_2 5), with mean values ranging from approximately 2 to 6 eV (Figure DL_2 6). This variation underscores the importance of averaging over multiple PFAS species to obtain representative performance metrics, as individual compounds may exhibit preferential interactions with specific surface sites or crystallographic facets.

Performance Analysis Across Computational Methods

Application of this scoring framework reveals distinct performance patterns across the three evaluation metrics. In static DFT calculations, several candidates demonstrate superior PFAS binding relative to the Ti_4O_7 benchmark (Figure DL_2 9), with A-series materials generally exhibiting stronger adsorption energies. However, the inclusion of explicit solvent effects in AIMD simulations significantly alters these rankings (Figure DL_2 11), with most candidates showing reduced binding scores due to solvent-mediated destabilization effects. This divergence highlights the critical importance of accounting for aqueous solvation in evaluating catalyst performance under realistic conditions.

The C–F bond activation scores (Figure DL_2 10) provide complementary insights into degradation potential. Notably, candidates that excel in binding may not necessarily facilitate efficient C–F bond cleavage, as these properties depend on distinct electronic and geometric factors. The ideal degradation catalyst must balance strong PFAS adsorption with effective bond activation—a trade-off that our multi-metric scoring framework quantitatively captures.

Implications for Catalyst Design

The comprehensive evaluation framework presented here reveals that no single candidate uniformly excels across all performance metrics. Strong surface binding, while essential for PFAS capture, does not guarantee facile C–F bond activation. Similarly, materials optimized for bond dissociation may suffer from weak adsorption, particularly under aqueous conditions where solvent competition is significant. This multi-dimensional performance landscape emphasizes the need for balanced catalyst design that considers both thermodynamic (binding) and kinetic (bond activation) factors, evaluated under both implicit and explicit solvation conditions to ensure robust performance predictions for practical applications.

PFAS Degradation Catalyst Structures

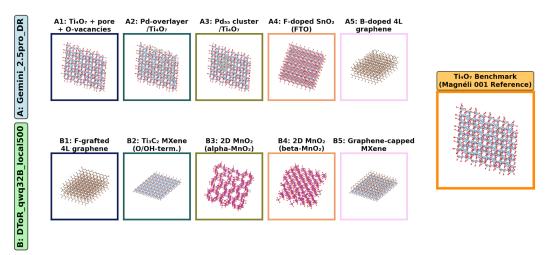


Figure DL_2 1: Structural configurations of candidate catalysts proposed by DR agents for PFAS degradation, with Ti₄O₇ [49] included as the domain benchmark reference.

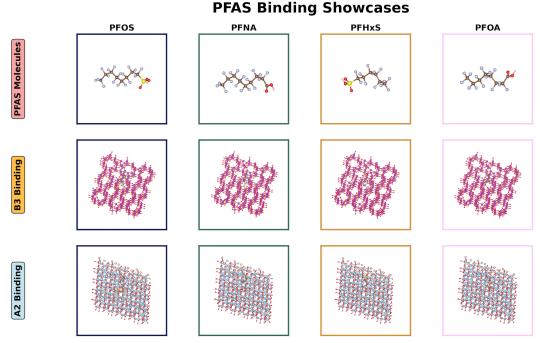


Figure DL_2 2: Representative PFAS compounds (PFOA, PFOS, PFNA, PFHxS) employed in degradation assessment, with corresponding optimized adsorption configurations on selected candidate surfaces.

PFAS Binding Energy (eV) Radar Chart

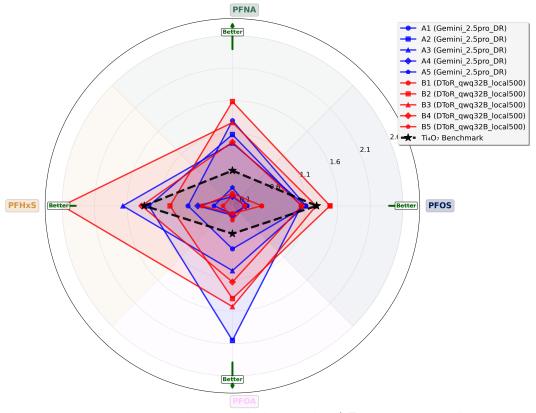


Figure DL_2 3: Radar plot comparing PFAS binding energies ($\Delta E_{\rm bind}$) across all candidate catalysts and benchmark Ti₄O₇ obtained from static DFT calculations with implicit solvent.

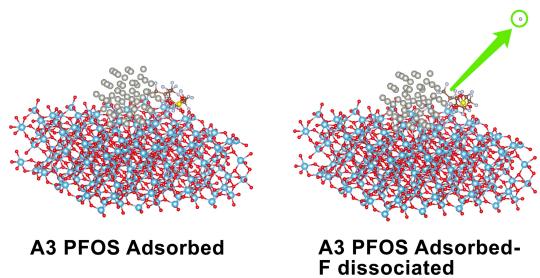


Figure DL_2 4: Illustration of C–F bond dissociation simulation protocol. Left: Initial configuration with PFOS adsorbed on A3 surface. Right: Final configuration with fluorine atom displaced to corner position (5 Å from each boundary) of the supercell.

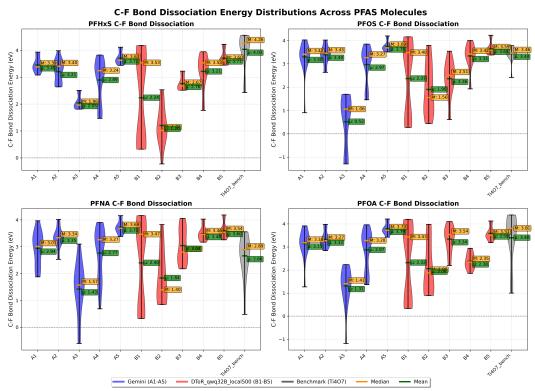


Figure DL_2 5: Violin plots showing distributions of C–F bond dissociation energies for different PFAS compounds (PFOA, PFOS, PFNA, PFHxS) across all candidate catalysts, with mean values indicated by horizontal bars.

PFAS Average C-F Bond Dissociation Energy (eV) Radar Chart

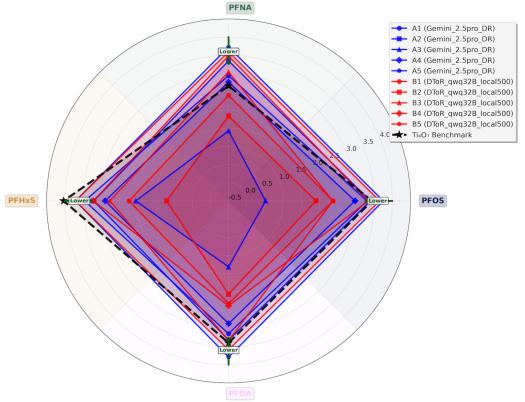


Figure DL_2 6: Radar plot comparing mean C–F bond dissociation energies across all candidate catalysts and benchmark ${\rm Ti_4O_7}$. Lower dissociation energies indicate more favorable C–F bond activation for degradation.

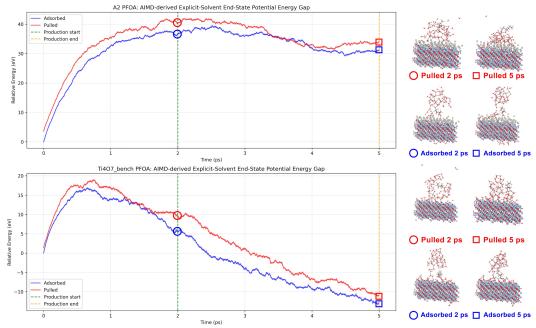


Figure DL_2 7: Representative AIMD trajectory showing system energy evolution under explicit solvent conditions. The green dashed line denotes the 2 ps equilibration period. Right panels display system configurations at highlighted time points along the 3 ps production trajectory.

AIMD-derived Explicit-Solvent End-State Potential Energy Gap (eV) Radar Chart

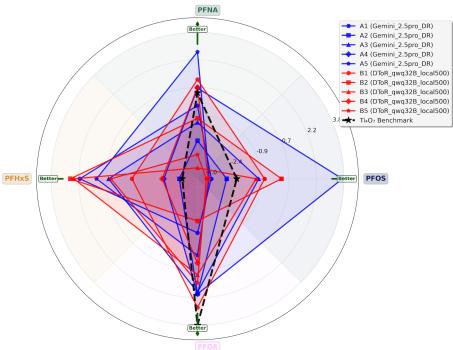


Figure DL_2 8: Radar plot comparing explicit-solvent end-state potential energy gaps ($\Delta E_{\rm pot}^{\rm ES}$) averaged over the 3 ps production period for all candidate catalysts and benchmark Ti₄O₇ from AIMD simulations.

DFT Static: PFAS Degradation Binding Scores

Implicit solvent, binding energy metric

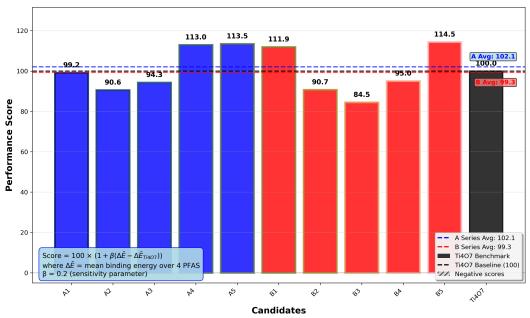


Figure DL_2 9: PFAS binding scores for all candidate catalysts from static DFT calculations, normalized to Ti₄O₇ baseline (100). Scores are calculated using the linear metric with sensitivity parameter $\beta=0.2$, where $\Delta \bar{E}=\overline{\Delta E}_{\text{PFAS}}$ represents the mean binding energy averaged over four PFAS compounds.

F Dissociation: PFAS Degradation Performance Scores

F atom dissociation energy, higher values indicate more stable C-F bonds

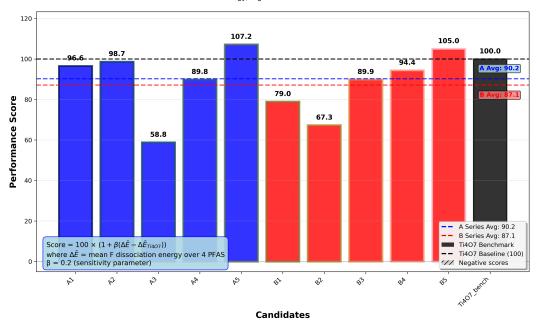


Figure DL_2 10: C–F bond activation scores for all candidate catalysts from static DFT calculations, normalized to Ti₄O₇ baseline (100). Scores are calculated using the linear metric with sensitivity parameter $\beta=0.2$, where $\Delta \bar{E}=\overline{E}_{\text{F-dis}}$ represents the mean C–F dissociation energy averaged over four PFAS compounds. Higher scores indicate lower dissociation barriers and more favorable degradation potential.

AIMD Explicit Solvent: PFAS Degradation Binding Scores

Explicit water molecules, captures solvent reorganization effects

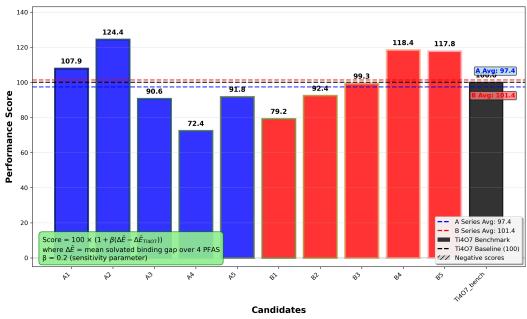


Figure DL_2 11: PFAS binding scores for all candidate catalysts from explicit-solvent AIMD simulations, normalized to Ti_4O_7 baseline (100). The scoring methodology remains identical to static DFT, with $\Delta \bar{E}$ representing the mean $\Delta E_{\text{pot}}^{\text{ES}}$ values averaged over four PFAS compounds from the 3 ps production period.

E.4 Battery Binder

For the third application: Battery Binder. The query as shown in D.1 is:

Battery Binder

Which advanced binder technologies are being developed to improve the \hookrightarrow performance and longevity of lithium-ion batteries?

The best commercial solution is: ChatGPT o4 mini high DR.

The best local solution is: DToR_qwq32B_local500.

The expert based on the two reports, provided 4 commercial DR report derived candidates (A1-A4), and 4 local DR report derived candidates (B1-B4), visualized in DL_3 1:

"latex

NCM811 cathode description

In order to create a test bench for lithium ion battery binder, we modeled a (104)-terminated slab of layered LiNi $_{0.8}$ Co $_{0.1}$ Mn $_{0.1}$ O $_2$ (NCM811) cathode in a periodic cell with a = 27.67 Å, b = 25.55 Å, c = 49.47 Å and α = β = 90°, γ = 72.51° (P1). The atomic slab spans \approx 12.1 Å along the surface normal, leaving \approx 37.4 Å of vacuum to suppress spurious interactions; the model contains 732 atoms (Li $_{174}$ Ni $_{129}$ Co $_{21}$ Mn $_{21}$ O $_{387}$). For structural relaxations, we used a large in-plane surface cell (tripled along b) and constrained the bottom \approx 1/3 of atoms to mimic bulk termination.

PVDF benchmark binder

For the benchmark binder, we built an all-trans PVDF 10-mer $(-CH_2-CF_2-)_{10}$, end-capped with H/F, pre-relaxed in vacuo, and then laid it flat on the pre-relaxed NCM811 (104) slab, constructed and constrained following our surface-build script. In the configuration, the chain backbone lies approximately parallel to the surface ($\sim 30^\circ$ off the b-axis), with a closest polymer–surface contact of ~ 2.0 Å, an in-plane backbone span of ≈ 23 Å, and a projected convex-hull footprint of ≈ 125 Å²; the polymer fragment contributes $C_{20}H_{21}F_{21}$, bringing the full model to 794 atoms.

A1 — PVPA (after binder adsorbed)

A poly(vinyl-phosphonic acid) 10-mer ($C_{20}H_{32}O_{30}P_{10}$) was adsorbed on the NCM811 (104) slab, using the same surface-build constraints as our base slab. The slab portion is 12.58 Å thick, leaving 36.88 Å of vacuum; the full system contains 824 atoms (slab = 732; polymer = 92). The PVPA backbone lies nearly parallel to the surface (\approx 34° off the b-axis), spans 20.59 Å in-plane, and covers a projected convex-hull footprint of 115.09 Ų. Phosphonate O atoms approach surface cations with a shortest O-Li contact of 2 Å, alongside O-O hydrogen-bond contacts to surface oxygen of 2.5–2.9 Å.

A2 — Alg-Ca (after binder adsorbed)

A Ca^{2+} -cross-linked alginate 5-mer (G–M blocks; carboxylates protonated except those coordinating Ca^{2+}) was laid approximately parallel to the NCM811(104) surface with the bottom oxide layers fixed. The combined model has 834 atoms (slab = 732; binder = 102, \approx $C_{30}H_{40}O_{31}Ca$). The slab portion is 12.20 Å thick with 37.27 Å vacuum. The polymer backbone aligns quasi-parallel to the surface (24.62° off **b**), spans 17.71 Šin-plane, and presents a convex-hull footprint of 92.64 Ų. Interfacial contacts are dominated by non-specific physisorption (nearest heavy-atom approach 3.23 Å; O_{binder} - O_{slab} 3.31–4.17 Å), while Ca^{2+} remains intrachain-coordinated (closest Ca- O_{slab} 7.13 Å).

A3 — Sulfonyl-polyimide (S-PI) (after binder adsorbed)

A tri-repeat S-PI fragment ($C_{42}H_{14}N_6O_{20}S_3$, 85 atoms) was constructed following the imide-dianhydride/diamino-sulfone motif and combined with the standard NCM811(104) slab (total 817 atoms; slab = 732, binder = 85). The polymer lies parallel to the surface, oriented 37.49° off the *b*-axis (principal in-plane span 25.71 Å) with a projected convex-hull footprint of 152.10 Å². This configuration represents the separated reference: the nearest polymer heavy atom is 17.42 Å above the top oxide layer (no interfacial contacts; shortest $O_{binder} \cdots O_{slab} = 19.17$ Å; $O_{binder} \cdots Ni/Co/Mn = 19.75$ Å).

A4 — PTFE (-CF₂-CF₂-)₁₀ (after binder adsorbed)

A perfluoroethylene decamer ($C_{20}F_{42}$, 62 atoms) is laid flat on the NCM811(104) surface to provide a non-reactive binder baseline; the combined model has 794 atoms (slab = 732; polymer = 62). The slab portion is 12.20 Å thick with 37.27 Å vacuum. The chain axis lies 39.61° off **b** within the **ab** plane, with a principal in-plane span of 25.14 Å and a projected convex-hull footprint of 80.55 Å². Interfacial contacts are dispersive: the nearest heavy-atom approach is 2.49 Å (F-Li); the shortest F-O_slab and F-Ni contacts are 2.84 Å and 3.43 Å, respectively, while C-O_slab \geq 3.63 Å and no metal-fluorine contact falls below 2.49 Å. These metrics set a PTFE benchmark for subsequent ΔE_a ds/A comparisons.

B1 — Schiff-base polymer (imine-linked poly(aldehyde-amine)) (after binder adsorbed)

An imine-linked poly(aldehyde-amine) fragment ($C_{40}H_{28}N_4$, 72 atoms) is adsorbed on the NCM811(104) slab; the combined model contains **804** atoms (slab = 732; binder = 72). The oxide portion is **13.07** Å thick with **36.39** Å vacuum. The polymer lies quasi-parallel to the surface, oriented **36.97**° off **b** within the **ab** plane; its principal in-plane span is **34.60** Å, and the projected convex-hull footprint is **168.51** Å². Interfacial contacts show a nearest heavy-atom approach of **2.57** Å (C-Li); the imine N approaches surface oxygen at **2.35** Å, while separations to cation sites remain larger (N-Li **3.68** Å; N-Ni **4.29** Å; N-Co **4.57** Å; N-Mn **4.76** Å).

B2 — lignin + CNT (after binder adsorbed)

A β -O-4 lignin dimer was positioned at the NCM811(104) interface and co-laid with an armchair (6,6) carbon nanotube segment (8 axial repeats; 192 C) oriented roughly parallel to the terrace; the tube's in-plane principal axis is rotated **36.86**° off **b** and spans **24.91** Å. The composite binder portion contains $\mathbf{C}_{234}\mathbf{H}_{80}\mathbf{O}_{46}$. The oxide portion spans **11.05** Å in thickness along the surface normal with **38.42** Å of vacuum retained. Lignin's phenolic/ether oxygens face the cation-terminated ridges, with the closest **O-(Ni/Co/Mn)** approach of **3.32** Å, consistent with physisorptive H-bond/ion-dipole contacts rather than direct chelation; the CNT lies flat to provide an extended π -contact, with a near-surface carbon convex-hull footprint of **84.71** Ų. Together, this layout tests cooperative anchoring: localized polar adhesion from lignin complemented by broad dispersive contact from the nanotube across the NCM811 terrace.

B3 — Lignin–CMC hybrid (β -O-4 lignin dimer + CMC monomer, Na⁺ exchanged) (after binder adsorbed)

The composite contains **876 atoms** (slab = 732; binder = 144 with $C_{50}H_{69}O_{24}Na_1$) and uses the same NCM811(104) platform and assembly recipe (lignin dimer H-bonded to a CMC unit, carboxylate presented to the terrace). The organic complex lies quasi-parallel to the surface, its principal in-plane axis oriented **38.18**° off **b**, with a **35.09** Å backbone span and a projected convex-hull footprint of **235.11** Å². The oxide portion is **11.09** Å thick with **38.38** Å of vacuum along **c**. Interfacial approach distances indicate mixed ionic/coordination contacts: nearest heavy-atom approach **2.25** Å (O_binder-Li), with **O_binder-Ni** = **2.99** Å, **O_binder-Co** = **4.09** Å, and **O_binder-Mn** = **4.29** Å; the exchanged **Na**⁺ sits remote from the terrace (**Na-O_slab** = **6.80** Å). During relaxations, the **bottom one-third of oxide atoms** are held fixed to emulate a semi-infinite substrate.

B4 — ZIF-8 face (Zn₄(2-mim)₈ cluster) (after binder adsorbed)

We approximate the MOF by a single sodalite-cage face—four Zn nodes bridged by eight 2-methylimidazolates—so the interface samples the correct local Zn–N₄/ μ_2 -imidazolate chemistry without imposing the artificial pore matching and long-range lattice constraints of a periodic ZIF on a finite oxide slab. The internal Zn–N network is kept near its intrinsic geometry via mild constraints, while only the interfacial coordinate is relaxed; one Zn is oriented toward a surface oxygen with an initial Zn–O set to ~2.10 Å to seed chemisorption and allow secondary N-Li/C-O_slab contacts to adjust around it. This charge-neutral ring thus captures the short-range bonding physics that governs adhesion at comparable cost to polymer fragments, enabling fair ranking by area-normalized adsorption energy ΔE_a ads/A across all candidates.

Battery Binder Adhesion Performance

The mechanical integrity of lithium-ion battery cathodes critically depends on binder adhesion strength, particularly under the chemically aggressive conditions of liquid electrolyte environments. To evaluate candidate binder materials proposed by the DR agents (Figure DL_3 1), we employed the same computational framework described for PFAS sensors, with binding energies normalized by the interfacial contact area A to account for varying polymer–cathode interaction footprints. The contact area was determined through convex hull analysis of atoms within 3.5 Å of the NCM811 surface, enabling direct comparison of adhesion strength per unit interfacial area.

Static DFT calculations reveal favorable adhesion for all candidates, with areal-normalized binding energies ranging from 0.015 to 0.111 eV/Å 2 (Figure DL_3 2). These values, equivalent to 0.24–1.78 J/m 2 , align well with established benchmarks for dry interface adhesion, such as graphene on SiO $_2$ (\sim 0.45 J/m 2) and metal–oxide contacts (\sim 2 J/m 2). Candidates B4 and B1 demonstrate the strongest adhesion, while the benchmark PVDF falls in the mid-range, suggesting potential improvements over current commercial binders.

However, the introduction of explicit battery electrolyte (EC:EMC:LiPF₆ = 9:9:2:2 at 1.25 g/cm³) through AIMD simulations dramatically alters the adhesion landscape (Figures DL_3 3 and DL_3 4). Following the same protocol of 2 ps equilibration and 3 ps production runs, most candidates exhibit negative $\Delta E_{\rm pot}^{\rm ES}/A$ values, indicating that the separated, electrolyte-wetted interface becomes energetically favorable over the bound state. This solvent-induced adhesion reversal reflects competitive adsorption of electrolyte molecules at the binder–cathode interface, disrupting direct polymer–surface interactions.

Notably, two exceptions emerge from this trend: candidate A3 (sulfonyl-polyimide) and B2 (lignin with CNT) maintain positive $\Delta E_{\rm pot}^{\rm ES}/A$ values even in explicit electrolyte, suggesting robust wet adhesion. This resilience likely arises from multidentate polar anchoring groups and extended π -conjugated contacts that effectively displace the first electrolyte layer. The performance scoring normalized to PVDF baseline (Figures DL_3 5 and DL_3 6) quantitatively captures this wet-interface reordering, where candidates optimized for dry adhesion may fail under operational electrolyte conditions while others demonstrate unexpected electrolyte compatibility. These findings underscore the critical importance of explicit-solvent validation in battery binder design, as interfacial energetics in liquid electrolyte environments cannot be reliably predicted from vacuum or implicit-solvent calculations alone.

Battery Binder Materials on NCM811 Cathode

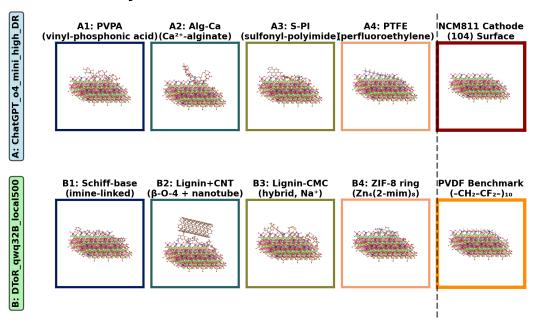


Figure DL_3 1: Structural configurations of battery binder candidates proposed by DR agents, with PVDF [50] included as the domain benchmark reference. Optimized adsorption configurations on the NCM811 cathode surface are shown for each candidate system.

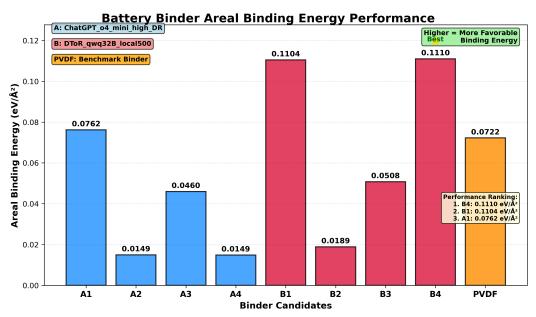


Figure DL_3 2: Comparison of areal-normalized binding energies ($\Delta E_{\rm bind}/A$) from static DFT calculations. Values are expressed in eV/Å², where A represents the interfacial contact area determined by convex hull analysis of atoms within 3.5 Å of the NCM811 surface.

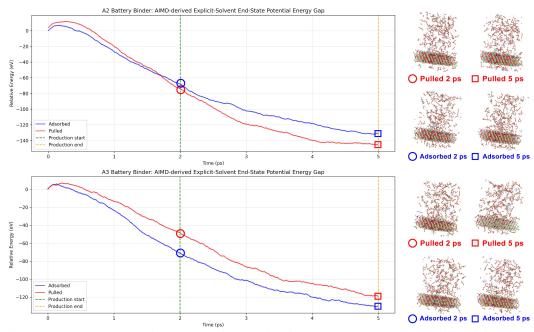


Figure DL_3 3: Representative AIMD trajectory showing system energy evolution in explicit battery electrolyte (EC:EMC:LiPF $_6$ = 9:9:2:2). The green dashed line denotes the 2 ps equilibration period. Right panels display system configurations at highlighted time points along the 3 ps production trajectory.

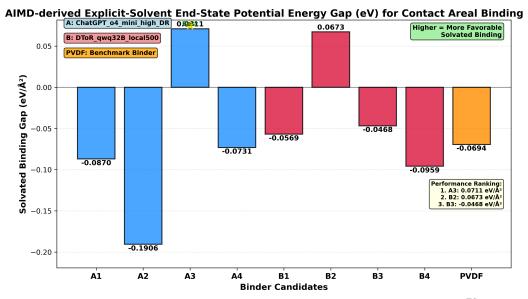


Figure DL_3 4: Areal-normalized explicit-solvent end-state potential energy gaps ($\Delta E_{\rm pot}^{\rm ES}/A$) from AIMD simulations. Negative values indicate unfavorable adhesion in the presence of electrolyte (Pulled state in explicit solvation stablizes the pulled state), reflecting competitive solvent adsorption at the binder–cathode interface.

DFT Static: Battery Binder Adhesion Performance Scores

Contact Areal Binding Energy (eV/Ų), static optimization

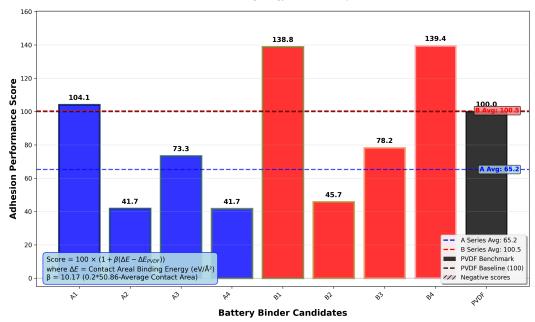


Figure DL_3 5: Adhesion performance scores for all binder candidates from static DFT calculations, normalized to PVDF baseline (100). Scores are calculated using the linear scaling metric with sensitivity parameter $\beta=10.17$, applied to areal-normalized binding energies.

AIMD Solvated Binding Gap: Battery Binder Adhesion Performance Scores

AIMD-derived Explicit-Solvent End-State Potential Energy Gap (eV) for Contact Areal Binding

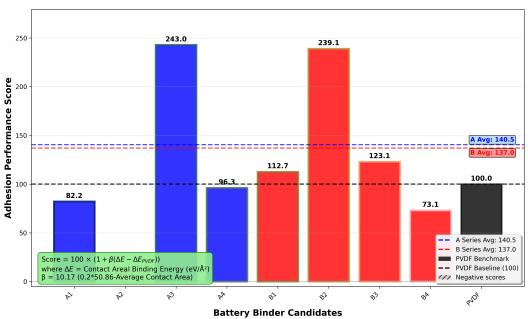


Figure DL_3 6: Adhesion performance scores for all binder candidates from explicit-electrolyte AIMD simulations, normalized to PVDF baseline (100). The scoring methodology remains identical to static DFT, with $\Delta E_{\rm bind}/A$ replaced by time-averaged $\Delta E_{\rm pot}^{\rm ES}/A$ from the 3 ps production period.

E.5 OER Stability

For the fourth application: OER Stability. The query as shown in D.1 is:

OER Stability

Identify the top nanostructured electrocatalyst materials for driving the oxygen

- → high-chloride, high-organic-load, or multi-ion streams-, for each the
- → material class/composition, comprehensively consider key performance metrics
- ightarrow surface-engineering strategies that confer corrosion resistance and

The best commercial solution is: **Perplexity_DR**.

The best local solution is: DToR_qwq32B_local500.

The expert based on the two reports, provided 4 commercial DR report derived candidates (A1-A4), and 4 local DR report derived candidates (B1-B4), visualized in DL 4 1:

Benchmark 110 IrO₂

For benchmark material, we created most commonly seen, a stoichiometric rutile IrO_2 (110) slab built as a 6 × 3 in-plane supercell along × [1 $\bar{1}0$], with cell parameters **a = 19.06** Å, **b = 19.11** Å, **c = 45.18** Å, $\alpha = \beta = \gamma = 90^{\circ}$ (P1), from which the underlying bulk metrics map to $a \approx 4.51$ Å and $c \approx 3.18$ Å. The slab comprises **five O-Ir-O trilayers** (thickness ≈ 15.18 Å) separated by ≈ 30.0 Å of vacuum along c, yielding symmetric, O-terminated (bridging-O) surfaces that expose the conventional row motif used for OER modeling (adjacent Ir_{cus}/O_{br} sites). The model is **strictly stoichiometric** (**Ir:O = 1:2; 180 Ir and 360 O**) and thus suitable as a clean, non-polar benchmark surface for computational electrocatalysis studies.

A1 — Ni-Fe layered double hydroxide (LDH) slab (basal)

A hexagonal brucite-derived Ni–Fe hydroxide (001) slab was constructed as a $\mathbf{6} \times \mathbf{6}$ in-plane supercell (basal lattice $a_0 \approx 3.126$ Å), with cell metrics $\mathbf{a} = \mathbf{b} = \mathbf{18.7560}$ Å, $\mathbf{c} = \mathbf{62.3777}$ Å, $\alpha = \beta = \mathbf{90^{\circ}}$, $\gamma = \mathbf{120^{\circ}}$ (P1). The model contains **three brucite-like** $[M(OH)_{-}6]$ octahedral sheets stacked along c (36 cation sites per sheet), with a Ni:Fe = 81:27 distribution (25% Fe on the cation sublattice; 9 Fe per sheet), and overall stoichiometry Ni_81Fe₂₇O₂₁₆H₂₁₆ (540 atoms). Both basal surfaces are symmetrically hydroxyl-terminated (H atoms constitute the outermost layers), yielding a stoichiometric, non-polar slab of thickness $\approx \mathbf{19.67}$ Å, separated by $\approx \mathbf{42.71}$ Å of vacuum within the c dimension; the inter-sheet metal-plane separations are 7.87 and 7.81 Å, consistent with LDH-like basal spacing. This geometry provides equivalent top/bottom (001) terminations exposing M–OH motifs suitable for benchmarking OER intermediates on LDH surfaces.

A2 — phosphate-intercalated Ni–Fe LDH (001)

Kept identical to **A1** in lattice and slab construction (6 × 6 hexagonal basal supercell; a = b = 18.756 Å, $\gamma = 120^{\circ}$, c = 62.3777 Å; P1) with symmetric hydroxyl-terminated basal surfaces, **A2 is derived from A1 by inserting protonated phosphate molecules (H_3PO_4)** into the interlayer galleries—**18 PO_4 units per supercell (9 per gallery)**, giving an overall composition **Ni_81Fe_{27}O_{288}H_{270}P_{18}** (684 atoms, Ni:Fe = 3:1). The intercalation shifts the metal-plane separations from \sim 7.87/7.81 Å (A1) to **8.89/8.90** Å, increasing the slab thickness to **22.22** Å and leaving \sim **40.16** Å of vacuum along c; the phosphorus centroids sit near the mid-planes of each gallery.

A3 — Ni single atoms on Sb-doped SnO_2(110)

A rutile $SnO_2(110)$ slab was built as a 4×8 in-plane supercell along $[1\bar{1}0]\times$ (cell a=26.8022 Å, b=25.4960 Å, c=57.0945 Å; $\alpha=\beta=\gamma=90^\circ;$ P1), comprising four O-Sn-O trilayers (slab thickness ≈ 12.72 Å) and ≈ 44.37 Å of vacuum. Eight Sb atoms substitute Sn in the subsurface cation plane (second layer from the top)—12.5 % within that plane, 3.2 % overall (Sn_244Sb₈). The top (110) surface carries four isolated Ni adatoms arranged in a 2×2 array (coverage ≈ 0.585 Ni nm⁻²), each adopting a square-pyramidal NiO_5 environment with four equatorial Ni-O = 1.922 Å and one apical Ni-O = 2.033 Å; the bottom surface is left clean and O-terminated. The full supercell stoichiometry is Sn_244Sb₈Ni₄O₅₀₈ (764 atoms).

A4 — Ce-doped, oxygen-deficient spinel NiCo_2O_4 slab (asymmetric OH termination)

Modeled in P1 with an oblique in-plane cell (a = 19.807 Å, b = 30.866 Å, γ = 108.715°; c = 41.150 Å), giving an in-plane area of 579.05 Ų (5.79 nm²); the atomic slab spans 14.87 Å along c and is separated by 26.28 Å of vacuum. The supercell composition is $Co_188Ce_8Ni_{92}O_{376}H_{96}$ (760 atoms), i.e., Ni: $Co \approx 1:2.04$ with Ce at 2.8 at.% of the cation sublattice. Relative to the ideal AB_2O_4 stoichiometry for 288 cations, the model contains 8 lattice-oxygen vacancies (x \approx 0.021 in Ni $Co_2O_4 - x$), corresponding—if evenly split between the two faces—to \approx 0.69 vac nm $^{-2}$ per surface. Ce dopants occupy two equivalent subsurface cation planes (four Ce per plane), located \sim 4.31 Å below the top and \sim 5.32 Å above the bottom surface, preserving an overall centrosymmetric placement. The top surface is Co-rich and partially hydroxylated (within the outermost 2 Å: 8 Co, 20 Co, 28 Co + Co

B1 — graphene-capped Ni–Fe LDH (001)

Retaining the A1 lattice and slab construction (6 × 6 hexagonal basal supercell; a = b = 18.756 Å, $\gamma = 120^{\circ}$, c = 70.6707 Å; P1), this model is derived from A1 by placing a **single-layer graphene overlayer** (128 C) above the top basal surface while keeping the LDH stoichiometry Ni_81Fe_{27}O_{216}H_{216} unchanged. The LDH comprises three brucite-like [M(OH)_6] sheets (Ni:Fe = 3:1) with metal-plane separations of 7.87 and 7.81 Å; both basal faces are hydroxylterminated, but only the **top face** is covered by graphene. The LDH portion spans 19.67 Å along c; including graphene, the occupied thickness is 24.67 Å, leaving 46.00 Å of vacuum. The graphene sheet is essentially planar (out-of-plane corrugation \leq 0.001 Å) and sits 5.00 Å above the outermost hydroxyl-H layer (5.97 Å above the topmost O layer), consistent with a physisorbed, non-perturbative interface that preserves the LDH surface structure; the bottom layers are identical to A1 and remain exposed/hydroxylated for benchmarking.

B2 — graphene-capped hematite (α -Fe_2O_3) (0001) with an aperture

A hexagonal α -Fe₂O₃(0001) slab was modeled in P1 with $\bf a=b=20.152~ \mathring{\bf A}, \gamma=120^\circ; c=60.339~ \mathring{\bf A}$ (in-plane area 351.70 $\mathring{\bf A}^2=3.517~ nm^2$), giving an atomistic thickness of 32.20 $\mathring{\bf A}$ separated by 28.14 $\mathring{\bf A}$ of vacuum. The slab is stoichiometric Fe₂O₃ in the interior (Fe₃84O₅₇₆) with both faces hydroxyl-terminated (H₉2 total) to render a non-polar model. A single sp²-carbon overlayer (C₁₁₅) caps only the top face and includes a circular \sim 6 $\mathring{\bf A}$ "window" (aperture) that leaves a patch of exposed surface Fe sites beneath the opening; elsewhere the carbon sheet remains graphene-like, with its centroid \sim 2.84 $\mathring{\bf A}$ above the outermost O plane (\approx 2.04 $\mathring{\bf A}$ above terminal H) and a small out-of-plane corrugation (Δ z \approx 0.91 $\mathring{\bf A}$). The bottom face remains hydroxylated and exposed for benchmarking.

B3 — BaCo 0.6Fe_{0.4}O₃ perovskite (001), BO₂||BaO asymmetric slab

In line with Branch B's chloride-resilience goal, a stoichiometric Ba(B)O₃ slab with mixed B-site cations (Ba₁₈₀Co₁₀₈Fe₇₂O₅₄₀; 900 atoms) was modeled in P1 with $\bf a=b=23.820~ \mathring{\bf A}, c=49.850~ \mathring{\bf A}, \alpha=\beta=\gamma=90^\circ$ (in-plane area 567.39 $\mathring{\bf A}^2=5.674~{\rm nm}^2$). The atomic slab spans 19.31 $\mathring{\bf A}$ along c and is separated by 30.54 $\mathring{\bf A}$ of vacuum. Terminations are asymmetric: the top face is a BO_2 layer that is Co-rich (Co:Fe = 27:9 \approx 3:1), while the bottom face is BaO-terminated; within $\pm 2~\mathring{\bf A}$ of the surfaces this corresponds to \sim 6.35 B nm⁻² (with \sim 12.69 O nm⁻²) at the top and \sim 6.35 Ba nm⁻² at the bottom. This construction furnishes a corrosion-tolerant, bare BO_2 active surface over a BaO support, consistent with the initial design brief to benchmark robust oxide surfaces against chloride attack.

B4 — coral-mimetic Ni–Fe LDH (001)

Keeping the MXene/LDH base identical to A1/B1 (6 × 6 hex basal supercell; a = b = 18.756 Å, $\gamma = 120^{\circ}$, c = 62.3777 Å; P1), the "coral" morphology was introduced by excising a single cylindrical pore (~8 Å diameter) through the **top two** hydroxide layers and hydroxylpassivating all newly exposed edges, as specified in the design note; both basal faces remain hydroxyl-terminated, with the **bottom face unchanged** from A1/B1. The resulting model is Ni_69Fe₂₅O₁₉₂H₂₁₆ (502 atoms), occupying 26.58 Å along c with 35.79 Å vacuum (in-plane area 304.66 Å²), yielding a non-polar, porous benchmark that preserves the original basal registry while exposing M–OH rim sites within the aperture.

Metal Site Choices Considering the complexity of each slab created based on DR report, we decide to select from each slab only representative sites for further study:

IrO_2

On $IrO_2(110)$ we target the canonical fivefold Ir terrace site (Ir_{5c}) on the top surface cation row. Its local geometry is square-pyramidal, with four coplanar surface O neighbors at 2.00 Å and a single subsurface apical O at 1.97 Å; the missing apical ligand points toward vacuum. Situated in the outermost O–Ir–O trilayer, this terrace Ir is the standard benchmark adsorption site on rutile-type oxides and is chosen to probe intrinsic, defect-free reactivity under well-defined coordination.

$\mathbf{A}_{\mathbf{1}}$

On A1 we probe two "regular" terrace cations on the **top basal (001) surface**—one Fe and one Ni located in the central, step-free region. Both exhibit near-ideal octahedral $M(OH)_6$ coordination with a 3-up/3-down O arrangement; their mean M–O distances are **2.09** Å (Fe) and **2.07** Å (Ni), and the in-plane metal–metal spacing reflects the basal lattice at \sim **3.13** Å. The Fe site sits in a Ni-rich cation shell (isolated Fe amid Ni), whereas the Ni site has a typical mixed neighborhood (\sim 4 Ni, 2 Fe), consistent with the \sim 3:1 Ni:Fe layer composition. These terrace positions were chosen as **representative**, **bulk-like references**—well coordinated and far from edges or defects—providing robust baselines for adsorption benchmarking on LDH surfaces.

A 2

In A2 we focus on two regular first-layer cations on the **top basal (001) surface**—one Fe and one Ni—above the phosphate-filled interlayer. Both sites are octahedral with 3 O above / 3 O below; the Fe exhibits Fe–O \approx 2.06/2.01 Å (above/below) and the Ni Ni–O \approx 2.10/2.07 Å. Relative to the subsurface phosphate, the Fe sits nearly on-axis over a PO₄ unit ($\Delta z \approx$ 4.30 Å, lateral offset \approx 0.71 Å), whereas the Ni is laterally off-axis from the nearest PO₄ ($\Delta z \approx$ 4.54 Å, lateral \approx 2.95 Å). These two surface metals are chosen as representative, symmetry-inequivalent terrace sites that sample distinct interlayer electrostatic environments (aligned vs off-axis to PO₄) while avoiding edges/defects—providing a solid, structure-based rationale for benchmarking metal-identity and phosphate-proximity effects.

A3

For A3, on the top SnO₂(110) terrace, the representative site comprises a single Ni atom anchored above a bridging-O row, adopting a square-pyramidal NiO₅ geometry with four surface O_{br} ligands in the equatorial plane (Ni–O ≈ 1.92 Å) and one apical O beneath (Ni–O ≈ 2.03 Å); the Ni sits $\sim\!0.19$ Å below the O_{br} plane and $\sim\!2.03$ Å above the apical O. The two nearest subsurface Sn cations on adjacent cation rows are each connected to the adatom via a Sn–O_{br}–Ni bridge (Sn–O ≈ 2.07 Å; Ni–Sn ≈ 3.12 Å). Situated on the upper, flat terrace away from step edges, this motif matches your original "Ni single-atom adjacent to Sn via bridging O" note and is selected as a representative active site because it captures the canonical O_{br}–M_{cus} environment of rutile(110) while providing quantitatively well-defined, low-coverage coordination ideal for benchmarking OER intermediates.

A4

For A4, we have choosed a total of 6 Co/Ni atoms. On the top Co–O(H)-terminated face of the Ce-doped, oxygen-deficient NiCo_2O_4 slab, we consider two Ce-proximal motifs and one Ce-remote control set. Ce-proximal: (i) a two-fold-coordinated Co sitting at the rim of a surface O-vacancy directly above the subsurface Ce layer, and (ii) a three-fold Co one O row back on the same terrace; alongside these, a three-fold Ni in the same Co-rich layer but still within the Ce-perturbed field captures the B-site identity contrast. Ce-remote (terrace center, far from dopants/vacancies): (iii) a two-/three-fold Co pair representative of the intact hydroxylated terrace, and (iv) a two-fold Ni surrounded by a largely unreconstructed O sublattice. Together these five local environments span coordination (2 vs. 3), cation identity (Co vs. Ni), and dopant/vacancy proximity (near-Ce vs. terrace-center) consistent with the A4 design (Ce in subsurface planes + engineered surface O-vacancies with hydroxyl passivation). This refelects the selection logic: Ce-adjacent rims/near-rim terrace sites plus Ce-remote terrace controls with 2- and 3-coordinate Co/Ni.

B1

On B1 (graphene-capped NiFe-LDH (001)), the representative first-layer surface Fe and Ni are the outermost cations in the top brucite-like sheet directly beneath the graphene overlayer. Each sits in a regular M(OH)_6 octahedron with nearly isotropic M–O bonds (Fe: 2.09 Å apical/2.09 Å apical; equatorial 2.09 \pm 0.00 Å; Ni: 2.08 Å apical/2.08 Å apical; equatorial 2.07 \pm 0.00 Å; octahedral distortion \approx 0.11 %), positioned 1.03 Å below the topmost O plane and 2.00 Å below terminal OH, with the graphene sheet 7.00 Å above the metal plane (5.97 Å above the surface O layer). The chosen Fe lies near the in-plane center (> 8.09 Å from periodic edges) and the Ni is also well inside the terrace (\approx 5.45 Å from edges), i.e., regular terrace sites rather than step/edge motifs—matching the "regular first-layer" description in the earlier manual note. These sites were selected to provide clean, defensible benchmarks for OER under Branch-B1's design goal (probing LDH activity beneath a protective graphene barrier): they are far from defects and intercalants, exhibit bulk-like octahedral geometry, and experience only physisorbed graphene influence, thus isolating intrinsic Fe/Ni–OH reactivity at the top basal surface.

B2

on B2, On the top (0001) face beneath the graphene aperture, four representative surface Fe sites were selected: two **rim** Fe located along the edge of the \sim 6 Å window where the sp² rim approaches the oxide, each making Fe–C contact (shortest Fe–C \approx 2.17 Å), with undercoordination (Fe–O CN \approx 3) and a slight outward relaxation (\approx 1.33 Å above the terminal O plane); and two **center-of-window** Fe positioned near the middle of the opening, far from the carbon edge (nearest C \geq 5.24 Å), retaining octahedral-like coordination (Fe–O CN \approx 6) and sitting \approx 1.19 Å below the top O/H plane. These two site families—carbon-contacted rim vs carbon-free aperture interior—follow the original design intent and provide a physically motivated contrast to quantify how the graphene shell perturbs intrinsic hematite reactivity.

B3

On the BO_2 -terminated top surface of B3 we sample four terrace cations: (i) one Fe and (ii) one Co sitting at the rim of a deliberate surface O-vacancy, each three-coordinated to lattice O (vacancy-edge sites, no apical O) with shortened in-plane M–O distances of \sim 1.90–1.97 Å; and, as references, (iii) one Fe and (iv) one Co on the same terrace but far from any defect (\geq 10.00 Å lateral separation), retaining the regular square-pyramidal BO_2 5 environment (four equatorial M–O \approx 1.95–2.02 Å, one apical \approx 2.05–2.10 Å). This choice—consistent with the original notes—targets the most contrastive local motifs (vacancy-edge undercoordination vs pristine terrace; Fe vs Co identity) while keeping all sites in the outermost BO_2 plane above the BaO underlayer.

B4

In B4's coral-type aperture, we target four representative rim sites: a Fe and a Ni in the **top metal plane** at the pore edge, and the corresponding **second-plane** Fe and Ni located directly beneath the rim (as in the earlier manual note). The two top-layer cations sit **1.88** Å (Fe) and **1.55** Å (Ni) below the outermost O plane and lie within **0.10–0.23** Å of the top-plane centroid; their in-plane separation is **3.21** Å. At the rim, Fe remains near-octahedral (6 O \leq **2.30** Å; $\langle d_{\text{M-O}} \rangle =$ **2.13** Å), whereas Ni is under-coordinated (4 O \leq **2.30** Å; $\langle d_{\text{M-O}} \rangle =$ **2.12** Å). The second-layer pair occupies the ring immediately below the aperture edge, **13.53–13.43** Å beneath the top O plane and within **0.09–0.02** Å of the second-plane centroid; both are under-coordinated (4 O \leq **2.30** Å; $\langle d_{\text{M-O}} \rangle =$ **2.12** Å (Ni), **2.06** Å (Fe)), with an in-plane separation of **3.21** Å. These sites were chosen to systematically span cation identity (Fe vs Ni), depth (surface vs subsurface), and coordination state at the pore edge, i.e., the structural factors most likely to govern *OH/*O/*OOH binding and chloride tolerance in the coral design.

OER Catalyst Stability Assessment

Discussion of OER stability. While recent advances in transition metal layered double hydroxides (LDH) and related materials have achieved excellent oxygen evolution reaction (OER) activity approaching theoretical limits, catalyst stability remains the critical bottleneck preventing commercial deployment of water splitting technologies. Operating potentials exceeding 1.5 V vs. RHE in alkaline conditions create thermodynamically favorable pathways for both metal dissolution and lattice oxygen loss, leading to rapid catalyst degradation. To address this fundamental challenge, we evaluated candidate catalysts proposed by the DR agents (Figure DL_4 1) using two complementary stability descriptors following established protocols from our previous work [51]: dissolution potential ($U_{\rm diss}$) and oxygen vacancy formation energy ($V_{\rm O}$).

Stability Descriptors and Physical Interpretation

The dissolution potential U_{diss} quantifies the thermodynamic driving force for metal atom leaching from the catalyst surface, calculated as:

$$U_{\rm diss}(M@{\rm slab}) = U_{\rm diss,bulk}^{\circ}(M^{z+}/M) + \frac{E_{\rm slab-vac}(M) + \mu_{\rm bulk}(M) - E_{\rm slab}}{ne}$$
 (6)

where $E_{\rm slab}$ and $E_{\rm slab-vac}(M)$ represent the DFT energies of the intact and metal-vacancy surfaces, respectively (Figure DL_4 2), $\mu_{\rm bulk}(M)$ is the bulk metal chemical potential, $U_{\rm diss,bulk}^{\circ}$ is the standard reduction potential, and n is the number of electrons transferred. Higher $U_{\rm diss}$ values indicate greater resistance to metal dissolution, with positive values suggesting thermodynamic stability under standard conditions.

The oxygen vacancy formation energy V_O captures the structural integrity of the oxide lattice during OER cycling. We following the following protocol for picking lattice oxygen atom to remove, oxygen atoms are identified within 3.0 Å of active metal sites and removed based on coordination number, with under-coordinated oxygen atoms preferentially selected (Figure DL_4 2). Higher V_O values indicate stronger metal—oxygen bonds and reduced susceptibility to lattice oxygen participation in the OER mechanism, which can lead to irreversible structural degradation.

Comparative Stability Analysis

Evaluation across representative surface sites (Figure DL_4 3) reveals significant variation in stability metrics among candidate catalysts. The dissolution potentials span from approximately $-0.4~\rm V$ to $+1.2~\rm V$ vs. SHE (Figure DL_4 4). Oxygen vacancy formation energies similarly exhibit broad distribution. In general, higher Vo suggesting enhanced lattice stability potentially arising from stronger metal—oxygen covalency or favorable coordination environments that stabilize the oxide framework.

Correlation Between Stability Metrics

A positive correlation emerges between $U_{\rm diss}$ and $V_{\rm O}$ across all 27 evaluated metal sites (Figure DL_4 6), indicating that resistance to metal dissolution and oxygen vacancy formation are intrinsically linked. This correlation suggests a unified electronic origin for both stability mechanisms: catalysts with stronger metal—oxygen bonds resist both metal leaching and oxygen loss. The observed relationship validates the use of either metric as a reasonable proxy for overall stability, though evaluating both provides more comprehensive assessment.

Normalized Stability Scoring

To enable quantitative comparison, we similarly applied the linear scoring framework with IrO₂ as the baseline reference (Score = 100), using sensitivity parameter $\beta = 0.2$ consistent with other applications. The U_{diss} -based scores (Figure DL_4 7). The V_{O} -based scores (Figure DL_4 8) show similar trends but not quite exactly the same rankings, highlighting the value of assessing both metrics.

The divergence between U_{diss} and V_O rankings for specific candidates reflects different stability-limiting mechanisms: some materials may excel at preventing metal dissolution while remaining vulnerable to oxygen loss, or vice versa. This complementary information guides rational catalyst design—for instance, candidates with high U_{diss} but moderate V_O might benefit from strategies to strengthen metal—oxygen bonds, while those with the opposite profile could require enhanced metal site anchoring.

OER Catalyst Structures

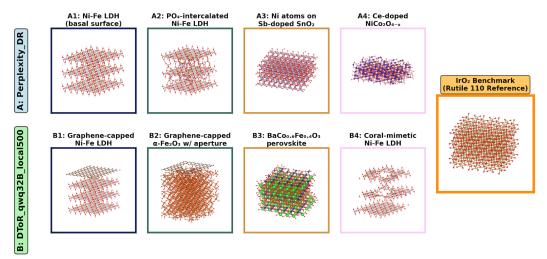


Figure DL_4 1: Structural configurations of OER catalyst candidates proposed by DR agents, with IrO_2 (110) included as the domain benchmark reference for stability assessment.

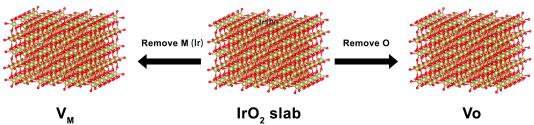


Figure DL_4 2: Schematic illustration of vacancy creation on the IrO $_2$ rutile (110) surface. Left: Metal vacancy (V $_M$) created by removing a surface metal atom for dissolution potential (U $_{diss}$) calculation. Right: Oxygen vacancy (V $_0$) created by removing a surface oxygen atom for vacancy formation energy calculation.

OER Catalyst Site Choices

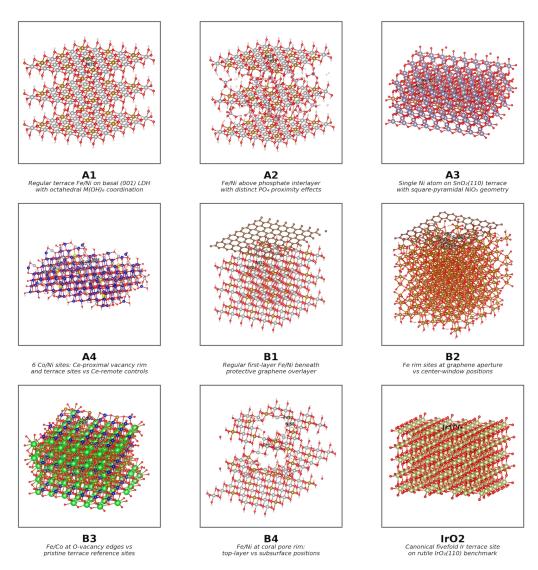


Figure DL_4 3: Matrix representation of selected metal sites for stability assessment across Aseries (A1–A4) and B-series (B1–B4) catalyst candidates. Representative surface metal atoms are highlighted for dissolution and vacancy formation calculations.

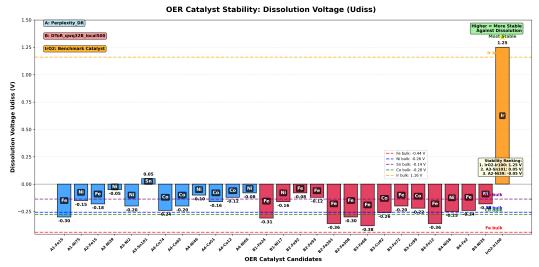


Figure DL_4 4: Dissolution potentials (U_{diss}) for all candidate catalysts at representative surface sites. Higher values indicate greater resistance to metal dissolution. Dashed horizontal lines denote standard bulk dissolution potentials for reference metals (Fe: -0.44 V, Co: -0.277 V, Ni: -0.257 V, Ir: +1.16 V vs. SHE).

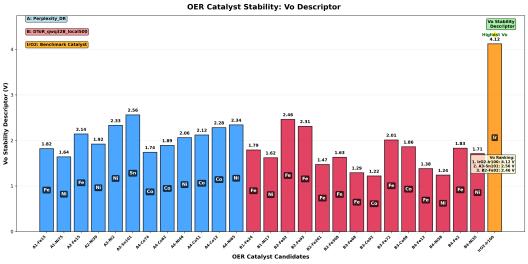


Figure DL_4 5: Oxygen vacancy formation energies $(V_{\rm O})$ for all candidate catalysts at representative surface sites. Higher formation energies indicate greater structural stability and resistance to oxygen loss during OER operation.

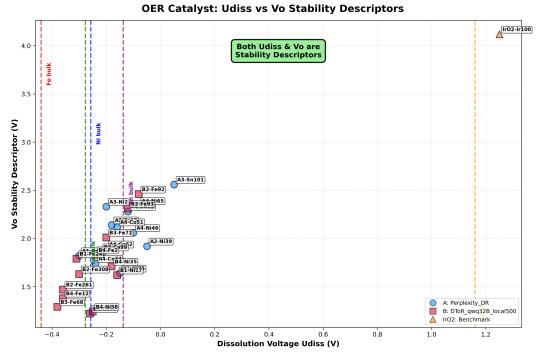


Figure DL_4 6: Correlation between dissolution potential (U_{diss}) and oxygen vacancy formation energy (V_O) across all 26 evaluated metal sites. The positive correlation indicates that catalysts resistant to metal dissolution also tend to resist oxygen vacancy formation.

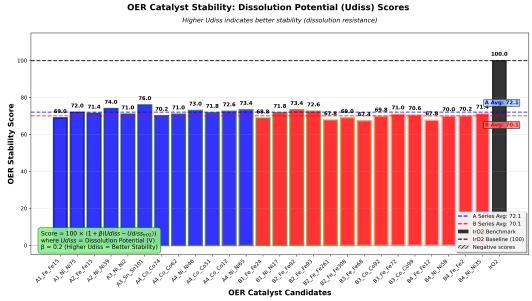


Figure DL_4 7: Stability scores based on dissolution potential (U_{diss}), normalized to IrO₂ baseline (100). Scores are calculated using the linear metric with sensitivity parameter $\beta=0.2$, where higher U_{diss} values yield higher scores.

OER Catalyst Stability: Oxygen Vacancy Formation Energy (Vo) Scores

Lower Vo indicates better OER activity (easier oxygen vacancy formation)

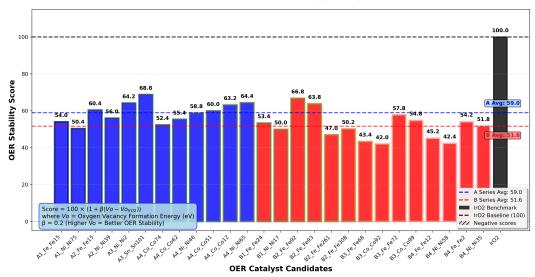


Figure DL_4 8: Stability scores based on oxygen vacancy formation energy (V_O), normalized to IrO₂ baseline (100). The scoring methodology employs the same linear scaling ($\beta=0.2$), with higher V_O values indicating superior stability.

E.6 CO₂ Sensor Probe

For the fifth application: CO₂ Sensor Probe. The query as shown in D.1 is:

CO2 Sensor Probe

Which two-dimensional materials-such as graphene derivatives, transition metal dichalcogenides, or MXenes-offer the highest CO_2 sensing performance in complex gas or aqueous environments, and how do they compare in terms of detection limit (ppm), selectivity against common interferents (e.g., O_2 , H_2O), response/recovery time, and long-term stability, including any functionalization or structural modifications that enhance these metrics? Think of novel candidates.

The best commercial solution is: **GlideLogic_DR**.

The best local solution is: **DToR_qwq32B_local500**.

The expert based on the two reports, provided 4 commercial DR report derived candidates (A1-A4), and 4 local DR report derived candidates (B1-B4), visualized in DL_5 1:

g-C₃N₄ description

For CO_2 binding, the benchmark substrate is a monolayer graphitic carbon nitride (g-C3N4) built as a 7×7 in-plane supercell of the primitive sheet (49 formula units; C147N196) in a hexagonal cell with a=b=33.50 Å, $\gamma=120^\circ$, and c=35.0 Å. All atoms lie in a single atomic plane at $z\approx5.0$ Å, leaving ~30 Å of vacuum to the next periodic image along c, i.e., no interlayer interactions. The in-plane network is a fully conjugated sp² C–N lattice: each C is three-coordinate, N sites are predominantly two- to three-coordinate ($\langle CN \rangle \approx 2.25$), and nearest-neighbor C–N distances cluster at 1.37 ± 0.06 Å. No heteroatoms, defects, terminations, or adsorbates are present (P1 by construction). This clean, stoichiometric, and strictly planar C_3N_4 slab (areal density ≈5.05 formula units nm⁻²) serves as our reference surface for CO_2 binding.

A1 (rGO-PEI)

Hexagonal P1 slab with a = b = 36.90 Å, γ = 120.00°, c = 45.00 Å (α = β = 90.00°). The model contains 605 atoms (C461H80O54N10). Atomic z positions span 4.63–16.63 Å, giving a slab thickness of 12.01 Å and 32.99 Å of vacuum along c. The in-plane area is 11.79 nm², yielding areal number densities of C 39.09 nm⁻², H 6.78 nm⁻², O 4.58 nm⁻², and N 0.85 nm⁻². Carbon heights resolve a dominant rGO basal plane centered at z \approx 7.39 \pm 0.64 Å (441 C) and a minor upper carbon domain at 12.09–15.35 Å (20 C). All N atoms reside in the upper domain (11.23–16.53 Å; $\langle z \rangle$ = 13.03 Å), consistent with an amine-rich polyethylenimine (PEI) branch moiety tethered above the rGO sheet; O atoms are concentrated below 10.65 Å, consistent with residual rGO oxygen functionalities. This rGO–PEI slab, isolated by ample vacuum, provides a structurally well-defined platform for site-specific CO₂ binding calculations.

A2 (phosphorene $|A|_2O_3$ cap with drilled through-pore)

Orthorhombic P1 cell (a = 26.36 Å, b = 27.26 Å, c = 27.08 Å; $\alpha = \beta = \gamma = 90.00^{\circ}$) comprising a phosphorene sublayer (P192) beneath an alumina cap (Al146O₂₂₀). Atomic z spans 3.95–17.44 Å, giving a slab thickness of 13.49 Å and 13.59 Å of vacuum. Phosphorene occupies z = 3.95–8.12 Å (thickness 4.17 Å; top P plane at 8.12 Å), while the cap exhibits O/Al = 1.51 (\approx Al₂O₃) with the Al framework at z = 8.83–15.83 Å (7.00 Å) and the O network extending 8.05–17.44 Å (9.39 Å), slightly overlapping the top P plane (Δ z = 0.07 Å) to form an interfacial contact. The in-plane area is 718.57 Å² (7.19 nm²). A single through-pore pierces the cap near (f_x, f_y) = (0.51, 0.50): the top aperture has radius 6.27 Å (\emptyset 12.54 Å), narrows to a bottleneck radius 5.11 Å (\emptyset 10.22 Å; section 82.02 Å², 11.42

A3 (MoS₂ monolayer with vacancies)

P1 cell with a = 27.58 Å, b = 22.29 Å, c = 28.13 Å ($\alpha = \beta = \gamma = 90.00^{\circ}$). The slab contains 200 atoms (Mo64S136; S/Mo = 2.13) and occupies z = 4.83–8.50 Å (thickness 3.67 Å) with 24.46 Å of vacuum along c. The in-plane area is 6.15 nm², giving areal number densities of Mo 10.41 nm² and S 22.13 nm². Layering is a single S–Mo–S tri-layer: the Mo plane is centered at z = 6.56 Å (spread 0.31 Å), bracketed by S planes at z = 5.01 Å and 8.12 Å, yielding an average S–S separation of 3.11 Å and Mo–S vertical offsets of 1.54 Å (bottom) and 1.57 Å (top). No heteroatoms, functional groups, or adsorbates are present. Five Mo atoms are removed to create Mo vacancies. This monolayer MoS² slab provides a simple, well-isolated reference for CO² adsorption benchmarks.

A4 ($Ti_3C_2T_x$, mixed terminations)

Hexagonal P1 slab with a = b = 30.67 Å, $\gamma = 120.00^{\circ}$, c = 51.70 Å ($\alpha = \beta = 90.00^{\circ}$); in-plane area 8.15 nm^2 . The model contains $Ti300C200O175F25H_{25}$, i.e., 100 Ti_3C_2 formula units with 200 surface sites. Terminations are asymmetric: the top face comprises 50.00

B1 (rGO-NH₂; same rGO basal slab as A1)

This model reuses the rGO lattice from A1 without change (a = b = 36.90 Å, γ = 120.00°, c = 45.00 Å; in-plane area 11.79 nm²), retaining the 441-C rippled basal plane centered at z = 7.29 \pm 0.67 Å. The slab contains 585 atoms (C441H70O54N20) with heights spanning 4.36–11.16 Å (thickness 6.80 Å) and 38.20 Å of vacuum. Oxygen functionalities lie over 4.36–10.64 Å, while primary amines are grafted exclusively on the upper face (N 8.37–10.38 Å, $\langle z \rangle$ = 9.40 \pm 0.63 Å; H up to 11.16 Å), giving an –NH₂ areal density of 1.70 nm⁻². Areal number densities are C 37.40 nm⁻², O 4.58 nm⁻², N 1.70 nm⁻², and H 5.94 nm⁻². With lattice vectors and basal carbon topology identical to A1, this slab isolates the effect of replacing the PEI overlayer with discrete –NH₂ terminations for CO₂ binding studies.

B2 (Ti₃C₂F₂, fully F-terminated; same construction protocol as A4)

Hexagonal P1 slab with a = b = 30.67 Å, γ = 120.00°, c = 51.70 Å (α = β = 90.00°); in-plane area 8.15 nm². The CIF contains Ti300C200F₂₀₀ (100 Ti₃C₂ formula units; 200 surface sites), i.e., both faces are 100.00

B3 (graphene | Ti3C2O2 heterojunction; same MXene cell as As A4/B2)

Hexagonal P1 slab with a=b=30.67 Å, $\gamma=120.00^\circ$, c=51.70 Å (in-plane area 8.15 nm²). The MXene sub-slab matches the A4/B2 Ti3C2 lattice and is fully *O-terminated on both faces (Ti300C200O200; 100 *O per face, 12.28 sites nm²). A single-layer graphene sheet (C338) is appended on the bottom side, giving a total composition Ti300C538O200. Atomic heights span z=7.31-19.41 Å (thickness 12.10 Å) with 39.60 Å of vacuum along c. The graphene forms a narrow plane at $z=7.79\pm0.31$ Å (7.31-8.45 Å), the bottom *O plane is at z=12.40 Å and the top *O plane at z=19.41 Å; the graphene–MXene van der Waals gap is 4.61 Å (center-to-plane; minimum 3.95 Å to *O), and the graphene–bottom Ti plane separation is 6.78 Å. The result is a well-isolated graphene|Ti3C2O2 heterojunction with symmetric *O terminations and a dense graphene areal carbon density of 41.49 nm².

B4 (Cu-inserted MoS₂; same MoS₂ host as A3)

Orthorhombic P1 cell with a = 27.58 Å, b = 22.29 Å, c = 28.13 Å ($\alpha = \beta = \gamma = 90.00^\circ$; in-plane area 6.15 nm²). The slab contains Mo₆₄Cu5S136 and occupies z = 4.80–8.32 Å (thickness 3.51 Å) with 24.61 Å of vacuum. The host S–Mo–S tri-layer is preserved ($\langle z \rangle$: S(bottom) = 5.03 Å, Mo = 6.59 Å, S(top) = 8.15 Å; S–S spacing 3.12 Å; Mo–S offsets 1.55/1.56 Å), as in A3. At A3's Mo vacancies, five Cu centers lie within the chalcogen sandwich at z = 6.55–6.89 Å ($\langle z \rangle$ = 6.70 Å), i.e., \sim 0.11 Å above the Mo plane and biased toward the top S layer; the areal Cu density is 0.81 nm⁻². Local coordination around Cu is consistent with insertion at chalcogen-vacancy–like sites: each Cu shows \sim 4–5 short Cu–S contacts at 2.28–2.38 Å (mean of first shell \approx 2.31 Å) with the next S neighbors near 2.95–3.06 Å, and \sim 2–3 Cu–Mo contacts at 2.78–3.26 Å. This configuration provides a well-defined model to probe how subsurface/embedded Cu at vacancy motifs modulates CO₂ binding relative to the pristine A3 monolayer.

CO₂ Sensor Performance Evaluation

The development of sensitive and selective CO_2 sensors addresses critical needs in environmental monitoring, industrial process control, and indoor air quality management. To evaluate the CO_2 sensing capabilities of candidates proposed by the DR agents (Figure DL_5 1), we quantified binding energies using a systematic grid-based sampling approach that captures the heterogeneity of surface adsorption sites, with g- C_3N_4 serving as the established benchmark reference [52].

The comprehensive sampling methodology employed a 5×5 grid across each surface unit cell, positioning CO_2 molecules at 25 distinct initial configurations to explore the complete adsorption landscape (Figure DL_5 2). At each grid point, CO_2 molecules were placed 1.8 Å above the highest surface atom within a 3.0 Å search radius, with initial linear configurations parallel to the surface. This systematic approach ensures unbiased sampling of potential binding sites, including atop atoms, bridge sites, hollow sites, and defect locations, providing statistically robust binding energy distributions rather than single-point estimates.

The resulting binding energy distributions reveal significant variability both within and between candidate systems (Figure DL_5 3). The violin plots demonstrate that surface heterogeneity plays a crucial role in CO_2 adsorption, with some candidates exhibiting narrow distributions indicating uniform binding sites, while others show broad distributions suggesting diverse adsorption environments. Mean binding energies range from approximately 0.15 to 0.90 eV across the evaluated systems, with A-series candidates generally demonstrating stronger CO_2 interactions compared to B-series materials.

Performance scoring using the linear transformation framework with sensitivity parameter $\beta=0.2$ (Figure DL_5 4) quantifies the relative sensing capabilities normalized to the g-C₃N₄ baseline. The distribution of scores reveals that commercial DR agent successfully identified materials spanning a range of binding strengths, from weak physisorption suitable for rapid-response sensors to stronger chemisorption appropriate for CO₂ capture applications.

The systematic evaluation highlights the importance of comprehensive surface sampling in sensor material assessment. Candidates with both high mean binding energies and narrow distributions, such as B4, represent ideal sensor materials offering consistent response across the entire surface. Conversely, materials with broad binding energy distributions may exhibit non-uniform sensor responses, potentially complicating calibration and reducing reliability. These findings underscore that effective ${\rm CO}_2$ sensor design requires not only optimizing average binding strength but also minimizing surface heterogeneity to ensure reproducible sensing performance.

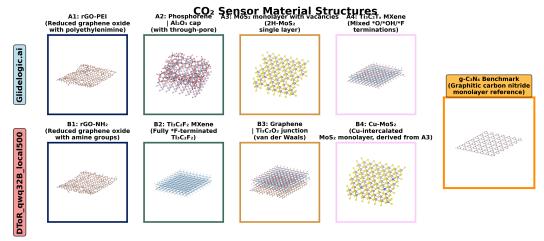


Figure DL_5 1: Structural configurations of CO_2 sensor candidates proposed by DR agents, with g-C₃N₄ [52] included as the domain benchmark reference.

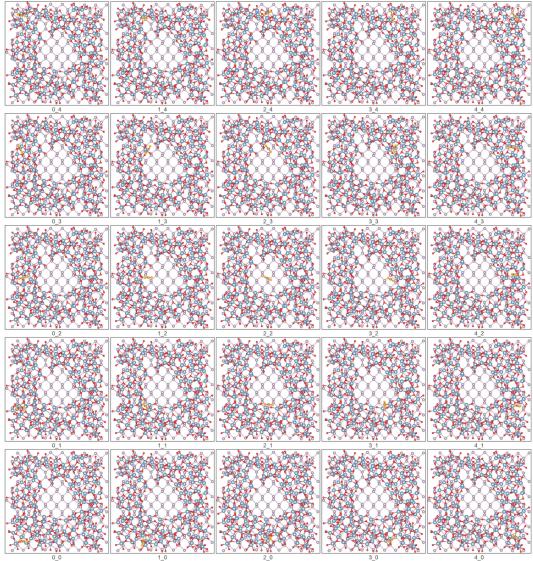


Figure DL_5 2: Top-view representation of the 5×5 grid sampling approach for CO₂ adsorption sites. Relaxed configurations show CO₂ molecules at 25 distinct initial positions on candidate A2, demonstrating comprehensive exploration of the binding landscape.

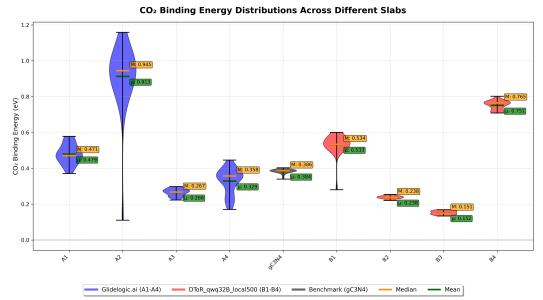


Figure DL_5 3: Distribution of CO_2 binding energies across all candidate systems visualized as violin plots. Mean (solid line) and median (dashed line) values are shown for each candidate, derived from 25 configurations per surface using the 5×5 grid sampling methodology.

CO2 Sensor: Binding Energy Scores

Higher binding energy indicates better CO2 sensing performance

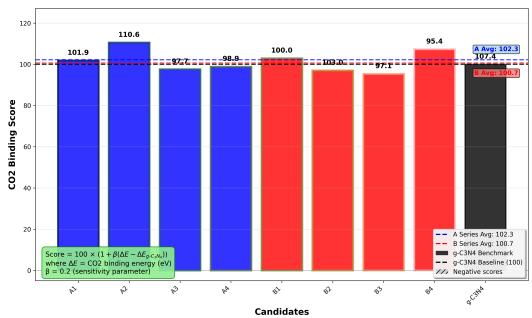


Figure DL_5 4: CO_2 sensing performance scores for all candidate systems from static DFT calculations, normalized to g- C_3N_4 baseline (100). Scores are calculated using the linear transformation with sensitivity parameter $\beta=0.2$, where higher binding energies yield higher scores.

F Power Consumption Estimation

Power estimation and LLM/RAG time decomposition. We convert runtime to energy using a calibrated two–level power model on our workstation (4×NVIDIA RTX A5500, Ubuntu). System power was monitored continuously during experiments using nvidia-smi for GPU power draw and powertop for system—level consumption. We acknowledge that our power measurement methodology is relatively basic; however, the core contribution of this work lies in demonstrating the feasibility of edge–device deployment, which already showcases significant democratization and decentralization advantages over commercial subscription—based solutions. The quantitative power estimates provided here serve primarily as reference points for relative comparisons rather than definitive consumption benchmarks. Across runs we *observed* a typical baseline draw during retrieval/IO of $P_{\rm base} \approx 76 \pm 3~\rm W$ (CPU, memory, disks, and idle GPUs). When the LLM generates tokens the GPU adds a model–dependent draw reported by nvidia-smi:

$$\Delta P_{\rm GPU}({\rm ds}14{\rm B}) = 282 \pm 4 \,{\rm W}, \quad \Delta P_{\rm GPU}({\rm ds}70{\rm B}) = 310 \pm 5 \,{\rm W}, \quad \Delta P_{\rm GPU}({\rm qwq}32{\rm B}) = 301 \pm 4 \,{\rm W}.$$

Hence $P_{\rm LLM}({\rm model}) = P_{\rm base} + \Delta P_{\rm GPU}({\rm model})$ is 358 W (ds14B), 386 W (ds70B), and 377 W (qwq32B). During RAG (vector search over the local store, parsing, and file IO) the GPU utilization remained near–idle, so we bill RAG at $P_{\rm RAG} = P_{\rm base} \approx 76~{\rm W.}^4$

Empirical split of runtime into LLM vs. RAG. Let $\bar{t}_{tot}(s, m, k)$ be the mean wall–clock time for strategy $s \in \{\text{Single}, \text{DToR}, \text{model } m \in \{\text{ds14B}, \text{qwq32B}, \text{ds70B}\}$, and local–RAG size $k \in \{0, 100, 500\}$. For k = 0 (no local vector store) we roughly estimated that all time is attributable to LLM generation:

$$t_{\rm LLM}(s,m,0) = \bar{t}_{\rm tot}(s,m,0), \qquad t_{\rm RAG}(s,m,0) = 0. \label{eq:tlm}$$

For k>0 we roughly split $\bar{t}_{tot}(s,m,k)$ into LLM and RAG by scaling the k=0 LLM time to reflect the additional *LLM* passes used in local–knowledge summarization and complementary–query generation. The scale factors were estimated from repeated single–node micro–benchmarks (profiling isolated summarization and query–generation calls with coarse-grained timing) and calibrated against observed aggregate runtime patterns across full runs. Given measurement noise and workload variability, we report representative values that *differed by strategy*:

$$\begin{split} & \text{Single:} \quad \alpha_{\text{ds}14\text{B}}^{\text{Single}} \approx 1.05, \quad \alpha_{\text{qwq}32\text{B}}^{\text{Single}} \approx 1.07, \quad \alpha_{\text{ds}70\text{B}}^{\text{Single}} \approx 1.09; \\ & \text{DToR:} \quad \alpha_{\text{ds}14\text{B}}^{\text{DToR}} \approx 1.10, \quad \alpha_{\text{qwq}32\text{B}}^{\text{DToR}} \approx 1.20, \quad \alpha_{\text{ds}70\text{B}}^{\text{DToR}} \approx 1.35. \end{split}$$

We then compute

$$t_{\text{LLM}}(s, m, k) = \alpha_m^s \bar{t}_{\text{tot}}(s, m, 0), \qquad t_{\text{RAG}}(s, m, k) = \bar{t}_{\text{tot}}(s, m, k) - t_{\text{LLM}}(s, m, k) \text{ (clipped at } \geq 0).$$

Two empirical patterns are worth noting. (i) In **Single** mode the increase in $t_{\rm LLM}$ from k=0 to k>0 is modest (ranging from +5% to +9% across models), and $t_{\rm RAG}$ exhibits a smaller difference than theoretically expected between k=100 and k=500 within measurement tolerance. (ii) In **DToR** mode, multi–hour runs accumulate more local–knowledge summarization and complementary–query steps; correspondingly $t_{\rm LLM}$ shows larger scaling (1.10/1.20/1.35 as above), and $t_{\rm RAG}$ increases with k due to heavier vector–store traffic and post–processing over longer context. When attribution is ambiguous we conservatively charge time to the LLM bucket.

Energy calculation. Per–report energy is the sum of LLM and RAG contributions:

$$E_{\rm LLM}(s,m,k) = \frac{P_{\rm LLM}(m)\,t_{\rm LLM}(s,m,k)}{1000}, \quad E_{\rm RAG}(s,m,k) = \frac{P_{\rm RAG}\,t_{\rm RAG}(s,m,k)}{1000}, \quad E_{\rm tot} = E_{\rm LLM} + E_{\rm RAG} \ (\rm kWh). \label{eq:ellm}$$

All results in Fig. 3d use $P_{\rm base}=76~{\rm W}, \Delta P_{\rm GPU}=\{282,301,310\}~{\rm W}$ (ds14B, qwq32B, ds70B), and the strategy–specific α values above. For example, substituting the measured $\{t_{\rm LLM}, t_{\rm RAG}\}$ of DToR_qwq32B_local500 yields $E_{\rm tot}=4.16~{\rm kWh}$ per completed report; configuration–wise LLM/RAG splits are provided alongside runtime heatmaps in the Appendix.

⁴Short web fetches last seconds relative to multi–hour runs and are neglected.