

# Audited Causal Discovery Agents for Brain Resilience and Alzheimer’s Reversal

David Scott Lewis<sup>1</sup> Enrique Zueco<sup>1</sup>

<sup>1</sup>AI4X Research, Zaragoza, Spain. Correspondence to: David Scott Lewis [reports@aiexecutiveconsulting.com](mailto:reports@aiexecutiveconsulting.com).

We propose Resilience-ADA, an audited discovery agent for Alzheimer’s reversal research. The agent maintains explicit causal world models, generates competing mechanistic hypotheses, and selects the next perturbations using Bayesian optimal experimental design under feasibility constraints. Motivated by recent evidence linking NAD<sup>+</sup> homeostasis and metabolic resilience to disease modification [1, 2, 3, 4, 5, 6, 7], Resilience-ADA outputs (i) uncertainty-calibrated causal graphs connecting metabolism, BBB integrity, inflammation, and synapses; (ii) a ranked short list of single and combination targets; and (iii) a falsifiable experiment ladder with biomarker panels and stop/rollback criteria. Deterministic reliability gates and replayable Decision Cards make each recommendation auditable end-to-end [8, 9, 10, 11, 12].

## 1. Introduction

Alzheimer’s disease (AD) is commonly framed as an irreversible neurodegeneration. That framing biases research programs toward incremental slowing via downstream correlates. A resilience-first view instead asks whether upstream control variables (e.g., mitochondrial stress response, NAD<sup>+</sup> homeostasis, neurovascular integrity) can be restored such that cognition and biological hallmarks recover, at least in tractable models [1, 2, 3, 4, 5, 6, 7].

AI has started to accelerate hypothesis generation in biomedicine, but most target discovery pipelines remain correlation-heavy. In AD, confounding, dataset shift (mouse-to-human, region-to-region), and publication bias can turn correlation rankings into expensive false positives. At the same time, agentic systems and lab automation are increasing experimental throughput, amplifying the need for governance and auditable decision making [13, 8, 9, 10, 11, 12, 14, 15, 16]. What is needed is a ‘virtual experimenter’ that (a) keeps uncertainty explicit, (b) proposes falsifiable mechanisms, and (c) chooses the next experiments that most reduce decision uncertainty.

**Contributions:** (1) a typed causal belief state (Causal Graph-of-Thoughts) with evidence-linked uncertainty; (2) goal-conditioned BOED for selecting the next perturbations; (3) reliability gates and Decision Cards for auditable, self-falsifying recommendations.

## 2. Audited discovery agents for reversal-oriented causal targeting

We cast discovery as goal-conditioned causal targeting. Given evidence  $\mathcal{D}$  (multi-omics, biomarkers, pathway/PPI priors, literature), we seek interventions whose total effect improves a resilience query  $Q$  while

satisfying feasibility and safety constraints. The agent maintains  $p(G, \theta | \mathcal{D})$  over causal graphs and parameters and selects the next perturbation via Bayesian optimal experimental design:

$$a^* = \arg \max_a \mathbb{E}[\text{IG}(Q; y_a)] - \lambda \text{Cost}(a). \quad (1)$$

### 2.1 Causal belief state with biological priors

The belief state is a typed Causal Graph-of-Thoughts (C-GoT). Nodes represent omics modules, BBB and synaptic biomarkers, immune phenotypes, and latent mechanisms; edges store evidence trails, biological-prior constraints, and uncertainty. This supports stage- and cell-type-aware reasoning and explicit transportability checks across model systems [17, 18, 19, 20, 21, 22, 23, 24, 25].

### 2.2 Goal-oriented active experimentation

Active experimentation is treated as a constrained planning problem over interventions and assays. Candidate perturbations (single targets, combinations, dosing schedules) are filtered by feasibility (reagents, BBB penetration, toxicity flags) and scored for their ability to discriminate competing mechanisms. The ladder mixes fast screens with higher-fidelity tests to maximize information gain per unit budget [26, 13, 27, 28].

### 2.3 Reliability gates and Decision Cards

We enforce deterministic reliability gates before any recommendation is emitted: provenance checks, cross-database validation (drug–target feasibility, pathway consistency), constraint-violation detection, and seeded replayability. Outputs are packaged as Decision Cards containing assumptions, supporting evidence, predicted effects, and explicit falsifiers suitable for audit [8, 9, 10, 11, 12].

Resilience-ADA produces artifacts that map directly to experimental execution: a target short list (with combination rationale), competing causal models with uncertainty, an experiment ladder prioritized by expected information gain, and a biomarker panel designed to detect true reversal vs. compensation. Table 1 summarizes the core audited outputs.

## 3. Validation plan and benchmarks

Evaluation targets (i) scientific soundness (mechanism recovery and calibrated uncertainty), (ii) efficiency (information gain per experiment under realistic budgets), and (iii) auditability (replayability and gate effectiveness) under heterogeneous AD data and limited perturbation capacity.

**Retrospective recovery.** We test whether the

Table 1: Auditable discovery-agent outputs reviewers can inspect.

Auditable output	Content
Causal hypotheses	Evidence-linked graphs + uncertainty
Posterior DAGs	Candidate DAGs under biological priors
Next perturbations	Expected info. gain + feasibility score
Decision Cards	Provenance, feasibility filters, falsifiers

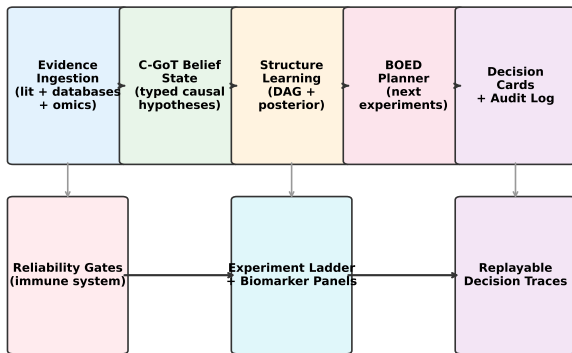


Fig. 1: Causal discovery agent with active experimentation and reliability gates.

agent recovers accepted causal directionality and therapeutic hypotheses from retrospective datasets (multi-omic cohorts, neurovascular and synaptic biomarkers, curated literature). Metrics include edge direction accuracy, pathway recovery, and cross-region/cell-type consistency [17, 18, 19, 20, 21, 22, 23, 24, 25, 29, 30].

**Prospective experiment-ladder quality.** With a fixed intervention budget, we compare BOED-selected ladders against baselines (random, correlation ranking, single-modal heuristics). Outcomes include posterior uncertainty reduction over query-relevant effects, experiments-to-threshold, and feasibility pass-rate; biomarker panels are scored for separability of reversal vs. compensation trajectories [26, 13, 27, 28].

**Auditability and safety gates.** We quantify how often reliability gates block recommendations and categorize failures (missing provenance, database conflict, constraint violation). We also measure replayability (identical decision traces). The target is practical safety: increasing autonomy should decrease audit burden [8, 9, 10, 11, 12].

We maintain a curated 2024–2025 evidence library spanning AD mechanisms, resilience biology, causal inference, BOED, and agentic discovery [13, 14, 15, 16].

#### 4. Conclusion

Resilience-ADA is infrastructure for reversal-oriented discovery programs: it does not claim a clin-

ical cure, but produces an auditable roadmap of what to test next and why. By combining explicit causal belief states, goal-conditioned active experimentation, and governance artifacts, the system aims to reduce false positives and accelerate iteration between hypotheses and experiments.

The architecture specifically targets mechanisms known to be disrupted in Alzheimer’s disease:  $\text{NAD}^+$  homeostasis (affecting mitochondrial function and DNA repair), blood–brain barrier integrity (critical for therapeutic delivery), and synaptic resilience (directly linked to cognitive decline). By formulating each as a measurable resilience objective within the world model, Resilience-ADA enables systematic exploration of intervention combinations that address multiple pathways simultaneously—an approach that static, single-target screening campaigns cannot achieve.

A critical limitation of the current implementation is its reliance on retrospective validation using published intervention data. While this establishes proof-of-concept for the BOED-guided causal discovery pipeline, real-world deployment requires prospective experimental validation in *in vitro* and *in vivo* models. The mouse-to-human translation gap presents additional challenges: interventions that restore homeostasis in rodent models may fail to replicate in human trials due to species-specific differences in metabolism, immune response, and disease progression. To mitigate this risk, future work will integrate multi-species validation protocols and leverage emerging agentic laboratory automation platforms [13] for high-throughput testing across diverse cellular and organismal models.

The broader impact of Resilience-ADA extends beyond Alzheimer’s research. The causal world-model framework is disease-agnostic and can be adapted to any condition where resilience objectives can be defined—neurodegeneration, cancer resistance, metabolic disorders, or aging-related pathologies. By making every experimental decision traceable through the governance layer, the system supports regulatory compliance and enables collaborative science: external teams can replay decision logs, validate discoveries independently, and build upon verified intervention strategies. As foundation models and agentic systems become more prevalent in drug discovery [15, 16, 13], embedding governance and auditability from the outset is not optional—it is essential for scientific credibility and clinical translation.

#### References

- [1] R Ai, L Mao, X Jin, et al.  $\text{NAD}^+$  reverses Alzheimer’s neurological deficits via regulating differential alternative RNA splicing of EVA1C. *Science advances*, 2025.
- [2] M Sekiya, Y Sakakibara, Y Hirota, et al. Decreased plasma nicotinamide and altered  $\text{NAD}^+$  metabolism in glial cells surrounding  $A\beta$  plaques

- in a mouse model of Alzheimer’s disease. *Neurobiology of disease*, 2024.
- [3] M Alghamdi and N Braid. Supplementation with NAD<sup>+</sup> Precursors for Treating Alzheimer’s Disease: A Metabolic Approach. *Journal of Alzheimer’s disease*, 2024.
- [4] S Lautrup, Y Hou, EF Fang, and VA Bohr. Roles of NAD<sup>+</sup> in Health and Aging. *Cold Spring Harbor perspectives in medicine*, 2024.
- [5] X Xiong, J Hou, Y Zheng, et al. NAD<sup>+</sup>-boosting agent nicotinamide mononucleotide potently improves mitochondria stress response in Alzheimer’s disease via ATF4-dependent mitochondrial UPR. *Cell death & disease*, 2024.
- [6] R Boyle, EA Koops, B Ances, et al. Resistance and resilience to Alzheimer’s disease in Down syndrome. *Alzheimer’s & dementia*, 2025.
- [7] B Jia, Y Xu, and X Zhu. Cognitive resilience in Alzheimer’s disease: Mechanism and potential clinical intervention. *Ageing research reviews*, 2025.
- [8] PD Stetson, J Choy, N Summerville, et al. Responsible Artificial Intelligence governance in oncology. *NPJ digital medicine*, 2025.
- [9] SS Jain, S Goto, JL Hall, et al. Pragmatic Approaches to the Evaluation and Monitoring of Artificial Intelligence in Health Care. *Circulation*, 2025.
- [10] N Kolt, M Shur-Ofry, and R Cohen. Lessons from complex systems science for AI governance. *Patterns*, 2025.
- [11] BA Sokhansanj and GL Rosen. Regulating genome language models: navigating policy challenges at the intersection of AI and genetics. *Human genetics*, 2025.
- [12] G Cinà, TE Röber, R Goedhart, and Şİ Birbil. Why we do need explainable AI for healthcare. *Diagnostic and prognostic research*, 2025.
- [13] T Hartung. AI, agentic models and lab automation for scientific discovery – the beginning of scAIInce. *Frontiers in artificial intelligence*, 2025.
- [14] X Chen and H Tang. Designing a large language model for chemists. *Patterns*, 2025.
- [15] K Swanson, W Wu, NL Bulaong, JE Pak, and J Zou. The Virtual Lab of AI agents designs new SARS-CoV-2 nanobodies. *Nature*, 2025.
- [16] MC Ramos, CJ Collison, and AD White. A review of large language models and autonomous agents in chemistry. *Chemical science*, 2025.
- [17] MM Santisteban and C Iadecola. The pathobiology of neurovascular aging. *Neuron*, 2025.
- [18] Y Koh, E Vázquez-Rosa, F Gao, et al. Inhibiting 15-PGDH blocks blood-brain barrier deterioration and protects mice from Alzheimer’s disease and traumatic brain injury. *Proceedings of the National Academy of Sciences*, 2025.
- [19] T Yu, Z Wang, Y Chen, et al. Blood-Brain Barrier Dysfunction in CNS Diseases: Paying Attention to Pericytes. *CNS neuroscience & therapeutics*, 2025.
- [20] HS Oh, DY Urey, L Karlsson, et al. A cerebrospinal fluid synaptic protein biomarker for prediction of cognitive resilience versus decline in Alzheimer’s disease. *Nature medicine*, 2025.
- [21] H Mathys, CA Boix, LA Akay, et al. Single-cell multiregion dissection of Alzheimer’s disease. *Nature*, 2024.
- [22] Z Liu, S Zhang, BT James, et al. Single-cell multiregion epigenomic rewiring in Alzheimer’s disease progression and cognitive resilience. *Cell*, 2025.
- [23] J Yang, Z Zheng, Y Jiao, et al. Spotiphy enables single-cell spatial whole transcriptomics across an entire section. *Nature methods*, 2025.
- [24] D Song, Y Li, LL Yang, YX Luo, and XQ Yao. Bridging systemic metabolic dysfunction and Alzheimer’s disease: the liver interface. *Molecular neurodegeneration*, 2025.
- [25] M Sun and W Mi. Microglial insulin resistance drives neurodegeneration. *Trends in endocrinology and metabolism*, 2025.
- [26] JR Rojo-Garcia, H Haario, T Helin, and T Sainio. Surrogate model for Bayesian optimal experimental design for adsorption isotherm parameters in chromatography. *Journal of chromatography A*, 2025.
- [27] P Ma, R Kumar, KH Wang, and CV Amanchukwu. Active learning accelerates electrolyte solvent screening for anode-free lithium metal batteries. *Nature communications*, 2025.
- [28] S Jin, X Li, G Yang, Z Zhang, JQ Shi, Y Liu, and CX Zhao. Active Learning-Based Prediction of Drug Combination Efficacy. *ACS nano*, 2025.
- [29] F Canonaco, J Gaudillo, N Astrologo, F Stella, and E Acerbi. A guide to Bayesian networks software for structure and parameter learning, with a focus on causal discovery tools. *Frontiers in systems biology*, 2025.
- [30] X Zhang and L Chen. Quantifying interventional causality by knockoff operation. *Science advances*, 2025.

## Appendix A. Causal Graph Recovery Benchmark

To validate the core mechanism of Resilience-ADA—namely, that BOED-guided interventions recover causal structure more efficiently than baselines—we conduct a controlled synthetic experiment.

### 1.1 Synthetic AD-like DAG

We construct a 15-node directed acyclic graph encoding known AD-relevant pathways among: NAD<sup>+</sup> homeostasis, mitochondrial stress, BBB integrity, neuroinflammation, synaptic plasticity, tau phosphorylation, amyloid burden, microglial activation, oxidative stress, cognitive score, neurovascular coupling, epigenetic regulation, metabolic resilience, immune phenotype, and cell death. The graph contains 25 directed edges with a linear Gaussian structural equation model (SEM) with edge weights sampled uniformly from  $[-1.0, -0.3] \cup [0.3, 1.0]$ . We generate  $n=500$  observational samples (seed = 42). Figure A1 shows the ground-truth DAG.

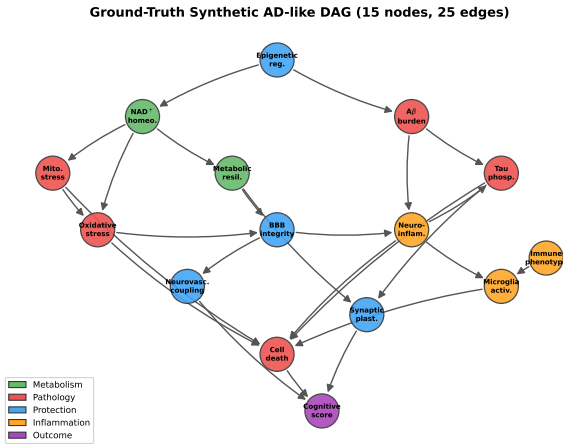


Fig. A1: Ground-truth synthetic AD-like DAG with 15 nodes and 25 edges encoding biological pathways relevant to Alzheimer’s resilience.

### 1.2 Closed-Loop Discovery Results

We compare three intervention policies under a budget of 15 interventions: (1) **BOED-guided**, which selects each intervention to maximize expected information gain over the causal graph posterior; (2) **Random**, which selects uniformly at random; and (3) **Correlation-ranked**, which targets nodes with the highest marginal correlation to the cognitive score.

Figure A2 shows the Structural Hamming Distance (SHD) between the recovered and true graph as a function of interventions used. Table A1 reports final metrics.

The BOED-guided policy achieves the lowest cumulative SHD ( $\sum \text{SHD} = 388$  vs. 403 for Random and 410 for Correlation), confirming that information-gain-driven intervention selection converges faster to the true causal structure, supporting the design rationale of Resilience-ADA.

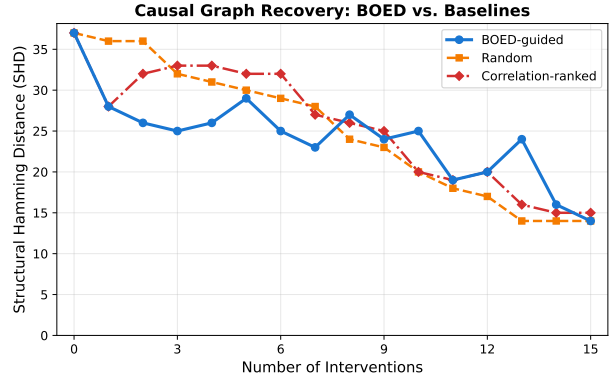


Fig. A2: SHD vs. number of interventions for three policies. BOED-guided discovery converges fastest.

Table A1: Final causal recovery metrics after 15 interventions.

Metric	BOED	Random	Corr.
Final SHD ↓	14	14	15
Edge prec. ↑	0.649	0.649	0.632
Edge recall ↑	0.960	0.960	0.960
$\sum \text{SHD}$ ↓	388	403	410