

Building an Atlas of Social Experiments to Link Studies, Reconcile Conflicts, and Bridge Gaps

Jiawei Zhang^{1,=} Honglin Bao^{1,=} Pengda Wang^{2,=} Alex Yan^{3,=}
Xiao Liu^{1,✉} James A. Evans¹

¹University of Chicago ²Rice University ³Yale University

⁼ Co-first authors: J.Z., H.B., P.W., and A.Y.

[✉] Correspondence: liuxiao@uchicago.edu

Abstract

Social and behavioral science runs thousands of experiments each year, yet their findings rarely accumulate into a coherent map of what is known, what conflicts, and what remains missing. We introduce EXATLAS, a framework for turning an archive of experiments into an atlas: a structured map in which studies link, conflict, or leave bridgeable gaps. Given a target study, EXATLAS searches for prior studies that are locally close in treatment and outcome space and asks whether their observed effects can be composed to predict the target effect. This yields three cases. If composition succeeds and agrees with the observed result, EXATLAS **links** the target to consistent prior evidence. If composition succeeds but disagrees, EXATLAS **reconciles** the conflict and proposes candidate moderators or higher-level theories that could explain it. If the composition fails, EXATLAS proposes **bridge** experiments to close the gap. We provide an error bound for composition under local smoothness of the treatment-effect surface. On held-out targets certified as locally supported, EXATLAS recovers effect direction in 98.6% of cases. Human evaluations further suggest that its proposed bridge experiments are plausible and exhibit connectedness, and that its conflict explanations are useful for theory generation. These results suggest that the archive of social experiments contains more latent structure than current practice extracts—and that making this structure explicit can guide both future theory and future experimentation.

1 Introduction

Social and behavioral science has long suffered what Watts called an “incoherency problem”: scholars examine related phenomena without integrating them, competing findings coexist without being adjudicated, and isolated studies accumulate without theories that explain how they relate (Watts, 2017). Research often proceeds one experiment at a time,

leaving the literature fragmented even after decades of productive research (Almaatouq et al., 2024).

Addressing this problem requires turning fragmented experimental evidence into a coherent research landscape. Such a landscape should make three things visible: which experiments speak to the same phenomenon, where they reinforce or contradict each other, and where gaps between them lie. Done well, it serves not as a summary but as an atlas—a map to guide the next theory and the next experiment. Figure 1 illustrates this vision: an experimental atlas links consistent experiments, reconciles conflicts, and bridges gaps in fragmented literature. We use *atlas* literally: a structured collection of local maps that together cover a territory, where each map is a region of treatment-outcome space in which experiments compose to predict each other’s effects, and where seams between maps are themselves objects of study.

The organization of an experimental landscape is inherently compositional. A target finding may bear little surface resemblance to any single prior study, yet be clearly situated within the literature when several studies are interpreted together. Consequently, matching experiments by surface similarity alone is insufficient. Directly prompting an LLM does not address the problem either. What is required is the ability to reason across many studies at once: to detect compositions, conflicts, and gaps, and to determine how the evidential landscape changes when studies are evaluated jointly rather than in isolation.

We introduce EXATLAS, a framework for constructing such an atlas from joint experimental evidence. At the core of EXATLAS is a composition step that asks whether a target experiment can be reconstructed from prior observed experiments. Consider a target study asking whether a 15-minute mindfulness break before an exam improves performance for first-generation college students. No single prior experiment may study this exact treatment-

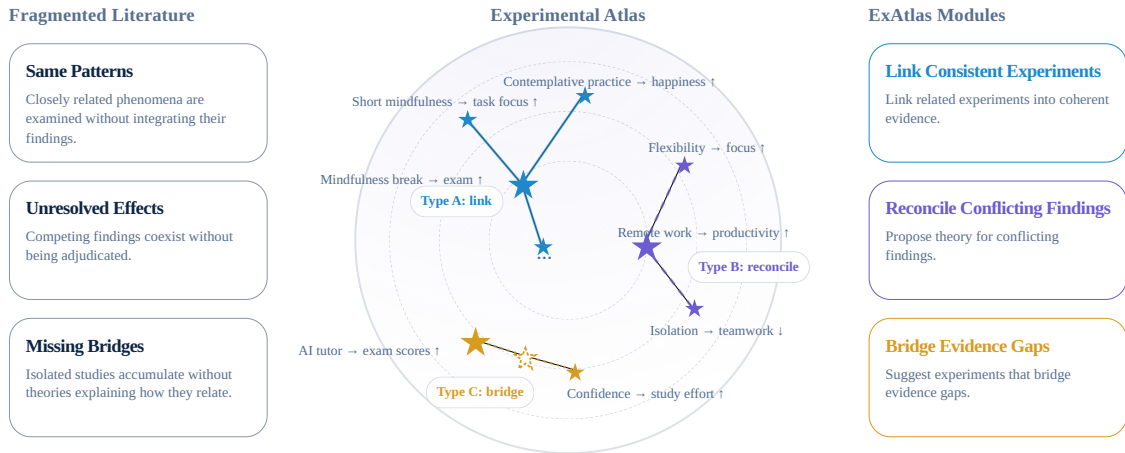


Figure 1: Illustration of how EXATLAS links, reconciles, and bridges experimental evidence.

outcome pair. But the archive may contain studies of shorter mindfulness interventions before tasks, of study-skills supports targeted at first-generation students, and of brief contemplative practices in academic settings. If the target sits inside the region these studies span—close in treatment and outcome meanings—its effect can be predicted by an appropriately weighted average of theirs.

The method is rejectable by design. Theoretically, the composition error decomposes into three terms: extrapolation (how far the target sits from its reconstruction), local model bias (how curved the treatment-effect surface is in the neighborhood), and unobserved residual (variation in factors no embedding captures). When all three are small, the composition is trustworthy. Empirically, we evaluate EXATLAS by holding out experiments from an archive and reconstructing their observed treatment effects from the remaining studies. The certification is the contribution, and the accuracy is its consequence. On targets that EXATLAS certifies as locally supported, the predicted direction matches the observed direction in 71 of 72 cases (98.61%).

Based on the composition result, EXATLAS identifies conflicts and gaps. When a target is composable from nearby experiments, but the composition disagrees with the observed effect, EXATLAS treats the disagreement as a theoretical opportunity. Prior work and the new result imply different effects despite being close in treatment-outcome space, and an LLM is prompted with both to propose a reconciliation that could explain the divergence. When a target cannot be composed, the gap itself is the finding. EXATLAS iteratively proposes bridge experiments until the target enters the composable region. In this way, the same machinery that pre-

dicts known effects also distinguishes two cases: Kuhnian anomalies (Kuhn, 1962), which are existing puzzles in the literature where experiments can be composed but yield conflicting predictions, and Lakatosian extensions (Lakatos, 1978), which mark the field’s next frontiers where existing experiments cannot yet be composed.

Human raters judge EXATLAS’s bridge suggestions more plausible and more connected to the source experiments than forced composition of existing experiments. This suggests that the bridge module can turn gaps in the archive into actionable experiment designs. For composable targets with conflicting evidence, qualitative expert evaluation further shows that EXATLAS can turn disagreements into candidate theoretical explanations. For example, EXATLAS surfaces a hidden 2×2 moderator structure in which the effect of candidate age changes direction across cells, and the LLM articulates a reconciliation that experts judge plausible and useful for interpreting the conflict.

Our contributions are as follows: (1) We introduce EXATLAS, a rejectable framework for composing experimental evidence, whose composition procedure turns experimental archives into an atlas of links, conflicts, and gaps. (2) We provide a theoretical error bound for composition under local smoothness, and show empirically that when EXATLAS certifies local support, composed effects recover the observed effect direction in 98.6% of cases. (3) We show that the atlas’s two failure modes (i.e., composition-with-conflict and non-composition) are themselves generative. The first yields theoretical proposals that human raters find substantive, the second yields bridge-experiment proposals that human raters find feasible, plausible,

and well connected to the source literature.

2 Related Work

LLMs for literature understanding and scientific discovery. A growing line of work asks LLMs to generate research ideas and hypotheses from the literature (Liu et al., 2026; Yang et al., 2024; Baek et al., 2025), with Si et al. (2025) finding that LLM-generated NLP ideas can match or exceed human experts in novelty. This direction has been extended to experiment design (Li et al., 2025b; Guo et al., 2026) and end-to-end research automation (Lu et al., 2026).

Most of these methods, however, retrieve literature shallowly, relying on keyword search and feeding titles or abstracts as background, which limits reasoning over relationships among studies. Some work organizes literature more structurally—Baek et al. (2025) draws on academic graphs and entity stores, and Li et al. (2025a) chains ideas across prior work—but the focus is on connecting and extrapolating existing ideas rather than surfacing the conflicts and gaps that often drive scientific progress (Kuhn, 1962). In contrast, EXATLAS constructs an experimental atlas that explicitly models how studies compose, conflict, or leave gaps, and uses these relations to guide theory generation and experiment design.

LLMs as experiment predictors. Recent work asks whether LLMs can replicate human experiments by simulating participants in social, economic, psychological, and management studies (Aher et al., 2023; Horton et al., 2023; Dillion et al., 2023; Cui et al., 2025), with further efforts training models for behavioral simulation (Binz et al., 2025; Cao et al., 2025). A countercurrent questions their reliability for human psychology and policy-relevant inference, particularly in high-stakes domains (Schröder et al., 2025; Hou et al., 2025).

A related line uses LLMs to predict experimental results directly (Hewitt et al., 2024; Lippert et al., 2024), treating the model as an outcome predictor that infers participant behavior or treatment effects from an experiment description. These predictions, however, rest on the model’s internal knowledge. In contrast, EXATLAS estimates a target effect from existing experiments, and only when those experiments can compose the target, making the evidential basis explicit and the method’s applicable scenarios observable.

3 The EXATLAS Framework

Given a target experiment and an archive of prior experiments, EXATLAS proceeds in two stages (Figure 2). First, it asks whether the target’s treatment and outcome can be composed from prior experiments. If so, it estimates the target’s effect by composing observed effects. Second, it uses the composition to organize the target: a match links the target into the atlas, a conflict triggers a meta/competing-theory proposal, and a failure to compose triggers bridge experiments.

3.1 Representing Experiments

Most experiments in social and behavioral science test causal relations between treatments and outcomes; we encode each study accordingly. Specifically, each experiment i is encoded as

$$e_i = (T_i, O_i, Z_i, \tau_i),$$

where T_i is the treatment description, O_i is the outcome or dependent-variable description, Z_i contains contextual information such as experiment settings, population, and control conditions, and τ_i is the observed treatment effect.

Because treatments and outcomes are often described briefly in papers, EXATLAS first uses an LLM (OpenAI o3 model) to enrich them into context-aware descriptions, for instance, enriching "creativity" to "employee creativity in the organization" by reading the content of the original paper. Human annotators verified 50 randomly sampled enrichments and found all of them accurate. The enriched descriptions incorporate information from Z_i , which helps make composition more accurate.

We then represent the experiments in a shared treatment–outcome space. Let \tilde{T}_i and \tilde{O}_i denote the enriched treatment and outcome descriptions. We define

$$x_i = [\text{emb}(\tilde{T}_i) \parallel \text{emb}(\tilde{O}_i) \parallel \text{emb}(\tilde{T}_i) \odot \text{emb}(\tilde{O}_i)],$$

where $\text{emb}(\cdot)$ denotes the text embedding function, \parallel denotes concatenation, and \odot denotes element-wise multiplication. The first two components encode the treatment and outcome semantics, while the interaction term captures their joint relation.

3.2 Assessing Composability

Given a target experiment e_t , EXATLAS first asks whether e_t can be composed from prior experiments in the archive. Because these experiments are fundamentally about causal relations between

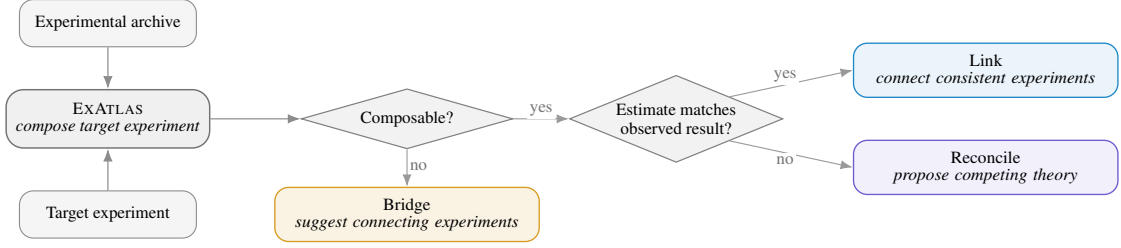


Figure 2: Overview of the EXATLAS framework. Given a target experiment and an archive of prior experiments, EXATLAS tries to compose the target’s treatment and outcome from prior experiments. If the target is composable, EXATLAS estimates its effect from prior observed effects. Consistent estimates link experiments, mismatches motivate theoretical reconciliations, and non-composable targets trigger connecting experiments.

treatments and outcomes, we evaluate composability in the shared treatment–outcome space defined above. The target experiment is *composable* when x_t can be closely approximated by a weighted combination of prior experiment representations.

Given a set of candidate source experiments $\mathcal{C}(t)$ and their reconstruction weights α_{tj} , we compute the distance between the target representation and its composed approximation:

$$r_t(\alpha) = \left\| x_t - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} x_j \right\|_2,$$

where $\|\cdot\|_2$ denotes the Euclidean norm. Because raw distances vary with archive density, we normalize the residual by the local scale:

$$s_t = \text{median}_{j \in \mathcal{C}(t)} \|x_t - x_j\|_2, \quad \rho_t(\alpha) = \frac{r_t(\alpha)}{s_t}.$$

We treat e_t as composable when $\rho_t(\alpha) \leq \lambda$, where λ is a predefined threshold. This makes composition rejectable. If the target is not locally supported, EXATLAS does not force a composition; instead, it treats the target as a gap in the atlas and uses the missing local support to guide follow-up experiment suggestions.

We obtain $\mathcal{C}(t)$ and α_{tj} through a local search because exhaustive search over all subsets is infeasible. EXATLAS instead forms a local candidate set of up to K nearby experiments and solves the constrained optimization problem:

$$\begin{aligned} \alpha_t^* &= \arg \min_{\alpha} r_t(\alpha), \\ \text{subject to } &\alpha_{tj} \geq 0, \quad \sum_{j \in \mathcal{C}(t)} \alpha_{tj} = 1. \end{aligned}$$

3.3 Estimating the Treatment Effect and Linking Consistent Experiments

If a target experiment e_t is composable, EXATLAS estimates its treatment effect by applying the same

weights to the observed effects of the source experiments: $\hat{\tau}_t = \sum_{j \in \mathcal{C}(t)} \alpha_{tj}^* \tau_j$. Here τ_j denotes the observed human treatment effect of source experiment j , while $\hat{\tau}_t$ is the reconstructed estimate for the target. When the target’s human result τ_t is available, it is used only for evaluation and atlas organization, not for estimating $\hat{\tau}_t$.

EXATLAS compares the direction of the composed estimate with the direction of the observed target effect. If $\text{sign}(\hat{\tau}_t) = \text{sign}(\tau_t)$, then the target experiment and its source experiments are treated as a coherent local region of the atlas. In this case, EXATLAS links the target to the source experiments, indicating that the target finding can be situated within the existing experimental landscape.

3.4 Reconciling Conflicts and Bridging Gaps

Not every target experiment that enters the atlas becomes a link. Besides linking consistent experiments, EXATLAS provides two diagnostic outputs to support theory building and experiment design.

The first occurs when a target experiment is composable, but the composed estimate conflicts with the observed target result. In this case, experiments that are close enough to reconstruct the target nevertheless imply different effects, signaling a local theoretical tension in the atlas. EXATLAS then uses an LLM to generate a reconciliation: an explanation for why the source experiments and the target experiment may diverge despite being compositionally related. The LLM is prompted with the target experiment, the source experiments, their observed effects, and the direction of the mismatch. This grounds the theory proposal in the local structure of the atlas rather than in the target experiment alone. The LLM then proposes an evidence-integration explanation, such as a meta-theory or moderator structure, that separates the conditions under which each effect obtains.

The second diagnostic output occurs when the target is not composable. In this case, there is a gap between the target experiment and the composable region of existing literature. EXATLAS uses an LLM to suggest connecting experiments that could bridge this gap, helping turn isolated regions of the literature into a more coherent experimental landscape. This process is iterative: an LLM proposes an experiment at each step, EXATLAS adds the proposed experiment as a hypothetical bridge to the experimental archive, and reruns the composability assessment, repeating the process until the target can be composed from existing and proposed experiments, or a step limit is reached.

We use OpenAI o3 to generate reconciliations and DeepSeek-R1 to propose connecting experiments, as both models performed better in their respective tasks; implementation details and prompts are provided in Appendices B and F, respectively.

4 Validating the Composition Method

Two checks support the method: a theoretical error bound shows that composition is effective when the target experiment is close to its composed approximation, and an empirical evaluation tests whether it recovers held-out experiments from the archive.

4.1 Theoretical Justification

When can a weighted composition of source effects recover the target effect? We show that this is possible when two conditions hold: the target is locally composable from prior experiments, and the treatment-effect surface is smooth in the neighborhood where composition occurs.

Let m_i denote the latent treatment-outcome mechanism representation of experiment e_i . We assume that the treatment effect can be written as

$$\tau_i = \mu(m_i) + \epsilon_i^\mu, \quad |\epsilon_i^\mu| \leq \delta_{\text{unobs}},$$

where μ is a smooth treatment-effect function, ϵ_i^μ captures residual variation due to factors not represented in m_i , and δ_{unobs} bounds the size of this unobserved residual.

The target experiment is locally composable if its latent representation can be approximated by a convex combination of source representations:

$$m_t = \sum_{j \in \mathcal{C}(t)} \alpha_{tj} m_j + \epsilon_t, \quad \alpha_{tj} \geq 0, \quad \sum_{j \in \mathcal{C}(t)} \alpha_{tj} = 1.$$

Here ϵ_t is the latent reconstruction residual. Small $\|\epsilon_t\|$ means that the target lies close to the local structure induced by the source experiments.

Let $\tau_t^{\text{comp}} = \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \tau_j$. Under local smoothness of the treatment-effect function μ , the error of composing source effects is bounded by

$$|\tau_t - \tau_t^{\text{comp}}| \leq L \|\epsilon_t\| + \frac{H}{2} \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \|m_j - \bar{m}_t\|^2 + \left| \epsilon_t^\mu - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \epsilon_j^\mu \right|,$$

where $\bar{m}_t = \sum_{j \in \mathcal{C}(t)} \alpha_{tj} m_j$. Each term is empirically diagnosable. The first term is extrapolation error: the target must be close to the composed source representation. The second term is local curvature error: composition is more accurate when the treatment-effect function is approximately linear in the local neighborhood. The third term is residual error from unobserved factors not captured by the latent representation and bounds the floor on composition error. Full assumptions and the proof are provided in Appendix A.

In practice, the latent representation m_i is not directly available to EXATLAS. It uses the treatment-outcome feature representation x_i as an empirical proxy for the latent mechanism representation. The normalized residual ρ_t estimates whether the target is close enough to the local source composition to be treated as composable.

4.2 Empirical Evaluation

Experimental setup. We evaluate composition on an archive of 360 experiments from five top management and psychology journals, collected by Cui et al. (2025). We hold out one experiment at a time as the target and use the remaining experiments as the source archive. Applying the composability assessment, EXATLAS identifies 72 experiments as composable. For these composable targets, we evaluate whether the composed effect recovers the reported human effect of the held-out target experiment. We use sign match to measure whether the composed effect has the same direction as the reported target effect, and use MSE, MAE, and Spearman’s ρ to assess effect-size error and rank preservation. Implementation details and metric definitions are in Appendix B and C.

Baselines. We compare EXATLAS against four baselines, with implementation details in Appendix D. (1) *Direct Prediction* asks the model to predict the target effect direction using only the target experiment description. (2) *Nearest-Experiment RAG* provides the model with the sin-

Method	Sign match \uparrow	MSE \downarrow	MAE \downarrow	Spearman ρ \uparrow
Direct Prediction	70.83%	6.4509	1.1941	0.3881
Nearest-Experiment RAG	79.17%	9.4243	1.1857	0.4668
Contributing-Experiment RAG	87.50%	6.3498	0.8209	0.6563
Synthetic Participants	83.10%	7.0761	1.6673	0.6168
EXATLAS Composition	98.61%	2.3477	0.7435	0.8474

Table 1: Performance comparison on predicting treatment effects for composable experiments.

gle most similar prior experiment retrieved by TF-IDF similarity, including its reported effect. (3) *Contributing-Experiment RAG* provides the model with the source experiments identified by EXATLAS as contributing to the target composition, including their reported effects. (4) *Synthetic Participants* uses the released LLM-as-participant replication results from Cui et al. (2025), which simulate participant responses to the original scenarios and survey questions and calculate treatment effects by aggregating responses. GPT-5.4 is used for the three prompting baselines, and Claude 3.5 Sonnet is used for the simulation baseline because it achieves the best performance in that study.

Results. Table 1 shows that EXATLAS achieves the strongest overall performance on predicting effects. It matches the reported effect direction for 98.61% of targets, substantially outperforming all the baselines, and also estimates the effect size more accurately. This indicates the effectiveness of the composition based on experiment embeddings.

We also observe that, by providing the contributing experiments identified by EXATLAS, Contributing-Experiment RAG outperforms the other prompting baselines and the simulation method. This shows that the experiments selected by EXATLAS contain useful evidential signal.

5 Evaluating EXATLAS-Guided Research Directions

5.1 Worked Example: A Hidden 2×2 in the Hiring Literature

Before turning to systematic evaluation, we present a real case (originally published by (Lee et al., 2015)) that illustrates how EXATLAS finds an anomaly and turns it into a candidate theoretical reconciliation. In a creativity-focused, competitively framed department, a new experiment found that higher candidate age increased hiring-recommendation strength. The algorithm surfaced a cluster of prior findings that, together with the new finding, form a hidden 2×2 crossing of departmental orientation (creativity vs. stability) and

climate (cooperative vs. competitive). As shown in Table 2, the age effect changes direction across cells. The contradiction does not lie in any single prior effect; it emerges from a moderator crossing that had been treated as separable in the literature.

	Cooperative	Competitive
Creativity	\downarrow (prior)	\uparrow (new)
Stability	\uparrow (prior)	\downarrow (prior)

Table 2: Sign of the age effect on hiring recommendation in the worked example. **Bold**: new finding.

The case shows a three-step division of labor between the algorithm, the LLM, and the expert. The algorithm surfaces a structural anomaly in the hidden 2×2 : the effect of candidate age changes direction across cells in a way that is not explained by either departmental orientation or climate alone. One human expert, blind to the algorithm, did not locate this cross-cell conflict; another expert proposed industry as a moderator, which was directionally relevant but remained at single-moderator resolution and left part of the disagreement unexplained. Based on the anomaly, the LLM articulated the full 2×2 as *Stereotype Fit-Context Alignment*, with cell-specific goals and a testable boundary condition: the reversal hinges on whether competition is framed as outcome accountability (favors age) or speed (favors youth). In *post hoc* review, the experts judged the framing both plausible and useful for interpreting the conflict. The remainder of this section systematically evaluates both diagnostic outputs of EXATLAS at scale: reconciliations when composition conflicts (§5.2) and bridge suggestions when composition fails (§5.3).

5.2 Human Evaluation: Reconciling Conflicts

We first evaluate the Reconcile module. The threshold λ used in §4.2 was selected to maximize prediction accuracy, so composable-with-mismatch cases, the Reconcile module’s input, are correspondingly scarce under this conservative setting (1/72). Here we relax the threshold to $\lambda' = 1.5\lambda$, admitting

targets whose composition is still well supported but no longer near-perfect. This yields 13 conflict cases, where the composed direction disagrees with the observed effect. For each conflict case, EXATLAS instructs an LLM to propose a meta/competing theory that could reconcile the conflicting evidence.

To evaluate these reconciliations, we recruited two Ph.D. students in social science from two US institutions, each with top-tier publications. Because reconciliations are difficult to evaluate with quantitative criteria, we use a qualitative expert-evaluation procedure. Each expert first examined the conflict case independently, without seeing the model output, and proposed their interpretation of how the findings might be reconciled. They were then shown the LLM-generated reconciliation and asked to compare it with their assessment.

In 8 of the 13 conflict cases, both experts and the LLM identified the same broad direction of reconciliation, while the LLM often supplied a more explicit mechanism or moderator framing than the experts initially articulated. The worked example in §5.1, the age \times creativity \times competition case, illustrates this pattern. The EXATLAS algorithm identifies the anomaly, the LLM reconciles the conflict by articulating the moderator structure, and the two experts ratify the framing in *post hoc* review. A second case involving racial-identity cues in the selection of a debate opponent follows the same pattern. One human expert wrote that “the key advantage of LLM feedback is that it gives a more detailed view of existing theory . . . human rater is more limited in this sense.”

Human experts do not always agree with the LLM. In one case, an experiment found that gender in a *non-empowering* coworker-assistance episode decreased self-perceived competence, echoing prior findings on *empowering* assistance but reversing the direction of the effect. The LLM proposed a threat-benefit meta-theory to reconcile this pattern. A human expert disagreed, writing that this reflected the “same mechanism, opposite types of interaction,” and therefore added “nothing new” Four additional cases follow this pattern. Whether a sign flip under a moderator-by-construction counts as a theoretical contribution is a methodological choice. These disagreements help clarify what the Reconcile module contributes. Reconcile does not determine whether a conflict should be treated as a new theory, a boundary condition, or a familiar mechanism instantiated in a new context. Rather, it renders that interpretive choice explicit by trans-

forming conflicting evidence into a theoretical proposal that experts can accept, revise, or reject. In this sense, Reconcile shifts the burden from implicit judgment to explicit adjudication, without claiming to settle the adjudication itself. We understand the framework as making precisely such choices legible. Taken together, the cases illustrate Reconcile’s role across modes: surfacing a cross-cell conflict in a hidden moderator structure that neither party would likely have located unaided, articulating the relevant prior literature more exhaustively than any single expert, and exposing a productive disagreement over what counts as a theoretical contribution.

5.3 Human Evaluation: Connecting Experiments

We then conduct a human evaluation to assess whether the connecting experiments proposed by EXATLAS are useful research directions for targets that are not locally composable from the existing experimental archive. These are the targets the composability assessment rejects: their normalized reconstruction residual ρ_t exceeds the threshold λ , so EXATLAS declines to compose and routes them to the Bridge module instead.

We compare two conditions on targets that EXATLAS identifies as non-composable. The first, Forced Composition, runs the composition algorithm on each target regardless of whether the composability assessment rejects it, returning for each target a list of contributing experiments from the existing archive. This represents a direct-composition strategy without a refusal or gap-detection mechanism. The second, EXATLAS Bridge, uses the proposed gap-bridging procedure: rather than forcing a composition under insufficient support, EXATLAS generates one or more connecting experiments intended to link the target to nearby experimental evidence, until the target becomes composable from existing and proposed experiments or a step limit is reached. This could also be viewed as the ablation study of the Bridge module.

Evaluation dimensions. Each experiment (contributing experiment in Forced Composition and connecting experiment in EXATLAS Bridge) is evaluated independently along three dimensions using a five-point Likert scale: **Feasibility**, **Plausibility**, and **Connectedness**. Feasibility captures whether the experiment can be empirically implemented with a clear design, including a manipulable independent variable, a measurable depen-

Dimension	Forced Composition	EXATLAS Bridge	t	df	p	95% CI
Feasibility	4.89 (0.31)	4.88 (0.33)	-0.465	396.78	.642	[-0.078, 0.048]
Plausibility	3.78 (0.94)	3.96 (0.75)	2.061	380.39	.040	[0.008, 0.342]
Connectedness	4.10 (1.03)	4.29 (0.75)	2.163	363.05	.031	[0.018, 0.372]

Table 3: Comparison of Forced Composition and EXATLAS Bridge scores across three evaluation dimensions. Values for each condition report means with standard deviations in parentheses. The table also reports Welch’s t -test statistics, degrees of freedom, p -values, and 95% confidence intervals for the mean difference.

dent variable, and a plausible data collection procedure. Plausibility assesses whether the proposed relationship between the independent and dependent variables is theoretically, empirically, or logically reasonable. Connectedness evaluates whether the suggested experiment meaningfully links or explains the original input experiments, beyond surface-level keyword overlap. Six raters (three Ph.D. students and three advanced undergraduates, all with backgrounds in organizational behavior or industrial-organizational psychology) evaluated 100 lists of contributing experiments per condition (Forced Composition vs. EXATLAS Bridge) on the three dimensions. Each experiment received at least two independent ratings. Full assignment and demographic details are in Appendix E.

Inter-Rater agreement. Inter-rater reliability was moderate to high across conditions (ICC(1,2) range 0.45–0.87), and within-one agreement exceeded 0.97 on every dimension (Table 4 in Appendix E). Feasibility ratings were near ceiling in both conditions, with 88–89% of ratings at the maximum score, likely attenuating ICC estimates for this dimension. We therefore focus substantive comparisons on Plausibility and Connectedness.

Comparison results. We conducted Welch’s independent-samples t -tests as appropriate to compare the two conditions across the three rating dimensions. Welch’s t -test was used because it does not assume equal variances across groups and is therefore more appropriate when the rating distributions differ in variance across methods. As shown in Table 3, the two methods received nearly identical Feasibility scores, with Forced Composition rated at 4.89, and EXATLAS Bridge rated at 4.88. This difference was not statistically significant, $t(396.78) = -0.47$, $p = .642$, suggesting that EXATLAS Bridge does not gain its advantage by proposing experiments that are harder or easier or clearer to implement. In contrast, EXATLAS Bridge performed better on the other two dimensions. For Plausibility, EXATLAS Bridge received higher ratings than Forced Composition, 3.96 ver-

sus 3.78, and this difference was statistically significant, $t(380.39) = 2.06$, $p = .040$. For Connectedness, EXATLAS Bridge also received higher ratings, 4.29 versus 4.10, with a statistically significant difference, $t(363.05) = 2.16$, $p = .031$.

Together, these results support the gap-bridging design of EXATLAS. EXATLAS Bridge proposes new experiments that help reduce the gap between the target and nearby experimental evidence. Importantly, the proposed experiments remain empirically grounded: compared with the published experiments retrieved by Forced Composition, they are equally feasible and more plausible.

6 Conclusion

This paper introduces EXATLAS, a rejectable framework for turning fragmented experimental archives into an atlas of links, conflicts, and gaps. Given a target experiment, it asks whether the existing archive of human experiments is close enough, in treatment and outcome, to compose a prediction. When it is, EXATLAS links consistent evidence or turns mismatches into theoretical reconciliation. When it is not, EXATLAS declines to compose and proposes bridge experiments intended to make composition possible. On the targets it certifies as locally supported, predicted directions match observed ones in 71 of 72 cases. On targets it declines, human raters judge the proposed bridge experiments to be feasible, plausible, and well connected to the target literature.

We do not claim that LLMs replace human participants or that composition based on text embeddings is universally accurate. We claim something narrower but more useful: that an experimental archive, read compositionally, supports a calibrated estimate of where its own evidence reaches, where it falls short, and where those limits point to meaningful research directions. If social science is to escape what Watts (2017) called the incoherency problem, it needs infrastructure for reading the archive as a map rather than a pile. EXATLAS is one such infrastructure.

Limitations

The strongest result applies only to locally supported targets. The 98.61% replication rate is measured on the 72-experiment subset where nearby existing experiments already match the target in treatment and outcome context. It should not be generalized to targets that require extrapolation beyond the archive. More fundamentally, EXATLAS cannot prove from text alone that an embedding space is mechanism-preserving. The paper therefore treats feature geometry as a candidate-retrieval and reconstruction device, not as a causal identification argument by itself. Composition is interpreted only when the residual is small and the target passes the composability assessment.

The archive is itself biased: what gets published, what gets coded as a clean treatment-outcome experiment, and what gets a coded effect direction all reflect editorial and disciplinary selection. Local composition inherits these biases; the atlas is dense where the literature is dense and silent where it is silent.

The meta-competing-theory reconciliation module is a generator, not a verifier. Section 5.2 shows that the LLM can articulate moderator structures that an expert ratifies, but ratification is not validation. The proper next step is to run the boundary-condition experiments that the module proposes.

Ethics Statement

The research ideas, bridge experiments, and theoretical reconciliations generated by LLMs may reflect biases present in their training data and may unintentionally resemble existing work without proper attribution. They should therefore be treated as candidate research directions rather than validated scientific claims, and should not be adopted for practical use without careful expert review, literature checking, and empirical validation. Because the framework operates on human-subject experimental evidence, researchers should also consider whether proposed experiments involve sensitive populations, identity-related manipulations, or other ethical risks before implementation. Any use of LLM-generated ideas or explanations should be disclosed transparently to preserve research integrity. In this paper, we only used AI to polish text and write code. All AI-generated content has been well examined by authors before adoption.

References

- Gati V. Aher, Rosa I. Arriaga, and Adam Tauman Kalai. 2023. **Using large language models to simulate multiple humans and replicate human subject studies.** In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 337–371. PMLR.
- Abdullah Almaatouq, Thomas L Griffiths, Jordan W Suchow, Mark E Whiting, James Evans, and Duncan J Watts. 2024. Beyond playing 20 questions with nature: Integrative experiment design in the social and behavioral sciences. *Behavioral and Brain Sciences*, 47:e33.
- Teresa M. Amabile. 1982. **Social psychology of creativity: A consensual assessment technique.** *Journal of Personality and Social Psychology*, 43(5):997–1013.
- Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. 2025. Researchagent: Iterative research idea generation over scientific literature with large language models. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6709–6738.
- Marcel Binz, Elif Akata, Matthias Bethge, Franziska Brändle, Fred Callaway, Julian Coda-Forno, Peter Dayan, Can Demircan, Maria K Eckstein, Noémi Éltető, and 1 others. 2025. A foundation model to predict and capture human cognition. *Nature*, 644(8078):1002–1009.
- Yong Cao, Haijiang Liu, Arnav Arora, Isabelle Augenstein, Paul Röttger, and Daniel Hershcovich. 2025. Specializing large language models to simulate survey response distributions for global populations. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3141–3154.
- Ziyan Cui, Ning Li, and Huaikang Zhou. 2025. **A large-scale replication of scenario-based experiments in psychology and management using large language models.** *Nature Computational Science*, 5:627–634.
- Douglas L. Dean, Jillian M. Hender, Thomas L. Rodgers, and Eric L. Santanen. 2006. **Identifying quality, novel, and creative ideas: Constructs and scales for idea evaluation.** *Journal of the Association for Information Systems*, 7(10).
- Danica Dillion, Niket Tandon, Yuling Gu, and Kurt Gray. 2023. **Can ai language models replace human participants?** *Trends in Cognitive Sciences*, 27(7):597–600.
- Yuwei Guo, Zihan Zhao, Xiaowei Liu, Xiangning Yu, and Deyu Zhou. 2026. **From script to stage: Automating experimental design for social simulations with llms.** *Preprint*, arXiv:2512.08935.

- Luke Hewitt, Ashwini Ashokkumar, Isaias Ghezae, and Robb Willer. 2024. [Predicting results of social science experiments using large language models](#). Preprint, arXiv:2404.11794.
- John J. Horton, Apostolos Filippas, and Benjamin S. Manning. 2023. [Large language models as simulated economic agents: What can we learn from homo silicus?](#) Working Paper 31122, National Bureau of Economic Research.
- Abe Bohan Hou, Hongru Du, Yichen Wang, Jingyu Zhang, Zixiao Wang, Paul Pu Liang, Daniel Khashabi, Lauren M Gardner, and Tianxing He. 2025. Can a society of generative agents simulate human behavior and inform public health policy? a case study on vaccine hesitancy. In *Second Conference on Language Modeling*.
- David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. [Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions](#). In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 169–182, Dublin, Ireland. Association for Computational Linguistics.
- Thomas Samuel Kuhn. 1962. *The Structure of Scientific Revolutions*. University of Chicago Press, Chicago.
- Imre Lakatos. 1978. *The methodology of scientific research programmes: Volume 1: Philosophical papers*, volume 1. Cambridge university press.
- Sun Young Lee, Marko Pitesa, Stefan Thau, and Madan M Pillutla. 2015. Discrimination in selection decisions: Integrating stereotype fit and interdependence theories. *Academy of Management Journal*, 58(3):789–812.
- Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang, Yuming Jiang, Yifei Xin, Ronghao Dang, Yu Rong, Deli Zhao, Tian Feng, and Lidong Bing. 2025a. [Chain of ideas: Revolutionizing research via novel idea development with LLM agents](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 8971–9004, Suzhou, China. Association for Computational Linguistics.
- Yu Li, Lehui Li, Lin Chen, Qingmin Liao, Fengli Xu, and Yong Li. 2025b. Agentextpt: Automating ai experiment design with llm-based resource retrieval agent. *arXiv preprint arXiv:2511.04921*.
- Steffen Lippert, Anna Dreber, Magnus Johannesson, Warren Tierney, Wilson Cyrus-Lai, Eric Luis Uhlmann, Thomas Pfeiffer, Emotion Expression Collaboration, and 1 others. 2024. Can large language models help predict results from a complex behavioural science study? *Royal Society Open Science*, 11(9).
- Xiao Liu, Haokun Liu, and Chenhao Tan. 2026. [Heuristic-based ideation for guiding llms toward structured creativity](#). In *ICLR Blogposts 2026*. <https://iclr-blogposts.github.io/2026/blog/2026/ideation-heuristics/>.
- Chris Lu, Cong Lu, Robert Tjarko Lange, Yutaro Yamada, Shengran Hu, Jakob Foerster, David Ha, and Jeff Clune. 2026. Towards end-to-end automation of ai research. *Nature*, 651(8107):914–919.
- Sarah Schröder, Thekla Morgenroth, Ulrike Kuhl, Valerie Vaquet, and Benjamin Paaßen. 2025. Large language models do not simulate human psychology. *arXiv preprint arXiv:2508.06950*.
- Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. 2025. Can llms generate novel research ideas? a large-scale human study with 100+ nlp researchers. In *International Conference on Learning Representations*, volume 2025, pages 94003–94092.
- Chris van der Lee, Albert Gatt, Emiel van Miltenburg, Sander Wubben, and Emiel Krahmer. 2019. [Best practices for the human evaluation of automatically generated text](#). In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 355–368, Tokyo, Japan. Association for Computational Linguistics.
- Duncan J Watts. 2017. Should social science be more solution-oriented? *Nature Human Behaviour*, 1(1):0015.
- Zonglin Yang, Xinya Du, Junxian Li, Jie Zheng, Soujanya Poria, and Erik Cambria. 2024. Large language models for automated open-domain scientific hypotheses discovery. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13545–13565.

A Details of the Theoretical Justification

A.1 Assumptions

We state the assumptions used in the theoretical analysis of EXATLAS.

A1. Source internal validity. For each source experiment e_i , random assignment makes the treatment and control groups comparable, so the difference in their average outcomes estimates the treatment effect τ_i . That is, treatment assignment $A_i \in \{0, 1\}$ is randomized:

$$A_i \perp (Y_i(1), Y_i(0)).$$

The observed outcome satisfies consistency:

$$Y_i = A_i Y_i(1) + (1 - A_i) Y_i(0).$$

Therefore, the source average treatment effect is identified by

$$\tau_i = \mathbb{E}[Y_i | A_i = 1] - \mathbb{E}[Y_i | A_i = 0].$$

A2. Approximate mechanism sufficiency. Each experiment e_i has a latent treatment–outcome representation

$$m_i = m(T_i, O_i, Z_i) \in \mathbb{R}^d,$$

which captures the treatment, outcome, and contextual features that are relevant for treatment-effect variation across experiments. We assume the treatment effect can be written as

$$\tau_i = \mu(m_i) + \epsilon_i^\mu, \quad |\epsilon_i^\mu| \leq \delta_{\text{unobs}},$$

where μ maps mechanism representations to treatment effects, and ϵ_i^μ captures residual variation from features not represented in m_i . The bound δ_{unobs} limits the contribution of such unobserved factors.

A3. Local smoothness. The treatment-effect function μ is locally smooth around the target and its selected source experiments. Specifically, we assume a local Lipschitz condition,

$$|\mu(m) - \mu(m')| \leq L\|m - m'\|,$$

and bounded local curvature,

$$\|\nabla^2 \mu(m)\|_{\text{op}} \leq H,$$

for experiments in this neighborhood. Intuitively, nearby experiments in mechanism space should have related treatment effects, and the effect surface should not change abruptly within the local neighborhood.

A4. Local support. The target representation lies approximately in the convex hull of nearby source representations:

$$m_t = \sum_{j \in \mathcal{C}(t)} \alpha_{tj} m_j + \epsilon_t,$$

$$\alpha_{tj} \geq 0, \quad \sum_{j \in \mathcal{C}(t)} \alpha_{tj} = 1.$$

The residual ϵ_t measures the degree to which the target lies outside local source support.

A5. Representation validity. The latent mechanism representation m_i is not observed directly. In practice, EXATLAS uses the text-based experiment representation x_i , constructed from the treatment and outcome descriptions, as an empirical proxy. We assume that this representation approximately preserves local neighborhood structure, so that nearby experiments and reconstruction weights in x -space are informative about the corresponding relationships in latent mechanism space.

A.2 Proof of the Error Bound

Let

$$\bar{m}_t = \sum_{j \in \mathcal{C}(t)} \alpha_{tj} m_j.$$

By local compositional support,

$$m_t = \bar{m}_t + \epsilon_t.$$

The target effect is

$$\tau_t = \mu(m_t) + \epsilon_t^\mu,$$

and the composed effect is

$$\tau_t^{\text{comp}} = \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \tau_j.$$

Using $\tau_j = \mu(m_j) + \epsilon_j^\mu$, we have

$$\tau_t^{\text{comp}} = \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(m_j) + \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \epsilon_j^\mu.$$

We decompose the approximation error:

$$\begin{aligned} |\tau_t - \tau_t^{\text{comp}}| &= \left| \mu(m_t) - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(m_j) \right. \\ &\quad \left. + \epsilon_t^\mu - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \epsilon_j^\mu \right| \\ &\leq \left| \mu(m_t) - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(m_j) \right| \\ &\quad + \left| \epsilon_t^\mu - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \epsilon_j^\mu \right|. \end{aligned}$$

For the mechanism component, add and subtract $\mu(\bar{m}_t)$:

$$\begin{aligned} &\left| \mu(m_t) - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(m_j) \right| \\ &\leq |\mu(m_t) - \mu(\bar{m}_t)| \\ &\quad + \left| \mu(\bar{m}_t) - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(m_j) \right|. \end{aligned}$$

The first term is controlled by local Lipschitz smoothness:

$$|\mu(m_t) - \mu(\bar{m}_t)| \leq L\|m_t - \bar{m}_t\| = L\|\epsilon_t\|.$$

We now bound the second term. For each source experiment e_j , apply Taylor’s theorem to $\mu(m_j)$ around \bar{m}_t :

$$\mu(m_j) = \mu(\bar{m}_t) + \nabla \mu(\bar{m}_t)^\top (m_j - \bar{m}_t) + R_j,$$

where the second-order remainder satisfies

$$|R_j| \leq \frac{H}{2} \|m_j - \bar{m}_t\|^2$$

because $\|\nabla^2 \mu(m)\|_{\text{op}} \leq H$ in the local neighborhood.

Taking the weighted sum gives

$$\begin{aligned} \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(m_j) &= \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(\bar{m}_t) \\ &\quad + \nabla \mu(\bar{m}_t)^\top \sum_{j \in \mathcal{C}(t)} \alpha_{tj} (m_j - \bar{m}_t) \\ &\quad + \sum_{j \in \mathcal{C}(t)} \alpha_{tj} R_j. \end{aligned}$$

Since the weights sum to one,

$$\sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(\bar{m}_t) = \mu(\bar{m}_t).$$

The first-order term cancels because

$$\begin{aligned} &\sum_{j \in \mathcal{C}(t)} \alpha_{tj} (m_j - \bar{m}_t) \\ &= \sum_{j \in \mathcal{C}(t)} \alpha_{tj} m_j - \bar{m}_t \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \\ &= \bar{m}_t - \bar{m}_t = 0. \end{aligned}$$

Thus,

$$\sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(m_j) = \mu(\bar{m}_t) + \sum_{j \in \mathcal{C}(t)} \alpha_{tj} R_j.$$

Therefore,

$$\begin{aligned} \left| \mu(\bar{m}_t) - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \mu(m_j) \right| &= \left| \sum_{j \in \mathcal{C}(t)} \alpha_{tj} R_j \right| \\ &\leq \sum_{j \in \mathcal{C}(t)} \alpha_{tj} |R_j| \\ &\leq \frac{H}{2} \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \|m_j - \bar{m}_t\|^2. \end{aligned}$$

Combining the two bounds gives

$$\begin{aligned} |\tau_t - \tau_t^{\text{comp}}| &\leq L \|\epsilon_t\| + \frac{H}{2} \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \|m_j - \bar{m}_t\|^2 \\ &\quad + \left| \epsilon_t^\mu - \sum_{j \in \mathcal{C}(t)} \alpha_{tj} \epsilon_j^\mu \right|. \end{aligned}$$

Since $|\epsilon_i^\mu| \leq \delta_{\text{unobs}}$ and the weights sum to one, the residual term is bounded by $2\delta_{\text{unobs}}$.

B Implementation Details of EXATLAS

In the composability assessment, the embeddings are computed with sentence-transformers/all-mpnet-base-v2, an MPNet-based sentence embedding model. For each target, EXATLAS first computes distances to all candidate source experiments in the representation space. It then forms a local candidate set by taking experiments whose distance is within 1.5 times the median candidate distance, capped at $K = 30$ nearest candidates for scalability.

Within this selected neighborhood, EXATLAS solves the nonnegative, unit-sum reconstruction problem, adding a ridge penalty $10^{-2} \|\alpha_t\|_2^2$ for numerical stability. The implementation uses OSQP when available and falls back to the default convex solver otherwise. If optimization fails to return weights, uniform weights over the selected neighbors are used as a conservative fallback.

We determine the threshold λ by maximizing $(1 - \text{scaled MSE}) \times \text{Coverage}$ to balance the accuracy and coverage of the method. Here, Coverage denotes the fraction of experiments that pass the composability assessment, i.e., the fraction with $\rho_t \leq \lambda$. And the resulted threshold is $\lambda = 0.462$.

We use the o3 model for the reconciliation step. When the composed direction derived from contributing experiments conflicts with the observed sign in human experiments, we prompt the model to propose a reconciliation (e.g., a meta theory) to account for the discrepancy (see Appendix F for the prompt).

For the bridging step, we follow a stepwise search process. At each step, we rerun the algorithm jointly over the existing experiments and the newly suggested ones, checking whether the ratio of isolated experiments decreases. We evaluated several advanced models for proposing connecting experiments and found that DeepSeek-R1 performed best (Figure 3). The first round connects nearly all of the isolated experiments¹ that can be connected; subsequent rounds yield diminishing returns. We attribute this to the nature of our dataset—the collection is incomplete, so some experiments are genuinely unlinkable. We thus show the results in the main paper based on one round of suggestions proposed by Deepseek-R1.

¹"isolate" means that they can neither contribute to other target experiments nor serve as a target themselves, thus fewer experiments than those rejected move on to the reconcile and bridge module.

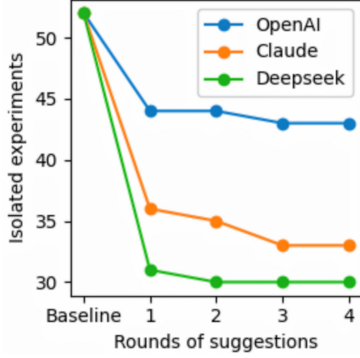


Figure 3: DeepSeek-R1 dominates the arena when it comes to suggesting new experiments to connect existing ones (vs. OpenAI o3 and Claude Opus 4.5).

C Evaluation Metrics for the Composition Experiment

We evaluate predicted treatment effects using four metrics: sign match, MSE, MAE, and Spearman’s ρ . Let $\hat{\tau}_t$ denote the predicted treatment effect for target experiment t , and let τ_t denote the reported human treatment effect.

We define sign match as the fraction of targets for which the predicted effect has the same direction as the reported human effect:

$$\text{SignMatch} = \frac{1}{M} \sum_{t=1}^M \mathbf{1}[\text{sign}(\hat{\tau}_t) = \text{sign}(\tau_t)],$$

where M is the number of evaluated target experiments.

Let $e_t = \hat{\tau}_t - \tau_t$ denote the prediction error. We report mean squared error:

$$\text{MSE} = \frac{1}{M} \sum_{t=1}^M e_t^2,$$

and mean absolute error:

$$\text{MAE} = \frac{1}{M} \sum_{t=1}^M |e_t|.$$

Finally, we report Spearman’s rank correlation between the predicted effects and the reported human effects. This metric evaluates whether the method preserves the relative ordering of treatment effects across target experiments.

D Baselines for the Composition Experiment

All prompting baselines ask an LLM to predict the direction of the target treatment effect.

Direct Prediction. The Direct Prediction baseline uses only the target experiment information,

including the treatment, outcome, and detailed conditions of the experiment.

Nearest-Experiment RAG. The Nearest-Experiment RAG baseline augments the target description with one retrieved prior experiment. The retrieved experiment is selected by TF-IDF cosine similarity over the experiment descriptions. The prompt includes the retrieved experiment’s reported effect, allowing the model to use a simple nearest-neighbor evidence signal.

Contributing-Experiment RAG. The Contributing-Experiment RAG baseline augments the target description with the contributing experiments identified by EXATLAS. Due to length limit, at most 20 contributing experiments are included, ordered by their weights.

Synthetic Participants. The Synthetic Participants baseline uses the released LLM-as-participant replication results from Cui et al. (2025), where an LLM simulates human participants given the questionnaires in the treatment and control conditions. For each target experiment, simulated responses are aggregated within treatment and control conditions.

E Human Evaluation of Connecting Experiments

We conducted a human evaluation to assess the quality of the generated experiment suggestions, following established practices for evaluating automatically generated text (Howcroft et al., 2020; van der Lee et al., 2019). We evaluated suggestions along three dimensions: Feasibility, Plausibility, and Connectedness, drawing on prior work on idea quality and subjective assessment of creative products (Dean et al., 2006; Amabile, 1982).

Six independent raters participated in the evaluation, including three Ph.D. students and three junior undergraduate students, all of whom had backgrounds in organizational behavior, industrial-organizational psychology, or closely related areas. The raters included two men and four women.

Raters were randomly assigned to one of two conditions. Three raters evaluated the 100 Forced Composition suggestions, and the remaining three raters evaluated the 100 EXATLAS Bridge suggestions. Within each condition, study assignments were organized such that every experiment was evaluated by at least two raters.

For each annotation task, raters were shown a experiment to evaluate (contributing experiment in Forced Composition and generated connecting experiment in EXATLAS Bridge) together with its source context. They then independently evaluated the experiment on a five-point Likert scale. Raters were instructed to assess each dimension independently, as a experiment could score highly on one dimension while scoring lower on another.

Feasibility. Feasibility captures the extent to which the experiment can be operationalized in practice. A highly feasible experiment is one for which a concrete empirical design, including data collection, manipulation or intervention, experimental procedure, and measurement strategy, can be readily specified and executed.

Item: How feasible is it to empirically implement the experiment?

- 1 – Not feasible: Cannot be realistically operationalized or tested.
- 2 – Low feasibility: Presents significant practical or methodological challenges.
- 3 – Moderate feasibility: Feasible with some effort or additional assumptions.
- 4 – High feasibility: Can be implemented using standard research methods.
- 5 – Very high feasibility: Easily operationalized with a clear and straightforward empirical design.

Plausibility. Plausibility captures the likelihood that a theoretically grounded, meaningful relationship exists between the independent and dependent variables. This dimension reflects whether the proposed relationship is supported by existing theoretical or empirical knowledge, as well as by the internal logic of the experiment.

Item: How plausible is the relationship between the independent and dependent variables of the experiment?

- 1 – Not plausible: The relationship is unlikely or lacks logical or theoretical support.
- 2 – Slightly plausible: The relationship has weak or unclear justification.
- 3 – Moderately plausible: The relationship has some reasonable basis, but remains uncertain.

- 4 – Very plausible: The relationship has strong theoretical or intuitive support.
- 5 – Highly plausible: The relationship is very likely and well supported by existing knowledge or logic.

Connectedness. Connectedness captures the degree to which the experiment helps explain, integrate, or bridge the original input experiments. A highly connected experiment provides a meaningful link between the input experiments and clarifies how they may be related, regardless of the direction or sign of the relationship.

Item: To what extent does the experiment help connect or explain the original input experiments?

- 1 – Not connected: Does not meaningfully relate to or explain the input experiments.
- 2 – Weakly connected: Only superficially relates to the input experiments.
- 3 – Moderately connected: Has some relevance, but limited explanatory value.
- 4 – Strongly connected: Clearly helps explain, extend, or relate the input experiments.
- 5 – Highly connected: Provides a strong and meaningful bridge that integrates or explains the input experiments.

F Prompts

This section documents the verbatim prompts used in our pipeline. Placeholders enclosed in curly braces (e.g., {iv}, {dv}) are dynamically substituted at runtime with the corresponding field values from each experimental record.

F.1 Bridge-Experiment Generation Prompt

This prompt asks the model to propose *bridge experiments* that logically connect a target experiment to an existing body of literature, without duplicating any experiment already known to the system.

You are a professor in management and psychology research domain. You mainly use experimental methods. Your task is to bridge the gap between a *target experiment* and the existing body of scientific literature to make the whole theory space more connected and self-contained.

[Target Experiment]:
"How does {IV} impact {DV}?"

Condition	Dimension	n_{pairs}	ICC(1,2)	Within-one agreement	Mean abs. diff	% at scale max
Forced Composition	Feasibility	100	0.45	1.00	0.14	0.89
	Plausibility	100	0.83	0.97	0.42	0.24
	Connectedness	100	0.87	0.97	0.42	0.45
EXATLAS Bridge	Feasibility	100	0.58	1.00	0.13	0.88
	Plausibility	100	0.75	0.97	0.39	0.24
	Connectedness	100	0.69	0.98	0.49	0.45

Table 4: Agreement statistics for Forced Composition and EXATLAS Bridge across three evaluation dimensions. ICC(1,2) measures inter-rater reliability, within-one agreement reports percentage within one point, mean absolute difference reports average score disagreement, and % at scale max reports percentage at the maximum scale value.

[Context]:

The target experiment currently lacks strong connections to established research or prior experimental evidence. Thus, it cannot be logically derived from the collection of existing body of literature. Relevant existing body of literature (separated by '; ') is listed below:
 "{literature}"

[Task]:

Propose bridge experiments that logically connect the target experiment to the existing body of literature, such that the derived synergy of prior evidence and the newly proposed experiments form a coherent and cumulative pathway toward linking, justifying and motivating the target experiment.

[Requirements for each proposed experiment]:

1. Try to state the proposed experiment in the **simple form of independent/dependent variables and their relations** (e.g., **Variable A positively/negatively impacts variable B**).
2. You should try to **avoid** introducing complex conditioning, mediation, or moderation relations between dependent and independent variables unless doing so is necessary to establish a meaningful link between experiments.
3. **Avoid** introducing theories/experiments that directly compete with or contradict the existing evidence or the derived synergy of existing evidence.
4. Ensure the proposed experiments are **logically coherent and empirically testable**.
5. Embed **concrete details** into each proposed experiment.
6. **Be concise and creative**.
7. You should select the number of proposed experiments that is **the most appropriate** for the connection to the space of the established experiments.

You should **avoid** proposing experiments that duplicate any of those in the given list below:

"{listofknown}"

Return only the proposed experiments.

If more than one experiment is proposed, **separate them with ";"**

F.2 Conflict-Reconciliation Prompts

A reconciliation prompt is used to probe whether a synthetic experimental finding contradicts prior literature, and, if so, to elicit a unified meta/competing-theoretical account.

I ran an experiment on humans and found:
 {Predicted result based on composition}

Prior literature suggests:
 {Results of contributing experiments}

Answer the following:

- Q1 - Consistency check:

Do my findings contradict the prior findings or the hypotheses implied by them? Answer Yes or No.

- Q2: If Yes, reconcile the conflict.

Source of discrepancy: Identify concrete reasons for the conflict (e.g., differences in context, moderator conditions, populations, operationalization of variables, model assumptions, or distinct variants within a broader meta-theory).

Unified explanation: Propose a single, coherent reconciliation (e.g., a meta theory) that can accommodate both my findings and the prior literature. Be explicit about the mechanisms or conditions under which each result holds.

Be creative. Skip Q2 if the answer to Q1 is No.