

AgentEconomist: An End-to-End Agentic System for Translating Economic Intuitions into Executable Computational Experiments

Anonymous ACL submission

Abstract

A long-standing challenge in economics lies not in the lack of intuition, but in the difficulty of translating intuitive insights into verifiable research. To address this challenge, we introduce **AgentEconomist**, an end-to-end interactive system designed to translate abstract intuitions into executable computational experiments. Grounded in a domain-specific knowledge base covering over 8,700 high-quality academic papers, the system employs a multi-agent architecture. Specifically, an Idea Development Agent and an Experimental Design Agent collaborate to formulate theoretically grounded hypotheses and draft experimental protocols. An Experimental Execution Agent then conducts the designed experiments, forming a closed-loop workflow that supports the rapid development, validation, and iterative refinement of economic intuitions. Through extensive experiments involving human expert evaluation and large language models (LLMs) as judges, we show that the system generates research ideas with stronger literature grounding and higher novelty and insight than state-of-the-art generic LLMs. Overall, AgentEconomist adopts a human-AI collaboration paradigm that enables researchers to focus on high-level intuitions, while delegating the labor-intensive processes of translation and computational execution to agents¹.

1 Introduction

Scientific research often begins with a spark of intuition, a pre-formal, qualitative understanding of how underlying mechanisms in the world function, which serves as the starting point of any formal ideas and hypotheses (Polanyi, 2009; Popper, 2005). In economics, such intuitions typically concern how agents respond to incentives (Becker, 1976), how institutions shape behavior (North,

¹<https://anonymous.4open.science/r/AgentEconomist-CB75/>

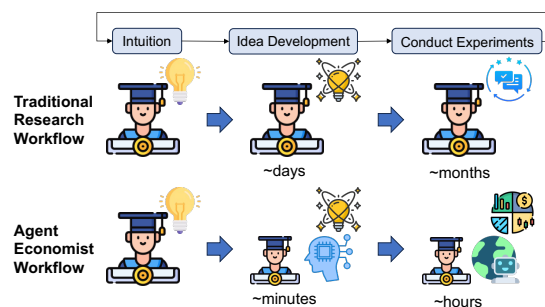


Figure 1: Workflow Comparison. Traditional research (top) suffers from long iteration cycles. In contrast, AgentEconomist (bottom) bridges intuition and execution, enabling rapid hypothesis verification via automated simulation.

1990), and how these interactions give rise to macro-level outcomes (Schelling, 2006). To move beyond mere conjecture, economists traditionally rely on theoretical derivation (Samuelson, 1948) or field experiments (Harrison and List, 2004; Dufo and Banerjee, 2011). However, as illustrated in Figure 1, this traditional workflow is inherently slow and resource-intensive, thereby limiting the scope of ideas a researcher can explore. Computational experiments, based on agent-based modeling (ABM), are a well-established methodology in economics and finance for simulation-based analysis (Tsfatsion, 2006; LeBaron, 2006; Farmer and Foley, 2009) and are frequently employed in teaching to convey complex economic mechanisms (Tisue et al., 2004; Epstein and Axtell, 1996; Railsback and Grimm, 2019). Yet, while such simulation environments provide a powerful medium for verifying economic insights, their inherent complexity often hinders their practical feasibility in research workflows (Axtell and Farmer, 2025).

This barrier between intuition and computational experiments manifests in three distinct dimensions in traditional modeling workflows. First, the process of idea development is often implicit. Novices struggle to systematize their intuition into con-

crete modeling choices because key methodological knowledge—such as selecting plausible mechanisms or assumptions—remains “tacit knowledge” locked in the minds of experts, rarely explicit in textbooks. Second, turning an idea into a computational experiment is itself challenging. Even when a researcher possesses a clear idea, translating it into an executable simulation experiments involves navigating complex codebases, creating a friction that slows down hypothesis testing. Third, and perhaps most critically, the research process lacks systematic experience accumulation. Scientific discovery is an iterative path of refinement, yet the rationale behind specific iterations, why a parameter was tweaked or a hypothesis rejected, is poorly organized alongside the simulation experiments. This loss of *epistemic context* forces researchers to rely on fragmented memory rather than a structured history of inquiry, making it difficult to learn from past failures or incrementally refine ideas.

Recent advances in leveraging LLM agents for autonomous knowledge discovery have attempted to streamline scientific workflows, yet they reveal substantial limitations when assisting researchers in such complex domains. First, existing efforts often target isolated stages, e.g., optimizing either hypothesis generation (Gottweis et al., 2025) or code execution (Swanson et al., 2025), rather than supporting the seamless workflow required to link literature to simulation. Second, holistic frameworks like *The AI Scientist* (Lu et al., 2024) tend to prioritize “outcome automation” (e.g., generating a final manuscript) over “process support”. By treating research as a black box to be automated, they neglect the critical need for *sense-making*: helping researchers structure their thinking and understand the “why” behind modeling decisions. Finally, current systems largely pursue generic applicability at the expense of domain depth, lacking the specific grounding required for economics, where research is deeply rooted in theoretical frameworks and institutional constraints.

To address these challenges, we focus on economics as a representative domain and introduce **AgentEconomist**, an end-to-end interactive research copilot for economic simulation. Built on top of AgentEconomy, a comprehensive agent-based economic simulator, AgentEconomist is designed to support the entire intuition-to-experiment workflow rather than replacing the researcher. Crucially, the system grounds idea development in a large-scale corpus of economic literature, lever-

aging over 8,700 academic papers from top-tier journals to make tacit theoretical knowledge explicit and accessible. The workflow is decomposed into three specialized roles. An *Idea Development Agent* supports literature-grounded sense-making by retrieving relevant economic studies and synthesizing mechanisms, assumptions, and variables. An *Experimental Design Agent* formalizes these ideas into testable hypotheses and executable experiment specifications. An *Experimental Execution Agent* operationalizes the design by running simulations and returning structured results. Additionally, A structured memory module preserves theoretical context, experimental decisions, and outcomes across iterations, enabling cumulative and context-aware reasoning.

Evaluating such open-ended research assistants presents a unique challenge, as standard benchmarks cannot capture the nuance of scientific reasoning. Therefore, we design a rigorous **mixed-methods evaluation protocol**. First, to assess the quality of Idea Generation, we employ a dual-evaluation framework where both advanced LLMs and human experts score generated hypotheses across eight distinct dimensions (e.g., economic soundness, novelty, and feasibility). Second, to validate the system’s utility in real-world workflows, we conduct a holistic user study. Through questionnaires and semi-structured interviews, we gather qualitative feedback on the researchers’ experience, focusing on how the system affects their cognitive load, sense of agency, and overall research efficiency. Across both evaluation protocols, our results show that AgentEconomist consistently outperforms strong baselines on the core dimensions of hypothesis novelty and literature grounding, while also producing hypotheses that are more readily operationalizable in simulation-based settings. The primary contributions of this work include:

- **A human-AI collaboration workflow for bridging the intuition-execution gap.** We conceptualize the translation of economic intuition into computational experiments as a collaborative process. By explicitly modeling the transition from implicit sense-making to operational execution, we provide a structured workflow that lowers the technical barriers for economic modeling.
- **An end-to-end system architecture for grounded experimentation.** We introduce *AgentEconomist*, a unified system that integrates

170 a retrieval-augmented knowledge base (8,700+
171 papers) with specialized agents. This design
172 ensures that generated hypotheses are theoretic-
173 ally sound and that designed experiments are
174 executable, forming a closed loop for iterative
175 discovery.

- 176 • **Empirical validation via mixed-methods as-**
177 **essment.** We conduct quantitative evalua-
178 tions across 8 quality dimensions and qualita-
179 tive user interviews, demonstrating that Agent-
180 Economist effectively produces higher-quality
181 research ideas and significantly streamlines the
182 research workflow, empowering users to explore
183 complex scenarios previously inaccessible to
184 novices.

185 2 Related Works

186 2.1 Task-Specific Scientific Assistants

187 Recent research has focused on augmenting spe-
188 cific stages of the research lifecycle. In the realm of
189 data management and deep research (Huang et al.,
190 2025), specific systems like *ChatPD* (Xu et al.,
191 2025) and *SciSciGPT* (Shao et al., 2025b) have
192 emerged to automate dataset discovery and litera-
193 ture analysis, though benchmarks like *DATASE-*
194 *TRESEARCH* (Li et al., 2025) reveal that cur-
195 rent agents still struggle significantly with out-of-
196 distribution dataset demands. For ideation, meth-
197 ods have evolved from iterative prompting (Got-
198 tweis et al., 2025) to multi-agent frameworks like
199 *VIRSCI* (Su et al., 2025), which simulate team col-
200 laboration to enhance idea novelty. In the evalu-
201 ation phase, approaches vary by technical archi-
202 tecture: *ReviewRL* (Zeng et al., 2025b) employs
203 reinforcement learning to optimize feedback qual-
204 ity, while *ReviewAgents* (Gao et al., 2025b) and
205 *DeepReview* (Zhu et al., 2025) utilize multi-agent
206 collaboration and chain-of-thought reasoning to
207 bridge the gap with human peer reviews. Crucially,
208 however, assessments like *IdeaBench* (Guo et al.,
209 2025) expose a paradox: while LLMs generate
210 highly novel ideas, they often lack practical feasibil-
211 ity. This highlights a limitation of such “open-loop”
212 tools: they produce concepts without the agency to
213 validate them. In contrast, AgentEconomist bridges
214 this gap by coupling literature-grounded ideation
215 with an execution toolbox, ensuring abstract intu-
216 itions are operationalized into testable simulations.

217 2.2 Autonomous Research Systems

218 A parallel stream of work aims to replicate the
219 full role of a human researcher through end-to-
220 end automation (Hu et al., 2025). Building on the
221 code-generation capabilities of *The AI Scientist* (Lu
222 et al., 2024), recent systems have achieved signifi-
223 cant breakthroughs: *CycleResearcher* (Weng et al.,
224 2024) introduced a closed-loop “research-review-
225 revise” mechanism, while *DeepScientist* (Weng
226 et al., 2025) incorporated Bayesian optimization for
227 long-term discovery. To support these autonomous
228 agents, infrastructure like *ToolUniverse* (Gao et al.,
229 2025a) has emerged to standardize tool usage. Fur-
230 thermore, recognizing the complexity of scientific
231 inquiry, frameworks like *OmniScientist* (Shao et al.,
232 2025a) and *MirrorMind* (Zeng et al., 2025a) have
233 expanded this vision by integrating diverse agents
234 to model professional research workflows and col-
235 laborative reasoning. However, critical gaps re-
236 main when applying these generalist systems to
237 economic inquiry. Existing models often function
238 as black boxes that prioritize outcome automation
239 over process support, denying novices the oppor-
240 tunity for sense-making. Moreover, they lack the
241 deep domain grounding required to map theoretical
242 constructs to high-dimensional simulation param-
243 eters. To address this, AgentEconomist positions
244 itself as an interactive co-pilot grounded in a spe-
245 cialized knowledge base of over 8,700 academic
246 papers. By coupling domain depth with rigorous
247 execution, it ensures that scientific discovery re-
248 mains a transparent, human-aligned process.

249 3 Method

250 3.1 Design Rationale: Human-Agent 251 Complementarity

252 Economic research involves heterogeneous cogni-
253 tive demands. Human researchers excel at forming
254 high-level intuitions, exercising normative judg-
255 ment, and deciding when an explanation or result
256 is sufficient. In contrast, transforming these intu-
257 itions into executable experiments requires system-
258 atic literature grounding, formal specification, and
259 reliable execution—tasks that are labor-intensive
260 and error-prone when performed manually.

261 AgentEconomist is designed around this comple-
262 mentarity. The system assigns abstraction, judg-
263 ment, and goal-setting to the human researcher,
264 while delegating literature-grounded sense-making,
265 experiment formalization, and execution manage-
266 ment to automated components. This separation

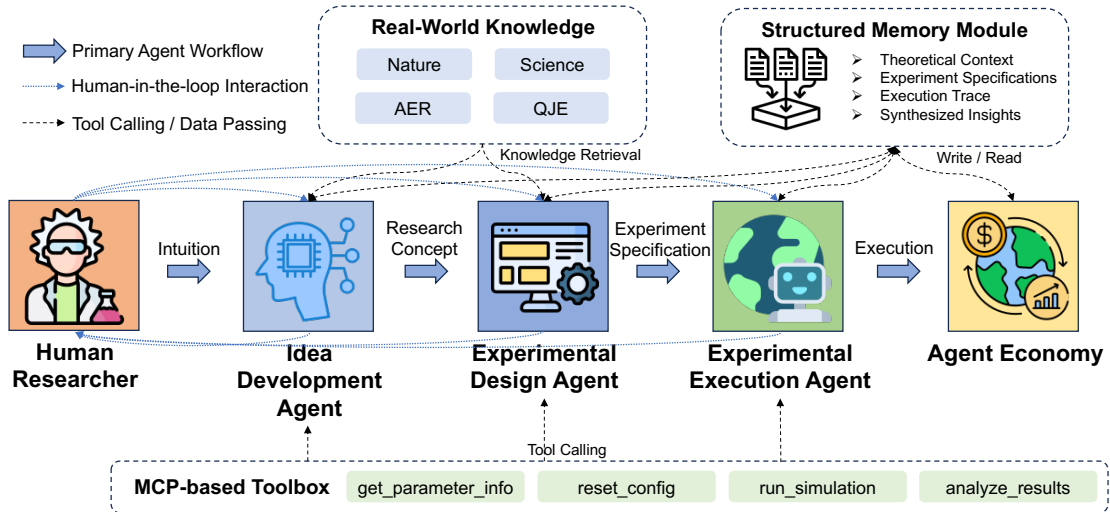


Figure 2: Overview of the AgentEconomist framework. The system decomposes the intuition-to-experiment workflow into literature-grounded ideation, experimental formalization, and execution. These stages are implemented through specialized agents that coordinate via a structured memory module and access the simulator through an MCP-based toolbox, with the human researcher in the loop.

allows each part of the workflow to be handled by the entity best suited for it, while preserving human control over research direction and interpretation.

3.2 The AgentEconomist Framework

Building on the design rationale above, we introduce the AgentEconomist framework shown in Figure 2. The framework operationalizes the intuition-to-experiment workflow through a modular agent architecture, supported by a shared infrastructure layer.

3.2.1 Agent Core

The cognitive workflow is decomposed into three specialized agents corresponding to key stages of economic research:

- **Idea Development Agent.** Responsible for literature-grounded sense-making. It translates human intuition into structured research concepts and candidate hypotheses by retrieving relevant economic literature via retrieval-augmented generation (RAG) and synthesizing key mechanisms and variables (Fig. 3b–c). Beyond initial ideation, this agent also supports iterative verification by checking whether emerging simulation results are consistent with the proposed hypotheses and by guiding hypothesis refinement across iterations (Fig. 3f).
- **Experimental Design Agent.** Responsible for formalization. Given a research concept, it translates abstract hypotheses into executable experi-

mental plans by specifying concrete simulation parameters, defining the number and structure of experiments (e.g., control–treatment or multi-condition setups), and selecting evaluation metrics that align with the simulation environment (Fig. 3d).

- **Experimental Execution Agent.** Responsible for operationalization. It executes experiments in AgentEconomy via the MCP-based toolbox, manages execution states, and returns structured results for downstream analysis and interpretation (Fig. 3e).

3.2.2 Agent Infrastructure

The agent core is supported by three foundational infrastructure components that ensure grounding, continuity, and reliable execution:

- **Real-World Knowledge Base.** To support literature-grounded reasoning, we construct a vector database containing over 8,700 academic papers from top-tier economics and interdisciplinary journals. Retrieved literature supports both idea development and experimental design by making relevant mechanisms, assumptions, and empirical precedents explicit.
- **Structured Memory Module.** The system maintains a persistent memory that records theoretical context, experiment specifications, execution traces, and synthesized outcomes. This enables

324	context-preserving iteration and cumulative reasoning across multiple research cycles.	371
325		372
326	• MCP-based Toolbox. AgentEconomist interacts with the simulator through a standardized toolbox following the Model Context Protocol (MCP). The toolbox abstracts low-level simulator APIs into semantic actions for parameter inspection, configuration, execution, and analysis, ensuring reliable and reproducible operation.	373
327		374
328		375
329		376
330		377
331		378
332		379
333		380
334	3.3 Simulation Substrate: AgentEconomy	381
335	To support the execution and evaluation of experiments specified by AgentEconomist, we develop AgentEconomy as a computational laboratory for rapid hypothesis verification. Unlike toy models that isolate specific sectors, AgentEconomy provides a comprehensive and flexible agent-based environment designed to capture the ripple effects of economic interventions across interacting components.	382
336		383
337		384
338		385
339		386
340		387
341		388
342		389
343	Comprehensive Economic Ecosystem. The environment models a closed-loop economic system comprising four core entity types: households, firms, a government, and a bank. These entities interact through two explicitly modeled markets grounded in real-world data:	390
344		391
345		392
346		393
347		394
348		395
349	• Labor Market: Matches household skills to firm job requirements, allowing wages and employment to emerge endogenously from supply–demand dynamics.	396
350		397
351		398
352		399
353		400
354		401
355		402
356		403
357		404
358		405
359		406
360		407
361		408
362		409
363		410
364		411
365		412
366		413
367		414
368		415
369		416
370		417
		418
		419
		420
		421
		422
		423
		424
		425
		426
		427
		428
		429
		430
		431
		432
		433
		434
		435
		436
		437
		438
		439
		440
		441
		442
		443
		444
		445
		446
		447
		448
		449
		450
		451
		452
		453
		454
		455
		456
		457
		458
		459
		460
		461
		462
		463
		464
		465
		466
		467
		468
		469
		470
		471
		472
		473
		474
		475
		476
		477
		478
		479
		480
		481
		482
		483
		484
		485
		486
		487
		488
		489
		490
		491
		492
		493
		494
		495
		496
		497
		498
		499
		500
		501
		502
		503
		504
		505
		506
		507
		508
		509
		510
		511
		512
		513
		514
		515
		516
		517
		518
		519
		520
		521
		522
		523
		524
		525
		526
		527
		528
		529
		530
		531
		532
		533
		534
		535
		536
		537
		538
		539
		540
		541
		542
		543
		544
		545
		546
		547
		548
		549
		550
		551
		552
		553
		554
		555
		556
		557
		558
		559
		560
		561
		562
		563
		564
		565
		566
		567
		568
		569
		570
		571
		572
		573
		574
		575
		576
		577
		578
		579
		580
		581
		582
		583
		584
		585
		586
		587
		588
		589
		590
		591
		592
		593
		594
		595
		596
		597
		598
		599
		600
		601
		602
		603
		604
		605
		606
		607
		608
		609
		610
		611
		612
		613
		614
		615
		616
		617
		618
		619
		620
		621
		622
		623
		624
		625
		626
		627
		628
		629
		630
		631
		632
		633
		634
		635
		636
		637
		638
		639
		640
		641
		642
		643
		644
		645
		646
		647
		648
		649
		650
		651
		652
		653
		654
		655
		656
		657
		658
		659
		660
		661
		662
		663
		664
		665
		666
		667
		668
		669
		670
		671
		672
		673
		674
		675
		676
		677
		678
		679
		680
		681
		682
		683
		684
		685
		686
		687
		688
		689
		690
		691
		692
		693
		694
		695
		696
		697
		698
		699
		700
		701
		702
		703
		704
		705
		706
		707
		708
		709
		710
		711
		712
		713
		714
		715
		716
		717
		718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829

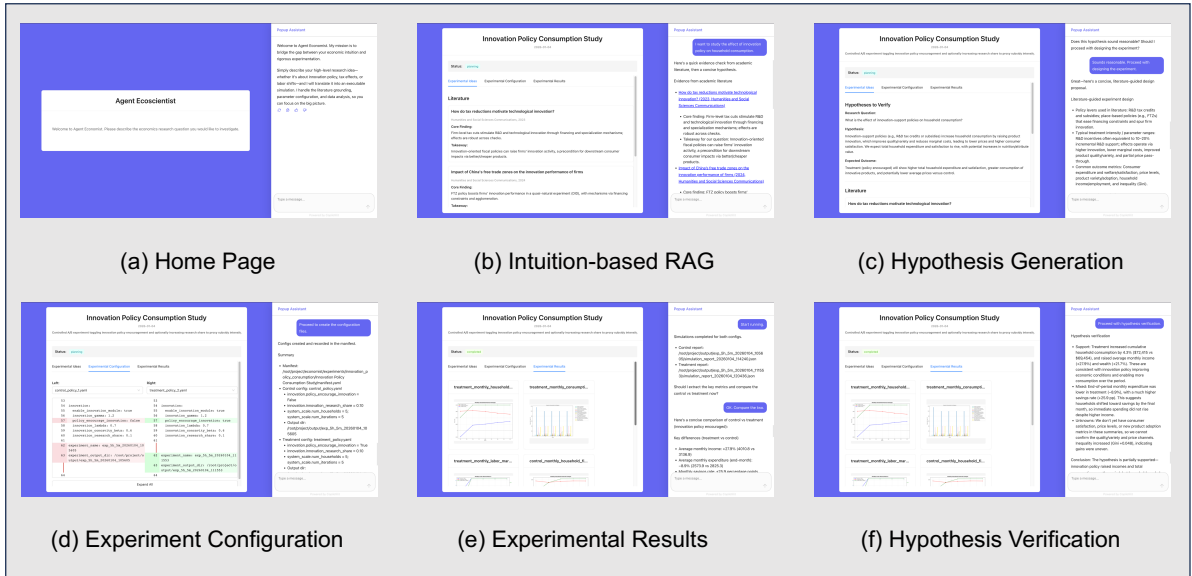


Figure 3: The interactive research workflow of AgentEconomist, illustrating the end-to-end process from intuition input and literature-grounded ideation to experiment configuration, execution, and hypothesis verification.

Accordingly, we adopt a mixed-methods evaluation design that combines controlled comparisons against strong LLM baselines with human-in-the-loop assessments based on real user interactions.

We first describe the experimental setup and data collection procedure (§4.1), then report quantitative results on hypothesis quality across multiple economic dimensions (§4.2), and finally present a qualitative analysis of user feedback to characterize perceived trust, workflow support, and usability limitations (§4.3).

4.1 Experimental Setup

Systems and Participants. We evaluate AgentEconomist against strong general-purpose LLM baselines commonly used as standalone research assistants. AgentEconomist is instantiated with GPT-5 as its backbone model and augmented with domain-specific retrieval, structured memory, and simulator-oriented tooling, as described in §3. Baselines include advanced LLMs such as GPT-5.2 and Gemini 3, used in a vanilla conversational setting without explicit memory or execution support.

We recruited 15 participants with backgrounds in economics, public policy, or related social science fields. Each participant interacted with the system by providing an economic intuition of their own choosing. All participants completed a structured questionnaire; however, due to incomplete submissions, only 14 participants provided full and usable interaction logs for hypothesis-level comparisons.

4.2 Idea Quality and Grounding (RQ1 & RQ2)

This section evaluates whether AgentEconomist improves the quality of generated hypotheses and their grounding in economic literature.

Evaluation Protocol. To isolate hypothesis-level quality from downstream execution variability, we conduct controlled comparisons at the hypothesis generation stage. Both AgentEconomist and baseline LLMs are prompted with identical user intuitions. Generated hypotheses are evaluated along eight dimensions capturing economic rigor, insight, and implementability, using both an LLM-based anonymous referee and human judgments.

Results. Figure 4 summarizes the evaluation results. Across both evaluation protocols, AgentEconomist shows clear advantages on the dimensions most central to the intuition-to-hypothesis translation task. In particular, Literature Grounding and Novelty & Insight exhibit the largest improvements. Under LLM-based judging, Literature Grounding increases from 3.36 to 4.93 and Novelty & Insight from 3.00 to 4.43. Human evaluations show consistent gains in the same dimensions, with Literature Grounding improving from 3.11 to 4.50 and Novelty & Insight from 3.12 to 4.05.

These results are especially significant given our evaluation focus. The primary challenge in economic research assistance is not producing fluent text, but translating vague intuitions into hypothe-

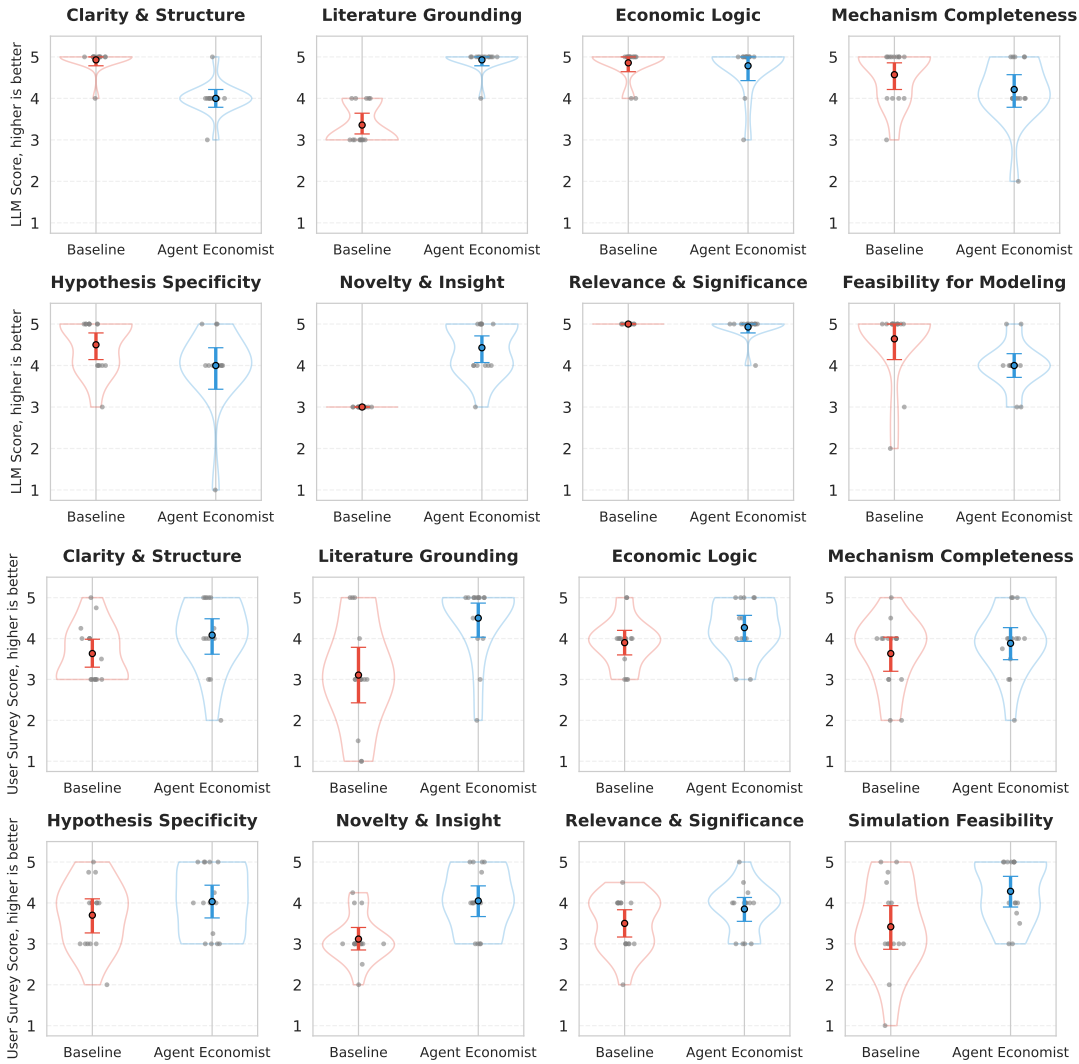


Figure 4: Distribution of hypothesis quality scores across eight evaluation dimensions. **Top:** LLM-based anonymous referee evaluation. **Bottom:** Human evaluation by study participants. Each subplot corresponds to one dimension. Across both judging protocols, AgentEconomist consistently outperforms the baseline, with the largest margins observed in *Literature Grounding* and *Novelty & Insight*.

478 ses that are both theoretically grounded and intel-
 479 lectually non-trivial. On these core criteria, Agen-
 480 tEconomist achieves substantial and consistent im-
 481 provements over strong baseline models.

482 **Human-LLM Evaluation Differences.** A sys-
 483 tematic divergence arises in the Clarity & Structure
 484 dimension. While the LLM-based referee slightly
 485 favors the baseline, human evaluators rate AgentE-
 486 conomist comparably or higher. We attribute this
 487 difference to distinct evaluation priorities. Baseline
 488 outputs tend to be shorter and rhetorically stream-
 489 lined, whereas AgentEconomist produces denser
 490 hypotheses that explicitly articulate mechanisms,
 491 assumptions, and literature connections. Although
 492 this additional structure may reduce surface-level
 493 clarity for automated judges, it aligns more closely

494 with human researchers’ needs when developing
 495 grounded and novel economic hypotheses.

4.3 Qualitative User Experience Analysis (RQ3) 496 497

498 To examine how users perceive AgentEconomist as
 499 an end-to-end research framework, we conduct a
 500 formative qualitative analysis of open-ended ques-
 501 tionnaire responses. The questionnaire consists
 502 of five prompts: Q1 (Perceived Advantages), Q2
 503 (Trust and Credibility), Q3 (Pain Points and Limita-
 504 tions), Q4 (Role Perception), and Q5 (Other Open
 505 Feedback).

506 Not all participants answered every prompt (Q1:
 507 $n = 8$, Q2-Q4: $n = 7$, Q5: $n = 3$). Accord-
 508 ingly, our analysis focuses on recurring themes
 509 supported by multiple responses rather than esti-

510 mating population-level prevalence.

511 **Method.** We apply an LLM-assisted grounded-
512 theory thematic analysis to aggregate responses
513 and extract emergent themes. Only patterns di-
514 rectly supported by verbatim evidence are reported,
515 and findings are interpreted as experience-level per-
516 ceptions rather than objective performance claims.

517 **Findings.** Participants consistently emphasized
518 grounded trust, attributing increased confidence to
519 literature-backed reasoning enabled by the RAG
520 module. A second dominant theme is operational-
521 ization support: users valued the system’s ability
522 to translate vague intuitions into simulator-aligned
523 experimental configurations. Third, respondents
524 highlighted mechanistic scaffolding, noting that
525 AgentEconomist more explicitly articulated causal
526 and behavioral chains than generic LLM outputs.
527 Together, these factors contributed to a perceived
528 role shift from a conversational assistant toward a
529 research assistant or collaborative expert spanning
530 hypothesis formulation and experiment setup.

531 **Pain Points and Implications.** Feedback also re-
532 vealed usability bottlenecks, including execution la-
533 tency, limited process transparency, and occasional
534 instruction-following drift over long interactions.
535 These observations suggest that future work should
536 prioritize improving responsiveness, exposing in-
537 terpretable intermediate states, and strengthening
538 long-horizon intent preservation.

539 4.4 Case Study

540 We illustrate the end-to-end use of AgentEconomist
541 with a representative interaction drawn from a real
542 user session. The user started from a high-level
543 intuition: *whether innovation-support policies in-
544 crease household consumption*. Given this un-
545 derspecified query, the system first retrieved and
546 summarized relevant economic literature, includ-
547 ing studies on R&D tax incentives, place-based
548 innovation policies, and production–consumption
549 linkages, grounding the inquiry in established evi-
550 dence.

551 Conditioned on this context, AgentEconomist
552 formulated a concrete hypothesis that innovation
553 policies raise household consumption through
554 firm innovation, price effects, and income chan-
555 nels. The system then designed a controlled
556 A/B experiment by toggling a policy parameter
557 (`innovation.policy_encourage_innovation`),
558 holding market structure constant, and specify-

559 ing treatment and control configurations. Key
560 parameters—including the number of households,
561 simulation horizon, and innovation research
562 share—were explicitly exposed and iteratively
563 adjusted by the user.

564 After executing the simulations, the system auto-
565 matically extracted and compared outcome metrics.
566 The treatment condition exhibited higher cumula-
567 tive household consumption (+4.3%), substantially
568 higher income (+27.9%) and wealth (+21.7%),
569 alongside a higher savings rate and slightly in-
570 creased inequality. AgentEconomist synthesized
571 these results into a structured verification report,
572 concluding that the hypothesis was partially sup-
573 ported and recommending further analysis.

574 The interaction terminated once the user was
575 satisfied with the explanation and next-step sugges-
576 tions. This case demonstrates that AgentEconomist
577 supports the full research workflow—from intu-
578 ition, to literature grounding, hypothesis genera-
579 tion, and experimental design and iteration—within
580 a unified and context-aware framework.

581 5 Conclusion

582 We presented AgentEconomist, an end-to-end agen-
583 tic system that supports economic research by trans-
584 lating abstract intuitions into executable computa-
585 tional experiments. The system adopts a human-
586 in-the-loop design that decomposes the intuition-
587 to-experiment workflow into literature-grounded
588 idea development, experiment formalization, and
589 execution, implemented by specialized agents. Em-
590 pirical evaluation shows that AgentEconomist im-
591 proves hypothesis quality on the dimensions most
592 critical to economic inquiry, particularly literature
593 grounding and novelty of insight, while remaining
594 effective in interactive research settings. Through-
595 out this process, the system remains controllable
596 executions and iterations explicitly mediated by
597 the user, mitigating potential risks associated with
598 autonomous behavior. More broadly, this work
599 highlights the importance of epistemic scaffolding
600 in scientific AI systems, where supporting theory-
601 grounded reasoning and iterative sense-making is
602 often more valuable than end-to-end automation.
603 We hope this perspective informs future work on
604 domain-grounded, interactive agents for scientific
605 discovery.

606 Limitations

607 This work has several limitations. First, our evaluation
608 focuses on the intuition-to-experiment workflow in an agent-based economic simulation, and
609 does not assess performance on real-world policy deployment or empirical data analysis. Second,
610 while the user study provides evidence of the system’s usefulness in interactive research settings,
611 the number of participants is limited and may not capture the full diversity of economic research practices.
612 Third, the effectiveness of AgentEconomist depends on the coverage and quality of the underlying literature corpus and simulation environment;
613 domains or research questions that fall outside these resources may be less well supported.
614 Finally, although the system maintains structured memory across interactions, long-horizon alignment and execution efficiency remain constrained
615 by current LLM capabilities and system latency. Future work may address these limitations by expanding empirical evaluation settings, incorporating
616 larger and more diverse user populations, and improving system robustness and scalability.
617
618
619
620
621
622
623
624
625
626
627
628

629 References

630 Robert L Axtell and J Doyne Farmer. 2025. Agent-based modeling in economics and finance: Past, present, and future. *Journal of Economic Literature*, 63(1):197–287.

631
632
633

634 Gary S Becker. 1976. *The economic approach to human behavior*, volume 803. University of Chicago press.

635

636 Esther Duflo and Abhijit Banerjee. 2011. *Poor economics*, volume 619. PublicAffairs New York.

637

638 Joshua M Epstein and Robert Axtell. 1996. *Growing artificial societies: social science from the bottom up*. Brookings Institution Press.

639
640

641 J Doyne Farmer and Duncan Foley. 2009. The economy needs agent-based modelling. *Nature*, 460(7256):685–686.

642
643

644 Shanghua Gao, Richard Zhu, Pengwei Sui, Zhenglun Kong, Sufian Aldogom, Yepeng Huang, Ayush Noori, Reza Shamji, Krishna Parvataneni, Theodoros Tsiligkaridis, and 1 others. 2025a. Democratizing ai scientists using tooluniverse. *arXiv preprint arXiv:2509.23426*.

645
646
647
648
649

650 Xian Gao, Jiacheng Ruan, Zongyun Zhang, Jingsheng Gao, Ting Liu, and Yuzhuo Fu. 2025b. Reviewagents: Bridging the gap between human and ai-generated paper reviews. *arXiv preprint arXiv:2503.08506*.

651
652
653

Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, and 1 others. 2025. Towards an ai co-scientist. *arXiv preprint arXiv:2502.18864*. 654
655
656
657
658

Sikun Guo, Amir Hassan Shariatmadari, Guangzhi Xiong, Albert Huang, Myles Kim, Corey M Williams, Stefan Bekiranov, and Aidong Zhang. 2025. Ideabench: Benchmarking large language models for research idea generation. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pages 5888–5899. 659
660
661
662
663
664
665

Glenn W Harrison and John A List. 2004. Field experiments. *Journal of Economic literature*, 42(4):1009–1055. 666
667
668

Ming Hu, Chenglong Ma, Wei Li, Wanghan Xu, Jiamin Wu, Jucheng Hu, Tianbin Li, Guohang Zhuang, Jiaqi Liu, Yingzhou Lu, and 1 others. 2025. A survey of scientific large language models: From data foundations to agent frontiers. *arXiv preprint arXiv:2508.21148*. 669
670
671
672
673
674

Yuxuan Huang, Yihang Chen, Haozheng Zhang, Kang Li, Huichi Zhou, Meng Fang, Linyi Yang, Xiaoguang Li, Lifeng Shang, Songcen Xu, and 1 others. 2025. Deep research agents: A systematic examination and roadmap. *arXiv preprint arXiv:2506.18096*. 675
676
677
678
679

Blake LeBaron. 2006. Agent-based computational finance. *Handbook of computational economics*, 2:1187–1233. 680
681
682

Keyu Li, Mohan Jiang, Dayuan Fu, Yunze Wu, Xiangkun Hu, Dequan Wang, and Pengfei Liu. 2025. Datasetresearch: Benchmarking agent systems for demand-driven dataset discovery. *arXiv preprint arXiv:2508.06960*. 683
684
685
686
687

Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. 2024. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*. 688
689
690
691

Douglass C North. 1990. *Institutions, institutional change and economic performance*. Cambridge university press. 692
693
694

Michael Polanyi. 2009. The tacit dimension. In *Knowledge in organisations*, pages 135–146. Routledge. 695
696

Karl Popper. 2005. *The logic of scientific discovery*. Routledge. 697
698

Steven F Railsback and Volker Grimm. 2019. *Agent-based and individual-based modeling: a practical introduction*. Princeton university press. 699
700
701

Paul Anthony Samuelson. 1948. Foundations of economic analysis. *Science and Society*, 13(1). 702
703

Thomas C Schelling. 2006. *Micromotives and macrobehavior*. WW Norton & Company. 704
705

706	Chenyang Shao, Dehao Huang, Yu Li, Keyu Zhao, Weiquan Lin, Yining Zhang, Qingbin Zeng, Zhiyu Chen, Tianxing Li, Yifei Huang, and 1 others. 2025a. Omniscientist: Toward a co-evolving ecosystem of human and ai scientists. <i>arXiv preprint arXiv:2511.16931</i> .	Minjun Zhu, Yixuan Weng, Linyi Yang, and Yue Zhang. 2025. Deepreview: Improving llm-based paper review with human-like deep thinking process. <i>arXiv preprint arXiv:2503.08569</i> .	762
707			763
708			764
709			765
710			
711	Erzhuo Shao, Yifang Wang, Yifan Qian, Zhenyu Pan, Han Liu, and Dashun Wang. 2025b. Sciscigt: advancing human-ai collaboration in the science of science. <i>Nature Computational Science</i> , pages 1–15.	A Experiment Details	766
712		A.1 Hypothesis Quality Dimensions	767
713		Participants were asked to evaluate each generated hypothesis along the following eight dimensions. For each dimension, participants were instructed to provide an independent score based on their own judgment, focusing on the hypothesis content rather than writing style or presentation.	768
714			769
715	Haoyang Su, Renqi Chen, Shixiang Tang, Zhenfei Yin, Xinzhe Zheng, Jinzhe Li, Biqing Qi, Qi Wu, Hui Li, Wanli Ouyang, and 1 others. 2025. Many heads are better than one: Improved scientific idea generation by a llm-based multi-agent system. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 28201–28240.		770
716			771
717			772
718			773
719			774
720			775
721			776
722			777
723	Kyle Swanson, Wesley Wu, Nash L Bulaong, John E Pak, and James Zou. 2025. The virtual lab of ai agents designs new sars-cov-2 nanobodies. <i>Nature</i> , 646(8085):716–723.		778
724			779
725			780
726			781
727	Leigh Tesfatsion. 2006. Agent-based computational economics: A constructive approach to economic theory. <i>Handbook of computational economics</i> , 2:831–880.		782
728			783
729			784
730			785
731	Seth Tisue, Uri Wilensky, and 1 others. 2004. Netlogo: A simple environment for modeling complexity. In <i>International conference on complex systems</i> , volume 21, pages 16–21. Boston, MA.		786
732			787
733			788
734			789
735	Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo Zhang, Jindong Wang, Yue Zhang, and Linyi Yang. 2024. Cyclereviewer: Improving automated research via automated review. <i>arXiv preprint arXiv:2411.00816</i> .		790
736			791
737			792
738			793
739			794
740	Yixuan Weng, Minjun Zhu, Qiuji Xie, Qiyao Sun, Zhen Lin, Sifan Liu, and Yue Zhang. 2025. Deepscientist: Advancing frontier-pushing scientific findings progressively. <i>arXiv preprint arXiv:2509.26603</i> .		795
741			796
742			797
743			798
744	Anjie Xu, Ruiqing Ding, and Leye Wang. 2025. Chatpd: An llm-driven paper-dataset networking system. In <i>Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2</i> , pages 5106–5116.		799
745			800
746			801
747			802
748			803
749	Qingbin Zeng, Bingbing Fan, Zhiyu Chen, Sijian Ren, Zhilun Zhou, Xuhua Zhang, Yuanyi Zhen, Fengli Xu, Yong Li, and Tie-Yan Liu. 2025a. Mirrormind: Empowering omniscientist with the expert perspectives and collective knowledge of human scientists. <i>arXiv preprint arXiv:2511.16997</i> .		804
750			805
751			
752			
753			
754			
755	Sihang Zeng, Kai Tian, Kaiyan Zhang, Yuru Wang, Junqi Gao, Runze Liu, Sa Yang, Jingxuan Li, Xinwei Long, Jiaheng Ma, and 1 others. 2025b. Reviewrl: Towards automated scientific review with rl. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 16942–16954.		
756			
757			
758			
759			
760			
761			

806	A.2 Scoring Protocol	research experience related to economic simulation	851
807	Each dimension is scored on a 5-point Likert scale:	or computational economics, ensuring sufficient	852
808	1 = Very poor, 3 = Acceptable, 5 = Excellent.	domain knowledge to meaningfully evaluate the	853
809	Judges are instructed to base scores solely on the	system.	854
810	hypothesis-generation content and to ignore writ-	Participants were compensated for their time at a	855
811	ing style, verbosity, or downstream analysis.	rate commensurate with their academic background	856
812	A.3 Qualitative User Experience Analysis	and local standards. We consider the payment to	857
813	Protocol	be adequate given the participants' demographic	858
814	This qualitative evaluation is designed as a for-	and level of expertise, and the study involved no	859
815	mative assessment to capture participants' subjec-	deceptive practices or sensitive data collection. All	860
816	tive experiences when interacting with AgentE-	participants provided informed consent prior to par-	861
817	conomist as a research framework. Participants	ticipation and were explicitly informed about how	862
818	are asked to provide free-text responses based on	their interaction data would be used for research	863
819	their own interaction trajectories, without assum-	and evaluation purposes.	864
820	ing complete task coverage or uniform interaction	B Prompts Demonstration	865
821	length.	In this section, we provide the prompts used in	866
822	Evaluation Dimensions. Open-ended questions	our evaluation pipeline. Each {...} placeholder	867
823	are organized around four aspects of user experi-	in the templates will be substituted with the cor-	868
824	ence:	responding study materials (e.g., hypothesis ex-	869
825	• Perceived Advantages: Key differences com-	cerpts, anonymized IDs, and aggregated free-text	870
826	pared to using a general-purpose LLM.	responses). For implementation details and exact	871
827	• Trust and Credibility: Factors influencing confi-	instantiations, please refer to our code release.	872
828	dence in generated hypotheses and experimental	Table 1 presents the LLM-based referee prompt	873
829	designs.	used for hypothesis-quality scoring. Table 2	874
830	• Pain Points and Limitations: Usability issues,	presents the grounded-theory prompt used for	875
831	missing features, or interaction difficulties.	LLM-assisted thematic analysis of open-ended	876
832	• Role Perception: How users conceptualize the	feedback.	877
833	system (e.g., search engine, research assistant,	C AI Assistance Disclosure	878
834	collaborative expert).	AI-based tools were used in the preparation of this	879
835	Response Format. Participants provide short	work for code development assistance and for gram-	880
836	free-text answers (approximately 100 words per	matical and stylistic editing of the manuscript. All	881
837	question). Responses are anonymized and aggre-	scientific content, experimental design, results, and	882
838	gated for analysis.	conclusions were conceived, implemented, and ver-	883
839	Analysis Protocol. Aggregated responses are an-	ified by the authors.	884
840	alyzed using an LLM-assisted grounded-theory		
841	workflow. The analysis extracts emergent themes		
842	through iterative coding and grouping, with final re-		
843	porting restricted to themes supported by multiple		
844	responses and direct textual evidence.		
845	A.4 Participant Background and		
846	Compensation		
847	All participants in the user study were doctoral		
848	students actively working in the area of economic		
849	modeling or simulation-based economic research.		
850	Each participant had at least six months of prior		

Role: Anonymous Economics Referee (LLM-as-a-Judge)

You are acting as an anonymous, neutral, and rigorous economics referee. Your task is to evaluate and compare the quality of **hypothesis-generation content** produced by two systems in response to the same economic question.

[Critical Scope Restriction]

Evaluate **ONLY the generated hypotheses themselves**. Do NOT consider:
- downstream analysis, simulations, policy discussion, extensions, or follow-up content.

[Evaluation Materials]

- Materials are provided as images or PDF pages.
- First {N} pages belong to System A; following {M} pages belong to System B.
- Output order does NOT imply priority. Ignore system identity.

[Evaluation Principles]

- 1) Content-only evaluation: ignore style, verbosity, tone, or rhetorical confidence.
- 2) Order invariance: do not favor earlier or later materials.
- 3) Model blindness: do not infer system identity or sophistication.
- 4) Referee-level standards: judge as if reviewing an economics paper's hypotheses.

[Dimensions: score each 1-5 (integer)]

- (1) Clarity & Structure
- (2) Literature Grounding & Factual Plausibility
- (3) Economic Logic / Soundness
- (4) Mechanism Completeness (at hypothesis level)
- (5) Hypothesis Specificity
- (6) Novelty & Insight
- (7) Relevance & Significance
- (8) Feasibility for Modeling or Simulation

[Required Output Format: Markdown Only]

- A dimension-wise score table comparing System A vs System B with brief justification.
- Overall assessment (200 words).
- Bias & scope compliance check:
 - a) evaluated hypotheses only? (Yes/No)
 - b) ignored identity and order? (Yes/No)
 - c) any dimensions hard to score due to insufficient detail?

Table 1: Prompt template for LLM-based hypothesis-quality judging (anonymous economics referee). Each {...} placeholder is instantiated with the corresponding evaluation materials.

Role: Social Science Researcher (Grounded Theory)

You are a social science researcher familiar with qualitative methods. Please use grounded theory (Grounded Theory) to conduct a systematic thematic analysis on the following interview materials.

[Methodological Constraint]

Remain highly sensitive to the data and avoid pre-set theoretical assumptions. Ensure conclusions are grounded in the interview materials themselves.

[Required Output Structure]

- Core themes
- Theme summary
- Verbatim example sentences

[Interview Materials]

{Paste aggregated anonymized responses here, e.g., P1:, P2:, ...}

Table 2: Prompt template for LLM-assisted grounded-theory thematic analysis of participants' open-ended feedback.