
AER-bench: A Comprehensive Benchmark for Automatic Economic Research with Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce **AER-bench**, a comprehensive benchmark for evaluating language
2 models on automatic economic research. Unlike existing benchmarks that focus
3 on isolated question-answering, AER-bench assesses models across the full re-
4 search pipeline through seven task types: Literature Review, Economic Knowledge,
5 Concept Identification, Paper Abstraction, Economic Analysis, Mathematical Mod-
6 eling, and Model Generation. Built from 191 top journal articles with automatically
7 extracted tasks and 900 manually annotated historical text segments, our bench-
8 mark employs a hybrid evaluation framework combining rule-based scoring and
9 LLM-as-judge assessment across multiple dimensions. We evaluate 16 state-of-the-
10 art language models and reveal three key findings: (1) **pure LLMs show unstable**
11 **performance on knowledge-intensive tasks**, suggesting the need for agentic sys-
12 tems with external knowledge access; (2) **models within the same family exhibit**
13 **consistent performance patterns**, indicating systematic capability differences
14 between model series; (3) **understanding of historical economic concepts is un-**
15 **stable across models**: several leading systems handle early-20th-century concepts
16 on par with recent ones, while weaker models drop 20%+ on early concepts—a
17 gap that emerges almost entirely on long source paragraphs. Our findings highlight
18 critical gaps in language models’ economic research capabilities and provide a
19 rigorous testbed for future development. The benchmark, code, and results are
20 available at GitHub to facilitate reproduction and potential data reconstruction.

21 1 Introduction

22 Economic research is a primary input to policy, market design, and aggregate welfare [9, 7, 22, 13, 1,
23 6, 3], yet producing each contribution still demands a labor-intensive pipeline: literature review, data
24 collection and cleaning, specification of theoretical or econometric models, derivation or estimation,
25 robustness checking, and result interpretation. A substantial share of this pipeline is mechanical and
26 repetitive, and a growing literature has begun to ask how generative AI could absorb it [10].

27 Frontier large language models are now plausible candidates for that role. Web-scale pretraining
28 has given them broad command of economic terminology, canonical models, and econometric
29 techniques [2, 23], and progress on multi-step reasoning, quantitative problem solving, and formal
30 mathematics [25, 11, 8] brings the symbolic and computational manipulations central to economic
31 analysis within reach. Building on foundational paradigms such as ReAct and Toolformer [26, 20],
32 agentic frameworks like OpenClaw [17], Claude Code [4], and Hermes [15] now package these
33 capabilities into reusable, tool-using workflows that browse, execute code, and iterate on intermediate
34 artifacts, while analogous systems have begun to drive scientific and algorithmic discovery end-
35 to-end [12, 16]. A first attempt has already appeared in economics itself: APE (Automatic Policy
36 Evaluation) [21], built on Claude Code, automates empirical policy research, yet on its own self-

37 reported scale still trails top journals by a wide margin—at once demonstrating the promise of
 38 automated economic research and underscoring the need for an independent, rigorous benchmark.

Table 1: Comparison of AER-bench with existing economic scenario benchmarks. AER-bench is the first benchmark designed to evaluate the pipeline of automatic economic research.

Benchmark	Data Source	Focus	Modality	# Items	# Task Types
EconGym [14]	Real-world data	Econ Behavior	Simulation	—	25+
GDPVal [18]	Industry reports	Econ Productivity	Text	1,320	9
EconLogicQA [19]	News articles	Econ Reasoning	Text	650	1
FinMR [5]	Web + Experts	Fin Reasoning	Text + Image	3,700	15
AER-bench (Ours)	Journals + Experts	Econ Research	Text	1,629	7

39 General-purpose academic benchmarks such as MMLU-Pro [24] include an economics subject, but
 40 each subject contributes only a few hundred multiple-choice items and cannot probe the research-
 41 grade capabilities we are interested in. Several economics-focused benchmarks¹ have therefore been
 42 proposed (Table 1). EconGym [14] is a scalable multi-agent simulation testbed in which LLMs play
 43 the role of households, firms, or governments across more than twenty real-data-driven economic
 44 tasks, evaluating decision behavior rather than research output. GDPVal [18] curates 1,320 tasks
 45 across nine GDP-relevant industries from industry reports to measure whether AI can substitute for
 46 economically valuable human labor, again outside the academic pipeline. EconLogicQA [19] distills
 47 650 multi-event sequencing questions from economic news articles, targeting the narrow capability
 48 of causal and temporal reasoning. FinMR [5] provides over 3,700 multimodal questions on financial-
 49 analyst-level chart and formula interpretation, likewise sourced from non-academic web material
 50 rather than economics research. Each benchmark is valuable within its own scope—behavioral
 51 simulation, productive substitution, event reasoning, and financial analysis, respectively—yet none
 52 draws its data from top-journal research articles, and consequently none evaluates the core capabilities
 53 that economic research relies on, such as domain knowledge, mathematical modeling, and economic
 54 analysis of empirical or theoretical results. This gap remains unaddressed.

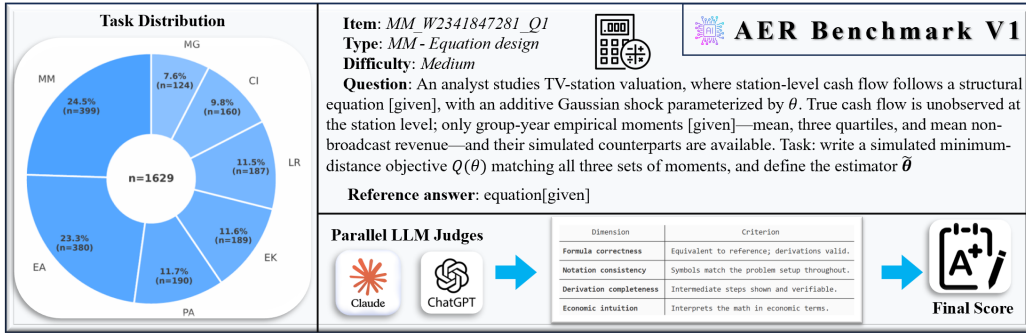


Figure 1: Overview of the AER benchmark. **Left:** a pie chart showing the distribution and proportion of tasks across the seven categories; **middle:** a representative task example; **bottom:** the evaluation pipeline, in which gpt-5.2 and claude-4.6-sonnet independently score each response against a shared rubric, and their judgments are aggregated into a final score.

55 To address this gap, we introduce **AER-bench**², the first comprehensive benchmark for systematically
 56 and objectively evaluating language models on automatic economic research. AER-bench com-
 57 bines two complementary data sources: (i) 191 articles drawn from top-5 economics journals—the
 58 *American Economic Review* and *Econometrica*—representing the modern frontier of the field; and
 59 (ii) roughly 900 text segments from pre-1950 economics literature, which support the analysis of
 60 how language models handle historical economic discourse. Tasks are constructed through two
 61 complementary pipelines (§2.2): an LLM-driven multi-round extraction pipeline that turns paper

¹We discuss only a few representative economics-related benchmarks here; additional benchmarks covering economic and financial scenarios exist but are omitted for brevity.

²An overview of our benchmark is shown in Figure 1, where the representative task example illustrates that our tasks differ fundamentally from those in previous benchmarks.

62 PDFs into self-contained items, and a manual annotation effort by economics doctoral students that
63 distills the historical segments into 160 concept-identification tasks across eight conceptual dimen-
64 sions. The result is a benchmark of 1,629 evaluation items organized into seven task types—EK,
65 MM, PA, LR, MG, EA, and CI—spanning capabilities from foundational knowledge recall and
66 symbolic manipulation to higher-order modeling and analysis (§2.1). To score these heterogeneous
67 tasks reliably, AER-bench adopts a hybrid evaluation framework (§2.3): rule-based scoring for EK
68 and CI, and multi-dimensional rubric-based LLM-as-judge scoring averaged over two independent
69 judges (claude-4.6-sonnet and gpt-5.2) for the remaining five tasks. We use AER-bench to
70 evaluate 16 frontier LLMs (§3); the full ranking and per-task breakdowns are reported in §4.

71 In summary, AER-bench is, to our knowledge, the first comprehensive benchmark specifically
72 designed to evaluate language models on the full pipeline of automatic economic research, combining
73 a workflow-aligned task taxonomy, top-journal and pre-1950 sources, and a hybrid rule-based
74 and dual-judge evaluation framework. Beyond the benchmark itself, the underlying construction
75 methodology—an LLM-driven multi-round extraction pipeline coupled with expert annotation—is
76 in itself a reusable asset for continuously extending the benchmark, constructing training data, and
77 adapting the approach to other research disciplines. Across 16 frontier LLMs we observe pervasive
78 shortcomings, most notably substantial instability on knowledge-intensive tasks, which we interpret
79 as evidence that agentic, retrieval-augmented architectures are a necessary next step toward reliable
80 automated economic research. Along the way, AER-bench surfaces several further phenomena worth
81 attention—including consistent within-family capability patterns and a pronounced historical concept
82 drift (a “modernity bias”) on pre-1950 economic concepts—which we examine in detail in §4.

83 2 Benchmark Design

84 To systematically evaluate the capability of Language Models in automating economic research, we
85 introduce AER-bench. Unlike generic benchmarks (e.g. MMLU) that provide only shallow coverage
86 of economic scenarios, AER-bench is grounded in the authentic research pipeline, assessing models
87 across both deterministic knowledge retrieval and open-ended generative analysis.

88 2.1 Task Taxonomy

89 AER-bench deconstructs the economic research workflow into seven distinct task categories. Each
90 task targets a specific phase of the research lifecycle and is evaluated along fine-grained dimensions:

- 91 • **Literature Review (LR):** Assesses the model’s ability to reconstruct the argument structure
92 of an introduction and appropriately cite prior literature. It is evaluated via a dual mechanism:
93 citation recall (capturing the core reference network) and a paper-specific writing rubric.
- 94 • **Economic Knowledge (EK):** A rigorous QA task comprising single-choice, true/false,
95 fill-in-the-blank, and ordering questions. It tests the model’s deep comprehension of paper-
96 specific methodologies, findings, and institutional contexts.
- 97 • **Concept Identification (CI):** A robust binary classification task requiring models to map
98 text chunks to eight structural economic properties: relevance, mathematical formalization,
99 uniqueness, stability, existence, equilibrium scene, time characteristic, and time horizon.
- 100 • **Paper Abstraction (PA):** The model must synthesize a concise abstract from a condensed
101 paper body. Performance is evaluated across four dimensions: coverage of atomic semantic
102 propositions, factual accuracy, conciseness, and academic style.
- 103 • **Economic Analysis (EA):** Given uninterpreted empirical or theoretical results, the model
104 generates a qualitative analysis. It is scored across six dimensions: accuracy, economic
105 interpretation, mechanism, theory connection, policy implication, and writing quality.
- 106 • **Mathematical Modeling (MM):** Given a textual economic scenario and notation con-
107 straints, the model must perform formal derivations. Outputs are judged on four dimensions:
108 correctness, notation consistency, derivation completeness, and economic intuition.
- 109 • **Model Generation (MG):** Evaluates formal construction of an economic model from a
110 de-symbolized narrative. Scoring uses a granular element rubric checking for presence and
111 correctness of sectors, objective functions, constraints, and equilibrium concepts.

112 **2.2 Data Collection & Generation Pipeline**

113 **Data Sources** To ensure the benchmark is both comprehensive, AER-bench combines an automated
 114 extraction pipeline for modern research with meticulous manual annotations for historical literature.
 115 The dataset is built upon two distinct corpora:

- 116 • *Modern Frontiers (191 articles)*: Sourced from public, Open Access³ versions of papers pub-
 117 lished in 2025 in top-tier journals (*American Economic Review*, *Econometrica*), guaranteeing
 118 that the benchmark challenges models with the most recent methodological innovations
 119 while strictly adhering to copyright and open-science protocols.
- 120 • *Historical Fragments (900 segments)*: Curated from multilingual text segments of canonical
 121 historical literature (1900–1950) to evaluate models on historical concept drift and early
 122 economic paradigms, which presents a challenging task due to archaic terminology and
 123 evolving theoretical frameworks.

	EK	MM	PA	LR	MG	EA
R1	Knowledge Point Extraction	Math Content Extraction	Paper Structure Parsing	Reference Extraction	Full Model Extraction	Result & Data Extraction
R2	Question Generation ✓	Problem Design ✓	Condensed Body Generation	Argument Analysis	Scenario Narrative	Analysis Excerpt & Task Build ✓
R3	—	—	Key Point Verification ✓	Condensed Body Generation	Rubric Construction ✓	—
R4	—	—	—	Rubric Construction ✓	—	—

Figure 2: Auto-acquisition pipelines for the six LLM-judged tasks. Each column is a task type and each row is a generation round; cell labels name the per-round operation. Blue cells consume the source PDF, gray cells operate from prior-round outputs alone (LLM-only), and the green check marks each task’s terminal round. Pipeline depth ranges from 2 to 4 rounds.

124 **Task Generation Pipeline** For the modern corpus, we employ a multi-round, LLM-assisted
 125 generation pipeline that begins with structural PDF parsing to extract raw text and equations, followed
 126 by multiple rounds of information extraction, desymbolization, and task construction with automated
 127 quality checks to ensure self-contained tasks free of structural errors⁴. Figure 2 summarizes the
 128 round-by-round operations across all six modern-corpus pipelines, illustrating which rounds consume
 129 the source PDF and where each pipeline terminates; pipeline depth ranges from two rounds (EK,
 130 MM, EA) to four rounds (LR), reflecting the differing structural complexity of each task. In contrast,
 131 the historical fragments used for the Concept Identification (CI) task were strictly manually annotated
 132 by expert annotators⁵ with advanced degrees in economics, who read each fragment and performed
 133 binary classification across eight predefined conceptual fields, providing gold-standard labels that
 134 form a rigorous testbed for evaluating AI comprehension of early, non-formalized economic concepts.

135 **2.3 Evaluation Framework**

136 AER-bench utilizes a hybrid evaluation framework that balances absolute objectivity with nuanced
 137 semantic assessment, adapting its scoring strategy to the nature of each task type.

³For articles not available under Open Access licenses from journals, we substitute with legally accessible working paper versions from NBER, arXiv, or authors’ personal homepages to ensure data legality.

⁴See Appendix G for detailed prompt templates, output formats, and scoring methods for each task.

⁵All annotators are economics doctoral students who have passed their qualifying examinations and language requirements. Moreover, we only select and retain data with consistent annotations to ensure the validity of the task answers.

138 **Rule-Based Scoring** For tasks with deterministic answers (EK and CI), we employ strict rule-based
 139 grading. CI uses categorical exact matching, while EK incorporates unidirectional word-boundary
 140 regular expressions for fill-in-the-blank items and normalized Kendall’s τ distance for ordering tasks.

141 **LLM-as-Judge** For open-ended generative tasks⁶ (MM, PA, LR, MG, EA), we adopt an LLM-
 142 as-Judge paradigm. To prevent holistic scoring biases, each judge model evaluates outputs across
 143 multiple isolated dimensions with detailed rubrics, producing a weighted dimension-level score per
 144 task that captures distinct aspects of quality. To further enhance robustness, we employ multiple
 145 judge models in parallel, mitigating individual model biases. Section 4.2 systematically compares the
 146 scoring differences between two parallel judges to validate this approach.

147 **Score Aggregation** For each item, dimension scores are first aggregated into a task score via
 148 configurable intra-task weights, then scores from multiple judges are arithmetically averaged to
 149 produce a consolidated item score. Overall benchmark scores are computed by weighted summation
 150 across task types, with specific weight values detailed in Section 3.

151 3 Experimental Setup

152 **Evaluated Models & Inference Details.** We evaluate 16 state-of-the-art Large Language Models
 153 to establish a comprehensive baseline for automated economic research. The cohort includes the
 154 OpenAI GPT family (gpt-5.4, gpt-5.2, gpt-5.1, gpt-5, gpt-5-mini, gpt-4o), the Anthropic
 155 Claude series (claude-4.6-opus, claude-4.6-sonnet, claude-3.7-sonnet), and the Google
 156 Gemini lineage (gemini-3.1-pro, gemini-3-flash, gemini-2.5-flash). Furthermore, we
 157 incorporate leading open-weight and regional models: kimi-k2.5, glm-5.1, doubao-seed-1.6,
 158 and qwen3-v1-32b-instruct. Note that all inferences are executed under a zero-shot setting.

Table 2: Dimension-level scoring weights for each task type. Dimensions are weighted according to their importance for evaluating model performance, with weights normalized to sum to 1 within each task. Detailed rubric definitions for each dimension are provided in Appendix E.

Task	EK	MM	PA	LR	MG	EA	CI
# Dims	4	4	4	2	Multiple	6	8
Weights	Equal	4:2:3:1	4:4:3:1	3:2	Equal	4:4:4:3:3:2	Equal

159 **Judge Configuration & Score Aggregation.** For the hybrid evaluation framework introduced in
 160 Section 2.3, we instantiate a programmatic evaluator for deterministic tasks (EK and CI). For the
 161 five open-ended generative tasks, we deploy claude-4.6-sonnet and gpt-5.2 as independent
 162 LLM-as-Judges to mitigate individual evaluator bias.

163 To formalize the multi-judge merging strategy outlined previously, let J be the set of judges evaluating
 164 model m on task t . The task-level score is computed as the arithmetic mean of the judges’ outputs:

$$S_{m,t} = \frac{1}{|J|} \sum_{j \in J} s_{m,t}^j \quad (1)$$

165 The overall benchmark score S_m is then calculated as the weighted sum of the task-level scores. In
 166 AER-bench, all seven task categories are assigned equal importance ($w_t = 1/7$):

$$S_m = \sum_{t \in T} w_t \cdot S_{m,t}, \quad \sum_{t \in T} w_t = 1 \quad (2)$$

167 Within each multi-dimensional task, sub-scores are aggregated according to predefined intra-task
 168 dimension weights, which reflect the relative importance of specific economic reasoning facets. These
 169 exact dimensional configurations are detailed in Table 2.

⁶Note that the Literature Review (LR) task employs both an LLM-as-Judge rubric and a rule-based citation recall score, which are combined via weighted aggregation.

170 **4 Results**

171 **4.1 Model Performance**

172 Table 3 presents the overall benchmark ranking across 16 evaluated language models. `gpt-5.4`
 173 achieves the highest overall score at 75.0%, excelling in Mathematical Modeling (92.8%), Model
 174 Generation (85.8%), and Economic Analysis (81.6%); however, a paired bootstrap test shows this
 175 lead is statistically indistinguishable from `claude-4.6-opus` (74.0%, $p = 0.092$), placing the two
 176 models in a statistical tie at the top. OpenAI models nonetheless exhibit notable version instabil-
 177 ity: `gpt-4o` scores only 54.0% with particularly weak Economic Knowledge (25.6%). In contrast,
 178 Anthropic’s Claude models demonstrate remarkably stable performance, with `claude-4.6-opus`
 179 (74.0%), `claude-3.7-sonnet` (73.7%), and `claude-4.6-sonnet` (73.4%) clustered within 0.6
 180 percentage points, and achieve the highest Concept Identification scores (72.7–72.9%), substan-
 181 tially outperforming `gpt-5.4` (58.9%). Google’s `gemini-3.1-pro` (71.6%) exhibits task-specific
 182 dominance, achieving state-of-the-art performance on Economic Knowledge (76.3%) and Paper
 183 Abstraction (74.3%) despite trailing in overall ranking, possibly reflecting targeted training on
 184 academic literature. `kimi-k2.5` (69.1%) demonstrates that open-weight models are approaching
 185 closed-source performance on domain-specific benchmarks. Critically, within-family performance
 186 generally correlates with model version (`gpt-5.4` > `gpt-5.2` > `gpt-5` > `gpt-5.1`), validating
 187 that AER-bench effectively captures incremental capability improvements. Per-task 95% bootstrap
 188 confidence intervals for all 16 models are reported in Appendix A (Table 4).

Table 3: Overall benchmark scores (in %) with item-level 95% bootstrap confidence intervals on the Overall column. Per-task columns show point estimates; bold values mark the best score in each column. Models are grouped by provider and sorted by Overall within each group.

Provider	Model	Overall	EK	MM	PA	LR	MG	EA	CI
OpenAI	<code>gpt-5.4</code>	75.0 [74.0, 76.0]	73.9	92.8	76.2	56.0	85.8	81.6	58.9
	<code>gpt-5.2</code>	73.0 [72.0, 74.0]	70.8	91.5	76.4	50.3	83.6	78.2	60.5
	<code>gpt-5</code>	71.0 [69.0, 72.0]	74.2	91.8	73.5	47.7	84.2	78.5	46.9
	<code>gpt-5.1</code>	70.0 [69.0, 71.0]	67.5	89.7	75.6	48.1	82.8	75.0	51.4
	<code>gpt-5-mini</code>	69.0 [68.1, 70.0]	72.2	90.8	71.1	48.1	79.2	76.8	45.1
	<code>gpt-4o</code>	54.0 [53.1, 55.0]	25.6	78.4	72.1	39.5	45.7	64.6	52.3
Anthropic	<code>claude-4.6-opus</code>	74.0 [73.1, 74.9]	71.5	89.4	76.1	52.5	78.1	77.7	72.7
	<code>claude-3.7-sonnet</code>	73.7 [72.8, 74.6]	67.0	88.6	76.8	50.7	81.7	78.2	73.0
	<code>claude-4.6-sonnet</code>	73.4 [72.5, 74.3]	66.7	88.8	76.3	50.6	81.8	76.9	72.9
Google	<code>gemini-3.1-pro</code>	71.6 [70.6, 72.5]	76.3	89.3	74.3	47.0	70.0	72.0	72.0
	<code>gemini-3-flash</code>	65.8 [64.9, 66.8]	74.8	86.5	77.0	46.2	62.4	72.3	41.6
	<code>gemini-2.5-flash</code>	64.9 [63.9, 65.9]	69.4	88.0	76.1	41.4	68.0	67.1	44.1
Others	<code>kimi-k2.5</code>	69.1 [68.0, 70.1]	68.9	87.4	73.5	47.3	71.9	76.3	58.3
	<code>glm-5.1</code>	65.5 [64.4, 66.6]	57.9	85.4	75.8	49.0	69.5	74.8	46.3
	<code>doubao-seed-1.6</code>	64.8 [63.8, 65.8]	66.0	84.9	74.7	42.2	61.5	72.7	51.4
	<code>qwen3-vl-32b-instruct</code>	61.3 [60.3, 62.3]	60.9	82.0	73.5	39.1	59.4	66.2	48.0

189 To preempt concerns that the ranking above might be an artifact of sampling noise, we additionally
 190 conduct pairwise paired-bootstrap significance tests on the Overall score for all $\binom{16}{2}$ model pairs
 191 (Appendix B, Table 6): every cross-tier comparison rejects equality at $p < 0.001$, while a small
 192 number of within-cluster pairs—most notably `gpt-5.4` vs. `claude-4.6-opus`—remain statistically
 193 tied, providing a conservative significance frame for every claim we make about model rankings.

194 Figure 3 illustrates the variation in task difficulty and discriminative power across the benchmark.
 195 Literature Review presents the lowest median score (~40%) with tight clustering, indicating that
 196 all models struggle uniformly with citation network reconstruction and argument structure recovery.
 197 Conversely, Mathematical Modeling exhibits high scores (median near 90%) with limited spread,
 198 suggesting that symbolic manipulation capabilities have largely converged across modern model
 199 families. Tasks like Economic Knowledge and Model Generation demonstrate the widest score
 200 ranges, serving as the primary differentiators between average and frontier AI systems in economic
 201 research contexts. See Appendix F for detailed performance analysis by task dimension.

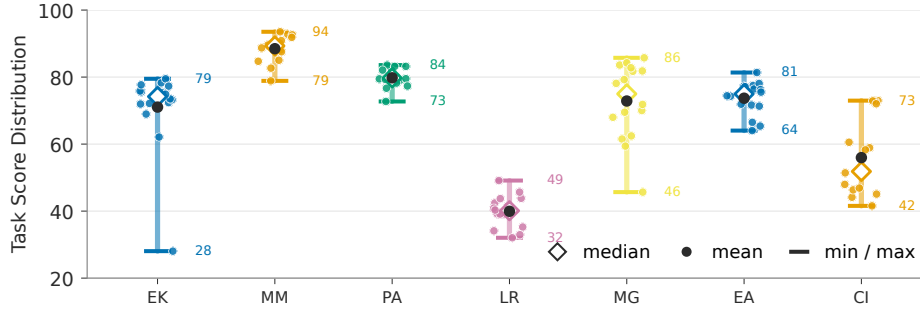


Figure 3: Task score distribution across all 16 evaluated models. Each vertical line represents the min-max range of model scores for a given task, with individual model scores shown as colored dots. White diamonds indicate median scores, and gray circles indicate mean scores. The figure reveals substantial performance variation across tasks.

202 4.2 LLM-as-a-Judge Quality

203 To validate the multi-judge evaluation framework detailed in Section 3, we analyze the scoring
 204 consistency between two independent judge models: `claude-4.6-sonnet` and `gpt-5.2`. Figure 4
 205 shows high overall agreement, with a mean absolute difference of 2.6 percentage points. This tight
 206 clustering confirms that our dimension-level, criteria-bound rubrics effectively constrain the subjectivity typical of LLM evaluators, producing reliable assessments across different judge architectures.

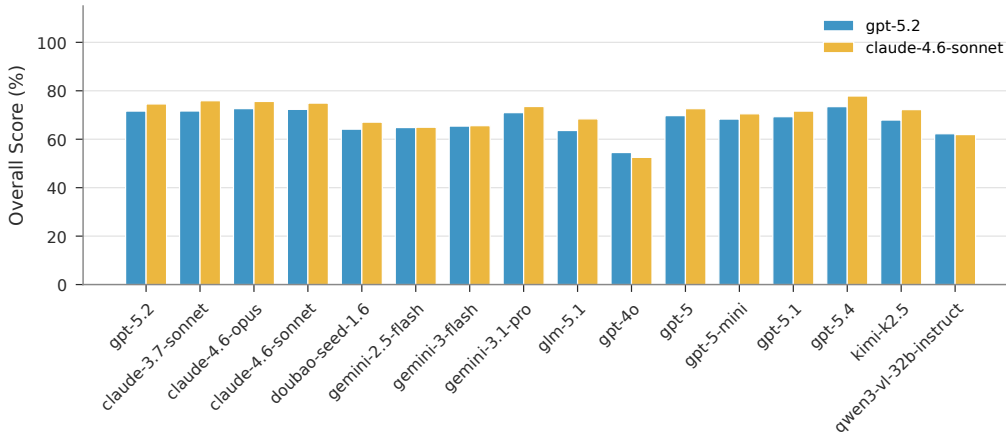


Figure 4: Judge agreement on overall benchmark scores across 16 models. Scores from two independent LLM judges (`claude-4.6-sonnet` and `gpt-5.2`) show strong consistency, with a mean absolute difference of 2.6 percentage points, validating the robustness of our rubric-based evaluation framework.

207 While overall agreement is strong, Figure 5 reveals systematic task-level scoring preferences between
 208 the judges. `claude-4.6-sonnet` assigns higher scores on Mathematical Modeling (+3.1 pp) and
 209 Model Generation (+13.0 pp), whereas `gpt-5.2` scores Paper Abstraction higher (-3.4 pp). To test
 210 whether these task-level offsets translate into a structural *self-preference bias*—each judge inflating
 211 scores for its own family—we computed the per-item Claude-minus-GPT score gap $\Delta_{m,i}$ for every
 212 examinee m , applied a paired Wilcoxon signed-rank test per model, and then compared the per-model
 213 mean Δ across families with a two-sided Mann-Whitney U test. Anthropic candidates do receive a
 214 larger “Claude-judge bonus” than OpenAI candidates (mean $\Delta = +2.90$ pp vs. $+1.41$ pp; “Other”
 215 families $+0.82$ pp), but the family-level difference falls short of conventional significance ($p = 0.095$,
 216 $n = 3$ vs. $n = 6$). The directional effect is real but small relative to the cross-tier gaps documented
 217 in Table 3, and arithmetic averaging across the two independent judges substantially attenuates it;
 218 per-model results and the full statistical analysis are reported in Appendix C.

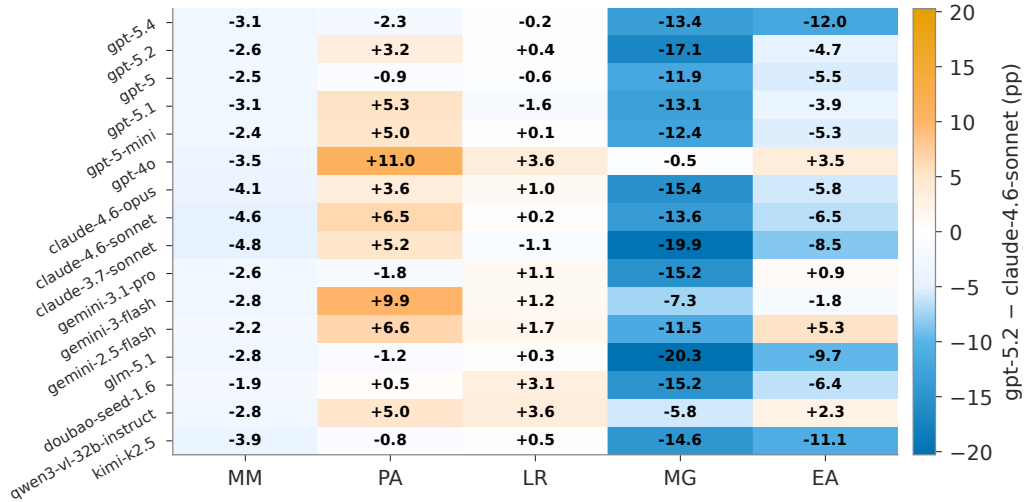


Figure 5: Task-level scoring differences of two judges. Heatmap shows the score difference (claude-4.6-sonnet minus gpt-5.2) across five generative tasks, with blue indicating higher scores from Claude and orange indicating higher scores from gpt-5.2. Models are grouped by provider and sorted by overall performance within each group.

220 4.3 Key Findings

221 Based on the quantitative results, we highlight three core findings.

222 **Finding 1: Performance Volatility in Knowledge-Intensive Tasks Demands Agentic Solutions.**

223 Pure language models exhibit severe instability when tasked with retrieving specific literature facts or
 224 constructing citation networks. This is evident in the extremely wide spread of Economic Knowledge
 225 (EK) scores (ranging from 28.1% to 79.5%) and the universally low performance in Literature Review
 226 (LR) shown in Figure 3. These results suggest that relying solely on an LLM’s parametric memory
 227 is insufficient for rigorous academic writing. Achieving expert-level reliability in these phases will
 228 likely require integrating pure LLMs with agentic systems.

229 **Finding 2: Distinct Family-Level Capability Fingerprints.**

230 Beyond version-level scaling, AER-
 231 bench reveals task-specific signatures that distinguish model families. All five GPT-5 variants occupy
 232 the top five positions on MM (90.9–93.5%), leading the best non-GPT-5 model by 3–4 points; the 3
 233 Claude models cluster tightly at 72.7–73.0% on CI—over 12 points above gpt-5.2—while staying
 234 within 0.5% overall, indicating a stable shared profile rather than checkpoint-level idiosyncrasy;
 235 and Gemini models trail GPT-5 and Claude by 5–15 points on LR, MG, and EA despite being
 236 competitive on EK and PA. These cross-family patterns are systematic rather than random, showing
 that AER-bench can diagnose the relative strengths of entire series, not merely rank isolated models.

237 **Finding 3: A Model-Stratified, Length-Conditional “Modernity Bias.”**

238 At face value, our longi-
 239 tudinal Concept Identification (CI) task suggests a pervasive “modernity bias”: aggregate scores are
 240 ~5 pp higher on 1930–1950 than on 1900–1929 chunks (Figure 6), led by the equilibrium_scene
 241 dimension. Two robustness checks (Appendix D, Tables 9–8) sharpen this picture. After partialling
 242 out source-chunk length (OLS, model FE, item-clustered SE), the Era₃ – Era₂ coefficient shrinks
 243 to +2.75 pp ($p = 0.44$); the gap lives entirely in the longest two length quartiles (Q3 +9.2, Q4
 244 +12.7 pp), while the shortest two reverse (Era₂ leads by 1.8–2.1 pp). A per-decade × per-model
 245 breakdown then reveals strong heterogeneity: the three Claude models, gemini-3.1-pro, gpt-5.2,
 246 and kimi-k2.5 score *higher* on 1900–1909 than on 1940–1950 (e.g., claude-4.6-sonnet 86.2 vs.
 247 64.2%), whereas gpt-5.4, gpt-5, gpt-5-mini, both Gemini Flash variants, and glm-5.1 show
 248 the canonical 20+ pp drop on the earliest decade. The honest characterisation is therefore a *model-*
 249 *stratified, length-conditional* bias: real for specific models on longer paragraphs, but closed—and in
 some cases inverted—by frontier Anthropic, Gemini-3.1, and GPT-5.2 systems. The pattern points to

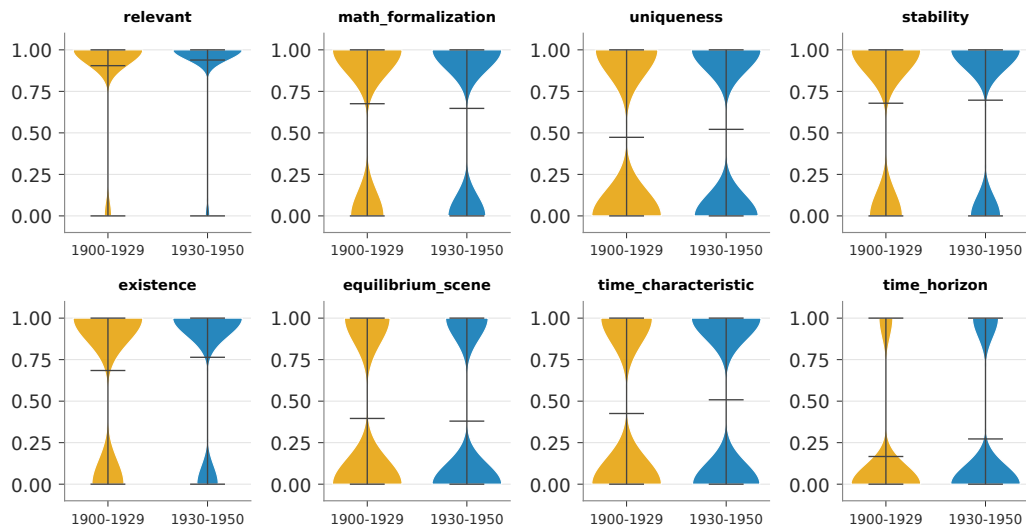


Figure 6: Concept identification performance across historical periods. Violin plots show the distribution of 16 model scores on each of the 8 concept dimensions, comparing publications from 1900–1929 (orange) versus 1930–1950 (blue). Models exhibit systematic performance gaps on earlier texts, particularly for `equilibrium_scene` and `existence` dimensions.

250 representational gaps in pre-training corpus composition rather than uniform conceptual drift, and
 251 reframes the CI task as a diagnostic for *which* models still carry this blind spot.

252 5 Limitations

253 While AER-bench establishes a rigorous evaluation framework, it presents several inherent limitations.
 254 First, the automated task generation pipeline is computationally expensive. Processing full-length
 255 papers through multi-turn interactions requires top-tier models, currently costing approximately \$6
 256 per paper. Consequently, we deliberately prioritized high-fidelity task extraction over massive scale.
 257 Second, the reliance on PDF parsing occasionally introduces formatting artifacts that can disrupt the
 258 extraction of complex mathematical notations. Third, although our dimension-level rubrics mitigate
 259 subjective biases, the LLM-as-a-Judge paradigm may still harbor residual architectural preferences.
 260 Finally, the modern corpus consists predominantly of English-language journals, potentially masking
 261 disparities in models’ multilingual economic reasoning capabilities. Additionally, the benchmark
 262 does not include a human expert baseline, a decision motivated by several practical and conceptual
 263 challenges discussed in the Appendix H.

264 6 Conclusion

265 This paper introduces AER-bench, the first comprehensive benchmark designed to evaluate LLMs
 266 across the full lifecycle of economic research. Moving beyond simple question-answering paradigms,
 267 AER-bench assesses 16 state-of-the-art models on a rigorous, multi-dimensional taxonomy spanning
 268 seven distinct tasks—from deterministic knowledge retrieval and conceptual identification to the
 269 generation of formal mathematical models and qualitative empirical analyses.

270 Our evaluations reveal that while frontier models excel in symbolic manipulation and text summariza-
 271 tion, they exhibit high volatility in knowledge-intensive tasks like citation network reconstruction.
 272 We also uncover a systematic “modernity bias” wherein models struggle to interpret early historical
 273 economic concepts. These findings highlight the necessity of transitioning from pure, parametrically-
 274 reliant LLMs to agentic, retrieval-augmented systems. Ultimately, AER-bench provides a robust and
 275 systematic evaluation yardstick that will guide the development of the next generation of AI systems
 276 capable of autonomously conducting, interpreting, and advancing rigorous economic research.

277 **References**

- 278 [1] Daron Acemoglu, Simon Johnson, and James A Robinson. The colonial origins of comparative
279 development: An empirical investigation. *American economic review*, 91(5):1369–1401, 2001.
- 280 [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
281 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4
282 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 283 [3] Joshua D Angrist and Jörn-Steffen Pischke. The credibility revolution in empirical economics:
284 How better research design is taking the con out of econometrics. *Journal of economic
285 perspectives*, 24(2):3–30, 2010.
- 286 [4] Anthropic. Claude code. <https://github.com/anthropics/claude-code>, 2025. Ac-
287 cessed: 2026-05.
- 288 [5] Shuangyan Deng, Haizhou Peng, Jiachen Xu, Rui Mao, Ciprian Doru Giurcaneanu, and Jiamou
289 Liu. *FinMR: A Knowledge-Intensive Multimodal Benchmark for Advanced Financial Reasoning*,
290 page 168–176. Association for Computing Machinery, New York, NY, USA, 2025. ISBN
291 9798400722202. URL <https://doi.org/10.1145/3768292.3770365>.
- 292 [6] Esther Duflo and Abhijit Banerjee. *Poor economics*, volume 619. PublicAffairs New York,
293 2011.
- 294 [7] Milton Friedman et al. The methodology of positive economics. *Essays in positive economics*,
295 3(3):145–178, 1953.
- 296 [8] Thomas Hubert, Rishi Mehta, Laurent Sartran, Miklós Z Horváth, Goran Žužić, Eric Wieser,
297 Aja Huang, Julian Schrittwieser, Yannick Schroecker, Hussain Masoom, et al. Olympiad-level
298 formal mathematical reasoning with reinforcement learning. *Nature*, pages 1–3, 2025.
- 299 [9] John Maynard Keynes. The general theory of employment. *The quarterly journal of economics*,
300 51(2):209–223, 1937.
- 301 [10] Anton Korinek. Generative ai for economic research: Use cases and implications for economists.
302 *Journal of Economic Literature*, 61(4):1281–1317, 2023.
- 303 [11] Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay
304 Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quan-
305 titative reasoning problems with language models. *Advances in neural information processing
306 systems*, 35:3843–3857, 2022.
- 307 [12] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scien-
308 tist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*,
309 2024.
- 310 [13] Robert E Lucas Jr. Econometric policy evaluation: A critique. In *Carnegie-Rochester conference
311 series on public policy*, volume 1, pages 19–46. North-Holland, 1976.
- 312 [14] Qirui Mi, Qipeng Yang, Zijun Fan, Wentian Fan, Heyang Ma, Chengdong Ma, Siyu Xia, Bo An,
313 Jun Wang, and Haifeng Zhang. Econgym: A scalable ai testbed with diverse economic tasks.
314 *arXiv preprint arXiv:2506.12110*, 2025.
- 315 [15] Nous Research. Hermes agent. <https://github.com/nousresearch/hermes-agent>,
316 2025. Accessed: 2026-05.
- 317 [16] Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt
318 Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian,
319 et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint
320 arXiv:2506.13131*, 2025.
- 321 [17] OpenClaw Contributors. OpenClaw: An open-source agentic framework. [https://github.
322 com/openclaw/openclaw](https://github.com/openclaw/openclaw), 2025. Accessed: 2026-05.

- 323 [18] Tejal Patwardhan, Rachel Dias, Elizabeth Proehl, Grace Kim, Michele Wang, Olivia Watkins,
324 Simón Posada Fishman, Marwan Aljubeh, Phoebe Thacker, Laurance Fauconnet, Natalie S. Kim,
325 Patrick Chao, Samuel Miserendino, Gildas Chabot, David Li, Michael Sharman, Alexandra Barr,
326 Amelia Glaese, and Jerry Tworek. Gdpval: Evaluating ai model performance on real-world
327 economically valuable tasks, 2025. URL <https://arxiv.org/abs/2510.04374>.
- 328 [19] Yinzhu Quan and Zefang Liu. EconLogicQA: A question-answering benchmark for evaluating
329 large language models in economic sequential reasoning. In Yaser Al-Onaizan, Mohit Bansal,
330 and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics:
331 EMNLP 2024*, pages 2273–2282, Miami, Florida, USA, November 2024. Association for
332 Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.125. URL [https://
333 aclanthology.org/2024.findings-emnlp.125/](https://aclanthology.org/2024.findings-emnlp.125/).
- 334 [20] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro,
335 Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models
336 can teach themselves to use tools. *Advances in neural information processing systems*, 36:
337 68539–68551, 2023.
- 338 [21] Social Catalyst Lab. APE: Automatic Policy Evaluation. [https://ape.socialcatalystlab.
339 org/](https://ape.socialcatalystlab.org/), 2025. Accessed: 2026-05.
- 340 [22] Robert M Solow. Technical change and the aggregate production function. *The review of
341 Economics and Statistics*, 39(3):312–320, 1957.
- 342 [23] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,
343 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly
344 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 345 [24] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo,
346 Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and
347 challenging multi-task language understanding benchmark. *Advances in Neural Information
348 Processing Systems*, 37:95266–95290, 2024.
- 349 [25] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
350 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.
351 *Advances in neural information processing systems*, 35:24824–24837, 2022.
- 352 [26] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
353 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,
354 2022.

355 **A Full Per-Task Ranking with Bootstrap Confidence Intervals**

356 To complement the main-text ranking (Table 3), we report in Table 4 the full per-task results for all
357 16 evaluated models, with item-level 95% bootstrap confidence intervals on both the Overall score
358 and every individual task. Confidence intervals are obtained from $B = 10,000$ bootstrap resamples
359 stratified within each task; the same resample indices are shared across models so that pairwise
360 comparisons preserve per-item correlation. For non-CI tasks, the per-item score is the arithmetic
361 mean of the `claude-4.6-sonnet` and `gpt-5.2` judges before resampling; the CI task uses the
362 rule-based judge.

363 **No single model dominates all dimensions.** While `gpt-5.4` attains the column-wise maximum on
364 five of the eight columns (Overall, MM, LR, MG, EA), it is decisively beaten on the remaining three:
365 `gemini-3.1-pro` leads on EK (76.3 vs. 73.9), `gemini-3-flash` leads on PA (77.0 vs. 76.2), and the
366 Anthropic family leads on CI by a wide margin (`claude-3.7-sonnet` 73.0; `claude-4.6-sonnet`
367 72.9; `claude-4.6-opus` 72.7) versus `gpt-5.4` at 58.9. The CI column in particular splits cleanly
368 along provider lines: every Anthropic model exceeds 72%, whereas the OpenAI family except `gpt-4o`
369 (52.3) sits between 45.1 and 60.5. Combined with the comparable ranking of `gemini-3.1-pro`
370 (72.0) on the same task, this pattern suggests that performance on early-20th-century economic-history
371 concepts depends more on pre-training corpus composition than on raw scale.

372 **Confidence-interval width tracks item count and judge type.** CI widths are far from uniform
373 across tasks. The two largest tasks, EA ($n = 380$) and MM ($n = 399$), produce the tightest
374 intervals: half-widths are typically ≤ 0.6 pp and ≤ 0.8 pp respectively (e.g., `gpt-5.4` on EA has
375 [81.0, 82.2]). The CI task is the opposite extreme, with half-widths near 5–6 pp on every model (e.g.,
376 `claude-4.6-opus` [67.3, 77.8]), reflecting both its smaller sample ($n = 160$) and the higher per-
377 item variance of multi-field rule-based scoring. LR ($n = 187$) and MG ($n = 124$) show intermediate
378 widths around 1–3 pp. These differences are quantitative reminders that gaps of 1–2 pp on CI should
379 not be over-interpreted, whereas similar-sized gaps on EA or MM are typically well outside the
380 bootstrap noise floor.

381 **Anthropic models occupy a tight cluster.** The three Claude models lie within 0.6 pp on Overall
382 (74.0, 73.7, 73.4) and their per-task CIs overlap extensively on every column. Their largest within-
383 family spread is on EK (66.7–71.5), where `claude-4.6-opus` clearly leads, and on MG (78.1–81.8),
384 where the two Sonnet variants slightly exceed Opus—suggesting the Sonnet line carries a marginal
385 advantage on procedural model construction. By contrast, the OpenAI line shows a near-monotonic
386 version-vs-score relationship on Overall (`gpt-5.4` > `gpt-5.2` > `gpt-5` > `gpt-5.1` > `gpt-5-mini`
387 > `gpt-4o`), with `gpt-4o` a clear outlier on EK at 25.6, more than 25 pp below every other model and
388 consistent with its older training cutoff.

389 **Where the rankings are statistically separable.** Reading the bootstrap intervals horizontally
390 provides an informal visual significance test on the Overall column: the top two rows (`gpt-5.4`
391 [74.0, 76.0] and `claude-4.6-opus` [73.1, 74.9]) overlap by roughly 1 pp, consistent with the paired-
392 bootstrap p -value of 0.092 reported in the main text. The next three rows (`claude-3.7-sonnet`,
393 `claude-4.6-sonnet`, `gpt-5.2`) form a tightly overlapping cluster between 72.5 and 74.6. After
394 `gemini-3.1-pro` the intervals start to separate cleanly, with `gpt-4o` ([53.1, 55.0]) sitting more than
395 6 pp below the next-lowest model. Per-task intervals tell the same story at finer resolution: many
396 cross-provider gaps on EK and CI are highly significant (no overlap at all), while LR and PA gaps are
397 typically within the noise floor.

Table 4: Per-task and overall scores (in %) with item-level 95% bootstrap confidence intervals ($B = 10,000$ resamples, stratified within task). Each cell stacks the point estimate above its $[ci_{lo}, ci_{hi}]$ interval. Bold marks the column-wise maximum. Models are grouped by provider and sorted by Overall within each group.

Provider	Model	Overall	EK	MM	PA	LR	MG	EA	CI
OpenAI	gpt-5.4	75.0 [74.0, 76.0]	73.9 [72.4, 75.4]	92.8 [92.0, 93.5]	76.2 [73.9, 78.5]	56.0 [54.7, 57.3]	85.8 [83.2, 88.1]	81.6 [81.0, 82.2]	58.9 [53.0, 64.8]
	gpt-5.2	73.0 [72.0, 74.0]	70.8 [69.3, 72.2]	91.5 [90.8, 92.3]	76.4 [74.7, 78.0]	50.3 [49.2, 51.3]	83.6 [81.2, 85.8]	78.2 [77.7, 78.6]	60.5 [54.5, 66.6]
	gpt-5	71.0 [69.9, 72.0]	74.2 [72.6, 75.8]	91.8 [91.0, 92.6]	73.5 [71.2, 75.6]	47.7 [46.9, 48.5]	84.2 [81.6, 86.7]	78.5 [78.0, 78.9]	46.9 [40.9, 52.9]
	gpt-5.1	70.0 [69.0, 71.0]	67.5 [65.9, 69.0]	89.7 [88.8, 90.6]	75.6 [73.6, 77.5]	48.1 [47.2, 49.0]	82.8 [80.6, 84.8]	75.0 [74.4, 75.6]	51.4 [45.5, 57.3]
	gpt-5-mini	69.0 [68.1, 70.0]	72.2 [70.7, 73.7]	90.8 [89.8, 91.7]	71.1 [69.0, 73.1]	48.1 [47.0, 49.2]	79.2 [76.8, 81.5]	76.8 [76.2, 77.3]	45.1 [39.3, 50.9]
	gpt-4o	54.0 [53.1, 55.0]	25.6 [23.2, 28.2]	78.4 [77.0, 79.8]	72.1 [70.9, 73.1]	39.5 [38.7, 40.4]	45.7 [43.7, 47.7]	64.6 [63.9, 65.4]	52.3 [46.3, 58.0]
Anthropic	claude-4.6-opus	74.0 [73.1, 74.9]	71.5 [70.0, 73.0]	89.4 [88.4, 90.5]	76.1 [74.3, 78.0]	52.5 [51.4, 53.6]	78.1 [75.4, 80.7]	77.7 [77.1, 78.2]	72.7 [67.3, 77.8]
	claude-3.7-sonnet	73.7 [72.8, 74.6]	67.0 [65.4, 68.5]	88.6 [87.6, 89.7]	76.8 [75.1, 78.4]	50.7 [49.5, 51.8]	81.7 [79.5, 83.8]	78.2 [77.6, 78.7]	73.0 [67.5, 78.2]
	claude-4.6-sonnet	73.4 [72.5, 74.3]	66.7 [65.1, 68.3]	88.8 [87.7, 89.8]	76.3 [74.5, 78.0]	50.6 [49.4, 51.6]	81.8 [79.4, 84.1]	76.9 [76.3, 77.5]	72.9 [67.6, 78.0]
Google	gemini-3.1-pro	71.6 [70.6, 72.5]	76.3 [75.0, 77.6]	89.3 [88.3, 90.3]	74.3 [71.5, 77.0]	47.0 [46.2, 47.9]	70.0 [67.3, 72.6]	72.0 [71.5, 72.6]	72.0 [67.1, 77.0]
	gemini-3-flash	65.8 [64.9, 66.8]	74.8 [73.4, 76.2]	86.5 [85.5, 87.6]	77.0 [75.7, 78.1]	46.2 [45.2, 47.2]	62.4 [60.4, 64.4]	72.3 [71.6, 72.9]	41.6 [35.9, 47.3]
	gemini-2.5-flash	64.9 [63.9, 65.9]	69.4 [67.9, 71.0]	88.0 [87.0, 89.1]	76.1 [74.0, 78.0]	41.4 [40.5, 42.3]	68.0 [65.4, 70.5]	67.1 [66.5, 67.7]	44.1 [38.4, 50.0]
Others	kimi-k2.5	69.1 [68.0, 70.1]	68.9 [67.3, 70.5]	87.4 [86.2, 88.5]	73.5 [70.8, 76.0]	47.3 [46.1, 48.4]	71.9 [68.9, 74.7]	76.3 [75.6, 77.0]	58.3 [52.5, 63.8]
	glm-5.1	65.5 [64.4, 66.6]	57.9 [55.3, 60.4]	85.4 [84.1, 86.7]	75.8 [73.3, 78.2]	49.0 [47.8, 50.2]	69.5 [66.8, 72.0]	74.8 [74.0, 75.6]	46.3 [40.6, 52.0]
	doubao-seed-1.6	64.8 [63.8, 65.8]	66.0 [64.4, 67.6]	84.9 [83.7, 86.2]	74.7 [72.2, 77.1]	42.2 [41.2, 43.3]	61.5 [58.8, 64.2]	72.7 [72.0, 73.4]	51.4 [45.9, 57.0]
	qwen3-vl-32b-instruct	61.3 [60.3, 62.3]	60.9 [59.3, 62.5]	82.0 [80.7, 83.4]	73.5 [71.3, 75.5]	39.1 [38.4, 40.0]	59.4 [57.6, 61.3]	66.2 [65.5, 66.9]	48.0 [42.3, 53.8]

398 **B Pairwise Paired-Bootstrap Significance on Overall Score**

399 The point-estimate ranking in Table 3 is informative but does not by itself establish whether two
 400 adjacent models differ in a statistically meaningful sense. Table 6 reports the two-sided p -value for
 401 each ordered pair (i, j) of models under a paired bootstrap on the Overall score. For each of the $B =$
 402 10,000 resamples we use the *same* per-task item indices to recompute every model’s Overall score, so
 403 that the difference $D_b = \widehat{O}_i^{(b)} - \widehat{O}_j^{(b)}$ inherits the per-item correlation between the two models. The
 404 reported p -value is the standard two-sided percentile, $\hat{p} = 2 \cdot \min(\Pr(D_b \geq 0), \Pr(D_b \leq 0))$, where
 405 the probabilities are taken over the bootstrap distribution. Bold cells flag $p < 0.05$; italic cells flag
 406 $0.05 \leq p < 0.10$; plain cells are non-significant. To keep the matrix readable, models are referenced
 407 by the integer IDs defined in Table 5; IDs follow the same provider-grouped ordering used throughout
 408 the paper.

Table 5: Model ID legend used in Table 6.

ID	Model	ID	Model
1	gpt-5.4	9	kimi-k2.5
2	claude-4.6-opus	10	gpt-5-mini
3	claude-3.7-sonnet	11	gemini-3-flash
4	claude-4.6-sonnet	12	glm-5.1
5	gpt-5.2	13	gemini-2.5-flash
6	gemini-3.1-pro	14	doubao-seed-1.6
7	gpt-5	15	qwen3-v1-32b-instruct
8	gpt-5.1	16	gpt-4o

409 **The top of the ranking is a statistical tie.** The most consequential pair for the headline claim is
 410 gpt-5.4 (ID 1) versus claude-4.6-opus (ID 2): their paired-bootstrap p -value is 0.092, which
 411 fails to reject equality at the 0.05 level and only barely crosses the 0.10 threshold. gpt-5.4 *is*
 412 significantly above the third-place claude-3.7-sonnet ($p_{1,3} = 0.020$) and claude-4.6-sonnet
 413 ($p_{1,4} = 0.010$), but its margin over claude-4.6-opus sits in the same regime as several within-
 414 cluster comparisons below. We therefore describe the leader as “statistically tied at the top” rather
 415 than as a clear winner in the main text.

Table 6: Pairwise paired-bootstrap two-sided p -values on the Overall benchmark score, $B = 10,000$ resamples with shared within-task item indices across models. **All entries are reported in units of 10^{-3}** , so a cell value of 92 means $p = 0.092$ and “< 1” means $p < 10^{-3}$. Rows and columns are indexed by the model IDs listed in Table 5. Bold: $p < 0.05$. Italic: $0.05 \leq p < 0.10$. Plain: $p \geq 0.10$. The matrix is symmetric ($p_{ij} = p_{ji}$); the diagonal is empty.

ID	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
1	—	92	20	10	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1
2	92	—	520	158	55	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1
3	20	520	—	512	164	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1
4	10	158	512	—	418	1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1
5	<1	55	164	418	—	3	3	<1	<1	<1	<1	<1	<1	<1	<1	<1
6	<1	<1	<1	1	3	—	331	6	<1	<1	<1	<1	<1	<1	<1	<1
7	<1	<1	<1	<1	3	331	—	85	3	<1	<1	<1	<1	<1	<1	<1
8	<1	<1	<1	<1	<1	6	85	—	47	68	<1	<1	<1	<1	<1	<1
9	<1	<1	<1	<1	<1	<1	3	47	—	989	<1	<1	<1	<1	<1	<1
10	<1	<1	<1	<1	<1	<1	<1	68	989	—	<1	<1	<1	<1	<1	<1
11	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	—	558	34	78	<1	<1
12	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	558	—	158	168	<1	<1
13	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	34	158	—	852	<1	<1
14	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	78	168	852	—	<1	<1
15	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	—	<1
16	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	<1	—

416 **A four-model cluster sits just below the leader.** IDs 2–5 (claude-4.6-opus,
417 claude-3.7-sonnet, claude-4.6-sonnet, gpt-5.2) are mutually indistinguishable: ev-
418 ery pairwise comparison among them yields $p \in [0.158, 0.520]$, with the only borderline case being
419 $p_{2,5} = 0.055$. The three Claude models in particular cannot be ordered statistically (all internal
420 p -values ≥ 0.158). Combined with the leader pair above, this means the top five rows of Table 3
421 should be read as a *cluster* rather than as a strict ranking.

422 **Lower-tier pairs that look close are also indistinguishable.** Several mid-table neighbours fail to
423 separate at conventional levels: $p_{6,7} = 0.331$ (gemini-3.1-pro vs. gpt-5), $p_{7,8} = 0.085$ (gpt-5
424 vs. gpt-5.1), $p_{9,10} = 0.989$ (kimi-k2.5 vs. gpt-5-mini), $p_{11,12} = 0.558$ (gemini-3-flash vs.
425 glm-5.1), and $p_{13,14} = 0.852$ (gemini-2.5-flash vs. doubao-seed-1.6). These five clusters
426 partition the leaderboard into five quasi-equivalence classes; differences within a class should not be
427 treated as evidence of model superiority.

428 **Cross-cluster gaps are robust.** By contrast, every cross-cluster comparison—e.g. any leader-
429 cluster model against any model below gemini-3.1-pro, or any mid-tier model against
430 qwen3-vl-32b-instruct (ID 15) or gpt-4o (ID 16)—yields $p < 0.001$. The lowest two models
431 are themselves strongly separated from each other ($p_{15,16} < 0.001$). Together these patterns indicate
432 that AER-bench resolves coarse capability tiers unambiguously while leaving the finest within-tier
433 ordering to be settled by additional data—a useful diagnostic for sample-size planning in future
434 benchmark extensions.

435 **C Judge Self-Preference Significance Analysis**

436 The two LLM judges used for the open-generation tasks (claude-4.6-sonnet and
 437 azure-gpt-5.2) are themselves candidates in the benchmark. Reviewers reasonably worry that this
 438 creates a structural conflict of interest: each judge could systematically inflate the scores of models
 439 from its own family. The main text (Section 4.2) summarises the result; this appendix documents the
 440 full procedure, the per-model statistics, and the limitations.

441 **Statistical setup.** For every examinee m and every item i on a non-CI task, we observe two
 442 normalised scores in $[0, 1]$, one from each judge. Define the per-item gap

$$\Delta_{m,i} = s_{m,i}^{\text{claude-judge}} - s_{m,i}^{\text{gpt-judge}},$$

443 so that $\Delta > 0$ means the Claude judge is more lenient on item i than the GPT judge. Stacking across
 444 the six non-CI tasks gives $n = 1,469$ paired observations per examinee. We run two complementary
 445 tests:

- 446 1. **Per-model paired Wilcoxon signed-rank test** on $\{\Delta_{m,i}\}_{i=1}^n$ (two-sided), to detect whether
 447 *any* systematic offset exists between the two judges for examinee m .
- 448 2. **Family-level Mann–Whitney U test** on the per-model means $\bar{\Delta}_m$, comparing the three
 449 Anthropic candidates against the six OpenAI candidates. This is the test that directly probes
 450 *in-family* bias: if the Claude judge favoured Anthropic candidates over OpenAI candidates,
 451 Anthropic models would have systematically larger $\bar{\Delta}_m$ than OpenAI models.

452 **Per-model results.** Table 7 reports $\bar{\Delta}_m$ and the Wilcoxon p -value for each of the 16 examinees.
 453 With $n = 1,469$ paired items the Wilcoxon test is highly powered, so almost every model shows a
 454 statistically significant non-zero gap; the magnitude rather than the p -value is the quantity of interest.

Table 7: Per-examinee mean Claude-minus-GPT judge gap $\bar{\Delta}_m$ (in percentage points) with paired Wilcoxon signed-rank p -values (two-sided, $n = 1,469$ items per row). Positive values mean the Claude judge scores examinee m higher than the GPT judge does. The footnote reports family-level mean gaps and the Mann–Whitney U test comparing Anthropic vs. OpenAI candidates.

Examinee	Family	n	$\bar{\Delta}$ (pp)	Wilcoxon p
claude-3.7-sonnet	Anthropic	1469	+3.96	< 0.001
claude-4.6-opus	Anthropic	1469	+2.41	< 0.001
claude-4.6-sonnet	Anthropic	1469	+2.32	< 0.001
gemini-2.5-flash	Google	1469	-1.43	< 0.001
gemini-3-flash	Google	1469	+0.05	0.486
gemini-3.1-pro	Google	1469	-0.01	< 0.001
gpt-4o	OpenAI	1469	-2.05	< 0.001
gpt-5	OpenAI	1469	+1.69	< 0.001
gpt-5-mini	OpenAI	1469	+1.60	< 0.001
gpt-5.1	OpenAI	1469	+1.39	< 0.001
gpt-5.2	OpenAI	1469	+2.00	< 0.001
gpt-5.4	OpenAI	1469	+3.85	< 0.001
doubao-seed-1.6	Other	1469	+1.56	< 0.001
glm-5.1	Other	1469	+3.57	< 0.001
kimi-k2.5	Other	1469	+3.53	< 0.001
qwen3-vl-32b-instruct	Other	1469	-1.52	< 0.001

Group means:
 Anthropic $\bar{\Delta} = +2.90$ pp ($n_g = 3$);
 OpenAI $\bar{\Delta} = +1.41$ pp ($n_g = 6$);
 Other $\bar{\Delta} = +0.82$ pp ($n_g = 7$).
 Mann–Whitney U test (Anthropic vs. OpenAI): $p = 0.095$.

455 **A directional but underpowered in-family effect.** The three Anthropic candidates receive an
456 average Claude-judge bonus of $\bar{\Delta} = +2.90$ pp, roughly twice the OpenAI candidates' $+1.41$ pp and
457 more than three times the "Other" families' $+0.82$ pp. The direction of the effect is exactly what the
458 self-preference hypothesis predicts. However, the Mann–Whitney U test on the per-model means
459 yields $p = 0.095$, falling short of the conventional 0.05 threshold: with only $n_g = 3$ Anthropic and
460 $n_g = 6$ OpenAI candidates the family-level test has limited power. We therefore *cannot* reject the
461 null of no in-family bias, but we equally *cannot* claim that the apparent ~ 1.5 pp Anthropic–OpenAI
462 gap is zero.

463 **Why the effect is small in practice.** Two facts bound the practical impact of this directional
464 bias on the benchmark rankings. First, the quantity that matters for the leaderboard is the *between-*
465 *family* component (~ 1.5 pp), not the global Claude-vs-GPT lenience offset ($+1.43$ pp averaged
466 across all 16 candidates), because the latter is absorbed by averaging the two judges rather than
467 picking either one. Second, the bias is far smaller than the leaderboard's gross dynamic range:
468 the top-to-bottom Overall gap is roughly 21 pp (gpt-5.4 75.0% vs. gpt-4o 54.0%), and the only
469 Anthropic–OpenAI comparisons whose ordering a ~ 1.5 pp shift could plausibly disturb already
470 sit inside the within-cluster ties that the pairwise paired-bootstrap analysis (Appendix B) flags
471 as statistically indistinguishable (most prominently gpt-5.4 vs. claude-4.6-opus, $p = 0.092$).
472 All Anthropic-vs-OpenAI pairs that *are* significant in Appendix B have point-estimate gaps that
473 substantially exceed the residual family-level bias, so no cross-tier ranking claim in the main text is
474 at risk of flipping under a 1.5 pp adjustment.

475 **What this analysis does not establish.** The two-judge protocol attenuates but does not eliminate
476 the conflict-of-interest concern: the test above is necessary but not sufficient. A truly independent
477 third-party judge—e.g. gemini-3.1-pro, which is a candidate but not currently a judge—would be
478 needed to fully decouple judging from candidate selection. We treat the addition of a third independent
479 judge as the single highest-value follow-up experiment for future versions of AER-bench, and discuss
480 it explicitly in the limitations section.

481 D Robustness Controls for the Historical Concept-Drift Finding

482 The main-text Finding 3 (Section 4) reframes our Concept Identification (CI) results as a *model-*
 483 *stratified, length-conditional* “modernity bias” rather than a uniform property of contemporary LLMs.
 484 This appendix documents the two robustness controls behind the reframing: an item-clustered OLS
 485 regression with chunk-length control, and a per-decade \times per-model breakdown of CI accuracy with
 486 bootstrap confidence intervals.

487 **Why these controls are needed.** The aggregate Era₃ – Era₂ gap of $\sim +5$ pp in Figure 6 can be
 488 confounded by at least three non-conceptual factors: (i) systematic length differences between eras;
 489 (ii) a 30-year era partition that may hide decade-level heterogeneity; and (iii) aggregation across 16
 490 models with potentially very different behaviours. We re-analyse the CI data along all three axes.

491 **Data and setup.** We recover each CI item’s publication year and source-chunk text from the
 492 benchmark’s raw CI files, use chunk character count as a length proxy, and bin items into five decades.
 493 Of 160 CI items, 146 carry a usable year + chunk pair; the decade breakdown is {10, 22, 27, 43, 44}
 494 items. Per-model bootstrap intervals on each decade use $B = 10,000$ item-level resamples.

495 **OLS with chunk-length control.** We estimate the panel regression

$$\text{score}_{m,i} = \alpha_m + \beta_{\text{Era}_3} \mathbb{1}[\text{Era}_3]_i + \beta_{\log \text{len}} \log(\text{len})_i + \beta_x \mathbb{1}[\text{Era}_3]_i \cdot \log(\text{len})_i + \varepsilon_{m,i},$$

496 with model fixed effects α_m , mean-centred $\log(\text{len})_i$, and item-clustered standard errors. We obtain
 497 $\hat{\beta}_{\text{Era}_3} = +0.0275$ ($p = 0.442$), $\hat{\beta}_{\log \text{len}} = -0.1502$ ($p < 10^{-5}$), $\hat{\beta}_x = +0.0572$ ($p = 0.319$). Two
 498 facts follow: (a) the unconditional +5 pp era gap shrinks to +2.75 pp and is no longer distinguishable
 499 from zero once length is controlled for; (b) longer paragraphs are uniformly harder (~ -15 pp per
 500 unit of \log -length), the dominant CI difficulty signal. The Era₃ \times length interaction is positive but not
 501 individually significant.

Table 8: Length-stratified CI accuracy by era. “Gap” is the per-quartile Era₃ – Era₂ difference in mean accuracy; the bottom row reports the OLS coefficients. After controlling for length, the Era₃ main effect is no longer distinguishable from zero.

Length quartile	Era2 acc. (%)	Era3 acc. (%)	Gap (pp)	n_{Era_2}	n_{Era_3}
Q1 (shortest)	69.1	67.3	-1.8	304	288
Q2	49.4	47.3	-2.1	240	336
Q3	48.3	57.5	+9.2	256	320
Q4 (longest)	42.3	55.0	+12.7	144	448

OLS with model fixed effects and item-clustered standard errors:

$$\beta_{\text{Era}_3} = +0.0275 \quad (p = 0.442);$$

$$\beta_{\log \text{len}} = -0.1502 \quad (p = 9.16 \times 10^{-6});$$

$$\beta_{\text{Era}_3 \times \log \text{len}} = +0.0572 \quad (p = 0.319).$$

502 **Length-quartile stratification.** Table 8 reveals a non-linear pattern the OLS coefficients hide: the
 503 per-quartile Era₃ – Era₂ gap is -1.8 pp on Q1, -2.1 pp on Q2, $+9.2$ pp on Q3, and $+12.7$ pp on Q4.
 504 The two shortest quartiles *reverse* the sign of the modernity bias, while the two longest account for
 505 essentially all of the aggregate +5 pp gap—consistent with representational gaps that surface only
 506 when models must reason over several inter-related sentences of pre-1930 prose.

507 **Per-decade, per-model heterogeneity.** Table 9 shows that aggregating across models obscures the
 508 dominant regularity. Along the 1900–1909 vs. 1940–1950 axis, models split into two clean groups:

- 509 • **Earliest-decade-strong** (claude-4.6-sonnet, claude-4.6-opus,
 510 claude-3.7-sonnet, gemini-3.1-pro, gpt-5.2, kimi-k2.5) score *higher* on
 511 1900–1909 than on 1940–1950, with old-to-recent drops between -3 and -22 pp (extreme:
 512 claude-4.6-sonnet 86.2% vs. 64.2%).
- 513 • **Modern-skewed** (gpt-5.4, gpt-5, gpt-5-mini, gemini-3-flash,
 514 gemini-2.5-flash, glm-5.1) show the canonical pattern: 19–41% on 1900–1909 rising
 515 roughly monotonically to 46–59% on 1940–1950.

Table 9: Per-decade CI accuracy (%) with item-level 95% bootstrap confidence intervals ($B = 10,000$). Models are ordered by their 1900–1909 score; the *Items* row gives the per-decade sample size, so intervals on the $n = 10$ column are correspondingly wide.

Model	1900-09	1910-19	1920-29	1930-39	1940-50
<i>Items</i>	10	22	27	43	44
claude-4.6-sonnet	86.2 [65.0, 100.0]	80.1 [66.5, 92.6]	78.2 [66.7, 88.4]	67.2 [57.6, 76.2]	64.2 [52.8, 75.6]
claude-4.6-opus	82.5 [62.5, 97.5]	83.5 [71.0, 94.3]	69.4 [54.6, 82.9]	69.5 [59.9, 78.5]	65.3 [54.0, 75.9]
claude-3.7-sonnet	76.2 [48.8, 97.5]	80.7 [66.5, 92.6]	79.2 [67.1, 89.4]	64.8 [54.4, 75.0]	67.3 [55.7, 77.6]
gemini-3.1-pro	81.2 [58.8, 96.2]	75.6 [59.1, 90.3]	75.0 [63.4, 86.1]	63.7 [54.1, 72.7]	72.2 [63.6, 80.1]
gpt-5.2	61.3 [32.5, 86.2]	66.5 [47.7, 83.5]	68.1 [53.2, 81.9]	57.0 [46.8, 67.2]	54.8 [44.0, 65.9]
kimi-k2.5	61.3 [33.8, 85.0]	51.7 [33.5, 69.9]	61.1 [46.8, 74.5]	58.4 [49.1, 67.4]	57.4 [47.4, 67.3]
gpt-5.4	41.2 [15.0, 68.8]	52.8 [35.2, 71.0]	56.9 [41.7, 71.8]	64.8 [55.8, 73.8]	59.4 [48.0, 70.5]
doubao-seed-1.6	66.2 [41.2, 87.5]	45.5 [27.8, 63.6]	54.2 [39.8, 68.5]	45.6 [36.3, 54.9]	51.7 [41.8, 61.9]
gpt-5.1	65.0 [40.0, 86.2]	51.1 [33.0, 69.3]	42.6 [28.2, 57.4]	52.3 [42.7, 61.9]	49.7 [38.9, 60.5]
gpt-4o	50.0 [22.5, 76.2]	44.9 [27.3, 63.6]	51.9 [37.0, 66.7]	57.3 [47.4, 66.6]	52.6 [41.8, 62.8]
qwen3-vl-32b-instruct	46.2 [20.0, 72.5]	42.6 [25.0, 60.8]	51.9 [37.0, 66.2]	50.0 [40.4, 59.6]	45.2 [35.2, 55.7]
gpt-5	31.2 [7.5, 57.5]	31.8 [15.9, 49.4]	45.8 [31.5, 60.6]	52.9 [42.1, 63.1]	56.2 [45.2, 67.6]
glm-5.1	31.2 [6.2, 57.5]	31.8 [15.9, 48.9]	45.4 [31.0, 59.7]	48.0 [37.8, 57.8]	55.7 [45.7, 65.9]
gemini-2.5-flash	30.0 [6.2, 56.2]	29.5 [14.2, 46.6]	43.1 [28.7, 57.4]	48.8 [38.1, 58.7]	52.8 [42.3, 63.4]
gpt-5-mini	18.8 [0.0, 38.8]	29.0 [14.2, 45.5]	47.2 [32.9, 62.0]	53.2 [43.3, 63.1]	46.9 [35.8, 57.7]
gemini-3-flash	18.8 [0.0, 37.5]	32.4 [15.9, 50.0]	40.3 [25.9, 54.2]	49.4 [38.7, 59.9]	46.3 [34.9, 57.1]

516 Mid-tier models (*doubao-seed-1.6*, *gpt-5.1*, *qwen3-vl-32b-instruct*, *gpt-4o*) sit between
517 the groups with comparatively flat profiles. Earliest-decade intervals are wide ($n = 10$), so the ~ 20
518 pp within-row swings should be read alongside the 1–2 pp swings on the densely-sampled 1940–1950
519 column.

520 **What the reframing buys us.** Together, the three controls support a much narrower claim than
521 “current AI systems exhibit historical concept drift”: the era main effect vanishes after length control
522 ($p = 0.442$); the descriptive gap concentrates in the longest two quartiles; and the per-decade structure
523 differs sharply by model, with frontier Anthropic / Gemini-3.1-pro / GPT-5.2 systems scoring *better*
524 on the earliest decade than on the most recent. We therefore frame the CI task as a *diagnostic for*
525 *which models still carry a representational blind spot on early-20th-century economic prose*, rather
526 than as evidence of a uniform property of contemporary LLMs. The same controls flag the two most
527 informative follow-ups: (i) inter-annotator agreement on the 1900–1909 subset, where bootstrap
528 intervals are widest, and (ii) controlled experiments holding chunk length fixed across eras.

529 E Detailed Rubric Definitions

530 This appendix provides detailed rubric definitions for all LLM-as-a-Judge tasks in AER-bench. Each
531 task employs multi-dimensional scoring to prevent holistic biases and ensure fine-grained evaluation
532 of model capabilities.

533 E.1 Mathematical Modeling (MM)

534 Mathematical Modeling evaluates the ability to derive equations, solve optimization problems, and
535 perform formal mathematical reasoning in economic contexts. Each question is scored across four
536 dimensions (0–10 points each, total 40 points):

537 **Correctness (Weight: 4)** Assesses whether formulas are mathematically equivalent to the reference
538 answer and derivation steps are valid.

- 539 • **9–10:** All formulas are mathematically equivalent to the reference; derivation steps are valid.
- 540 • **6–8:** Minor errors (sign, index, simplification) but core structure is correct.
- 541 • **3–5:** Significant errors but partial understanding of the approach is evident.
- 542 • **0–2:** Fundamentally incorrect or mostly missing.

543 **Notation Consistency (Weight: 2)** Evaluates whether symbols match the problem setup and are
544 used consistently throughout the derivation.

- 545 • **9–10:** All symbols match the problem setup and are used consistently.
- 546 • **6–8:** Minor inconsistencies that do not affect readability.
- 547 • **3–5:** Some symbols undefined or used inconsistently.
- 548 • **0–2:** Major notation confusion or deviation from given symbols.

549 **Completeness (Weight: 3)** Checks whether all intermediate steps are shown and a reader can
550 verify each step.

- 551 • **9–10:** All intermediate steps shown; a reader can verify each step.
- 552 • **6–8:** Key steps present but some intermediate algebra omitted.
- 553 • **3–5:** Major steps missing; jumps to conclusions without justification.
- 554 • **0–2:** Only final answer given, or derivation is fragmentary.

555 **Economic Intuition (Weight: 1)** Assesses whether the candidate provides clear interpretation of
556 results in economic terms and explains what the mathematics means.

- 557 • **9–10:** Clear interpretation of results in economic terms; explains what the math means.
- 558 • **6–8:** Some economic interpretation provided.
- 559 • **3–5:** Minimal economic context; mostly mechanical derivation.
- 560 • **0–2:** No economic interpretation at all.

561 E.2 Paper Abstraction (PA)

562 Paper Abstraction evaluates the ability to synthesize a condensed paper body (approximately 30%
563 of original length, excluding the abstract) into a coherent abstract that captures key contributions,
564 methods, and findings. Scoring is based on four dimensions (0–10 points each, total 40 points):

565 **Coverage (Weight: 4)** Measures whether all or nearly all reference key points (atomic semantic
566 propositions extracted from the original abstract) are mentioned or clearly implied.

- 567 • **9–10:** All or nearly all reference key points are mentioned or clearly implied.
- 568 • **6–8:** Most key points covered; 1–2 minor omissions.

- 569 • **3–5:** Several important points missing.
- 570 • **0–2:** Only 1–2 points covered, or major elements absent.

571 **Accuracy (Weight: 4)** Assesses factual correctness of all stated claims per the paper content.

- 572 • **9–10:** All stated claims are factually correct per the paper content.
- 573 • **6–8:** Minor inaccuracies that do not fundamentally mislead.
- 574 • **3–5:** Some factual errors or misrepresentations of results.
- 575 • **0–2:** Major factual errors that distort the paper’s findings.

576 **Conciseness (Weight: 3)** Evaluates whether the abstract stays within the specified word limit with
577 no redundancy.

- 578 • **9–10:** Within word limit; no redundancy; every sentence adds value.
- 579 • **6–8:** Mostly concise; minor redundancy or slightly over limit.
- 580 • **3–5:** Noticeably wordy or exceeds word limit significantly.
- 581 • **0–2:** Excessively long or highly redundant.

582 **Academic Style (Weight: 1)** Judges whether the abstract reads like a published abstract in a top
583 economics journal.

- 584 • **9–10:** Reads like a published abstract in a top economics journal.
- 585 • **6–8:** Generally academic but with some informal phrasing.
- 586 • **3–5:** Mixed register; some casual language or poor structure.
- 587 • **0–2:** Not academic in style.

588 **E.3 Literature Review (LR)**

589 Literature Review evaluates the ability to reconstruct an introduction section that aligns with the
590 paper’s citation network and argumentation structure. This task employs a hybrid scoring approach
591 combining citation recall and writing quality.

592 **Citation Score (Weight: 3)** Measures weighted recall of citations across three tiers:

- 593 • **Tier 1 (weight 0.8):** Core citations indispensable to the argument, typically from top journals
594 or highly influential sources.
- 595 • **Tier 2 (weight 0.2):** Supplementary citations from reputable academic sources.
- 596 • **Tier 3 (weight 0.0):** Optional citations with lower relevance or authority; not scored.

597 The citation score is computed as: $\text{citation_score} = 0.8 \times \text{Tier1_recall} + 0.2 \times \text{Tier2_recall}$.

598 **Rubric Score (Weight: 2)** Assesses writing quality across paper-specific dimensions tailored to
599 the introduction’s argumentation structure. Common dimensions include:

- 600 • **Research question clarity:** Is the core question clearly stated?
- 601 • **Literature positioning:** Does the introduction position the paper within existing literature?
- 602 • **Contribution statement:** Are the paper’s contributions explicitly articulated?
- 603 • **Methodological preview:** Is the research approach previewed appropriately?
- 604 • **Logical flow:** Does the introduction follow a coherent argumentative structure?

605 Each dimension is scored 0–10, with weights summing to 100. The rubric is dynamically constructed
606 for each paper based on its specific introduction structure.

607 **E.4 Model Generation (MG)**

608 Model Generation evaluates the ability to construct a complete formal economic model from a
609 desymbolized narrative description. Unlike other tasks, MG uses a flexible element-based rubric
610 where the number and type of required elements vary by question.

611 **Element Rubric (Flexible Dimensions)** Each model is decomposed into constituent elements
612 (agents, objective functions, constraints, equilibrium concepts, first-order conditions). Each element
613 is scored based on presence and correctness:

- 614 • **Full (100%)**: Element is present, correctly formalized, uses proper mathematical notation,
615 and is economically meaningful.
- 616 • **Substantial (70%)**: Element is present and mostly correct, but has minor issues (e.g.,
617 missing a variable, slight notation inconsistency, or incomplete specification).
- 618 • **Partial (40%)**: Element is mentioned or attempted, but has significant issues (e.g., wrong
619 functional form, missing key variables, or conceptually flawed formalization).
- 620 • **Minimal (15%)**: Element is only vaguely referenced without proper formalization.
- 621 • **Missing (0%)**: Element is not present at all.

622 The total score is the sum of all element scores, normalized to 0–100. The maximum possible score
623 and element weights are determined by the reference model’s complexity. For example, a model with
624 2 agents, 2 objective functions, 3 constraints, 1 equilibrium concept, and 2 FOCs would have 10
625 elements, each contributing proportionally to the total score.

626 **E.5 Economic Analysis (EA)**

627 Economic Analysis evaluates the ability to interpret empirical or theoretical results with substantive
628 economic reasoning. Given a purely factual data description (e.g., regression coefficients, significance
629 levels, trends), the model must produce a professional analysis paragraph. Scoring is based on six
630 dimensions (0–10 points each), with applicability determined per question:

631 **Accuracy (Weight: 2)** Assesses whether numerical descriptions and trend characterizations are
632 factually correct.

- 633 • **9–10**: All numerical and directional claims are accurate; no fabrication.
- 634 • **6–8**: Minor inaccuracies that do not fundamentally mislead.
- 635 • **3–5**: Some factual errors in data description.
- 636 • **0–2**: Major errors or fabricated numbers.

637 **Economic Interpretation (Weight: 2)** Evaluates whether the analysis provides substantive eco-
638 nomic meaning (e.g., magnitude interpretation, heterogeneity discussion).

- 639 • **9–10**: Rich interpretation of economic significance and practical implications.
- 640 • **6–8**: Reasonable interpretation provided.
- 641 • **3–5**: Superficial or generic interpretation.
- 642 • **0–2**: No substantive interpretation.

643 **Mechanism (Weight: 2)** Assesses whether economic channels or causal pathways are identified.

- 644 • **9–10**: Clear identification of driving mechanisms.
- 645 • **6–8**: Some mechanisms discussed but not fully developed.
- 646 • **3–5**: Vague or incomplete mechanism discussion.
- 647 • **0–2**: No mechanism identified.

648 **Theory Connection (Weight: 1.5)** Evaluates whether results are connected to theoretical predic-
649 tions or prior work.

- 650 • **9–10:** Strong theoretical grounding with specific references.
- 651 • **6–8:** Some theoretical context provided.
- 652 • **3–5:** Weak or generic theory connection.
- 653 • **0–2:** No theory connection.

654 **Policy Implication (Weight: 1.5)** Assesses whether policy implications are discussed appropriately
655 (only applicable when results support policy discussion).

- 656 • **9–10:** Nuanced, well-grounded policy discussion.
- 657 • **6–8:** Reasonable policy points raised.
- 658 • **3–5:** Superficial policy discussion.
- 659 • **0–2:** No policy discussion.
- 660 • **N/A:** Not applicable for this result.

661 **Writing Quality (Weight: 1)** Judges whether the analysis is well-structured and clearly written.

- 662 • **9–10:** Publication-ready prose; logical flow; precise language.
- 663 • **6–8:** Clear and organized; minor style issues.
- 664 • **3–5:** Readable but poorly organized or imprecise.
- 665 • **0–2:** Disorganized or unclear.

666 **Note:** Not all dimensions are applicable to every question. For example, a purely theoretical result may
667 not have policy implications. The normalized score is computed as: $\text{normalized_score} = \frac{\text{total}}{\text{max_possible}}$,
668 where max_possible counts only applicable dimensions (each worth 10 points).

669 **E.6 Judge Model Configuration**

670 All LLM-as-a-Judge evaluations in AER-bench use the following judge models:

- 671 • **claude-4.6-sonnet** (Anthropic)
- 672 • **gpt-5.2** (OpenAI, accessed via Azure)

673 For each generative task, both judges independently score all dimensions. The final item score is
674 the arithmetic mean of the two judges' scores. Section 4.2 presents a systematic comparison of
675 inter-judge agreement, demonstrating a mean absolute difference of 2.6 percentage points across all
676 models and tasks, validating the robustness of this multi-judge approach.

677 **F Detailed Performance Analysis by Task Dimension**

678 This appendix provides a comprehensive breakdown of model performance across all seven task
 679 dimensions in AER-bench. Figure 7 presents a visual summary of the overall ranking table (Table 3),
 680 showing task-level performance for all evaluated models grouped by provider. Each subsection then
 681 presents a heatmap visualization showing dimension-level scores for all evaluated models, followed
 682 by detailed analysis of performance patterns, strengths, and weaknesses.

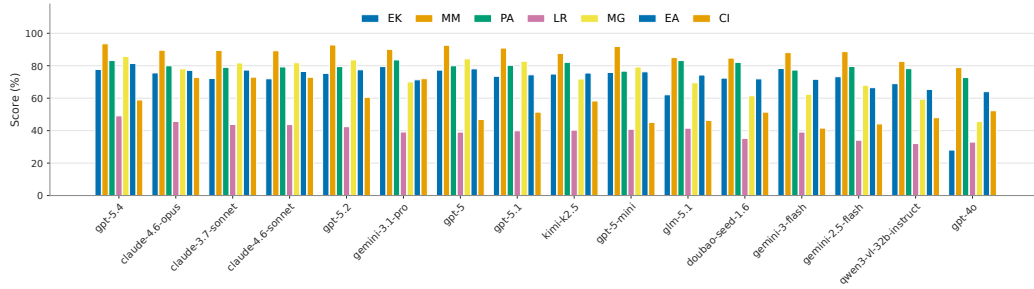


Figure 7: Task-level performance across all evaluated models, grouped by provider. This figure provides a visual complement to Table 3, showing the distribution of scores across the seven task dimensions (EK, MM, PA, LR, MG, EA, CI) for each model.

683 **F.1 Economic Knowledge (EK)**

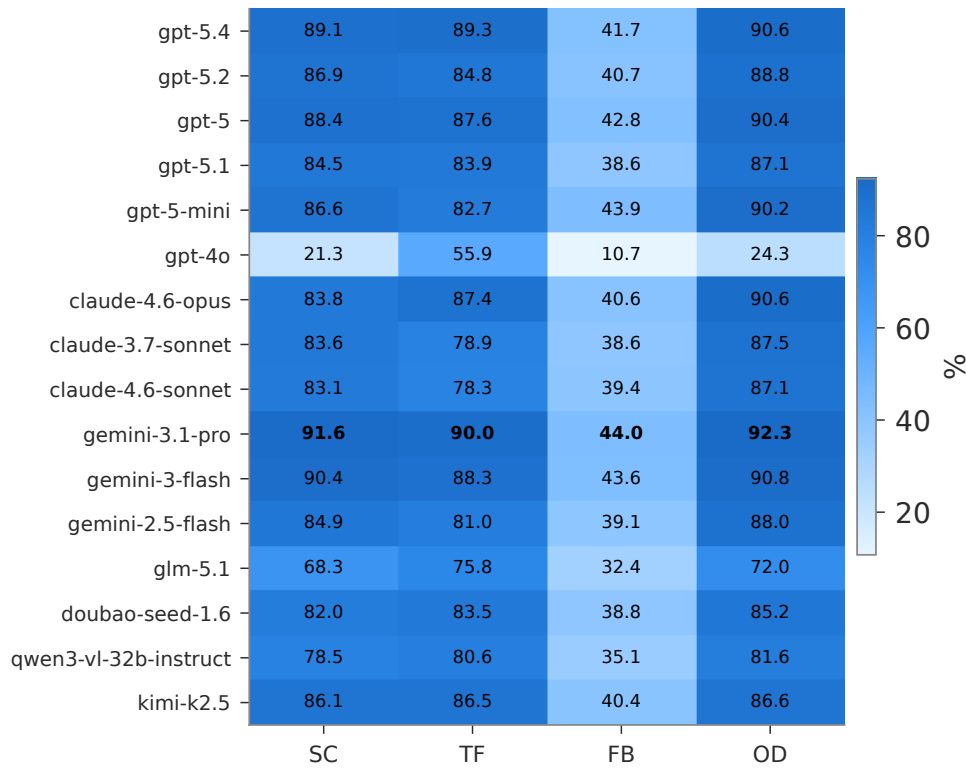


Figure 8: Performance breakdown across Economic Knowledge (EK) question types: Single Choice (SC), True/False (TF), Fill-in-Blank (FB), and Ordering (OD). Values represent percentage scores (0–100).

684 The Economic Knowledge task evaluates models across four question formats: single choice,
 685 true/false, fill-in-blank, and ordering. As shown in Figure 8, performance varies significantly
 686 across these formats. Google’s gemini-3.1-pro achieves the highest overall EK score (79.5%),
 687 demonstrating particularly strong performance on single-choice questions. OpenAI’s GPT-5 series
 688 models show consistent performance across all four formats, with gpt-5.4 and gpt-5 both scoring
 689 above 77%.

690 Notably, fill-in-blank questions prove most challenging across all models, with average scores
 691 approximately 40–50 percentage points lower than single-choice questions (e.g., gemini-3.1-pro
 692 scores 91.6% on single choice but only 44.0% on fill-in-blank). This suggests that open-ended
 693 knowledge recall is fundamentally more difficult than recognition-based tasks. Ordering questions
 694 generally yield the highest scores across models (up to 92.3%), likely because they test relational
 695 reasoning among familiar economic concepts. True/false questions also show high performance,
 696 though slightly below ordering.

697 Among the lower-performing models, gpt-4o shows a marked weakness in this dimension (28.1%
 698 overall), particularly struggling with fill-in-blank and ordering tasks. This indicates that older model
 699 generations lack the economic domain knowledge depth required for AER-bench tasks.

700 **F.2 Mathematical Modeling (MM)**

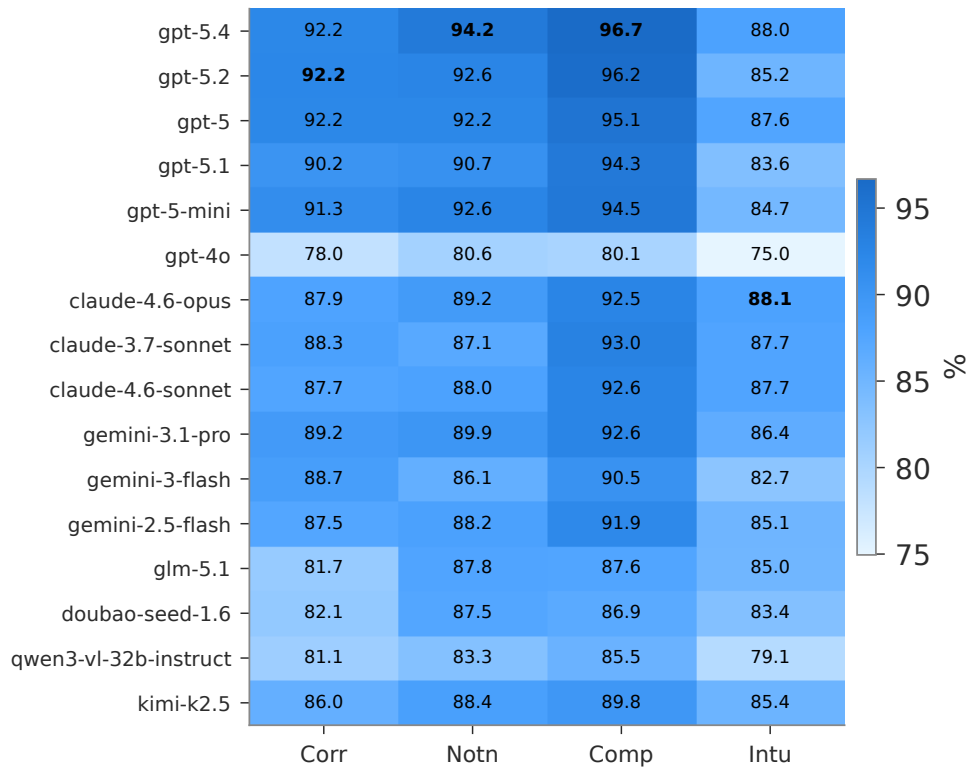


Figure 9: Performance breakdown across Mathematical Modeling (MM) dimensions: Correctness (Corr), Notation (Notn), Completeness (Comp), and Economic Intuition (Intu). Values represent percentage scores (0–100).

701 Mathematical Modeling represents the highest-scoring task across all models, with most frontier
 702 models exceeding 85% overall performance. Figure 9 reveals that gpt-5.4 achieves exceptional
 703 performance (93.5%), with near-perfect scores across all four evaluation dimensions: correctness,
 704 notation, completeness, and economic intuition.

705 The heatmap shows remarkably consistent performance across dimensions for top-tier models.
 706 Correctness scores are universally high (>90%) for GPT-5 series models (gpt-5.4, gpt-5.2, gpt-5

707 all at 92.2%, gpt-5-mini at 91.3%, gpt-5.1 at 90.2%), indicating strong mathematical reasoning
 708 capabilities. Claude 4.6 models achieve slightly lower but still strong correctness scores (87.7–87.9%).
 709 Notation scores, which assess proper use of mathematical symbols and formatting, are similarly high,
 710 suggesting models have learned standard economic notation conventions from training data.

711 Completeness and economic intuition dimensions show slightly more variation. Economic intuition,
 712 which evaluates whether models provide meaningful economic interpretation alongside mathematical
 713 derivations, proves most challenging. Models like gemini-3-flash and qwen3-vl-32b-instruct
 714 show a 5–8 point gap between correctness and intuition scores, indicating they can solve equations
 715 correctly but struggle to explain economic meaning.

716 The uniformly high MM scores across models suggest that mathematical manipulation and symbolic
 717 reasoning are well-developed capabilities in current LLMs, likely benefiting from extensive training
 718 on mathematical and scientific texts.

719 F.3 Paper Abstraction (PA)

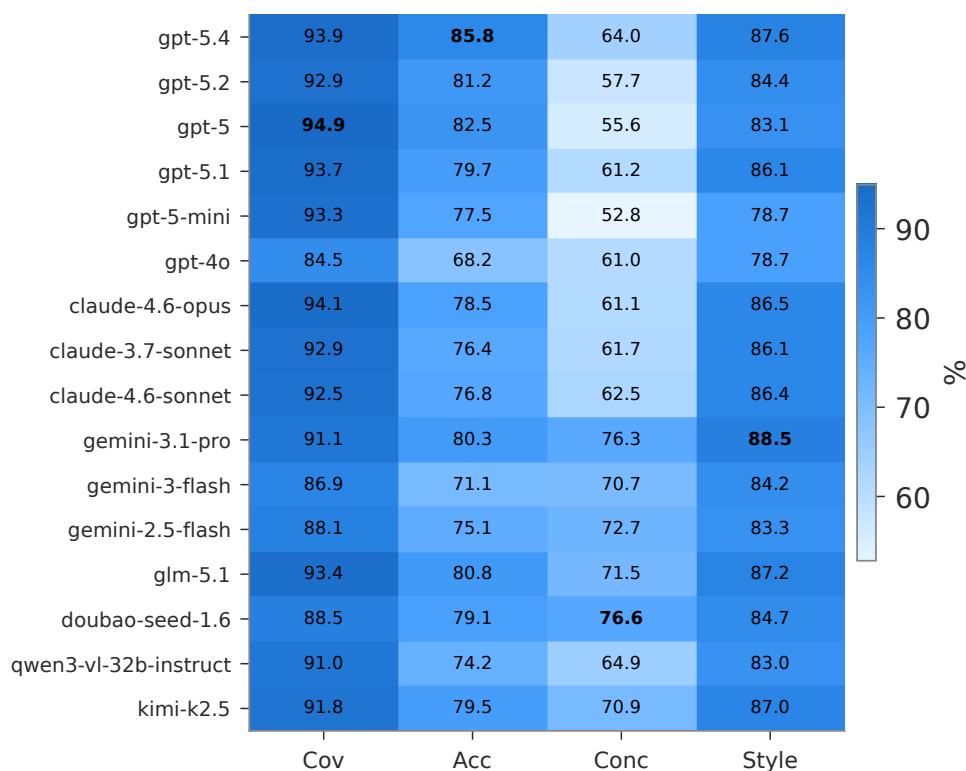


Figure 10: Performance breakdown across Paper Abstraction (PA) dimensions: Coverage (Cov), Accuracy (Acc), Conciseness (Conc), and Academic Style (Style). Values represent percentage scores (0–100).

720 Paper Abstraction tasks require models to summarize and analyze economic research papers across
 721 four dimensions: coverage, accuracy, conciseness, and academic style. Figure 10 shows that
 722 gemini-3.1-pro leads this task (83.6%), with particularly strong performance in coverage and
 723 accuracy dimensions.

724 Coverage scores, measuring whether summaries capture all key contributions, show the widest per-
 725 formance spread among dimensions. Top models (gpt-5, gpt-5.4, claude-4.6-opus, gpt-5.1,
 726 glm-5.1) achieve 91–95% coverage, while mid-tier models range from 86–91%. This suggests that
 727 identifying salient information from lengthy academic papers requires sophisticated comprehension
 728 capabilities.

729 Accuracy scores are generally high across all models, indicating that when models do extract infor-
 730 mation, they rarely introduce factual errors. However, conciseness proves challenging—many models
 731 produce verbose summaries that include excessive detail. The heatmap reveals that Claude mod-
 732 els (claude-4.6-opus, claude-3.7-sonnet, claude-4.6-sonnet) show particularly balanced
 733 performance across all four dimensions, suggesting well-calibrated summarization capabilities.

734 Academic style scores, evaluating whether summaries use appropriate scholarly language and struc-
 735 ture, correlate strongly with overall model capability. Frontier models consistently score above 80%,
 736 while smaller or older models (gpt-4o, qwen3-vl-32b-instruct) drop below 70%, indicating
 737 that academic writing conventions require substantial model capacity to master.

738 F.4 Literature Review (LR)

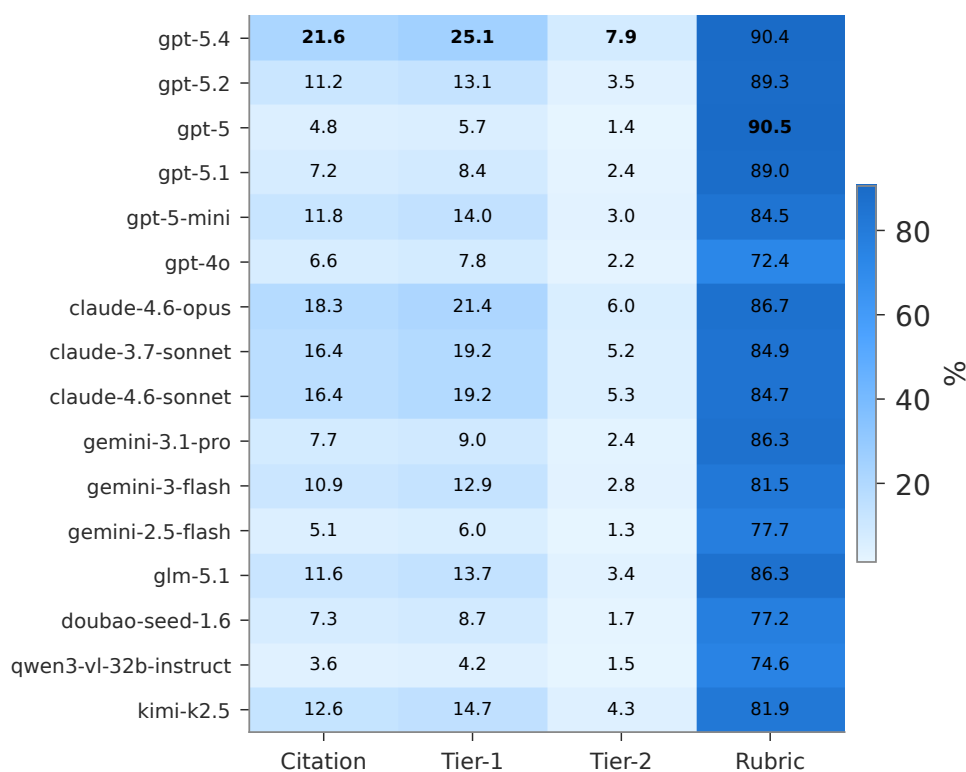


Figure 11: Performance breakdown across Literature Review (LR) dimensions: Citation Score (Cite) and Rubric Score (Rubric). Values represent percentage scores (0–100).

739 Literature Review is the most challenging task in AER-bench, with the highest-performing model
 740 (gpt-5.4) achieving only 49.1%. Figure 11 decomposes performance into citation score (accuracy
 741 of cited papers) and rubric score (quality of synthesis and analysis). The heatmap reveals that both
 742 dimensions are extremely difficult, with most models scoring below 45% on each.

743 Citation scores measure whether models correctly identify and cite relevant papers from a provided
 744 reference list. The extremely low scores indicate that models struggle with precise bibliographic
 745 matching and relevance assessment. Even the best-performing model, gpt-5.4, achieves only
 746 21.6% citation accuracy, followed by claude-4.6-opus at 18.3%. Most models score below 15%,
 747 suggesting that fine-grained citation retrieval and bibliographic matching remain a fundamental
 748 weakness for current LLMs.

749 Rubric scores, evaluating the quality of literature synthesis, thematic organization, and critical
 750 analysis, are substantially higher. Top models achieve 85–90% rubric scores (e.g., gpt-5 at 90.5%,
 751 gpt-5.4 at 90.4%), indicating that models can produce well-structured literature review prose when
 752 the task does not require precise citation matching. The stark contrast between citation and rubric

753 scores reveals that the primary bottleneck in literature review is not analytical writing quality but
 754 rather the ability to accurately identify and attribute relevant sources.

755 Interestingly, the performance gap between frontier and mid-tier models is smaller in LR than in
 756 other tasks. This suggests that literature review capabilities scale less predictably with model size and
 757 training, possibly because the task requires specialized academic reasoning patterns underrepresented
 758 in general training data.

759 **F.5 Method Generation (MG)**

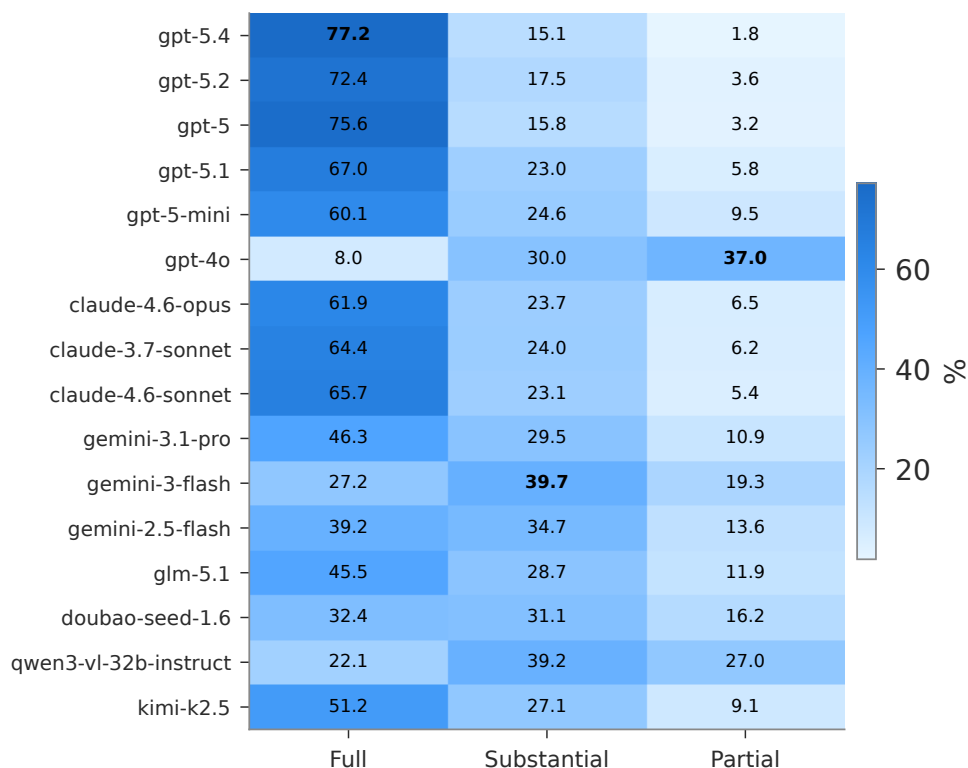


Figure 12: Performance breakdown across Method Generation (MG) rubric elements. Due to the flexible rubric structure, dimension labels vary by specific task item. Values represent percentage scores (0–100).

760 Method Generation tasks evaluate models’ ability to design novel research methodologies, including
 761 identification strategies, data collection plans, and empirical specifications. Figure 12 shows perform-
 762 ance across multiple rubric elements, which vary by specific task item due to the flexible evaluation
 763 structure.

764 gpt-5.4 leads this task (85.8%), followed closely by gpt-5 (84.2%) and gpt-5.1 (82.8%). The
 765 heatmap reveals that top models achieve consistently high scores (80–90%) across most rubric
 766 elements, indicating strong methodological reasoning. However, certain elements—particularly those
 767 requiring creative problem-solving or domain-specific knowledge of econometric techniques—show
 768 greater performance variation.

769 Claude models demonstrate balanced performance across rubric elements, with claude-3.7-sonnet
 770 achieving 81.7%. Gemini models show more uneven performance, with gemini-3.1-pro scoring
 771 only 70.0% overall despite strong performance in other tasks. This suggests that method generation
 772 requires different capabilities than comprehension or analysis tasks.

773 Lower-tier models (doubao-seed-1.6, qwen3-vl-32b-instruct) score below 65%, indicating
 774 that generating methodologically sound research designs requires substantial reasoning capacity.

775 The task demands understanding of causal inference principles, data requirements, and econometric
 776 techniques—knowledge that appears concentrated in frontier models.

777 **F.6 Economic Analysis (EA)**

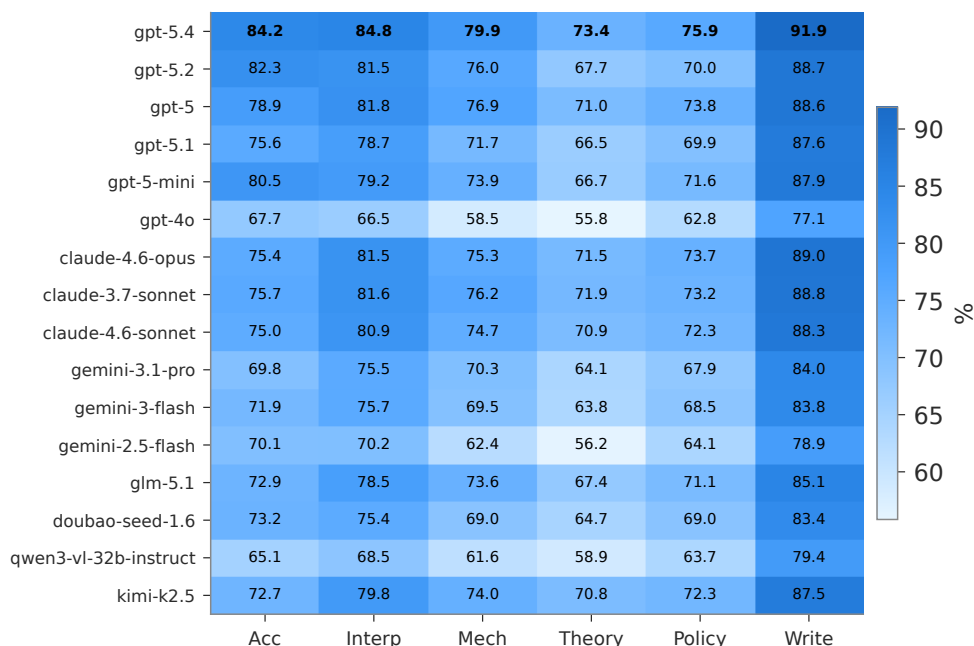


Figure 13: Performance breakdown across Economic Analysis (EA) dimensions: Accuracy (Acc), Economic Interpretation (Interp), Mechanism (Mech), Theory Connection (Theory), Policy Implication (Policy), and Writing Quality (Write). Values represent percentage scores (0–100).

778 Economic Analysis tasks require models to interpret regression results, explain economic mechanisms,
 779 connect findings to theory, and discuss policy implications. Figure 13 shows performance across six
 780 dimensions, revealing substantial variation in model capabilities.

781 gpt-5.4 achieves the highest overall EA score (81.4%), with particularly strong performance in
 782 accuracy and writing quality dimensions. The heatmap shows that accuracy scores—measuring
 783 correct interpretation of statistical results—range from 65–84% across models, with frontier models
 784 (gpt-5.4 at 84.2%, gpt-5.2 at 82.3%, gpt-5-mini at 80.5%) scoring above 80%, indicating solid
 785 quantitative reasoning. However, deeper analytical dimensions show more variation.

786 Economic interpretation and mechanism explanation prove moderately challenging, with most models
 787 scoring 70–85%. These dimensions require understanding causal pathways and economic logic
 788 beyond statistical patterns. Theory connection is the most challenging EA dimension, with scores
 789 ranging from 56–73%. Even top models like gpt-5.4 achieve only 73.4% on theory connection,
 790 suggesting that linking empirical findings to established theoretical frameworks requires sophisticated
 791 domain knowledge.

792 Policy implication scores are slightly higher than theory connection (ranging 63–76%), but still
 793 represent a significant challenge. This dimension requires extrapolating from research findings to
 794 real-world applications, considering practical constraints and normative judgments.

795 Writing quality scores are uniformly high (>75%) across frontier models, indicating that models can
 796 produce well-structured, coherent analytical text. The strong correlation between writing quality and
 797 overall EA performance suggests that clear communication and analytical depth are closely linked
 798 capabilities.

799 **E.7 Concept Identification (CI)**

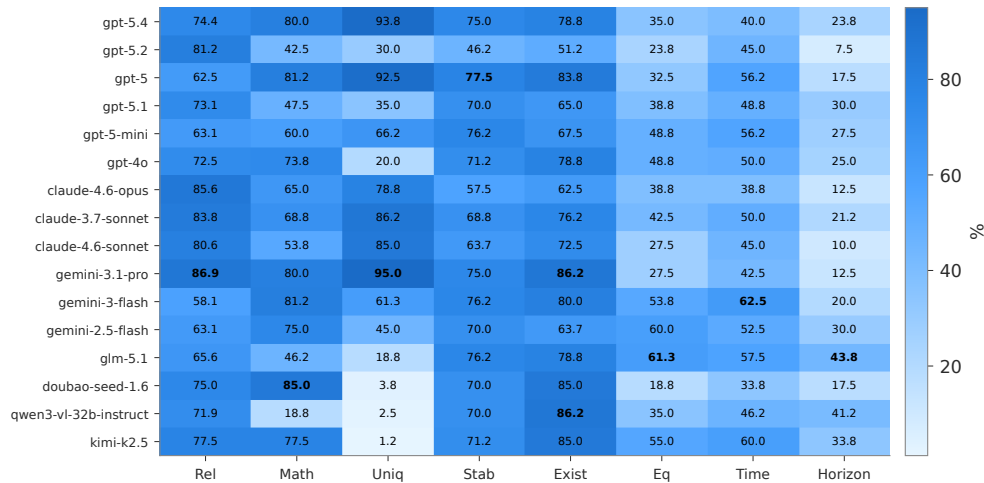


Figure 14: Performance breakdown across Concept Identification (CI) dimensions: Relevant, Math Formalization, Uniqueness, Stability, Existence, Equilibrium Scene, Time Characteristic, and Time Horizon. Values represent percentage scores (0–100).

800 Concept Identification tasks evaluate models’ ability to identify and classify economic concepts
 801 from historical AER papers, assessing understanding across eight dimensions. Figure 14 reveals that
 802 this task shows the highest performance variance across models, with `claude-3.7-sonnet` leading
 803 (73.0%) while `gpt-5` scores only 46.9%.

804 The heatmap shows that Claude models (`claude-4.6-opus`, `claude-3.7-sonnet`,
 805 `claude-4.6-sonnet`) consistently outperform other providers on CI tasks, achieving 72.7–
 806 73.0% overall scores. However, their performance varies dramatically across individual dimensions,
 807 ranging from as low as 10% (time horizon for `claude-4.6-sonnet`) to as high as 86.2%
 808 (uniqueness for `claude-3.7-sonnet`). This suggests that Anthropic’s training or architecture may
 809 be particularly well-suited for certain aspects of conceptual classification and historical document
 810 understanding, while struggling with temporal reasoning dimensions.

811 Relevance scores are generally the highest across models (ranging 58–87%), indicating that iden-
 812 tifying whether a concept is economically relevant is relatively straightforward. However, other
 813 dimensions show much greater variance and difficulty. Math formalization scores range from 18.8%
 814 to 85.0%, with some models excelling (`doubao-seed-1.6`: 85.0%, `gpt-5`: 81.2%) while others
 815 struggle (`qwen3-vl-32b-instruct`: 18.8%). More nuanced dimensions—uniqueness, stability,
 816 existence—also prove challenging, with uniqueness showing extreme variance (1.2–95.0%) across
 817 models.

818 The equilibrium scene, time characteristic, and time horizon dimensions show the greatest perfor-
 819 mance spread. These require understanding subtle distinctions in how economic concepts are defined
 820 and applied across different historical contexts. The low scores suggest that temporal reasoning and
 821 historical concept evolution are particularly difficult for current LLMs.

822 Interestingly, GPT-5 series models show unusually weak CI performance relative to their strong
 823 performance on other tasks. This may indicate that the training data or optimization objectives
 824 for these models prioritize contemporary language patterns over historical economic terminology,
 825 creating a systematic weakness in this dimension.

826 G Task Generation Methodology

827 This appendix provides comprehensive documentation of the multi-round generation pipeline used
828 to construct evaluation items for each task dimension in AER-bench. All tasks except Concept
829 Identification (CI) employ LLM-based generation with carefully designed prompts and validation
830 steps. CI items are manually annotated from historical AER papers and are therefore not covered in
831 this appendix.

832 For each task, we describe: (1) the complete multi-round generation workflow, (2) key design rationale,
833 (3) full prompt templates with all critical instructions, (4) complete output format specifications, and
834 (5) scoring rules and evaluation criteria. All generation uses gpt-4o as the extraction model.

835 G.1 Economic Knowledge (EK)

836 G.1.1 Task Overview

837 Economic Knowledge (EK) is a Level 1 (foundational capability) task that evaluates models’ depth of
838 understanding of economic literature. The core evaluation goal is to test whether a model can answer
839 questions about specific papers’ model setups, core findings, and methodological choices. Each
840 paper produces 20 objective questions: 5 single-choice, 5 true/false, 5 fill-in-blank, and 5 ordering
841 questions. The task depends only on PDF input (no supplementary materials required) and uses a
842 2-round generation pipeline.

843 G.1.2 Generation Workflow

844 EK employs a two-round pipeline to ensure balanced knowledge coverage and question diversity:

- 845 1. **Round 1 (PDF input):** Extract 12–15 paper-specific knowledge points spanning six cate-
846 gories: model setup, findings, methods, mechanisms, contributions, and institutional context.
847 Each knowledge point is tagged as *quantitative* (involving specific numbers) or *conceptual*
848 (involving qualitative understanding), with a target of at least 50% conceptual points to
849 avoid over-extraction of numerical facts.
- 850 2. **Round 2 (text input):** Based on the structured knowledge point list from Round 1, generate
851 20 objective questions: 5 single-choice, 5 true/false, 5 fill-in-blank, and 5 ordering questions.
852 Each question references the source knowledge point and includes difficulty tags.

853 G.1.3 Design Rationale

854 Separating knowledge extraction from question generation prevents two common failure modes: (1)
855 knowledge point omission when simultaneously reading the paper and generating 20 questions, and
856 (2) uneven distribution across question types or content categories. The structured knowledge point
857 list serves as a checklist to ensure comprehensive coverage.

858 **Critical Design Constraint** All questions are designed to test whether an AI system or economics
859 researcher **knows this specific paper**. Questions provide ONLY the first author’s surname + year + a
860 short topic phrase (e.g., “Autor et al. (2013), a paper on Chinese import competition and U.S. labor
861 markets”). The full paper title, journal name, and co-author details do NOT appear in the stem. The
862 answer must require knowledge of the paper’s specific details that cannot be guessed from the topic
863 phrase alone.

864 G.1.4 Round 1 Prompt Template

865 You are a professor of economics. Please read the attached PDF paper
866 in full and extract ****paper-specific knowledge points**** that can be
867 used for exam questions.

868
869 **### Critical design constraint**

870
871 The exam questions you will later create are designed to test whether
872 an AI system or an economics researcher ****knows this specific paper****.
873 Therefore:

874
875 - ****ALL knowledge points must be about this paper's specific content****:
876 its model setup, assumptions, variables, findings, methodology,
877 mechanisms, contributions, or institutional context.
878 - The questions will provide **ONLY** the ****first author's surname + year**
879 **+ a short topic phrase**** (e.g., "Autor et al. (2013), a paper on
880 Chinese import competition and U.S. labor markets"). The full paper
881 title, journal name, and co-author details will **NOT** appear in the
882 stem. The ****answer must require knowledge of the paper's specific**
883 **details**** that cannot be guessed from the topic phrase alone.
884 - A well-read economist who has studied this paper (or can search for
885 it) should be able to answer; one who has not should find it difficult.
886
887 **### Knowledge point categories**
888
889 Extract knowledge points spanning as many of the following categories
890 as the paper supports:
891
892 - 'model_setup': Specific assumptions, agent definitions, variable
893 choices, functional forms, or information structures in the paper's
894 model
895 - 'finding': Core empirical or theoretical results -- coefficient
896 directions, magnitudes, significance levels, comparative statics,
897 or proposition statements
898 - 'method': Identification strategy, estimation approach, instrument
899 choice, data sources, or sample construction decisions and their
900 rationale
901 - 'mechanism': Economic mechanisms or causal channels that the paper
902 identifies, tests, or discusses
903 - 'contribution': The paper's stated contributions, how it differs
904 from or extends prior work
905 - 'institution': Institutional background, policy details, or market
906 features specific to the paper's setting
907
908 **### CRITICAL: Content type balance**
909
910 You **MUST** extract a balanced mix of ****quantitative**** and ****conceptual****
911 knowledge points:
912
913 - ****Quantitative**** knowledge points: involve specific numbers,
914 coefficients, percentages, counts, or magnitudes from the paper
915 (e.g., "The estimated effect is -0.37", "The sample includes 126 DMAs")
916 - ****Conceptual**** knowledge points: involve model assumptions, mechanism
917 descriptions, methodology choices, variable definitions, institutional
918 features, or qualitative findings (e.g., "The paper uses a shift-share
919 instrument", "The model assumes risk-neutral agents")
920
921 ****Target: at least 50% of knowledge points must be conceptual.**** Do not
922 over-extract numerical facts at the expense of conceptual understanding.

923 **G.1.5 Round 1 Output Format**

```

924 {
925   "paper_id": "<paper ID>",
926   "paper_title": "<paper title>",
927   "knowledge_points": [
928     {
929       "kp_id": "KP01",
930       "category": "model_setup | finding | method | mechanism |
931         contribution | institution",
932       "content_type": "quantitative | conceptual",
933       "statement": "<specific knowledge point>",
934       "paper_context": "<2-3 sentence description stating what aspect
935         of the paper this addresses, the specific detail
936         tested, and where it appears>"

```

```
937     }
938   ]
939 }
```

940 **G.1.6 Round 2 Prompt Template**

```
941 You are a professor of economics. Below is a list of paper-specific
942 knowledge points extracted from an economics paper. Based on these
943 knowledge points, generate **20 objective questions** to test whether
944 an examinee knows this specific paper.
945
946 ### Question type requirements
947
948 Generate exactly:
949 - 5 single-choice questions (4 options each, labeled A/B/C/D)
950 - 5 true/false questions
951 - 5 fill-in-the-blank questions
952 - 5 ordering questions (4 items to order)
953
954 ### CRITICAL: Content type balance
955
956 Ensure that across all 20 questions, at least 50% test **conceptual**
957 knowledge (not just numerical recall). Use the content_type tags from
958 the knowledge points to guide this balance.
959
960 ### Mandatory stem format
961
962 Every question stem MUST begin with the paper reference in this exact
963 format:
964
965 "In [FirstAuthor] et al. ([Year]), a paper on [brief topic phrase], ..."
966
967 Example: "In Autor et al. (2013), a paper on Chinese import competition
968 and U.S. labor markets, what identification strategy is used?"
969
970 Do NOT include:
971 - Full paper title
972 - Journal name
973 - Full author list beyond first author
974
975 ### Question design principles
976
977 1. **Paper-specific**: Every question must require knowledge of this
978    specific paper's content. Avoid generic economics questions.
979 2. **Unambiguous**: Each question must have a single, clearly correct
980    answer based on the paper.
981 3. **Difficulty distribution**: Aim for a mix of medium (60%) and hard
982    (40%) questions. Easy questions that can be guessed should be avoided.
983 4. **Source tracking**: Each question must reference its source
984    knowledge point via source_kp field.
985
986 ### Single-choice question requirements
987
988 - Provide exactly 4 options (A, B, C, D)
989 - All distractors must be plausible but clearly incorrect based on the
990   paper
991 - Avoid "all of the above" or "none of the above" options
992 - Distractors should not be trivially eliminable
993
994 ### True/false question requirements
995
996 - Statement must be specific enough that truth value is unambiguous
997 - Avoid statements that are "technically true but misleading" or vice
998   versa
999 - Explanation must cite specific paper content
```

```

1000
1001 ### Fill-in-the-blank question requirements
1002
1003 - Blank should test a specific term, number, or concept
1004 - Use "___" to mark the blank position
1005 - Provide canonical answer and list of acceptable synonyms/variants
1006 - Acceptable answers should account for: capitalization differences,
1007   singular/plural forms, common abbreviations
1008
1009 ### Ordering question requirements
1010
1011 - Provide exactly 4 items to be ordered
1012 - Specify the ordering criterion clearly in the stem (e.g., "from
1013   smallest to largest effect size", "in chronological order of
1014   implementation")
1015 - Items must have an unambiguous correct ordering based on the paper
1016 - Answer field contains the items in correct order

```

1017 **G.1.7 Round 2 Output Format**

```

1018 {
1019   "paper_id": "<paper ID>",
1020   "questions": [
1021     {
1022       "id": "EK_<paper_id>_SC01",
1023       "type": "single_choice",
1024       "content_type": "conceptual | quantitative",
1025       "source_kp": "KP03",
1026       "difficulty": "medium | hard",
1027       "knowledge_tags": ["<knowledge category>"],
1028       "stem": "<question stem starting with paper reference>",
1029       "options": {
1030         "A": "<option A>",
1031         "B": "<option B>",
1032         "C": "<option C>",
1033         "D": "<option D>"
1034       },
1035       "answer": "B",
1036       "explanation": "<explanation referencing specific paper content>"
1037     },
1038     {
1039       "id": "EK_<paper_id>_TF01",
1040       "type": "true_false",
1041       "content_type": "conceptual | quantitative",
1042       "source_kp": "KP07",
1043       "difficulty": "medium | hard",
1044       "knowledge_tags": ["<knowledge category>"],
1045       "stem": "<statement about the paper to judge true/false>",
1046       "answer": true,
1047       "explanation": "<explanation referencing specific paper content>"
1048     },
1049     {
1050       "id": "EK_<paper_id>_FB01",
1051       "type": "fill_in_blank",
1052       "content_type": "conceptual | quantitative",
1053       "source_kp": "KP05",
1054       "difficulty": "medium | hard",
1055       "knowledge_tags": ["<knowledge category>"],
1056       "stem": "<fill-in stem with ___ referencing paper content>",
1057       "answer": "<canonical answer>",
1058       "acceptable_answers": ["<synonym 1>", "<synonym 2>"],
1059       "explanation": "<explanation referencing specific paper content>"
1060     },
1061     {
1062       "id": "EK_<paper_id>_OD01",

```

```

1063     "type": "ordering",
1064     "content_type": "conceptual | quantitative",
1065     "source_kp": "KP08",
1066     "difficulty": "medium | hard",
1067     "knowledge_tags": ["<knowledge category>"],
1068     "stem": "<ordering stem stating the criterion, referencing paper>",
1069     "items": ["<item A>", "<item B>", "<item C>", "<item D>"],
1070     "answer": ["<item B>", "<item A>", "<item D>", "<item C>"],
1071     "explanation": "<explanation referencing specific paper content>"
1072   }
1073 ]
1074 }

```

1075 G.1.8 Scoring Rules

Question Type	Fully Correct	Partially Correct	Incorrect
Single Choice (5 questions)	1 point	–	0 points
True/False (5 questions)	1 point	–	0 points
Fill-in-Blank (5 questions)	1 point (case-insensitive, or in acceptable_answers)	–	0 points
Ordering (5 questions)	1 point (exact match)	0.5 points (Kendall $\tau \geq 0.5$)	0 points

1076 Total score = sum of individual question scores / number of questions.

1077 G.2 Mathematical Modeling (MM)

1078 G.2.1 Task Overview

1079 Mathematical Modeling (MM) is a Level 1 (foundational capability) task that evaluates models'
1080 ability to write reasonable model equations, constraints, or complete key derivation steps given an
1081 economic scenario. The core evaluation goal is to test formula-level mathematical skills. Each paper
1082 produces 3 questions spanning three types: equation design, formula derivation, and equilibrium
1083 solving. The task depends only on PDF input and uses a 2-round generation pipeline.

1084 G.2.2 Generation Workflow

1085 MM uses a two-round pipeline to systematically extract and select mathematical content:

- 1086 **1. Round 1 (PDF input):** Systematically extract all mathematical modeling content from
1087 the paper, including model setup (agents, variables, parameters), key equations (objective
1088 functions, constraints, FOCs), propositions and proofs, and derivation chains. For empirical
1089 papers, extract regression equations and identification assumptions. Output includes a
1090 `suitable_for_mm` flag.
- 1091 **2. Round 2 (text input):** Based on Round 1's structured mathematical content, design 3
1092 self-contained mathematical modeling questions spanning three types: equation design,
1093 formula derivation, and equilibrium solving. Each question includes scenario description,
1094 notation conventions, known conditions, and standard answer in LaTeX.

1095 G.2.3 Design Rationale

1096 Economic papers scatter mathematical content across main text, appendices, and footnotes. Round
1097 1 performs exhaustive extraction and builds derivation chains, allowing Round 2 to select the
1098 most suitable elements for exam questions without missing key formulas or choosing overly trivial
1099 derivations.

1100 **LaTeX Escaping in JSON** Since output is JSON containing LaTeX strings, all backslashes must
1101 be escaped: `\frac{a}{b}` becomes `\\frac{a}{b}` in JSON. Every backslash `\` in LaTeX must be
1102 written as `\\` in the JSON string. Double backslashes in LaTeX (e.g., `\\` for line breaks) must be
1103 `\\\\` in JSON.

1104 G.2.4 Round 1 Prompt Template

1105 You are an expert in mathematical economics. Read the full attached
1106 PDF paper and systematically extract the paper's mathematical modeling
1107 content.

1108

1109 ### Extraction requirements

1110

1111 Identify and record the following:

1112

- 1113 1. **Model setup**: Definitions of all economic agents, variables, and
1114 parameters
- 1115 2. **Key equations**: Objective functions, constraints, equilibrium
1116 conditions, first-order conditions, etc.
- 1117 3. **Propositions and proofs**: Propositions, Theorems, and Lemmas in
1118 the paper and their core derivation steps
- 1119 4. **Derivation chains**: Logical relations between formulas (which
1120 formula follows from which prior conditions)

1121

1122 ### Handling non-theoretical papers

1123

1124 If the paper is empirical:

- 1125 - Extract regression / estimation equations (e.g. DID, IV, structural
1126 models, etc.)
- 1127 - Extract mathematical formulations of identification assumptions
- 1128 - Extract statistical or asymptotic properties of the model

1129

1130 If the paper contains almost no mathematical derivations suitable for
1131 exam items (e.g. only reports regression results, lacking
1132 reconstructible equations and derivation chains):

- 1133 - Set 'suitable_for_mm = false'
- 1134 - Provide 'reason_not_suitable'
- 1135 - Do not invent mathematical models or derivation chains

1136

1137 ### LaTeX in JSON -- escaping rules

1138

1139 Since the output is JSON containing LaTeX strings, you MUST follow
1140 these escaping rules:

- 1141 - Every backslash '\' in LaTeX must be written as '\\\' in the JSON
1142 string
- 1143 - For example: '\frac{a}{b}' must be written as '\\frac{a}{b}'
- 1144 - Newlines in LaTeX should be '\\n' in the JSON string
- 1145 - Double backslashes in LaTeX (e.g., '\\\' for line breaks) must be
1146 '\\\\\' in JSON
- 1147 - Test: the JSON must be parseable by a standard JSON parser

1148 G.2.5 Round 1 Output Format

1149 {

```
1150 "paper_id": "<paper ID>",
1151 "paper_title": "<paper title>",
1152 "suitable_for_mm": true,
1153 "reason_not_suitable": "",
1154 "paper_type": "theoretical | empirical | mixed",
1155 "model_elements": {
1156   "agents": [
1157     {"name": "<agent>", "variables": ["<variable (LaTeX)>"],
1158     "parameters": ["<parameter>"]}
1159   ],
1160   "key_equations": [
1161     {
1162       "eq_id": "EQ01",
1163       "label": "<equation label, e.g. objective function, budget
1164       constraint, FOC>",
1165       "latex": "<LaTeX formula>",
```

```

1166     "context": "<where the equation appears and its meaning
1167             (1 sentence)>",
1168     "depends_on": ["EQ00"]
1169   }
1170 ],
1171 "propositions": [
1172   {
1173     "prop_id": "P01",
1174     "statement": "<verbal statement of the proposition>",
1175     "key_result_latex": "<LaTeX expression of the core result>",
1176     "proof_sketch": "<key proof steps (2-3 steps)>",
1177     "depends_on": ["EQ01", "EQ03"]
1178   }
1179 ]
1180 },
1181 "candidate_questions": [
1182   {
1183     "suggested_type": "equation_design | derivation |
1184                     equilibrium_solving",
1185     "target_element": "EQ03 | P01",
1186     "why_suitable": "<why this element is suitable as an exam item
1187                   (1 sentence)>",
1188     "estimated_difficulty": "medium | hard"
1189   }
1190 ]
1191 }

```

1192 G.2.6 Round 2 Prompt Template

```

1193 You are an expert in writing mathematical economics exam questions.
1194 Below is extracted mathematical modeling content from an economics
1195 paper. Based on this content, design 3 mathematical modeling
1196 questions.
1197
1198 ### Question type requirements
1199
1200 Design exactly 3 questions, one from each type:
1201 1. Equation design: Given a scenario, write the appropriate
1202    objective function, constraint, or equilibrium condition
1203 2. Formula derivation: Derive a key result from given premises
1204 3. Equilibrium solving: Characterize equilibrium or solve for
1205    optimal choices
1206
1207 ### Self-containment principle
1208
1209 Each question must be completely self-contained:
1210 - Provide all necessary background and notation in the question stem
1211 - Define all variables and parameters used
1212 - State all assumptions and known conditions
1213 - The examinee should NOT need to refer to the original paper
1214
1215 ### Question structure
1216
1217 Each question should include:
1218 1. Scenario description: Economic context and setup (2-3 sentences)
1219 2. Notation conventions: List of all symbols with their meanings
1220 3. Known conditions: What is given or assumed
1221 4. Task: What the examinee needs to derive or solve
1222 5. Standard answer: Complete solution in LaTeX
1223
1224 ### LaTeX escaping reminder
1225
1226 Remember: every '\' in LaTeX must be '\\' in JSON output.
1227
1228 ### Difficulty calibration

```

1229
 1230 - **Medium**: Requires applying standard techniques (FOC, envelope
 1231 theorem, budget constraint substitution)
 1232 - **Hard**: Requires multi-step reasoning, non-obvious substitutions,
 1233 or careful handling of corner cases
 1234
 1235 Target: 1-2 medium questions, 1-2 hard questions.

1236 G.2.7 Round 2 Output Format

```

1237 {
1238   "paper_id": "<paper ID>",
1239   "questions": [
1240     {
1241       "id": "MM_<paper_id>_Q01",
1242       "type": "equation_design | derivation | equilibrium_solving",
1243       "difficulty": "medium | hard",
1244       "scenario": "<economic scenario description (2-3 sentences)>",
1245       "notation": {
1246         "<symbol>": "<meaning>",
1247         "x": "consumption quantity",
1248         "p": "price"
1249       },
1250       "known_conditions": [
1251         "<condition 1>",
1252         "<condition 2>"
1253       ],
1254       "task": "<what the examinee needs to do>",
1255       "standard_answer": {
1256         "latex": "<complete LaTeX solution>",
1257         "explanation": "<step-by-step explanation>"
1258       },
1259       "rubric": {
1260         "correctness": "<what constitutes a correct answer>",
1261         "key_elements": ["<element 1>", "<element 2>"]
1262       }
1263     }
1264   ]
1265 }

```

1266 G.2.8 Scoring Method

1267 MM uses a hybrid scoring approach combining structured comparison and LLM-as-Judge:

1268 Scoring Dimensions (0–10 points each, total 40)

- 1269 1. **Formula correctness** (0–10): All formulas are mathematically equivalent to the reference;
 1270 derivation steps are valid. 9–10: fully correct; 6–8: minor errors but core structure correct;
 1271 3–5: significant errors but partial understanding evident; 0–2: fundamentally incorrect.
- 1272 2. **Notation consistency** (0–10): All symbols match the problem setup and are used consis-
 1273 tently. 9–10: perfect consistency; 6–8: minor inconsistencies; 3–5: some symbols undefined
 1274 or inconsistent; 0–2: major notation confusion.
- 1275 3. **Derivation completeness** (0–10): All intermediate steps shown; a reader can verify each
 1276 step. 9–10: all steps present; 6–8: key steps present but some algebra omitted; 3–5: major
 1277 steps missing; 0–2: only final answer or fragmentary.
- 1278 4. **Economic intuition** (0–10): Clear interpretation of results in economic terms; explains
 1279 what the math means. 9–10: excellent economic interpretation; 6–8: some interpretation
 1280 provided; 3–5: minimal economic context; 0–2: no economic interpretation.

1281 Normalized score = total / 40.

1282 **G.3 Paper Abstraction (PA)**

1283 **G.3.1 Task Overview**

1284 Paper Abstraction (PA) is a Level 2 (comprehensive capability) task that evaluates models' ability
1285 to write abstracts after reading condensed body text. The core evaluation goal is to test whether
1286 generated abstracts cover key points from the original abstract. Each paper produces 1 evaluation
1287 item containing condensed body text (approximately 30% of original length, excluding abstract) and
1288 5–10 atomic semantic propositions as evaluation criteria. The task depends only on PDF input and
1289 uses a 3-round generation pipeline.

1290 **G.3.2 Generation Workflow**

1291 PA employs a three-round pipeline to produce condensed text and atomic evaluation criteria:

- 1292 1. **Round 1 (PDF input)**: Extract the original abstract verbatim, identify section structure, and
1293 summarize core arguments for each major section (1–2 sentences per section). Also extract
1294 2–4 claimed contributions.
- 1295 2. **Round 2 (PDF + Round 1)**: Generate condensed body text at approximately 30% of the
1296 original length, excluding the abstract. Use Round 1's section structure to guide selective
1297 retention of core arguments, key findings, and methodological details while removing
1298 redundant exposition.
- 1299 3. **Round 3 (text input)**: Decompose the original abstract into 5–10 atomic semantic proposi-
1300 tions (key points) that serve as evaluation criteria. Cross-validate these propositions against
1301 Round 1's section arguments to ensure completeness.

1302 **G.3.3 Design Rationale**

1303 The three-round separation ensures: (1) Round 1 provides structural guidance for Round 2's conden-
1304 sation, (2) Round 2 can refer back to the PDF to avoid incorrect deletions of critical arguments, and
1305 (3) Round 3 focuses purely on atomic proposition decomposition without needing PDF access.

1306 **Examinee Input** Examinees receive only the condensed body text (approximately 30% of original
1307 length, excluding abstract) in plain text format. This design reduces token cost, supports models
1308 without PDF input, and focuses evaluation on “distilling abstracts from core information” rather than
1309 “processing long documents.”

1310 **G.3.4 Round 1 Prompt Template**

```
1311 You are an expert in economics research methods. Please read the
1312 attached PDF paper and complete the following information extraction.
1313
1314 ### Task 1: Extract the original abstract
1315
1316 Locate and extract the paper's abstract text in full (verbatim copy;
1317 do not paraphrase).
1318
1319 ### Task 2: Paper structure and core arguments by section
1320
1321 Identify the paper's section structure and summarize the core argument
1322 in 1-2 sentences for each major section.
1323
1324 ### Task 3: Identify the paper's core contributions
1325
1326 List 2-4 main contributions the paper explicitly claims (typically at
1327 the end of the introduction or in the conclusion).
```

1328 **G.3.5 Round 1 Output Format**

```
1329 {
1330   "paper_id": "<paper ID>",
```

```

1331 "paper_title": "<paper title>",
1332 "original_abstract": "<full original abstract, verbatim>",
1333 "sections": [
1334   {
1335     "section_number": "1",
1336     "section_title": "<section title>",
1337     "core_argument": "<core argument (1-2 sentences)>"
1338   }
1339 ],
1340 "claimed_contributions": ["<contribution 1>", "<contribution 2>"],
1341 "methodology_summary": "<one-sentence overview of methods>",
1342 "key_findings_summary": "<one-sentence overview of findings>"
1343 }

```

1344 G.3.6 Round 2 Prompt Template

1345 You are an academic editor. Below is the section structure extracted
1346 from an economics paper. Please re-read the attached PDF source and
1347 compress the body of the paper into a **condensed** version of
1348 approximately 30% of the original length.

1349
1350 **Compression goal**

1351
1352 Produce a **plain-text condensed body** whose length is approximately
1353 30% of the original paper body (excluding references and appendices).
1354 This text will be the sole input for examinees, who will write an
1355 abstract based on it.

1356
1357 **Compression principles**

- 1358
- 1359 1. **Must not include the original abstract** -- that is what examinees
1360 need to produce
- 1361 2. **Preserve the paper's core argumentative chain**: research question
1362 -> methods -> key findings -> conclusions
- 1363 3. **Organize in section order**, retaining the 2-4 most critical
1364 sentences of core argument per section
- 1365 4. **Retain key data and quantitative results** (e.g., core regression
1366 coefficients, effect sizes, statistical significance)
- 1367 5. **Retain core methodological description** (e.g., identification
1368 strategy, model specification, data sources)
- 1369 6. **Omit**: detailed literature review, secondary variants of
1370 robustness checks, appendix material, redundant exposition
- 1371 7. **Do not paraphrase** -- prefer direct quotations or minimally
1372 shortened phrases from the original
- 1373 8. Use **English** and maintain an academic register
- 1374 9. Prefix each section paragraph with '[Section N: Title]' for clarity

1375 G.3.7 Round 3 Prompt Template

1376 You are an expert in academic writing. Below is the original abstract
1377 of an economics paper. Decompose it into **5-10 atomic semantic**
1378 **propositions** (key points) that will serve as evaluation criteria for
1379 examinee-generated abstracts.

1380
1381 **Decomposition principles**

- 1382
- 1383 1. **Atomic**: Each key point should express ONE distinct claim or fact
- 1384 2. **Verifiable**: Each point should be checkable against candidate text
- 1385 3. **Complete coverage**: Together, the points should cover all major
1386 elements of the original abstract
- 1387 4. **Category balance**: Include points from multiple categories:
1388 research question, method, finding, contribution

1389
1390 **Key point categories**

1391

1392 - 'research_question': What economic question the paper addresses
 1393 - 'method': Identification strategy, data, or modeling approach
 1394 - 'finding': Core empirical or theoretical results
 1395 - 'contribution': How the paper advances the literature
 1396
 1397 ### must_appear flag
 1398
 1399 Mark a key point as 'must_appear: true' if it represents a core
 1400 contribution or finding that MUST be mentioned in any reasonable
 1401 abstract of this paper. Mark as 'false' if it is a supporting detail
 1402 that could reasonably be omitted in a concise abstract.
 1403
 1404 Target: 60-70% of key points should be must_appear: true.

1405 G.3.8 Round 3 Output Format

```
1406 {
1407   "paper_id": "<paper ID>",
1408   "key_points": [
1409     {
1410       "kp_id": "KP01",
1411       "statement": "<atomic proposition>",
1412       "category": "research_question | method | finding | contribution",
1413       "must_appear": true
1414     }
1415   ]
1416 }
```

1417 G.3.9 Scoring Method

1418 PA uses LLM-as-judge evaluation across four dimensions (0-10 points each):

- 1419 1. **Coverage** (0-10): Are all or nearly all reference key points mentioned or clearly implied?
 1420 9-10: all key points covered; 6-8: most covered with 1-2 minor omissions; 3-5: several
 1421 important points missing; 0-2: only 1-2 points covered.
- 1422 2. **Accuracy** (0-10): Are all stated claims factually correct per the paper content? 9-10: all
 1423 correct; 6-8: minor inaccuracies; 3-5: some factual errors; 0-2: major errors.
- 1424 3. **Conciseness** (0-10): Is the abstract within word limit with no redundancy? 9-10: within
 1425 limit, every sentence adds value; 6-8: mostly concise with minor redundancy; 3-5: notice-
 1426 ably wordy or exceeds limit; 0-2: excessively long.
- 1427 4. **Academic style** (0-10): Does it read like a published abstract in a top economics journal?
 1428 9-10: publication-ready; 6-8: generally academic with some informal phrasing; 3-5: mixed
 1429 register; 0-2: not academic.

1430 G.4 Literature Review (LR)

1431 G.4.1 Task Overview

1432 Literature Review (LR) is a Level 2 (comprehensive capability) task that evaluates models' ability to
 1433 write introductions with citation networks consistent with the original paper. The core evaluation
 1434 goal is to test whether models can produce literature reviews that align with the paper's actual citation
 1435 structure. Each paper produces 1 evaluation item containing condensed body text (approximately 50%
 1436 of original length, excluding introduction and references), three-tier citation classification, citation
 1437 role annotations, and a paper-specific writing rubric. The task depends only on PDF input and uses a
 1438 4-round generation pipeline.

1439 G.4.2 Generation Workflow

1440 LR uses a four-round pipeline to construct citation-aware evaluation items:

- 1441 1. **Round 1 (PDF input):** Extract the complete reference list from the paper, including author
 1442 names, titles, publication venues, and years. Tag each reference by publication type (journal
 1443 article, book, working paper, etc.).
- 1444 2. **Round 2 (PDF + Round 1):** Analyze the introduction section: identify introduction
 1445 boundaries, extract the argumentation structure, and annotate each citation’s role in the
 1446 argument. Classify citations into three tiers based on journal quality and relevance to the
 1447 paper’s core argument.
- 1448 3. **Round 3 (PDF + Round 2):** Generate condensed body text at approximately 50% of original
 1449 length, excluding the introduction and references. Use Round 2’s introduction boundaries to
 1450 ensure precise exclusion.
- 1451 4. **Round 4 (text input):** Construct a paper-specific writing rubric based on the argumentation
 1452 structure identified in Round 2. The rubric captures the paper’s specific rhetorical strategy
 1453 rather than applying a generic template.

1454 G.4.3 Design Rationale

1455 The four-round separation ensures: (1) clean reference list extraction before citation role annotation,
 1456 (2) precise introduction exclusion using Round 2 boundaries, and (3) customized rubric construction
 1457 based on the paper’s specific argumentation strategy.

1458 **Citation Tier System** Citations in the introduction are classified into three tiers based on journal
 1459 quality and relevance:

- 1460 • **Tier 1 (Must-cite):** Core references indispensable to the argument, typically from top
 1461 journals (AER, Econometrica, QJE, JPE, REStud) or high-impact sources, but not limited to
 1462 these—any high-impact work critical to the argument qualifies.
- 1463 • **Tier 2 (Should-cite):** Relevant formal academic publications from reputable but non-top
 1464 journals, academic books, or well-known working papers (NBER, CEPR).
- 1465 • **Tier 3 (Optional):** Lower-relevance or less authoritative sources (obscure journals, general
 1466 preprints, technical reports).

1467 **Scoring Formula** Weighted citation recall: $0.8 \times \text{Tier1_recall} + 0.2 \times \text{Tier2_recall}$ (Tier 3 not
 1468 scored). Final score: 50% citation recall + 50% writing rubric.

1469 G.4.4 Round 1 Prompt Template

1470 You are an expert in academic reference management. Read the attached
 1471 PDF paper and extract the **complete reference list** from the
 1472 References / Bibliography section.

1473
 1474 **### Extraction requirements**
 1475
 1476 For each reference, extract:
 1477 - Authors (standard format)
 1478 - Year
 1479 - Title
 1480 - Journal / publisher
 1481 - Publication type (venue_type)
 1482 - A standardized citation_key

1483
 1484 **### venue_type categories**

1485
 1486 Label each item by its publication source:
 1487 - ‘top5’: economics Top 5 journals (AER, Econometrica, QJE, JPE, REStud)
 1488 - ‘top_field’: widely recognized top field journals (e.g., JF, JFE,
 1489 RFS, JME, JEEA, AEJ series, JHE, JDE, JUE, JPubE, JLE, RAND, EJ,
 1490 REE, etc.)
 1491 - ‘good_journal’: well-known but non-top-tier formal journals
 1492 - ‘book’: scholarly monographs or book chapters
 1493 - ‘conference’: conference papers

- 1494 - 'working_paper_top': working papers from prominent institutions
- 1495 (NBER, CEPR, World Bank, IMF, Fed, etc.)
- 1496 - 'working_paper_other': working papers from other institutions
- 1497 - 'preprint': preprints (arXiv, SSRN, etc.)
- 1498 - 'other': reports, technical documents, news, web pages, etc.

1499 **G.4.5 Round 2 Prompt Template**

1500 You are an expert in economics research writing. Re-read the attached
 1501 PDF paper and analyze the **introduction section** in detail.

1502
 1503 **### Task 1: Identify introduction boundaries**

1504
 1505 Precisely identify where the introduction begins and ends (page numbers
 1506 and section numbers).

1507
 1508 **### Task 2: Extract argumentation structure**

1509
 1510 Identify the major argumentative modules in the introduction (e.g.,
 1511 motivation, literature gap, methodology preview, contribution summary).
 1512 For each module, provide a 2-3 sentence summary.

1513
 1514 **### Task 3: Annotate each citation's role**

1515
 1516 For every citation that appears in the introduction, record:

- 1517 - citation_key (matching Round 1 format)
- 1518 - The argumentative role it plays (e.g., "establishes empirical
 1519 regularity", "represents prior theoretical approach", "identifies
 1520 literature gap")
- 1521 - Tier classification (1/2/3) based on journal quality and relevance

1522
 1523 **### Tier classification criteria**

1524
 1525 ****Tier 1 (Must-cite)**:** Core references indispensable to the argument.
 1526 Typically from top journals (AER, Econometrica, QJE, JPE, REStud) or
 1527 high-impact sources, but not limited to these---any high-impact work
 1528 critical to the argument qualifies.

1529
 1530 ****Tier 2 (Should-cite)**:** Relevant formal academic publications from
 1531 reputable but non-top journals, academic books, or well-known working
 1532 papers (NBER, CEPR).

1533
 1534 ****Tier 3 (Optional)**:** Lower-relevance or less authoritative sources.

1535
 1536 **### IMPORTANT: Tier 1 and Tier 2 count requirements**

1537
 1538 After classification, verify:

- 1539 - Tier 1 should have ≥ 10 citations
- 1540 - Tier 2 should have ≥ 15 citations

1541
 1542 If counts are insufficient, re-examine the introduction and promote
 1543 relevant citations from lower tiers.

1544 **G.4.6 Round 3 Prompt Template**

1545 You are an academic editor. Below is the introduction boundary
 1546 information extracted from an economics paper. Re-read the attached
 1547 PDF and compress the body of the paper (EXCLUDING the introduction
 1548 and references) into approximately 50% of the original length.

1549
 1550 **### Compression goal**

1551
 1552 Produce plain-text condensed body at ~50% of original length, excluding
 1553 introduction and references. This will be the examinee's input for
 1554 writing an introduction.

1555
 1556 **### Compression principles**
 1557
 1558 1. ****Must exclude the entire introduction section**** (use Round 2
 1559 boundaries)
 1560 2. ****Must exclude the references section****
 1561 3. ****Preserve core argumentative chain****: methods -> results ->
 1562 robustness -> conclusion
 1563 4. ****Retain key quantitative results**** and methodological details
 1564 5. ****Omit****: detailed proofs, secondary robustness checks, appendix
 1565 material
 1566 6. Use section prefixes '[Section N: Title]' for clarity

1567 **G.4.7 Round 4 Prompt Template**

1568 You are an expert in economics research writing. Based on the
 1569 introduction analysis from Round 2, construct a ****paper-specific**
 1570 writing rubric** that captures this paper's actual rhetorical strategy.

1571 **### Rubric construction principles**
 1572
 1573
 1574 1. ****Paper-specific dimensions****: Identify 4-6 dimensions based on what
 1575 THIS paper's introduction actually does (not a generic template)
 1576 2. ****Concrete evaluation criteria****: For each dimension, list specific
 1577 elements to check
 1578 3. ****Weight allocation****: Assign weights summing to 100 based on
 1579 importance
 1580 4. ****Tier-classified citation list****: Aggregate citation_keys by tier
 1581 from Round 2

1582 **### Example dimensions (adapt to paper)**
 1583
 1584
 1585 - Motivation and research question clarity
 1586 - Literature gap identification
 1587 - Methodology preview
 1588 - Contribution summary
 1589 - Institutional context (if applicable)
 1590 - Theoretical framework setup (if applicable)

1591 **### Output requirements**
 1592
 1593
 1594 - writing_rubric: array of dimension objects with name, description,
 1595 key_elements_to_check, weight
 1596 - gold_citations_tiered: object with tier_1, tier_2, tier_3 arrays
 1597 - scoring_formula: explicit formula for combining citation and rubric
 1598 scores

1599 **G.4.8 Scoring Method**

1600 LR uses dual-track scoring:

1601 **Citation Recall (50% of final score)**

1602 tier1_recall = |tier1_gold pred_set| / |tier1_gold|
 1603 tier2_recall = |tier2_gold pred_set| / |tier2_gold|
 1604 weighted_recall = 0.8 × tier1_recall + 0.2 × tier2_recall
 1605 citation_score = weighted_recall × 100

1606 Design rationale: Use Recall (not F1) because introduction writing is open-ended; examinees may
 1607 reasonably cite works outside the gold set. Tier 1 (80% weight) represents indispensable core
 1608 references; Tier 2 (20% weight) represents quality supplementary references. Tier 3 not scored.

1609 **Writing Rubric (50% of final score)** LLM-as-judge evaluates the introduction against paper-
1610 specific rubric dimensions. Each dimension scored 0–10, weighted by dimension weights, normalized
1611 to 0–100 scale.

1612 **Final Score**

1613 $\text{final_score} = 0.5 \times \text{citation_score} + 0.5 \times \text{rubric_score}$

1614 **G.5 Method Generation (MG)**

1615 **G.5.1 Task Overview**

1616 Method Generation (MG) is a Level 2 (comprehensive capability) task that evaluates models' ability
1617 to construct economic models from natural language scenarios. The core evaluation goal is to
1618 test whether models can clearly specify economic agents, optimization problems, constraints, and
1619 equilibrium/game mechanisms. Each paper produces 1 evaluation item containing a desymbolized
1620 scenario narrative, standard model structure, and element-counting rubric. The task depends only on
1621 PDF input and uses a 3-round generation pipeline.

1622 **G.5.2 Generation Workflow**

1623 MG employs a three-round pipeline to transform formal models into natural language scenarios:

- 1624 1. **Round 1 (PDF input):** Fully extract the paper's core economic model, including all agents
1625 and their choice variables, objective functions (with optimization direction), all constraints
1626 (budget, resource, IC, PC, feasibility), equilibrium concept, key FOCs, and core mechanisms.
1627 For empirical papers, extract the theoretical framework or structural model if present.
- 1628 2. **Round 2 (PDF + Round 1):** Rewrite Round 1's formal model into a desymbolized natural
1629 language scenario description. Remove all mathematical notation and variable symbols,
1630 replacing them with plain economic intuition and verbal descriptions of relationships.
- 1631 3. **Round 3 (text input):** Construct an element-counting rubric based on Round 1's model
1632 structure. Dynamically allocate points across dimensions (agents, objectives, constraints,
1633 equilibrium characterization) based on the actual number of elements in the model.

1634 **G.5.3 Design Rationale**

1635 Round 1 focuses purely on faithful extraction without rewriting. Round 2 can then generate targeted
1636 natural language descriptions while referring back to the PDF to ensure economic intuition accuracy.
1637 Round 3 is purely computational—allocating rubric points based on element counts—and requires no
1638 PDF access.

1639 **G.5.4 Round 1 Prompt Template**

1640 You are an economic theory expert. Read the attached PDF paper and
1641 **fully extract** the paper's core economic model.

1642

1643 **### Extraction requirements**

1644

1645 Strictly extract the following from the paper as written (do not omit
1646 any agent or constraint):

1647

- 1648 1. **Economic agents**: All agents in the model and their choice
1649 variables
- 1650 2. **Objective functions**: Each agent's optimization objective
1651 (including maximization/minimization direction)
- 1652 3. **Constraints**: All constraints each agent faces (budget, resource,
1653 incentive compatibility, participation, feasibility, etc.)
- 1654 4. **Equilibrium concept**: The equilibrium / solution concept used
- 1655 5. **Key first-order conditions**: FOCs / characterization results as
1656 derived or stated
- 1657 6. **Core mechanism**: Economic forces that drive the model's results

1658
1659 **### Handling non-theoretical papers**
1660
1661 If the paper is empirical:
1662 - Extract the theoretical framework or structural model it relies on
1663 (if any)
1664 - Extract the economic logic behind the identification strategy
1665 - If the paper has no theoretical model at all, state that the paper
1666 is unsuitable for the MG task
1667
1668 When 'suitable_for_mg = false':
1669 - Provide 'reason_not_suitable', briefly explaining why
1670 - Other model fields may be empty arrays / empty objects (do not
1671 fabricate a model)

1672 **G.5.5 Round 2 Prompt Template**

1673 You are an economics textbook author. Below is the full model structure
1674 extracted from an economics paper. Rewrite it as a ****desymbolized**
1675 economic scenario narrative** for examinees to model from.
1676
1677 **### IMPORTANT: Check suitability first**
1678
1679 If the Round 1 output has 'suitable_for_mg = false', output a minimal
1680 JSON indicating the paper is unsuitable and stop.
1681
1682 **### Desymbolization requirements**
1683
1684 1. ****Remove ALL mathematical symbols and variable names****: No x, y, p,
1685 w, , , etc.
1686 2. ****Use plain economic language****: "consumption quantity" not "x",
1687 "wage rate" not "w"
1688 3. ****Describe relationships verbally****: "utility increases with
1689 consumption" not "U(x) is increasing"
1690 4. ****Preserve economic intuition****: The scenario must convey the same
1691 economic trade-offs as the formal model
1692 5. ****Self-contained****: Examinee should NOT need the original paper
1693
1694 **### Scenario structure**
1695
1696 1. ****Economic context**** (2-3 sentences): What is the economic
1697 environment?
1698 2. ****Agents and their goals**** (1 sentence per agent): Who are the
1699 decision-makers and what do they want?
1700 3. ****Constraints and trade-offs**** (2-3 sentences): What limits their
1701 choices?
1702 4. ****Equilibrium/outcome concept**** (1 sentence): What constitutes a
1703 solution?
1704
1705 **### Examinee task description**
1706
1707 After the scenario, provide clear instructions: "Based on the scenario
1708 above, construct a formal economic model that captures the key
1709 trade-offs and decision problems. Your model should specify: (1) all
1710 economic agents and their choice variables, (2) each agent's objective
1711 function, (3) all relevant constraints, (4) the equilibrium concept."

1712 **G.5.6 Round 3 Prompt Template**

1713 You are an exam rubric designer. Based on the model structure extracted
1714 in Round 1, construct an ****element-counting rubric**** that dynamically
1715 allocates points based on the actual number of model elements.
1716
1717 **### Rubric construction algorithm**
1718

- 1719 Total points: 100
 1720
 1721 1. **Agent identification** (15-25 points): Allocate based on number
 1722 of agents. If 1 agent: 15 pts. If 2 agents: 20 pts. If 3+ agents:
 1723 25 pts.
 1724
 1725 2. **Objective specification** (20-30 points): Allocate based on
 1726 complexity. Simple single-period utility: 20 pts. Multi-period or
 1727 strategic objectives: 25-30 pts.
 1728
 1729 3. **Constraint completeness** (30-40 points): Allocate proportionally
 1730 to number of constraints. Each constraint type worth 5-10 pts
 1731 depending on complexity.
 1732
 1733 4. **Equilibrium characterization** (15-20 points): Standard
 1734 equilibrium concept: 15 pts. Complex equilibrium (e.g., perfect
 1735 Bayesian): 20 pts.
 1736
 1737 **Quality levels for each element**
 1738
 1739 - **correct** (full weight): Element present, correctly formalized,
 1740 economically meaningful
 1741 - **substantial** (70% weight): Mostly correct with minor issues
 1742 - **partial** (40% weight): Attempted but significant issues
 1743 - **minimal** (15% weight): Vaguely referenced without proper
 1744 formalization
 1745 - **missing** (0 points): Not present

1746 G.5.7 Scoring Method

1747 MG uses element-counting rubric (total 100 points) with LLM-assisted validation of element cor-
 1748 rectness. Points are dynamically allocated based on model complexity. For each rubric element,
 1749 LLM judges quality level (correct/substantial/partial/minimal/missing) and awards corresponding
 1750 percentage of per-element points. Total score is sum of all element scores.

1751 G.6 Economic Analysis (EA)

1752 G.6.1 Task Overview

1753 Economic Analysis (EA) is a Level 2 (comprehensive capability) task that evaluates models' ability
 1754 to write professional, multi-dimensional economic analysis given data results or model outputs. The
 1755 core evaluation goal is to test whether models can produce substantive economic interpretations of
 1756 empirical or theoretical results. Each paper produces 1–2 evaluation items (one per core result),
 1757 each containing a purely factual data description, a reference analysis paragraph, applicable scoring
 1758 dimensions, and examinee instructions. The task depends only on PDF input and uses a 2-round
 1759 generation pipeline.

1760 G.6.2 Generation Workflow

1761 EA uses a two-round pipeline to separate data description from analytical writing:

- 1762 1. **Round 1 (PDF input)**: Scan all tables and figures in the paper, label their type and
 1763 importance, and select 1–2 core results that best represent the paper's contribution (typically
 1764 main regression tables or core theoretical result figures). For each selected result, extract
 1765 a context paragraph (150–250 words covering research question, data, methodology, and
 1766 what the result represents) and a purely factual data description (200–400 words) without
 1767 any interpretation or analysis.
- 1768 2. **Round 2 (PDF + Round 1)**: For each result identified in Round 1, locate the corresponding
 1769 analysis paragraph(s) in the paper's main text (400–600 words). Extract the author's original
 1770 analysis as the gold standard. Identify which of the six evaluation dimensions apply to this
 1771 specific result.

1772 **G.6.3 Design Rationale**

1773 The core challenge is cleanly separating “result description” from “analysis paragraph.” Round 1
1774 identifies core results and extracts pure data descriptions (no interpretation), then Round 2 locates
1775 the author’s analysis of these results as the gold standard. Single-round generation tends to blur the
1776 boundary between data description and analytical content.

1777 **G.6.4 Round 1 Prompt Template**

1778 You are an expert in empirical economics research. Read the attached
1779 PDF paper, identify the paper’s **core empirical/theoretical results**,
1780 and extract a purely factual data description.
1781
1782 **### Task 1: Scan all results in the paper**
1783
1784 List all Tables and Figures in the paper (including titles). Label
1785 the type and importance of each:
1786 - Type: main_regression | robustness | descriptive | heterogeneity |
1787 mechanism | theory_result | other
1788 - Importance: core | supporting | auxiliary
1789
1790 **### Task 2: Select core results**
1791
1792 Choose **1-2** results that best represent the paper’s core
1793 contribution. Selection criteria:
1794 - Prefer main regression tables (e.g., Table 1-3) or core theoretical-
1795 result figures
1796 - Avoid robustness checks, descriptive statistics, and other auxiliary
1797 tables
1798 - Selected results should allow substantive analysis
1799
1800 **### Task 3: Write a rich context paragraph AND extract a purely**
1801 **factual data description**
1802
1803 **Context paragraph** (field: ‘context’, 150-250 words): For each
1804 selected result, write a comprehensive background paragraph that
1805 enables a reader who has NEVER seen the paper to understand the result.
1806 This paragraph **MUST** include ALL of the following:
1807 - **Research question & motivation** (2-3 sentences)
1808 - **Data & sample** (1-2 sentences)
1809 - **Methodology overview** (2-3 sentences)
1810 - **What the result represents** (1 sentence)
1811
1812 **CRITICAL**: The context must NOT contain any interpretation,
1813 mechanism discussion, or analytical conclusions.
1814
1815 **Data description** (field: ‘data_description’, 200-400 words):
1816 Describe the data/numbers themselves in plain text:
1817 - You may include: coefficient values, standard errors, significance
1818 levels, sample size, R-squared, direction of trends, etc.
1819 - Do NOT include: causal interpretation, mechanism explanation, policy
1820 implications, links to theory
1821
1822 **### Handling theoretical papers**
1823
1824 If the paper is purely theoretical:
1825 - ‘selected_results’ should focus on core propositions, theorems, or
1826 comparative statics results
1827 - ‘data_description’ should state the formal result in plain language
1828 (what the theorem says) without explaining why it holds
1829 - ‘result_type’ should be ‘theoretical_result’

1830 G.6.5 Round 2 Prompt Template

1831 You are an expert in economics research methods. Re-read the attached
1832 PDF paper and complete the following tasks.

1833
1834 **### Task 1: Locate and extract a FOCUSED reference analysis**
1835

1836 For each 'selected_result' from Round 1, find the passage(s) in the
1837 main text where the author ****directly analyzes/interprets that result****.

1838
1839 ****CRITICAL length constraint****: The reference_analysis must be
1840 ****400-600 words****. This means:

- 1841 - Extract only the ****most essential analytical paragraphs**** that
1842 directly interpret the selected result
- 1843 - Do NOT include extended discussions that cover other results, lengthy
1844 institutional background, or multi-page analyses
- 1845 - If the author's discussion spans multiple pages, ****select and**
1846 **condense**** to the 400-600 word core
- 1847 - Prefer the paragraphs immediately following the result presentation

1848
1849 **### Task 2: Tag applicable scoring dimensions**
1850

1851 For each result, label which of the following 6 dimensions are
1852 ****actually supported by that result****. Use EXACTLY these dimension IDs:

- 1853
1854 - 'accuracy': Can the result's numbers/patterns be meaningfully checked
1855 for accuracy?
- 1856 - 'economic_interpretation': Does the result support substantive
1857 economic interpretation (magnitude, heterogeneity)?
- 1858 - 'mechanism': Does the result support discussion of economic mechanisms
1859 or causal channels?
- 1860 - 'theory_connection': Can the result be meaningfully connected to
1861 theoretical predictions or prior literature?
- 1862 - 'policy_implication': Does the result have implications for policy?
- 1863 - 'writing_quality': Always applicable -- assesses logical coherence
1864 and professional language

1865
1866 **### Task 3: Build the complete task**
1867

1868 Integrate data descriptions from Round 1 and reference analysis from
1869 this round into the final evaluation item.

1870 G.6.6 Evaluation Dimensions

1871 Six dimensions are used consistently across Round 2 annotation and scoring prompts:

- 1872 • **Accuracy** (accuracy): Correctness of numerical descriptions, trends, and statistical claims
- 1873 • **Economic interpretation** (economic_interpretation): Substantive interpretation of
1874 magnitudes, heterogeneity, or practical significance
- 1875 • **Mechanism** (mechanism): Explanation of economic mechanisms or causal logic behind
1876 results
- 1877 • **Theory connection** (theory_connection): Linkage to theoretical predictions or prior
1878 literature
- 1879 • **Policy implication** (policy_implication): Discussion of policy relevance (only when
1880 results support it)
- 1881 • **Writing quality** (writing_quality): Logical clarity and professional language

1882 Not every dimension applies to every result. writing_quality is always applicable; accuracy is
1883 almost always applicable. For theoretical results, accuracy refers to whether the formal statement is
1884 correctly characterized, and mechanism refers to the economic intuition behind the result.

1885 **G.6.7 Examinee Input**

1886 Examinees receive: (1) a context paragraph describing the research question, data, and methodology
1887 (150–250 words), and (2) a purely factual data description (e.g., “Table 3 shows that the coefficient on
1888 X is -0.37 (SE = 0.12), significant at the 1% level”) without any analytical interpretation. They are
1889 asked to write a 300–500 word analysis paragraph interpreting the result across applicable dimensions.

1890 **G.6.8 Scoring Method**

1891 EA uses multi-dimension LLM-as-judge scoring (0–10 points per applicable dimension):

- 1892 1. **Accuracy** (0–10): 9–10: all factual claims match data, no errors; 6–8: minor inaccuracies,
1893 core facts correct; 3–5: some factual errors affecting interpretation; 0–2: major factual
1894 errors.
- 1895 2. **Economic interpretation** (0–10): 9–10: insightful interpretation with proper economic
1896 reasoning; 6–8: adequate interpretation, some depth lacking; 3–5: superficial, mostly restates
1897 numbers; 0–2: no meaningful interpretation.
- 1898 3. **Mechanism** (0–10): 9–10: clear identification of driving mechanisms; 6–8: some mech-
1899 anisms discussed but not fully developed; 3–5: vague or incomplete; 0–2: no mechanism
1900 identified.
- 1901 4. **Theory connection** (0–10): 9–10: strong theoretical grounding with specific references;
1902 6–8: some theoretical context; 3–5: weak or generic theory connection; 0–2: no theory
1903 connection.
- 1904 5. **Policy implication** (0–10): 9–10: nuanced, well-grounded policy discussion; 6–8: reason-
1905 able policy points; 3–5: superficial; 0–2: no policy discussion.
- 1906 6. **Writing quality** (0–10): 9–10: publication-ready prose, logical flow, precise language;
1907 6–8: clear and organized, minor style issues; 3–5: readable but poorly organized; 0–2:
1908 disorganized or unclear.

1909 Normalized score = total points / maximum possible points (counting only applicable dimensions,
1910 each worth 10 points max).

1911 **G.7 Summary**

1912 Table 10 summarizes the generation pipeline characteristics across all six tasks.

Task	Level	Rounds	Items/Paper	Scoring	Key Challenge
EK	L1	2	20 questions	Rule-based	Knowledge coverage balance
MM	L1	2	3 questions	LLM-as-Judge	Self-containment of questions
PA	L2	3	1 item	LLM-as-Judge	Clean abstract/body separation
LR	L2	4	1 item	Hybrid	Citation tier classification
MG	L2	3	1 item	LLM-as-Judge	Desymbolization fidelity
EA	L2	2	1–2 items	LLM-as-Judge	Data/analysis boundary

Table 10: Summary of task generation pipeline characteristics.

1913 All six tasks employ multi-round generation pipelines with clear separation of concerns: early rounds
1914 perform exhaustive extraction or structural analysis with PDF access, while later rounds focus on
1915 question design, rubric construction, or validation using structured intermediate outputs. This design
1916 ensures high-quality, diverse evaluation items while maintaining consistency and reproducibility
1917 across the benchmark.

1918 **H Why AER-bench Omits a Human Expert Baseline**

1919 A natural question when presenting a benchmark like AER-bench is how large language models
1920 compare to human experts. After careful consideration, we have chosen not to collect a human
1921 performance baseline. This appendix details the reasons behind that decision.

1922 **Economic expertise is deeply fragmented.** Modern economics is an enormously broad discipline,
1923 spanning macroeconomics, microeconomic theory, econometrics, labor, development, industrial
1924 organization, international trade, finance, and many other fields. A typical economics PhD student or
1925 even an experienced researcher usually possesses deep knowledge in only one or two of these areas.
1926 AER-bench draws tasks from 191 top-5 journal articles covering the full spectrum of contemporary
1927 economics. Consequently, no single human expert can be expected to perform competently across
1928 all seven task types and all subfields. Assembling a panel of specialists who could collectively
1929 cover the entire benchmark would not only be logistically daunting but would also introduce severe
1930 inter-annotator variability, making a single “human baseline” number difficult to interpret.

1931 **Human experts routinely rely on tools.** In real-world research, economists seldom work from
1932 memory alone. They consult literature databases, use statistical software, and increasingly interact
1933 with large language models themselves. If we were to allow human participants to use such tools
1934 in order to simulate realistic research conditions, the resulting “human” baseline would in fact be a
1935 hybrid human–AI system. That would largely defeat the purpose: our aim is to measure how far pure
1936 LLMs are from performing end-to-end economic research, not to compare them against a human–AI
1937 combination that is not under controlled conditions.

1938 **Restricting human experts to a closed-book setting is unrealistic and difficult to enforce.** Con-
1939 versely, if we were to forbid any external aids, we would be testing an artificial skill—unaided
1940 recall of paper details—that is not how economists actually work. Moreover, ensuring that remotely
1941 participating experts do not quickly search for answers is practically impossible. The resulting data
1942 would be noisy and likely unrepresentative of genuine human expertise.

1943 **The benchmark already embeds a form of expert knowledge.** It is worth noting that AER-bench
1944 already incorporates human expertise in two critical ways. First, all Concept Identification tasks are
1945 manually annotated by economics doctoral students, providing gold-standard labels. Second, the
1946 LLM-as-a-Judge rubrics are designed and weighted to reflect what a domain expert would value in an
1947 answer. These design choices ensure that the evaluation criteria are grounded in expert judgment,
1948 even without a separate human performance baseline.

1949 **Conclusion.** Given these considerations—the breadth of expertise required, the confounding effect
1950 of tool use, the impracticality of pure closed-book testing, and the expert knowledge already embedded
1951 in the benchmark design—we believe that adding a human baseline would not strengthen the
1952 actionable insights that AER-bench provides. We instead focus on comparisons across models, which
1953 already reveal consistent performance hierarchies and systematic weaknesses. We encourage future
1954 work to explore targeted human–model comparisons on narrower subsets of the benchmark where
1955 domain specialists can be meaningfully recruited.

1956 **NeurIPS Paper Checklist**

1957 **1. Claims**

1958 Question: Do the main claims made in the abstract and introduction accurately reflect the
1959 paper’s contributions and scope?

1960 Answer: [Yes]

1961 Justification: The abstract and introduction accurately describe AER-bench as a comprehen-
1962 sive benchmark for evaluating LLMs on economic research tasks, with claims supported by
1963 experimental results in Section 4.

1964 Guidelines:

- 1965 • The answer [N/A] means that the abstract and introduction do not include the claims
1966 made in the paper.
- 1967 • The abstract and/or introduction should clearly state the claims made, including the
1968 contributions made in the paper and important assumptions and limitations. A [No] or
1969 [N/A] answer to this question will not be perceived well by the reviewers.
- 1970 • The claims made should match theoretical and experimental results, and reflect how
1971 much the results can be expected to generalize to other settings.
- 1972 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1973 are not attained by the paper.

1974 **2. Limitations**

1975 Question: Does the paper discuss the limitations of the work performed by the authors?

1976 Answer: [Yes]

1977 Justification: Limitations are discussed in a dedicated section, including LLM-as-a-Judge
1978 subjectivity, English-only scope, and temporal bias in historical concept identification.

1979 Guidelines:

- 1980 • The answer [N/A] means that the paper has no limitation while the answer [No] means
1981 that the paper has limitations, but those are not discussed in the paper.
- 1982 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 1983 • The paper should point out any strong assumptions and how robust the results are to
1984 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1985 model well-specification, asymptotic approximations only holding locally). The authors
1986 should reflect on how these assumptions might be violated in practice and what the
1987 implications would be.
- 1988 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1989 only tested on a few datasets or with a few runs. In general, empirical results often
1990 depend on implicit assumptions, which should be articulated.
- 1991 • The authors should reflect on the factors that influence the performance of the approach.
1992 For example, a facial recognition algorithm may perform poorly when image resolution
1993 is low or images are taken in low lighting. Or a speech-to-text system might not be
1994 used reliably to provide closed captions for online lectures because it fails to handle
1995 technical jargon.
- 1996 • The authors should discuss the computational efficiency of the proposed algorithms
1997 and how they scale with dataset size.
- 1998 • If applicable, the authors should discuss possible limitations of their approach to
1999 address problems of privacy and fairness.
- 2000 • While the authors might fear that complete honesty about limitations might be used by
2001 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
2002 limitations that aren’t acknowledged in the paper. The authors should use their best
2003 judgment and recognize that individual actions in favor of transparency play an impor-
2004 tant role in developing norms that preserve the integrity of the community. Reviewers
2005 will be specifically instructed to not penalize honesty concerning limitations.

2006 **3. Theory assumptions and proofs**

2007 Question: For each theoretical result, does the paper provide the full set of assumptions and
2008 a complete (and correct) proof?

2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062

Answer: [N/A]

Justification: This paper introduces a benchmark dataset and does not include theoretical results requiring formal proofs.

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

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

Answer: [Yes]

Justification: Section 3 provides complete details on evaluation protocols, scoring methods, and model configurations. All rubrics are detailed in Appendix E.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

2063 Question: Does the paper provide open access to the data and code, with sufficient instruc-
2064 tions to faithfully reproduce the main experimental results, as described in supplemental
2065 material?

2066 Answer: [Yes]

2067 Justification: The benchmark dataset, evaluation scripts, and complete documentation will
2068 be released under an open-source license upon publication.

2069 Guidelines:

- 2070 • The answer [N/A] means that paper does not include experiments requiring code.
- 2071 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
2072 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 2073 • While we encourage the release of code and data, we understand that this might not
2074 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
2075 including code, unless this is central to the contribution (e.g., for a new open-source
2076 benchmark).
- 2077 • The instructions should contain the exact command and environment needed to run to
2078 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
2079 neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 2080 • The authors should provide instructions on data access and preparation, including how
2081 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 2082 • The authors should provide scripts to reproduce all experimental results for the new
2083 proposed method and baselines. If only a subset of experiments are reproducible, they
2084 should state which ones are omitted from the script and why.
- 2085 • At submission time, to preserve anonymity, the authors should release anonymized
2086 versions (if applicable).
- 2087 • Providing as much information as possible in supplemental material (appended to the
2088 paper) is recommended, but including URLs to data and code is permitted.

2089 6. Experimental setting/details

2090 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
2091 rameters, how they were chosen, type of optimizer) necessary to understand the results?

2092 Answer: [Yes]

2093 Justification: Section 3 specifies all evaluation details including model APIs, temperature
2094 settings, judge model configurations, and scoring weight distributions.

2095 Guidelines:

- 2096 • The answer [N/A] means that the paper does not include experiments.
- 2097 • The experimental setting should be presented in the core of the paper to a level of detail
2098 that is necessary to appreciate the results and make sense of them.
- 2099 • The full details can be provided either with the code, in appendix, or as supplemental
2100 material.

2101 7. Experiment statistical significance

2102 Question: Does the paper report error bars suitably and correctly defined or other appropriate
2103 information about the statistical significance of the experiments?

2104 Answer: [Yes]

2105 Justification: Each task contains approximately 200 items, providing sufficient sample size
2106 for statistically significant results. We report mean absolute differences between judges (2.6
2107 pp) and task-level score distributions to demonstrate result robustness.

2108 Guidelines:

- 2109 • The answer [N/A] means that the paper does not include experiments.
- 2110 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
2111 intervals, or statistical significance tests, at least for the experiments that support the
2112 main claims of the paper.

- 2113 • The factors of variability that the error bars are capturing should be clearly stated (for
2114 example, train/test split, initialization, random drawing of some parameter, or overall
2115 run with given experimental conditions).
- 2116 • The method for calculating the error bars should be explained (closed form formula,
2117 call to a library function, bootstrap, etc.)
- 2118 • The assumptions made should be given (e.g., Normally distributed errors).
- 2119 • It should be clear whether the error bar is the standard deviation or the standard error
2120 of the mean.
- 2121 • It is OK to report 1-sigma error bars, but one should state it. The authors should
2122 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
2123 of Normality of errors is not verified.
- 2124 • For asymmetric distributions, the authors should be careful not to show in tables or
2125 figures symmetric error bars that would yield results that are out of range (e.g., negative
2126 error rates).
- 2127 • If error bars are reported in tables or plots, the authors should explain in the text how
2128 they were calculated and reference the corresponding figures or tables in the text.

2129 8. Experiments compute resources

2130 Question: For each experiment, does the paper provide sufficient information on the com-
2131 puter resources (type of compute workers, memory, time of execution) needed to reproduce
2132 the experiments?

2133 Answer: [Yes]

2134 Justification: The paper specifies that evaluations use commercial API endpoints (OpenAI,
2135 Anthropic, Google) with standard inference configurations, requiring no specialized compute
2136 infrastructure.

2137 Guidelines:

- 2138 • The answer [N/A] means that the paper does not include experiments.
- 2139 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
2140 or cloud provider, including relevant memory and storage.
- 2141 • The paper should provide the amount of compute required for each of the individual
2142 experimental runs as well as estimate the total compute.
- 2143 • The paper should disclose whether the full research project required more compute
2144 than the experiments reported in the paper (e.g., preliminary or failed experiments that
2145 didn't make it into the paper).

2146 9. Code of ethics

2147 Question: Does the research conducted in the paper conform, in every respect, with the
2148 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

2149 Answer: [Yes]

2150 Justification: The research adheres to the NeurIPS Code of Ethics, using only publicly
2151 available academic papers and transparent evaluation methodologies.

2152 Guidelines:

- 2153 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
2154 Ethics.
- 2155 • If the authors answer [No], they should explain the special circumstances that require a
2156 deviation from the Code of Ethics.
- 2157 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
2158 eration due to laws or regulations in their jurisdiction).

2159 10. Broader impacts

2160 Question: Does the paper discuss both potential positive societal impacts and negative
2161 societal impacts of the work performed?

2162 Answer: [Yes]

2163 Justification: The paper discusses how AER-bench can guide responsible development of
2164 AI systems for economic research, while acknowledging risks of over-reliance on automated
2165 evaluation.

2166 Guidelines:

- 2167 • The answer [N/A] means that there is no societal impact of the work performed.
- 2168 • If the authors answer [N/A] or [No], they should explain why their work has no societal
2169 impact or why the paper does not address societal impact.
- 2170 • Examples of negative societal impacts include potential malicious or unintended uses
2171 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
2172 (e.g., deployment of technologies that could make decisions that unfairly impact specific
2173 groups), privacy considerations, and security considerations.
- 2174 • The conference expects that many papers will be foundational research and not tied to
2175 particular applications, let alone deployments. However, if there is a direct path to any
2176 negative applications, the authors should point it out.
- 2177 • The authors should consider possible harms that could arise when the technology is
2178 being used as intended and functioning correctly, harms that could arise when the
2179 technology is being used as intended but gives incorrect results, and harms following
2180 from (intentional or unintentional) misuse of the technology.
- 2181 • If there are negative societal impacts, the authors could also discuss possible mitigation
2182 strategies (e.g., gated release of models, providing defenses in addition to attacks,
2183 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
2184 feedback over time, improving the efficiency and accessibility of ML).

2185 11. Safeguards

2186 Question: Does the paper describe safeguards that have been put in place for responsible
2187 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
2188 image generators, or scraped datasets)?

2189 Answer: [N/A]

2190 Justification: The benchmark consists of curated academic papers and evaluation protocols,
2191 posing no misuse risks. No pre-trained models are released.

2192 Guidelines:

- 2193 • The answer [N/A] means that the paper poses no such risks.
- 2194 • Released models that have a high risk for misuse or dual-use should be released with
2195 necessary safeguards to allow for controlled use of the model, for example by requiring
2196 that users adhere to usage guidelines or restrictions to access the model or implementing
2197 safety filters.
- 2198 • Datasets that have been scraped from the Internet could pose safety risks. The authors
2199 should describe how they avoided releasing unsafe images.
- 2200 • We recognize that providing effective safeguards is challenging, and many papers do
2201 not require this, but we encourage authors to take this into account and make a best
2202 faith effort.

2203 12. Licenses for existing assets

2204 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
2205 the paper, properly credited and are the license and terms of use explicitly mentioned and
2206 properly respected?

2207 Answer: [Yes]

2208 Justification: All source papers are from the American Economic Review, properly cited,
2209 and used in accordance with academic fair use principles for research purposes.

2210 Guidelines:

- 2211 • The answer [N/A] means that the paper does not use existing assets.
- 2212 • The authors should cite the original paper that produced the code package or dataset.
- 2213 • The authors should state which version of the asset is used and, if possible, include a
2214 URL.

- 2215
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - 2216 • For scraped data from a particular source (e.g., website), the copyright and terms of
2217 service of that source should be provided.
 - 2218 • If assets are released, the license, copyright information, and terms of use in the
2219 package should be provided. For popular datasets, paperswithcode.com/datasets
2220 has curated licenses for some datasets. Their licensing guide can help determine the
2221 license of a dataset.
 - 2222 • For existing datasets that are re-packaged, both the original license and the license of
2223 the derived asset (if it has changed) should be provided.
 - 2224 • If this information is not available online, the authors are encouraged to reach out to
2225 the asset's creators.

2226 13. **New assets**

2227 Question: Are new assets introduced in the paper well documented and is the documentation
2228 provided alongside the assets?

2229 Answer: [Yes]

2230 Justification: The AER-bench dataset includes comprehensive documentation of task def-
2231 initions, rubrics, data sources, and evaluation protocols in the paper and supplemental
2232 materials.

2233 Guidelines:

- 2234 • The answer [N/A] means that the paper does not release new assets.
- 2235 • Researchers should communicate the details of the dataset/code/model as part of their
2236 submissions via structured templates. This includes details about training, license,
2237 limitations, etc.
- 2238 • The paper should discuss whether and how consent was obtained from people whose
2239 asset is used.
- 2240 • At submission time, remember to anonymize your assets (if applicable). You can either
2241 create an anonymized URL or include an anonymized zip file.

2242 14. **Crowdsourcing and research with human subjects**

2243 Question: For crowdsourcing experiments and research with human subjects, does the paper
2244 include the full text of instructions given to participants and screenshots, if applicable, as
2245 well as details about compensation (if any)?

2246 Answer: [N/A]

2247 Justification: The benchmark construction involved expert annotation by the authors, not
2248 crowdsourcing or human subjects research requiring IRB approval.

2249 Guidelines:

- 2250 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
2251 with human subjects.
- 2252 • Including this information in the supplemental material is fine, but if the main contribu-
2253 tion of the paper involves human subjects, then as much detail as possible should be
2254 included in the main paper.
- 2255 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
2256 or other labor should be paid at least the minimum wage in the country of the data
2257 collector.

2258 15. **Institutional review board (IRB) approvals or equivalent for research with human 2259 subjects**

2260 Question: Does the paper describe potential risks incurred by study participants, whether
2261 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
2262 approvals (or an equivalent approval/review based on the requirements of your country or
2263 institution) were obtained?

2264 Answer: [N/A]

2265 Justification: The paper does not involve human subjects research. All data is derived from
2266 publicly available academic publications.

2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291

Guidelines:

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

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: LLMs are the central focus of this benchmark. The paper evaluates 16 LLMs and uses LLM-as-a-Judge for scoring generative tasks, with full methodological details in Section 3.

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.