

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 SYNTHWORLDS: CONTROLLED PARALLEL WORLDS FOR DISENTANGLING REASONING AND KNOWLEDGE IN LANGUAGE MODELS

006 **Anonymous authors**

007 Paper under double-blind review

010 ABSTRACT

013 Evaluating the reasoning ability of language models (LMs) is complicated by
014 their extensive parametric world knowledge, where benchmark performance of-
015 ten reflects factual recall rather than genuine reasoning. Existing datasets and
016 approaches (e.g., temporal filtering, paraphrasing, adversarial substitution) cannot
017 cleanly separate the two. We present SYNTHWORLDS, a framework that disentan-
018 gles task reasoning complexity from factual knowledge. In SYNTHWORLDS, we
019 construct parallel corpora representing two worlds with identical interconnected
020 structure: a real-mapped world, where models may exploit parametric knowledge,
021 and a synthetic-mapped world, where such knowledge is meaningless. On top of
022 these corpora, we design two mirrored tasks as case studies: multi-hop question
023 answering and page navigation, which maintain equal reasoning difficulty across
024 worlds. Experiments in parametric-only (e.g., closed-book QA) and knowledge-
025 augmented (e.g., retrieval-augmented) LM settings reveal a persistent *knowledge*
026 *advantage gap*, defined as the performance boost models gain from memorized
027 parametric world knowledge. Knowledge acquisition and integration mechanisms
028 reduce but do not eliminate this gap, highlighting opportunities for system im-
029 provements. Fully automatic and scalable, SYNTHWORLDS provides a controlled
030 environment for evaluating LMs in ways that were previously challenging, en-
031 abling precise and testable comparisons of reasoning and memorization.

032 1 INTRODUCTION

034 Language model (LM) agents are increasingly ex-
035 pected to autonomously complete complex tasks that
036 require retrieve new information, reason over it, and
037 synthesize novel insights. These capabilities un-
038 derpin emerging applications such as web navi-
039 gation, where agents need to traverse linked infor-
040 mation to locate relevant content (Ning et al., 2025);
041 personal health insights, where they must connect
042 medical data with external resources to inform ad-
043 vice (Heydari et al., 2025); and scientific discovery,
044 where it is necessary to integrate findings scattered
045 across research articles to form new hypotheses (Ya-
046 mada et al., 2025). Success in these settings requires
047 operating over richly structured knowledge environ-
048 ments, navigating interlinked documents, resolving
049 indirect references, and integrating evidence spread
across multiple sources.

050 Yet, as LMs continue to be trained on massive web
051 corpora (often with undisclosed training data), it re-
052 mains unclear to what extent their performance re-
053 flects genuine reasoning versus the reciting of mem-
054 orized knowledge (Carlini et al., 2023; Wu et al.,

Multi-hop QA

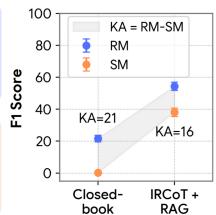
Q: Who is the namesake of
the city where **Ryan
Reynolds** was born?

A: **George Vancouver**

Q: Who is the namesake
of the city where **Derek
Veylan** was born?

A: **Altheon Metronis**

GPT-5-mini Performance



Page Navigation

Navigate from **Geoffrey
Hinton** to **Ryan Reynolds**

Navigate from **Caleb Ardent**
to **Derek Veylan**

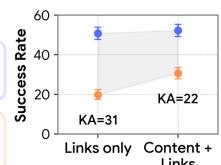


Figure 1: **Controlled experiments from SYNTHWORLDS corpora.** We measure the *knowledge advantage gap* (KA) as the performance difference between parallel tasks mapped to **real-world** (RM) and **synthetic** (SM) entities. Retrieval and page content boosts performance but the gap persists.

054 2024). Many benchmark tasks depend on *factual world knowledge* models likely encountered during
 055 training (Sainz et al., 2023; Xu et al., 2024b; Zhou et al., 2023). This undermines two goals: scientifically,
 056 it prevents isolating reasoning ability (i.e., functional linguistic competence) from memorization
 057 (i.e., formal linguistic competence) (Mahowald et al., 2024; Lu et al., 2024); practically, it
 058 limits confidence in deploying systems to novel environments (i.e., scientific discovery).

059 To distinguish reasoning from reciting, researchers have explored several strategies. One approach
 060 is manual curation of “*clean*” evaluation sets, which provides novelty but is costly, difficult to scale,
 061 and requires continual updates. For example, ToolQA (Zhuang et al., 2023), a benchmark released
 062 in 2023 to distinguish between questions answerable from an LM’s internal knowledge and those
 063 requiring external information, included GSM8K questions derived from “*error cases made by Chat-
 064 GPT*” at the time. However, subsequent work has shown that newer LMs may already memorize
 065 many of these answers (Zhang et al., 2024; Mirzadeh et al., 2025). Another approach, synthetic
 066 dataset generation (Huang et al., 2025; Hsieh et al., 2024), promises scalability, but often involves
 067 using existing content directly (e.g., novels) and thereby results in parametric knowledge leakage
 068 or relies on overly simplistic templates (e.g., “*The job of David is a farmer. The hobby of David is
 069 birdwatching.*”), limiting their ability to probe reasoning in realistic, richly interconnected settings.

070 Crucially, evaluations based *only* on synthetic unseen tasks still leave open questions about per-
 071 formance. Success demonstrates reasoning in isolation, but it does not reveal how much models
 072 typically rely on prior knowledge as a scaffold. Failure, on the other hand, is ambiguous: the rea-
 073 soning chain underlying the task may be too difficult for models to succeed, or the model may
 074 simply lack the background knowledge it usually exploits. Without controlling both task difficulty
 075 and requirements for parametric knowledge, such evaluations leave the contributions of reasoning
 076 and memorization entangled.

077 To address this, we introduce SYNTHWORLDS, a framework for disentangling reasoning from fac-
 078 tual knowledge. Parallel synthetic corpora are constructed to represent different *worlds* that replicate
 079 the structure and complexity of real-world information ecosystems. One corpus is mapped to **real-
 080 world** entities (e.g. *Geoffrey Hinton*), while the other is mapped to **synthetic** entities (e.g. *Caleb
 081 Ardent*), thereby obscuring the usefulness of parametric knowledge. This design allows us to quanti-
 082 fy the *knowledge advantage gap* (i.e., the performance difference between real-mapped [RM] and
 083 synthetic-mapped [SM] settings) and to evaluate how knowledge acquisition methods (e.g., provid-
 084 ing page content, retrieval-augmented generation) and integration strategies (e.g., chain-of-thought
 085 prompting, agentic reasoning) impact this gap (Fig. 1). The gap clarifies to what extent models rely
 086 on reasoning versus recall, and whether augmentation substitutes for or amplifies prior knowledge.

087 To support comparisons at scale, SYNTHWORLDS automatically generates parallel corpora from
 088 triplet facts in a knowledge graph (§3). To obscure factual knowledge, entities are renamed with
 089 surface-form-consistent transformations that preserve both type and name-derivation consistency
 090 before rendering facts into documents (Agarwal et al., 2021; Josifoski et al., 2023). Specifically,
 091 people receive person names (Geoffrey Hinton → Caleb Ardent), cities receive city names (Toronto
 092 → Metrovale), and derived names maintain consistency (University of Toronto → University of
 093 Metrovale, not University of Grandvale). This process yields corpora with identical reasoning struc-
 094 tures while removing familiarity with entity-specific facts, resulting in coherent worlds where tasks
 095 require reasoning over complex documents under controlled relevance of parametric knowledge.

096 To demonstrate the utility of our SYNTHWORLDS framework, we generate two parallel corpora de-
 097 rived from Wikidata: SYNTHWORLD-RM and SYNTHWORLD-SM (§4). On top of each corpus, we
 098 construct two reasoning-intensive tasks as case studies: multi-hop question answering (QA) (Trivedi
 099 et al., 2022; Ho et al., 2020) and page navigation (West & Leskovec, 2012) with fine-grained control
 100 over difficulty.

101 In our experiments, we evaluate LMs on these tasks to quantify the knowledge advantage gap, first
 102 in settings where models rely only on parametric knowledge (closed-book QA for multi-hop reason-
 103 ing and page names only for navigation), and then under conditions where knowledge augmentation
 104 (retrieval for QA, access to page contents for navigation) and integration strategies (e.g., chain-of-
 105 thought prompting) are provided (§5). Across both tasks, we find clear performance gaps between
 106 real-mapped and synthetic-mapped settings. While knowledge integration improves performance in
 107 both cases (and in some instances narrows the gap), the gap persists. This persistence highlights

108 opportunities for future work to design more effective knowledge integration schemes and to systematically study system behavior when models encounter novel environments (§6). **We contribute:**
 109
 110

111 1. **A scalable framework for generating rich, interconnected corpora and tasks** that disentangle
 112 and task reasoning difficulty from parametric knowledge.
 113 2. **Two parallel corpora with corresponding task datasets.** We instantiate the SYNTH-
 114 WORLDs framework with SYNTHWORLD-RM and SYNTHWORLD-SM, paired at the document,
 115 fact, and task levels to enable controlled evaluation. Each corpus contains 6,920 documents cov-
 116 ering 161K facts, along with 1.2K multi-hop QA and 1K page navigation instances. To support
 117 future research, we release these resources publicly.¹
 118 3. **An empirical analysis of LMs across parametric-only and knowledge-augmented settings**
 119 using our parallel datasets to quantify the knowledge advantage gap, which prior setups do not
 120 fully isolate. Our analysis reveals persistent shortcomings even with knowledge augmentation.

121 2 RELATED WORK

122 **Human Curated Data for Reasoning Evaluation.** As LM capabilities continue to improve and be-
 123 come widely deployed, researchers have relied on manually curated benchmarks to evaluate reasoning
 124 in settings not already covered by training data (Kazemi et al., 2025; Wei et al., 2025; Hendrycks
 125 et al., 2021; Cobbe et al., 2021; Bean et al., 2024; Srivastava et al., 2023; Tang & Yang, 2024; SU
 126 et al., 2025). These benchmarks are effective when first released but grow less informative over
 127 time as time passes. For example, MuSiQue (Trivedi et al., 2022), released in 2021 as a multi-hop
 128 QA benchmark, was originally designed to contain questions that models could not answer without
 129 the reference text. Despite this intent, it is still used across many evaluations today (Li et al., 2024;
 130 Zhang et al., 2025; Gutiérrez et al., 2025), even though current LMs (e.g., Llama-3.3-70B) achieve
 131 over 26% F1 score on these questions without any documents (Gutiérrez et al., 2025). This makes
 132 it difficult to assess whether improved performance reflects genuine advances in reasoning and re-
 133 trieval capabilities that would be informative of systems deployed in unseen environments. As a
 134 result, researchers must continually spend effort to construct new datasets and tasks (Gu et al., 2024;
 135 Tang & Yang, 2024; Monteiro et al., 2024; Bai et al., 2025). These efforts require substantial ex-
 136 pertise, grow increasingly complex as models advance, and are slow and costly to scale. In contrast,
 137 SYNTHWORLDs introduces a scalable framework to construct complex text data and associated
 138 reasoning tasks, reducing the manual curation burden while maintaining evaluation quality.

139 **Synthetic/Perturbed Data for Reasoning Generalization.** Given the resources needed to build
 140 high-quality human-generated data, researchers have developed methods to compose synthetic data
 141 or introduce perturbations to evaluate the reasoning generalization of LMs (Huang et al., 2025; Wu
 142 et al., 2024; Levy et al., 2024; Hsieh et al., 2024; Gu et al., 2025). These approaches reveal im-
 143 portant weaknesses when LMs are tested outside familiar conditions or over long contexts, but they
 144 do not disentangle reasoning ability from reliance on parametric factual knowledge. Other efforts
 145 address this separation more directly, for example by focusing on real-time factual updates (Kasai
 146 et al., 2023; Vu et al., 2024) or by generating synthetic text (Gong et al., 2025; Allen-Zhu & Li,
 147 2024; Monea et al., 2024). However, such work typically targets narrow aspects of knowledge or
 148 simplifies away the complexity and interconnectedness of real-world corpora, making it difficult to
 149 generalize findings to realistic scenarios (e.g., web navigation). Our work complements these lines
 150 by isolating the independent impacts of LM reasoning and parametric factual knowledge on task
 151 performance. Through controllable parallel dataset construction, we enable precise measurement of
 152 the *knowledge advantage gap* across common LM settings (e.g., in-context learning, RAG, agentic
 153 workflows), analysis on how different forms of knowledge augmentation influence this gap, and
 154 direct comparisons of model behavior in *novel* versus *familiar* settings.

155 3 SYNTHWORLDs: PARALLEL CORPORA FOR CONTROLLED EVALUATION

156 The main idea of SYNTHWORLDs is to construct parallel corpora and tasks that describe two worlds:
 157 one grounded in real-world entities, where factual knowledge encoded in language models’ parame-
 158

161 ¹Dataset and evaluation code available at <https://anonymous.4open.science/r/synthworld-experiments-CE26/>.

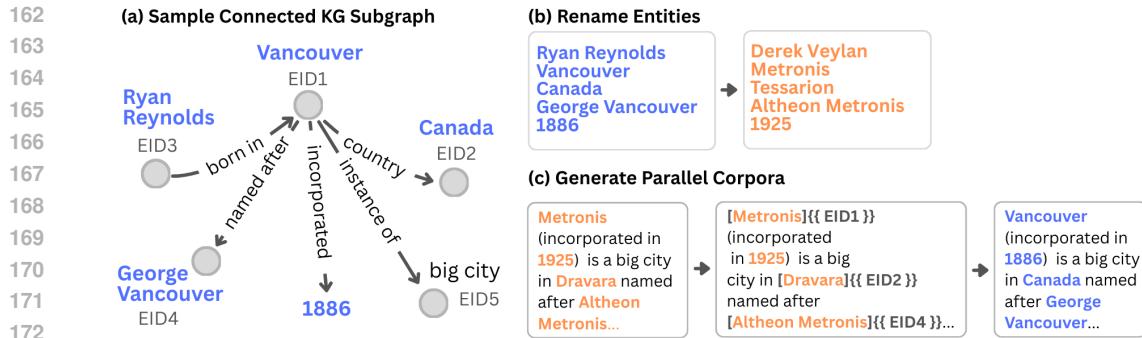


Figure 2: **Overview of SYNTHWORLDS Corpora Construction (Toy Example).** A connected subgraph is sampled from a large knowledge base (a). To obscure factual knowledge, entity labels are renamed from real-world labels (real-mapped) to synthetic name (synth-mapped) (b). From synth-mapped triplets, we generate synth-mapped documents. These documents are converted to real-mapped documents through additional LM steps with symbolic references (c). The final output is two parallel corpora: one **real-mapped**, one **synth-mapped**. Using the corpora, we construct parallel reasoning tasks (§4.1).

ters is potentially useful, and another built from synthetic entities, where such knowledge is deliberately uninformative. We define **factual knowledge** as entity-specific world knowledge tied to named entities (e.g., “*Barack Obama served as U.S. President from 2009 to 2017*”). In contrast, **domain-general knowledge** is not tied to named entities (e.g., arithmetic, physical laws, or the concept of an election or a university).² This distinction ensures that tasks maintain equivalent reasoning demands while preventing solutions that rely solely on recalling memorized entity facts. The reasoning preserved includes *commonsense* (e.g., hospitals have doctors), *compositional* (e.g., if a university has a medical school and medical schools train doctors, then the university trains doctors), *logical* (e.g., the parent of the parent of X is X ’s grandparent), and *temporal* reasoning.

Quantifying Parametric Knowledge in Reasoning Tasks. Constructing parallel corpora and tasks enables us to formally quantify the contribution of parametric knowledge. For a task, let P_R denote performance on the corpus with real-world entities (where parametric knowledge is useful), and P_S denote performance on the corpus with synthetic entities (where parametric knowledge is uninformative). We define the *knowledge advantage* gap as $KA = P_R - P_S$, quantifying the contribution of parametric knowledge to task performance. We further distinguish between two settings: the baseline case, where models rely only on their parametric knowledge ($KA^{\text{base}} = P_R^{\text{base}} - P_S^{\text{base}}$), and the augmented case, where models are provided with external knowledge acquisition and integration strategies ($KA^{\text{ext}} = P_R^{\text{ext}} - P_S^{\text{ext}}$). Examples include RAG for multi-hop QA and reading page content during agentic page navigation, though agents may employ other exploration and knowledge integration strategies in novel environments. In the baseline setting, P_S^{base} is expected to be near random (e.g., 0% accuracy for multi-hop QA; random walk performance for page navigation) since parametric knowledge is uninformative, so KA^{base} reflects the pure contribution of parametric memory. Additionally, with $KA^{\text{base}} - KA^{\text{ext}}$, we quantify how much the knowledge advantage closes when allowing external knowledge integration.

Framework Goals. To fairly measure KA, SYNTHWORLDS corpora and tasks are constructed with four core goals: (1) **emulate real-world complexity** by capturing the structure, interconnections, and both factual consistency (facts are mutually coherent) and semantic consistency, where semantic consistency requires that surface forms remain compatible with the entity’s ontological type. For example, university names remain university-like (University of Toronto → University of Grandvalle, not → Grandvalle Bank) and libraries remain library-like (Central Library → Oakwood Public Library, not → Central Stadium). This ensures surface form artifacts do not differ between real and synthetic variants, making performance on SYNTHWORLDS informative of reasoning in realistic tasks; (2) **enable parallel real- and synthetic-entity variants** to disentangle reasoning and factual knowledge; (3) **precisely control task difficulty** to support observations across levels of task com-

²Practically, we define named entities as proper nouns (i.e., capitalized) in common usage (Wikidata contributors, 2025) or recognized by NER models (e.g., the common noun *actor* vs. the named entity *Ryan Reynolds*).

plexity; and (4) **be fully automatic** such that new SYNTHWORLDS corpora and tasks can be readily constructed to continually provide novel evaluation data (guarding against evaluation corpora being included in pre- and post-training datasets).

Obscuring Factual Knowledge in Synthetically Generated Corpora. Similar to Wikipedia documents, SYNTHWORLDS’ corpora consist of documents about a specific entity with references to other entities in the corpus. The pipeline operates in three stages (Fig 2): (1) universe construction, (2) surface-form perturbation of named entities and timestamps, and (3) document generation.

First, to ensure the world is factually consistent, the pipeline samples a universe of connected triplet facts (i.e., subject → relation → object) from an existing (and assumed to be consistent) knowledge base (Fig 2a). Next, to remove parametric knowledge while maintaining consistency, entities are systematically renamed while preserving type information and context (e.g., ensuring that the rename of *Vancouver* is still a city named after *George Vancouver* the person) (Fig 2b). Finally, based on the synth-mapped facts (using the knowledge graph structure and new synthetic names), we generate documents using LMs, following prior work on generating documents from knowledge graph facts (Fig 2c) (Agarwal et al., 2021; Josifoski et al., 2023). Specifically, we first generate documents in the synth-mapped universe consistent with the triplets. We then insert symbolic references to entities in the text. Finally, we map these references to real-mapped labels, converting each synthetic document into its real-mapped counterpart.

The pipeline outputs two parallel corpora derived from a shared set of knowledge graph triplets: one mapped to real-world entities and the other to synthetic entities. Both corpora preserve identical sentence structures and world-consistent facts, differing only in their surface-form labels. For space, we include details of SYNTHWORLDS’ generation framework in Appendix A.

4 SYNTHWORLDS-RM AND SYNTHWORLDS-SM CORPORA AND TASKS

Using the SYNTHWORLDS framework (§3; Appendix A), we construct two parallel corpora and tasks: SYNTHWORLD-RM consisting of real-mapped entities and SYNTHWORLD-SM containing synthetic named entities. For space, we include dataset construction details in Appendix B which instantiates the framework on Wikidata. Our specific Wikidata pipeline can cheaply generate new datasets through natural stochasticity (sampling of renames) and by varying hyperparameters (e.g., entity types, seed nodes, knowledge graph sampling procedures), preventing SYNTHWORLDS from being overfit in evaluations.

Pages	Tokens	Facts	Entity Types	Relation Types	Avg Degree	Density	# Mhop QA	# Nav Pairs
6,290	~1.5M	161K	956	354	14.6	0.23%	1.2K	1K

Table 1: Summary Statistics for SYNTHWORLD-RM and SYNTHWORLD-SM.

Dataset Statistics. Table 1 summarizes our dataset. SYNTHWORLD-RM/SM each contain 6290 documents and over 1.5M tokens in total. The hyperlink graph is sparse, with an edge density of 0.23%. Its degree distribution is heavy-tailed: most pages have only a few links, while a small number act as hubs with disproportionately many incoming or outgoing connections. Both characteristics mirror the structure of real-world information networks such as the Web or Wikipedia (Adamic & Huberman, 2000; Kumar et al., 2000). Additional figures/tables (including cost of constructing our datasets) and qualitative examples of the dataset are provided in Appendix B.7 and B.8.

4.1 CASE STUDIES: PARALLEL TASKS WITH CONTROLLABLE DIFFICULTY

Given SYNTHWORLD-RM/SM corpora, we construct two tasks as case studies to evaluate LM reasoning: multi-hop QA and page navigation.

Multi-hop QA. Multi-hop questions are questions which require reasoning across multiple sources of evidence (Fig 3). For constructing these questions, we follow MuSiQue (Trivedi et al., 2022) and construct multi-hop questions through single-hop question composition (Fig. 3b). We build each multi-hop question using a specific graph motif composed of triplets, where each triplet corresponds



Figure 3: **Multi-hop QA Construction.** Subgraphs matching reasoning motifs are sampled with constraints to ensure uniqueness, diversity, and multi-hop reasoning (a). From their triplet facts, we generate synth-mapped single-hop questions (b), which are composed into a synth-mapped multi-hop question (c). Using the synth-to-real entity mapping, we replace synth names with real names (d). The final output is parallel sets of **real-mapped** and **synth-mapped** multi-hop questions.

to one single-hop question that can be composed into the final multi-hop question. This graph motif indicates a specific multi-hop reasoning structure. Table 3 summarizes all motifs used in our dataset.

Specifically, given the facts used to generate the synth-mapped documents, we first construct a global fact graph G_{facts} where nodes represent entities and edges represent facts, with each fact annotated by the page where it occurs. The fact graph structure G_{facts} is identical for both the synth-mapped and real-mapped corpora. From this graph, we sample subgraphs $S \subseteq G_{\text{facts}}$ that match desired reasoning motifs, ensuring that each reasoning step draws from a different page.

Next, we use an LM to generate a single-hop question for each unique triplet $(u, r, v) \in S$, where u and v denote the subject and object entities, respectively, and r denotes the relation between them. We start with the synth-mapped entities to generate single-hop questions. For automatic quality validation, we verify that the subject entity is mentioned in the corresponding question. We prompt a LM to compose a multi-hop question from the single-hop questions. We ensure that root entities in the subgraph are mentioned in the question while all bridge entities (non-root and non-leaf) are not mentioned in the question text. Finally, to create a parallel task, we remap the entity names in both the question and answer.

This approach allows us to control task difficulty through different reasoning motifs while maintaining task parallelism by using the same sampled subgraph S across both corpora. Specifically, difficulty is determined by two factors: (1) the number of hops (or equivalently, the number of decomposed single-hop questions needed to answer the question), and (2) the structural complexity of the question’s motif pattern. Table 3 illustrates examples of reasoning motifs and resulting questions. In our dataset, motifs range from 2 decomposed questions (motif A) to 4 (motifs D, E, F). Motifs that contain others as subgraphs are strictly more difficult (e.g., D > B > A and F > C, E > C) since they require an additional reasoning hop. Additional details on multi-hop QA construction and ensuring task quality and diversity are included in Appendix B.4 with prompts in B.10.

Page Navigation. In page navigation, an agent is asked to navigate from a source to target page (e.g., navigate from *Geoffrey Hinton* to *Ryan Reynolds*) using only the hyperlinks on the page. This task is broadly related to web navigation and agentic reasoning. At each page, agents must formulate hypotheses (e.g., "the link to *University of Toronto* might lead closer to *Ryan Reynolds* since both are Canadian"), evaluate alternative decisions, and integrate information learned from prior steps (Yao et al., 2023; Wang et al., 2025). Pages that are more difficult to navigate (i.e., requiring more steps and presenting more choices at each step) further increase the demands on reasoning.

We treat the symbolic references created during document generation as hyperlinks to other pages. From this, we construct a document graph $G_{\text{doc}} = (V_{\text{doc}}, E_{\text{doc}})$ where nodes V_{doc} are documents centered around specific entities and edges $(u, v) \in E_{\text{doc}}$ indicate a hyperlink from document u to document v . Note that this graph structure is identical for both the synth-mapped and real-mapped corpora, preserving task parallelism. Creating a page navigation task simply requires specifying a source and target page. To measure and control for difficulty, we use the expected random walk distance (i.e., expected number of steps for a random walk) between two nodes as a proxy for task difficulty and sample node pairs according to different distance buckets (Chandra et al., 1989). Higher values indicate that more intermediate decisions are required at each step, as the agent must navigate through a longer chain of choices to reach the goal.

324 **Task Statistics.** In total, we construct 1,200 parallel multi-hop questions spanning six reasoning
 325 structures, as well as 1,000 parallel page-navigation pairs organized into five difficulty buckets
 326 (random-walk distances of 50–1K, 1K–10K, 10K–100K, 100K–1M, and 1M–10M).
 327

328 5 EXPERIMENTS 329

330 To study the *knowledge advantage gap*, we evaluate models on SYNTHWORLD-RM/SM, in multi-
 331 hop QA and page navigation. We evaluate [six](#) models: GPT-5-mini ([OpenAI, 2025](#)) (reasoning
 332 effort set to medium), Gemini-2.0-Flash ([Gemini Team, 2025](#)), [gpt-oss-20B](#), [gpt-oss-120B](#) ([Ope-
 333 nAI et al., 2025](#)), [Kimi-K2-Instruct](#) ([AI & the Kimi Team, 2025a](#)), and [Kimi-K2-Thinking](#) ([AI &
 334 the Kimi Team, 2025b](#)), enabling observations across model families, model sizes, and instruct vs.
 335 [thinking models](#). Additional experiment details and evaluation prompts are in [Appendix C](#).
 336

337 **Multi-hop QA Baselines.** We evaluate three primary baselines: (1) *Closed-book*, where the model
 338 has no access to documents and answers directly from its parametric knowledge (KA^{base}); (2) *One-
 339 step RAG*, where the model retrieves supporting documents once before answering (KA^{RAG}); and
 340 (3) *IRCoT + RAG* ([Trivedi et al., 2022](#)), which interleaves retrieval with chain-of-thought reasoning,
 341 enabling iterative reasoning and retrieval steps($KA^{CoT+RAG}$). For retrieval, we use the HippoRAG
 342 2 retriever, designed for factual, multi-hop contexts ([Gutiérrez et al., 2024](#)).
 343

344 In addition, we include a *Reading Comprehension* condition in which the model is given all gold (2–
 345 4 documents depending on graph motif, examples in [Table 3](#)) and additional distractor documents,
 346 equaling 10 total. This condition serves two interpretations: (i) it provides an upper bound when
 347 retrieval is not a bottleneck, and (ii) it separates the inherent difficulty of the reasoning task from
 348 the challenge of retrieving relevant evidence in unfamiliar settings. All baseline prompts for QA are
 349 included in [Appendix C.1](#).
 350

351 **Page Navigation Baselines.** Page navigation tests an agent’s ability to plan and reason over a
 352 linked knowledge environment. For page navigation, we follow the design of existing tool-use
 353 agents ([Yang et al., 2024](#); [Gu et al., 2025](#)) and evaluate an agent equipped with two function-
 354 calling tools: `click_link`, which allows the agent to click any link on the current page, and
 355 `backtrack`, which allows the agent to return to a previously visited page. To address our navi-
 356 gation research questions, we evaluate the agent under two observation conditions: (1) *Links Only*,
 357 where the agent observes only the set of outgoing links on each page (KA^{base}); and (2) *Content*
 358 + *Links*, where the agent observes both the outgoing links and the full page text ($KA^{content}$). We
 359 include all prompts for agentic navigation in [Appendix C.2](#).
 360

361 The *Links Only* condition isolates the contribution of parametric knowledge and semantic familiarity,
 362 since navigation must rely entirely on recognizing entities in link text. The *Content + Links* condi-
 363 tion tests whether access to textual content can compensate for the absence of parametric knowledge
 364 by providing additional evidence for navigation decisions. In both settings, the agent is limited to
 365 a maximum of 30 steps. This cap is well above the distribution of shortest path lengths (median
 366 5, maximum 11), ensuring all tasks remain solvable while avoiding unbounded exploration. In our
 367 subsequent results, we observe that this bound is sufficient for meaningful exploration.
 368

369 **Metrics.** For all multi-hop QA experiments, we report token-based F1 scores for task performance
 370 following prior work ([Trivedi et al., 2022](#)). Following HippoRag 2 ([Gutiérrez et al., 2025](#)), we also
 371 report recall@5 for RAG baselines to evaluate retrieval quality. For page navigation, we report the
 372 success rate of reaching the target page.
 373

374 6 RESULTS AND DISCUSSION 375

376 We show results across task buckets for multi-hop QA and page navigation in [Figures 4](#) and [5](#). We
 377 report aggregated results for all task instances in [Table 4](#) and [5](#) in the Appendix.
 378

379 **RQ1: What is the knowledge advantage gap when relying solely on parametric knowledge?**
 380 In multi-hop QA, across models, we observe the baseline performance in RM, $P_R^{base} \approx 20$, indicat-
 381 ing that SYNTHWORLD-RM presents questions that LMs *can* answer using parametric knowledge
 382 ([Table 4](#); Closed-book, RM). In contrast, the near-zero P_S^{base} validates that SYNTHWORLD-SM
 383

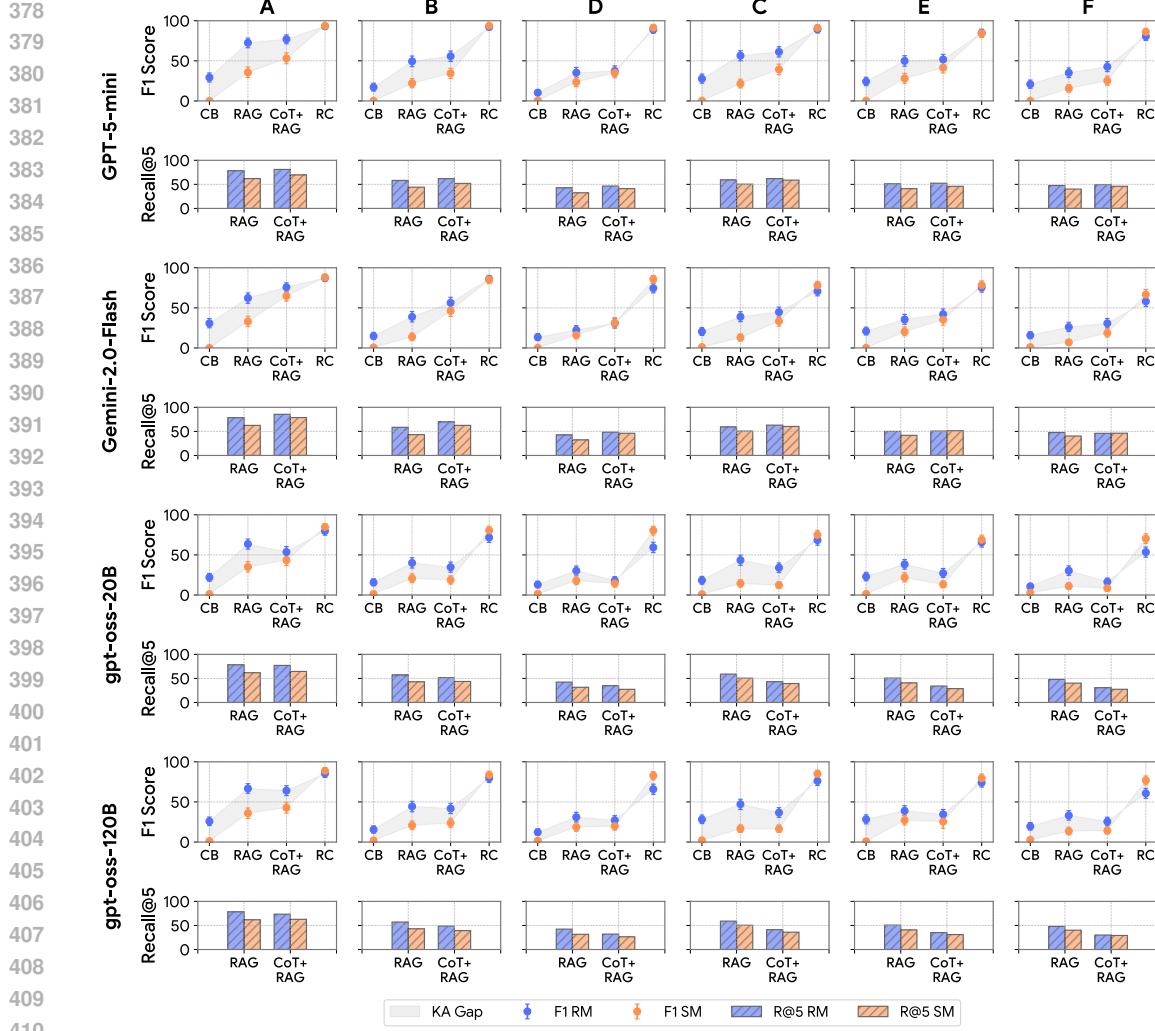


Figure 4: **Multi-hop QA Results by Reasoning Motifs.** We report F1 scores on SYNTHWORLD-RM (**RM**) and SYNTHWORLD-SM (**SM**), along with the *knowledge advantage gap* ($KA = F1_{RM} - F1_{SM}$). Settings: CB = Closed-book, RAG = One-step RAG, CoT+RAG = IRCoT + RAG, RC = Reading Comprehension. We show Recall@5 for RAG baselines (by construction, CB has recall = 0 and RC has recall = 1). IRCoT + RAG substantially reduces the KA gap compared to the CB baseline, primarily due to improved retrieval. Example questions for each motif are given in Table 3.

questions *cannot* be solved with parametric knowledge alone (Table 4; Closed-book, SM). Overall, $KA_{base} \approx 20$ (Table 4; Closed-book, KA). As task difficulty increases, P_R^{base} decreases as expected, while P_S^{base} remains at 0, showing that the gap would be even wider if we restricted evaluation to easier QA tasks (Fig. 4; CB left to right). In the reading comprehension setting, performance is equalized or even stronger in the SM cases because LMs are not distracted by parametric knowledge that could interfere with grounding its reasoning in the content (Monea et al., 2024).

For page navigation, we find a larger gap for GPT-5-mini and Kimi-K2 models ($KA^{base} \approx 30$) than for Gemini-2.0-Flash and gpt-oss models ($KA^{base} \approx 20$) (Table 5; Links Only, KA). This suggests the first set of models are better able to leverage parametric knowledge to locate the target page. Across difficulty levels, performance drops for both RM and SM tasks, but the gap persists. At the easiest difficulty, the gap narrows slightly, as models in SM can exploit the structure and semantics of hyperlinks to achieve modest success.

RQ2: To what extent does knowledge augmentation help close the gap? Knowledge augmentation with One-step RAG improves absolute performance across both RM and SM tasks. However,

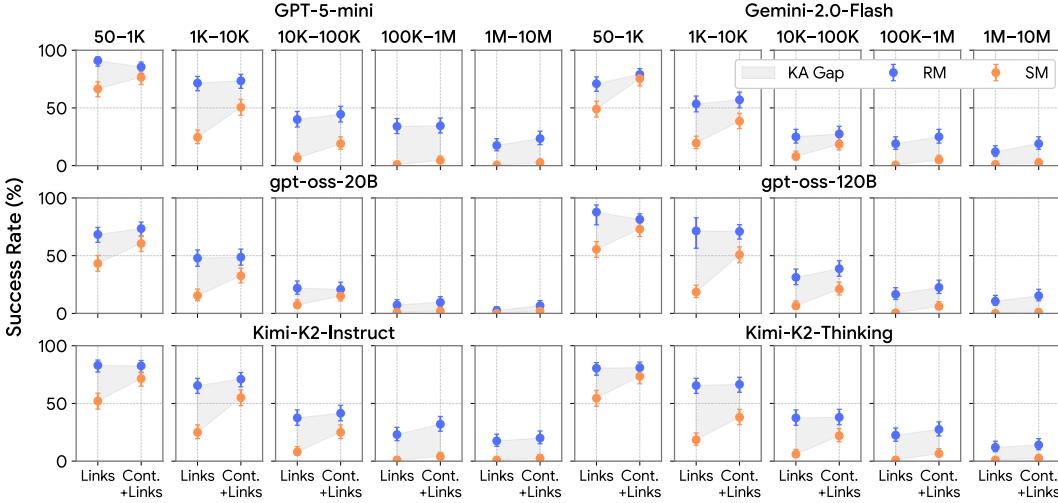


Figure 5: **Page Navigation Results by Difficulty (i.e., Expected Random Walk Distance).** We report success rate on SYNTHWORLD-RM (RM) and SYNTHWORLD-SM (SM) and the *knowledge advantage gap* (KA = $\text{Success}_{\text{RM}} - \text{Success}_{\text{SM}}$). Models consistently perform better on real-mapped corpora, especially in harder navigation tasks, indicating that parametric knowledge enables shortcuts. Page content (*Content + Links* vs. *Links Only*) benefits models more on synth-mapped corpora, narrowing the gap and showing its value in novel environments.

the knowledge advantage does not shrink; in fact, it widens. Specifically, $\text{KA}^{\text{base}} - \text{KA}^{\text{RAG}} = -4.0$ for GPT-5-mini and -1.3 for Gemini-2.0-Flash and similarly for other models (Table 4; Closed-book – One-step RAG), a pattern consistent across multiple difficulty levels (Fig. 4; A, B, C, F). This suggests that while One-step RAG benefits both RM and SM, it disproportionately benefits RM and reinforces models’ reliance on parametric knowledge. Meanwhile, IRCoT + RAG reduces the gap. Overall, $\text{KA}^{\text{base}} - \text{KA}^{\text{IRCoT+RAG}}$ is positive for both models, 5.2 for GPT-5-mini and 10.3 for Gemini 2.0-Flash (Table 4; Closed-book – IRCoT + RAG). We observe the gap closing across reasoning motifs (Fig. 4), indicating that interleaving retrieval with reasoning better aligns knowledge integration with task demands.³

To further probe this effect, we compare with the reading comprehension setting (i.e., perfect recall by construction). Triangulating reading comprehension F1-scores with F1-scores and retrieval recall from One-step RAG and IRCoT + RAG (Fig. 4; rows 2 and 4), we can infer that retrieval quality is a main driver of observed performance gaps. Retrieval performance improves slightly with IRCoT in both RM and SM, but retrieval in SM remains consistently lower than in RM. HippoRAG 2 uses an LM (GPT-5-mini in our experiments) to separately index RM and SM corpora and to retrieve documents given the input query. Given this setup, our results suggest that LM-based retrievers may not generalize well in novel environments, raising questions about the robustness of LM-indexed retrieval pipelines.

With respect to page navigation, across all task instance, we observe granting the agent access to page content improves performance, yielding differences of $\text{KA}^{\text{base}} - \text{KA}^{\text{content}} = 9.3$ and 7.0 for GPT-5-mini and Gemini-2.0-Flash, respectively (Table 5; Links Only – Content + Links). The performance gap narrows most on simpler navigation pairs (Fig. 5), though it remains present on more difficult ones.

To potentially explain the knowledge advantage gap, we analyze agent behavior by measuring how often externalized reasoning traces mention entities not observed during page navigation. For example, when tasked with navigating to the Brussels metropolitan area, a model trace included the statement: “*Ghent is in Belgium and likely links to Belgian geography or Brussels-related pages.*” We count the mentions of *Belgium* and *Belgian* as external, since they had not appeared in any previously visited page. In the SM setting, this rate is 0 by construction (and confirmed empirically).

³We note IRCoT + RAG does not improve absolute performance compared to One-step RAG for gpt-oss models as they struggle to follow the IRCoT prompt format. These results point to the importance of nuanced studies into the impact of knowledge integration.

486 Meanwhile, in the RM setting, we observe frequent reliance on external knowledge: under the *Links*
 487 *Only* condition, at least one external entity is mentioned in 48% of steps for GPT-5-mini and 60%
 488 for Gemini-2.0-Flash. Expanding access to *Content + Links* reduces these rates to 35% and 15%,
 489 respectively. Without page content, RM models tend to fall back on stored factual knowledge. In
 490 contrast, SM-like settings (where information is novel) offer only limited scope for fallback. This
 491 points to an opportunity to design agentic systems that both remain effective and efficiently acquire
 492 the necessary background knowledge.

493 **Insights enabled by SYNTHWORLDS.** The parallelism of SYNTHWORLDS enables controlled
 494 comparisons that isolate different aspects of model behavior. For example, it can allow us to ask
 495 when models take longer reasoning paths in the absence of recall or whether (and under what con-
 496 ditions) error types shift. It also makes it possible to investigate which system-level factors (such
 497 as retrieval quality in QA) and which core LM capabilities (as measured by reasoning or agentic
 498 benchmarks) lead to narrower or wider knowledge advantage gaps.

499 In our experiments, we studied knowledge integration through retrieval, both in single-step RAG
 500 and when interleaved with chain-of-thought or agentic workflows. These methods improved perfor-
 501 mance but did not fully eliminate the knowledge advantage gap. In QA, we see that it is a problem
 502 about knowledge acquisition (i.e., obtaining all the relevant documents), but additional thinking
 503 (e.g., CoT) can help. Meanwhile, in page navigation, even when models have the same content
 504 available, there is a gap as factual knowledge enables shortcuts. Beyond our case study results,
 505 SYNTHWORLDS allows researchers to examine alternative integration schemes. For example, in
 506 page navigation, what if models are integrated with retrieval to better plan their navigation? To what
 507 extent do long-context methods, where models must synthesize and retain relevant information with-
 508 out retrieval (Hsieh et al., 2024), or multi-agent workflows (Du et al., 2024), where group discussion
 509 and feedback shape integration, can help with knowledge augmentation?

510 **Future work and extending SYNTHWORLDS.** Our current work only scratches the surface of
 511 these possibilities. A limitation is that our experiments were conducted on the specific SYNTHWORLDS
 512 corpora and task designs we introduced, which may restrict the generality of our findings.
 513 These choices do not cover the full space of “constructed worlds” (or tasks) that could be defined
 514 by different relation types, connective structures, or contexts. Altering the way the corpora is con-
 515 structed could lead to different outcomes. Nonetheless, because SYNTHWORLDS is fully automatic,
 516 inexpensive, and flexible given any input knowledge base, we can generate alternate parallel corpora
 517 and probe these questions more broadly (see Appendix B.6 for an expanded discussion). Future work
 518 could impose targeted constraints on graph construction to highlight particular reasoning challenges
 519 or examine how parametric knowledge interacts with different underlying knowledge bases.

520 In general, our framework requires a high-quality knowledge graph to encode complex relation-
 521 ships and ensure factual consistency in synthesized text. Therefore, as modern extraction methods
 522 make knowledge graph construction from text increasingly feasible (Sainz et al., 2024; Xu et al.,
 523 2024a), graphs can be constructed directly from unstructured text (e.g., Wikipedia). Once built with
 524 consistent, accurate facts, the SYNTHWORLDS framework follows naturally. Likewise, SYNTHWORLDS
 525 is not limited to one knowledge graph or domain. The core idea, i.e., sampling a subgraph,
 526 renaming entities, and constructing parallel corpora and tasks, generalizes to any knowledge graph.
 527 The main requirement is understanding which entities should be renamed, which depends on the
 528 domain. For example, in mathematics, we could create parallel worlds with different notation sys-
 529 tems (RM: $x, y, f(x)$; SM: $\alpha, \beta, \phi(\alpha)$). Similarly, for code generation, we can consistently rename
 530 entire libraries (e.g., numpy/pandas) and function calls for the SM variant. By supporting con-
 531 trolled studies of reasoning, memory, and adaptation across varied settings, SYNTHWORLDS lays
 532 the groundwork for developing LM systems that are more robust and generalizable.

533 7 CONCLUSION

535 We present SYNTHWORLDS, a framework for disentangling the role of parametric knowledge in
 536 LM reasoning and retrieval. By constructing parallel corpora and tasks with controllable difficulty,
 537 SYNTHWORLDS reveals persistent performance gaps even when models have access to retrieval or
 538 page content. These findings highlight opportunities for advancing reasoning in novel environments
 539 and position SYNTHWORLDS as a scalable testbed for developing methods that generalize beyond
 540 reliance on parametric knowledge.

540 8 REPRODUCIBILITY
541

542 We provide full details of dataset construction, experimental setup, hyperparameters, and prompts
543 in the Appendix ensuring that our dataset and results could be reproduced. The dataset used in our
544 experiments is included in the supplementary material and will be publicly released. The code and
545 dataset for running all experiments is available at <https://anonymous.4open.science/r/synthworld-experiments-CE26/>.
546

547
548 REFERENCES
549

550 Lada A. Adamic and Bernardo A. Huberman. Power-law distribution of the world wide web.
551 *Science*, 287(5461):2115–2115, 2000. doi: 10.1126/science.287.5461.2115a. URL <https://www.science.org/doi/abs/10.1126/science.287.5461.2115a>.
552

553 Oshin Agarwal, Heming Ge, Siamak Shakeri, and Rami Al-Rfou. Knowledge graph based synthetic
554 corpus generation for knowledge-enhanced language model pre-training. In Kristina Toutanova,
555 Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cot-
556 terrell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of*
557 *the North American Chapter of the Association for Computational Linguistics: Human Lan-*
558 *guage Technologies*, pp. 3554–3565, Online, June 2021. Association for Computational Linguis-
559 *tics*. doi: 10.18653/v1/2021.naacl-main.278. URL <https://aclanthology.org/2021.naacl-main.278/>.
560

561 Moonshot AI and the Kimi Team. Kimi-k2: Open agentic intelligence (kimi-k2-instruct model
562 card). <https://moonshotai.github.io/Kimi-K2/>, 2025a. Open-source mixture-of-
563 experts (MoE) LLM with 1T parameters and 128K context length.
564

565 Moonshot AI and the Kimi Team. Kimi k2 thinking. <https://moonshotai.github.io/Kimi-K2/thinking.html>, 2025b. Documentation / model card for the Kimi-K2 model.
566

568 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: part 3.1, knowledge storage and
569 extraction. In *Proceedings of the 41st International Conference on Machine Learning*, ICML’24.
570 JMLR.org, 2024.
571

572 Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng
573 Xu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. LongBench v2: Towards deeper under-
574 standing and reasoning on realistic long-context multitasks. In Wanxiang Che, Joyce Nabende,
575 Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meet-*
576 *ing of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3639–3664,
577 Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-
578 251-0. doi: 10.18653/v1/2025.acl-long.183. URL <https://aclanthology.org/2025.acl-long.183/>.
579

580 Andrew Michael Bean, Simeon Hellsten, Harry Mayne, Jabez Magomere, Ethan A Chi, Ryan An-
581 drew Chi, Scott A. Hale, and Hannah Rose Kirk. LINGOLY: A benchmark of olympiad-level
582 linguistic reasoning puzzles in low resource and extinct languages. In *The Thirty-eighth Con-*
583 *ference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL
584 <https://openreview.net/forum?id=cLga8GStdK>.
585

586 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan
587 Zhang. Quantifying memorization across neural language models. In *The Eleventh International*
588 *Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=TatRHT_1cK.
589

590 A. K. Chandra, P. Raghavan, W. L. Ruzzo, and R. Smolensky. The electrical resistance of a graph
591 captures its commute and cover times. In *Proceedings of the Twenty-First Annual ACM Sym-*
592 *posium on Theory of Computing*, STOC ’89, pp. 574–586, New York, NY, USA, 1989. As-
593 sociation for Computing Machinery. ISBN 0897913078. doi: 10.1145/73007.73062. URL
594 <https://doi.org/10.1145/73007.73062>.
595

594 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
 595 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
 596 Schulman. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.

597

598 Wikipedia contributors. Wikipedia: Popular pages. https://en.wikipedia.org/wiki/Wikipedia:Popular_pages, 2025. Accessed: 2025-09-12.

599

600

601 Fred J. Damerau. A technique for computer detection and correction of spelling errors. *Commun. ACM*, 7(3):171–176, March 1964. ISSN 0001-0782. doi: 10.1145/363958.363994. URL <https://doi.org/10.1145/363958.363994>.

602

603

604 Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving
 605 factuality and reasoning in language models through multiagent debate. In *Proceedings of the*
 606 *41st International Conference on Machine Learning*, ICML’24. JMLR.org, 2024.

607

608 Gemini Team. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long
 609 context, and next generation agentic capabilities. Technical Report arXiv:2507.06261, Google
 610 DeepMind, 2025. URL <https://doi.org/10.48550/arXiv.2507.06261>. Version 4.

611

612 Albert Gong, Kamilė Stankevičiūtė, Chao Wan, Anmol Kabra, Raphael Thesmar, Johann Lee, Julius
 613 Klenke, Carla P Gomes, and Kilian Q Weinberger. Phantomwiki: On-demand datasets for rea-
 614 soning and retrieval evaluation. In *Forty-second International Conference on Machine Learning*,
 615 2025. URL <https://openreview.net/forum?id=DIZItj8ueN>.

616

617 Ken Gu, Ruoxi Shang, Ruien Jiang, Keying Kuang, Richard-John Lin, Donghe Lyu, Yue Mao,
 618 Youran Pan, Teng Wu, Jiaqian Yu, Yikun Zhang, Tianmai M. Zhang, Lanyi Zhu, Mike A Merrill,
 619 Jeffrey Heer, and Tim Althoff. BLADE: Benchmarking language model agents for data-
 620 driven science. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the*
 621 *Association for Computational Linguistics: EMNLP 2024*, pp. 13936–13971, Miami, Florida,
 622 USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.
 623 findings-emnlp.815. URL [https://aclanthology.org/2024.findings-emnlp.815/](https://aclanthology.org/2024.findings-emnlp.815).

624

625 Ken Gu, Zhihan Zhang, Kate Lin, Yuwei Zhang, Akshay Paruchuri, Hong Yu, Mehran Kazemi,
 626 Kumar Ayush, A. Ali Heydari, Maxwell A. Xu, Girish Narayanswamy, Yun Liu, Ming-Zher Poh,
 627 Yuzhe Yang, Mark Malhotra, Shwetak Patel, Hamid Palangi, Xuhai Xu, Daniel McDuff, Tim
 628 Althoff, and Xin Liu. Radar: Benchmarking language models on imperfect tabular data, 2025.
 629 URL <https://arxiv.org/abs/2506.08249>.

630

631 Bernal Jimenez Gutierrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. HippoRAG: Neuro-
 632 biologically inspired long-term memory for large language models. In *The Thirty-eighth Annual*
 633 *Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=hkujvAPVsg>.

634

635 Bernal Jiménez Gutiérrez, Yiheng Shu, Weijian Qi, Sizhe Zhou, and Yu Su. From RAG to mem-
 636 ory: Non-parametric continual learning for large language models. In *Forty-second International*
 637 *Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=LWH8yn4HS2>.

638

639 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
 640 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset.
 641 In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks*
 642 *Track (Round 2)*, 2021. URL <https://openreview.net/forum?id=7Bywt2mQsCe>.

643

644 Lucas Torroba Hennigen, Zejiang Shen, Aniruddha Nrusimha, Bernhard Gapp, David Sontag, and
 645 Yoon Kim. Towards verifiable text generation with symbolic references. In *First Conference on*
 646 *Language Modeling*, 2024. URL <https://openreview.net/forum?id=fib9qidCpY>.

647

A. Ali Heydari, Ken Gu, Vidya Srinivas, Hong Yu, Zhihan Zhang, Yuwei Zhang, Akshay Paruchuri,
 648 Qian He, Hamid Palangi, Nova Hammerquist, Ahmed A. Metwally, Brent Winslow, Yubin Kim,
 649 Kumar Ayush, Yuzhe Yang, Girish Narayanswamy, Maxwell A. Xu, Jake Garrison, Amy Arem-
 650 nto Lee, Jenny Vafeiadou, Ben Graef, Isaac R. Galatzer-Levy, Erik Schenck, Andrew Barakat,

648 Javier Perez, Jacqueline Shreibati, John Hernandez, Anthony Z. Faranesh, Javier L. Prieto, Con-
 649 nor Heneghan, Yun Liu, Jiening Zhan, Mark Malhotra, Shwetak Patel, Tim Althoff, Xin Liu,
 650 Daniel McDuff, and Xuhai "Orson" Xu. The anatomy of a personal health agent, 2025. URL
 651 <https://arxiv.org/abs/2508.20148>.

652 Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-
 653 hop QA dataset for comprehensive evaluation of reasoning steps. In Donia Scott, Nuria Bel,
 654 and Chengqing Zong (eds.), *Proceedings of the 28th International Conference on Computational
 655 Linguistics*, pp. 6609–6625, Barcelona, Spain (Online), December 2020. International Com-
 656 mittee on Computational Linguistics. doi: 10.18653/v1/2020.coling-main.580. URL <https://aclanthology.org/2020.coling-main.580/>.

657 Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, and
 658 Boris Ginsburg. RULER: What's the real context size of your long-context language models? In
 659 *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=kI0Bbc76Sy>.

660 Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai,
 661 Hui Yuan, Runzhe Wang, Yue Wu, Ming Yin, Shange Tang, Yangsibo Huang, Chi Jin, Xinyun
 662 Chen, Chiyuan Zhang, and Mengdi Wang. MATH-perturb: Benchmarking LLMs' math reason-
 663 ing abilities against hard perturbations. In *Forty-second International Conference on Machine
 664 Learning*, 2025. URL <https://openreview.net/forum?id=OZy70UggXr>.

665 Martin Josifoski, Marija Sakota, Maxime Peyrard, and Robert West. Exploiting asymmetry for
 666 synthetic training data generation: SynthIE and the case of information extraction. In Houda
 667 Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical
 668 Methods in Natural Language Processing*, pp. 1555–1574, Singapore, December 2023. Asso-
 669 ciation for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.96. URL <https://aclanthology.org/2023.emnlp-main.96/>.

670 Jungo Kasai, Keisuke Sakaguchi, yoichi takahashi, Ronan Le Bras, Akari Asai, Xinyan Velocity Yu,
 671 Dragomir Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui. Realtime QA: What's the answer
 672 right now? In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and
 673 Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=HfKOIPCvsv>.

674 Mehran Kazemi, Bahare Fatemi, Hritik Bansal, John Palowitch, Chrysovalantis Anastasiou, San-
 675 ket Vaibhav Mehta, Lalit K Jain, Virginia Aglietti, Disha Jindal, Peter Chen, Nishanth Dikkala,
 676 Gladys Tyen, Xin Liu, Uri Shalit, Silvia Chiappa, Kate Olszewska, Yi Tay, Vinh Q. Tran, Quoc V
 677 Le, and Orhan Firat. BIG-bench extra hard. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova,
 678 and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association
 679 for Computational Linguistics (Volume 1: Long Papers)*, pp. 26473–26501, Vienna, Austria, July
 680 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/
 681 2025.acl-long.1285. URL <https://aclanthology.org/2025.acl-long.1285/>.

682 Ravi Kumar, Prabhakar Raghavan, Sridhar Rajagopalan, D. Sivakumar, Andrew Tompkins, and Eli
 683 Upfal. The web as a graph. In *Proceedings of the Nineteenth ACM SIGMOD-SIGACT-SIGART
 684 Symposium on Principles of Database Systems*, PODS '00, pp. 1–10, New York, NY, USA, 2000.
 685 Association for Computing Machinery. ISBN 158113214X. doi: 10.1145/335168.335170. URL
 686 <https://doi.org/10.1145/335168.335170>.

687 Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catan-
 688 zaro, and Wei Ping. NV-embed: Improved techniques for training LLMs as generalist embedding
 689 models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL
 690 <https://openreview.net/forum?id=lgSYLSsDRe>.

691 Mosh Levy, Alon Jacoby, and Yoav Goldberg. Same task, more tokens: the impact of input length
 692 on the reasoning performance of large language models. In Lun-Wei Ku, Andre Martins, and
 693 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-
 694 putational Linguistics (Volume 1: Long Papers)*, pp. 15339–15353, Bangkok, Thailand, August
 695 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.818. URL
 696 <https://aclanthology.org/2024.acl-long.818/>.

702 Zhuowan Li, Cheng Li, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. Retrieval aug-
 703 mented generation or long-context LLMs? a comprehensive study and hybrid approach. In
 704 Franck Dernoncourt, Daniel Preoțiuc-Pietro, and Anastasia Shimorina (eds.), *Proceedings of*
 705 *the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*,
 706 pp. 881–893, Miami, Florida, US, November 2024. Association for Computational Linguistics.
 707 doi: 10.18653/v1/2024.emnlp-industry.66. URL <https://aclanthology.org/2024.emnlp-industry.66/>.

708

709 Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Harish Tayyar Madabushi, and Iryna Gurevych.
 710 Are emergent abilities in large language models just in-context learning? In Lun-Wei Ku, Andre
 711 Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association*
 712 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 5098–5139, Bangkok, Thailand,
 713 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.279.
 714 URL <https://aclanthology.org/2024.acl-long.279/>.

715

716 Kyle Mahowald, Anna A. Ivanova, Idan A. Blank, Nancy Kanwisher, Joshua B. Tenenbaum, and
 717 Evelina Fedorenko. Dissociating language and thought in large language models. *Trends in*
 718 *Cognitive Sciences*, 28(6):517–540, 2024. ISSN 1364-6613. doi: <https://doi.org/10.1016/j.tics.2024.01.011>. URL <https://www.sciencedirect.com/science/article/pii/S1364661324000275>.

719

720 Seyed Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and
 721 Mehrdad Farajtabar. GSM-symbolic: Understanding the limitations of mathematical reasoning in
 722 large language models. In *The Thirteenth International Conference on Learning Representations*,
 723 2025. URL <https://openreview.net/forum?id=AjXkRZIvjb>.

724

725 Giovanni Monea, Maxime Peyrard, Martin Josifoski, Vishrav Chaudhary, Jason Eisner, Emre Kici-
 726 man, Hamid Palangi, Barun Patra, and Robert West. A glitch in the matrix? locating and
 727 detecting language model grounding with fakepedia. In Lun-Wei Ku, Andre Martins, and
 728 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-*
 729 *putational Linguistics (Volume 1: Long Papers)*, pp. 6828–6844, Bangkok, Thailand, August
 730 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.369. URL
 731 <https://aclanthology.org/2024.acl-long.369/>.

732

733 Joao Monteiro, Pierre-Andre Noel, Étienne Marcotte, Sai Rajeswar, Valentina Zantedeschi, David
 734 Vazquez, Nicolas Chapados, Christopher Pal, and Perouz Taslakian. RepliQA: A question-
 735 answering dataset for benchmarking LLMs on unseen reference content. In *The Thirty-eight*
 736 *Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
 737 URL <https://openreview.net/forum?id=4diKTLmg2y>.

738

739 Liangbo Ning, Ziran Liang, Zhuohang Jiang, Haohao Qu, Yujuan Ding, Wenqi Fan, Xiao-yong Wei,
 740 Shanru Lin, Hui Liu, Philip S. Yu, and Qing Li. A survey of webagents: Towards next-generation
 741 ai agents for web automation with large foundation models. In *Proceedings of the 31st ACM*
 742 *SIGKDD Conference on Knowledge Discovery and Data Mining V.2*, KDD ’25, pp. 6140–6150,
 743 New York, NY, USA, 2025. Association for Computing Machinery. ISBN 9798400714542. doi:
 10.1145/3711896.3736555. URL <https://doi.org/10.1145/3711896.3736555>.

744

745 OpenAI. Gpt-5 mini. <https://openai.com/gpt-5>, 2025. Variant of the GPT-5 model,
 746 launched August 7, 2025.

747

748 OpenAI, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Ar-
 749 bus, Rahul K. Arora, Yu Bai, Bowen Baker, Haiming Bao, Boaz Barak, Ally Bennett, Tyler
 750 Bertao, Nivedita Brett, Eugene Brevdo, Greg Brockman, Sébastien Bubeck, Che Chang, Kai
 751 Chen, Mark Chen, Enoch Cheung, Aidan Clark, Dan Cook, Marat Dukhan, Casey Dvorak, Kevin
 752 Fives, Vlad Fomenko, Timur Garipov, Kristian Georgiev, Mia Glaese, Tarun Gogineni, Adam
 753 Goucher, Lukas Gross, Katia Gil Guzman, John Hallman, Jackie Hehir, Johannes Heidecke, Alec
 754 Helyar, Haitang Hu, Romain Huet, Jacob Huh, Saachi Jain, Zach Johnson, Chris Koch, Irina
 755 Kofman, Dominik Kundel, Jason Kwon, Volodymyr Kyrylov, Elaine Ya Le, Guillaume Leclerc,
 James Park Lennon, Scott Lessans, Mario Lezcano-Casado, Yuanzhi Li, Zhuohan Li, Ji Lin,
 Jordan Liss, Lily Liu, Jiancheng Liu, Kevin Lu, Chris Lu, Zoran Martinovic, Lindsay McCal-
 lum, Josh McGrath, Scott McKinney, Aidan McLaughlin, Song Mei, Steve Mostovoy, Tong Mu,

756 Gideon Myles, Alexander Neitz, Alex Nichol, Jakub Pachocki, Alex Paino, Dana Palmie, Ash-
 757 ley Pantuliano, Giambattista Parascandolo, Jongsoo Park, Leher Pathak, Carolina Paz, Ludovic
 758 Peran, Dmitry Pimenov, Michelle Pokrass, Elizabeth Proehl, Huida Qiu, Gaby Raila, Filippo
 759 Raso, Hongyu Ren, Kimmy Richardson, David Robinson, Bob Rotsted, Hadi Salman, Suvansh
 760 Sanjeev, Max Schwarzer, D. Sculley, Harshit Sikchi, Kendal Simon, Karan Singhal, Yang Song,
 761 Dane Stuckey, Zhiqing Sun, Philippe Tillet, Sam Toizer, Foivos Tsimpourlas, Nikhil Vyas, Eric
 762 Wallace, Xin Wang, Miles Wang, Olivia Watkins, Kevin Weil, Amy Wendling, Kevin Whinnery,
 763 Cedric Whitney, Hannah Wong, Lin Yang, Yu Yang, Michihiro Yasunaga, Kristen Ying, Wojciech
 764 Zaremba, Wenting Zhan, Cyril Zhang, Brian Zhang, Eddie Zhang, and Shengjia Zhao. gpt-oss-
 765 120b & gpt-oss-20b model card, 2025. URL <https://arxiv.org/abs/2508.10925>.
 766
 767 Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and
 768 Eneko Agirre. NLP evaluation in trouble: On the need to measure LLM data contamination for
 769 each benchmark. In *The 2023 Conference on Empirical Methods in Natural Language Processing*,
 770 2023. URL <https://openreview.net/forum?id=KivNpBsFAS>.
 771
 772 Oscar Sainz, Iker García-Ferrero, Rodrigo Agerri, Oier Lopez de Lacalle, German Rigau, and Eneko
 773 Agirre. GoLLIE: Annotation guidelines improve zero-shot information-extraction. In *The Twelfth
 774 International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Y3wpuxd7u9>.
 775
 776 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam
 777 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska,
 778 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Ko-
 779 curek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda
 780 Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Johan
 781 Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew
 782 La, Andrew Kyle Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh
 783 Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum,
 784 Arul Menezes, Arun Kirubarajan, Asher Mollokandov, Ashish Sabharwal, Austin Herrick, Avia
 785 Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej
 786 Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno
 787 Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinon, Cameron Diao, Cameron
 788 Dour, Catherine Stinson, Cedrick Argueta, Cesar Ferri, Chandan Singh, Charles Rathkopf, Chen-
 789 lin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Christopher Waites, Christian Voigt,
 790 Christopher D Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro,
 791 Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks,
 792 Dan Kilman, Dan Roth, C. Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí
 793 González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa,
 794 David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin,
 795 Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar
 796 Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova,
 797 Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, El-
 798 lie Pavlick, Emanuele Rodolà, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang,
 799 Ethan A Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozh-
 800 skii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chol-
 801 let, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germàn Kruszewski,
 802 Giambattista Parascandolo, Giorgio Mariani, Gloria Xinyue Wang, Gonzalo Jaimovich-Lopez,
 803 Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Ha-
 804 jishirzi, Harsh Mehta, Hayden Bogar, Henry Francis Anthony Shevlin, Hinrich Schuetze, Hi-
 805 romu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack
 806 Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B Si-
 807 mon, James Koppel, James Zheng, James Zou, Jan Kocon, Jana Thompson, Janelle Wingfield,
 808 Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski,
 809 Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Je-
 sujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller,
 John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-
 Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule,
 Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina

810 Ignatyeva, Katja Markert, Kaustubh Dhole, Kevin Gimpel, Kevin Omondi, Kory Wallace Math-
 811 ewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson,
 812 Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-
 813 Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis
 814 Oliveros-Colón, Luke Metz, Lütfi Kerem Senel, Maarten Bosma, Maarten Sap, Maartje Ter Ho-
 815 eve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco
 816 Maru, Maria Jose Ramirez-Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin
 817 Potthast, Matthew L Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova,
 818 Melody Arnaud, Melvin McElrath, Michael Andrew Yee, Michael Cohen, Michael Gu, Michael
 819 Ivanitskiy, Michael Starritt, Michael Strube, Michał Świdrowski, Michele Bevilacqua, Michi-
 820 hiro Yasunaga, Mihir Kale, Mike Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, Mo Ti-
 821 wari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun
 822 Peng, Nathan Andrew Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas
 823 Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Ni-
 824 tish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang,
 825 Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth
 826 Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy
 827 Liang, Peter W Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush
 828 Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade,
 829 Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm
 830 Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan
 831 Sikand, Roman Novak, Roman Sitelew, Ronan Le Bras, Rosanne Liu, Rowan Jacobs, Rui Zhang,
 832 Russ Salakhutdinov, Ryan Andrew Chi, Seungjae Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan
 833 Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wise-
 834 man, Samuel Gruetter, Samuel R. Bowman, Samuel Stern Schoenholz, Sanghyun Han, Sanjeev
 835 Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebas-
 836 tian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank
 837 Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar,
 838 Shubham Toshniwal, Shyam Upadhyay, Shyamolima Shammie Debnath, Siamak Shakeri, Simon
 839 Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene,
 840 Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie
 841 Lin, Stephen Prasad, Steven Piantadosi, Stuart Shieber, Summer Misherghi, Svetlana Kiritchenko,
 842 Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsunori Hashimoto, Te-Lin
 843 Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo
 844 Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj,
 845 Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas
 846 Raunak, Vinay Venkatesh Ramasesh, vinay uday prabhu, Vishakh Padmakumar, Vivek Sriku-
 847 mar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong,
 848 Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song,
 849 Yasaman Bahri, Yejin Choi, Yichi Yang, Sophie Hao, Yifu Chen, Yonatan Belinkov, Yu Hou,
 850 Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang,
 851 and Ziyi Wu. Beyond the imitation game: Quantifying and extrapolating the capabilities of lan-
 852 guage models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL
 853 <https://openreview.net/forum?id=uyTL5Bvosj>. Featured Certification.

854 Hongjin SU, Howard Yen, Mengzhou Xia, Weijia Shi, Niklas Muennighoff, Han yu Wang, Liu
 855 Haisu, Quan Shi, Zachary S Siegel, Michael Tang, Ruoxi Sun, Jinsung Yoon, Sercan O Arik,
 856 Danqi Chen, and Tao Yu. BRIGHT: A realistic and challenging benchmark for reasoning-
 857 intensive retrieval. In *The Thirteenth International Conference on Learning Representations*,
 858 2025. URL <https://openreview.net/forum?id=ykuc5q381b>.

859 Yixuan Tang and Yi Yang. Multihop-RAG: Benchmarking retrieval-augmented generation for multi-
 860 hop queries. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=t4eB3zYWBK>.

861 Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. MuSiQue:
 862 Multihop questions via single-hop question composition. *Transactions of the Association for
 863 Computational Linguistics*, 10:539–554, 2022. doi: 10.1162/tacl_a_00475. URL <https://aclanthology.org/2022.tacl-1.31/>.

864 Denny Vrandečić. Wikidata: a new platform for collaborative data collection. In *Proceedings of*
 865 *the 21st International Conference on World Wide Web*, WWW '12 Companion, pp. 1063–1064,
 866 New York, NY, USA, 2012. Association for Computing Machinery. ISBN 9781450312301. doi:
 867 10.1145/2187980.2188242. URL <https://doi.org/10.1145/2187980.2188242>.
 868

869 Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan
 870 Sung, Denny Zhou, Quoc Le, and Thang Luong. FreshLLMs: Refreshing large language mod-
 871 els with search engine augmentation. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar
 872 (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 13697–13720,
 873 Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/
 874 v1/2024.findings-acl.813. URL <https://aclanthology.org/2024.findings-acl.813>.
 875

876 Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and Graham Neubig. Agent workflow mem-
 877 ory. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=NTAhi2JEEE>.
 878

879 Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won
 880 Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecmp: A simple yet
 881 challenging benchmark for browsing agents, 2025. URL <https://arxiv.org/abs/2504.12516>.
 882

883 Robert West and Jure Leskovec. Human wayfinding in information networks. In *Proceedings of*
 884 *the 21st International Conference on World Wide Web*, WWW '12, pp. 619–628, New York,
 885 NY, USA, 2012. Association for Computing Machinery. ISBN 9781450312295. doi: 10.1145/
 886 2187836.2187920. URL <https://doi.org/10.1145/2187836.2187920>.
 887

888 Wikidata contributors. Help:label. <https://www.wikidata.org/wiki/Help:Label>,
 889 2025. Accessed: 2025-09-22.
 890

891 Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim,
 892 Jacob Andreas, and Yoon Kim. Reasoning or reciting? exploring the capabilities and limita-
 893 tions of language models through counterfactual tasks. In Kevin Duh, Helena Gomez, and Steven
 894 Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the As-
 895 sociation for Computational Linguistics: Human Language Technologies (Volume 1: Long Pa-
 896 pers)*, pp. 1819–1862, Mexico City, Mexico, June 2024. Association for Computational Linguis-
 897 tics. doi: 10.18653/v1/2024.naacl-long.102. URL <https://aclanthology.org/2024.naacl-long.102>.
 898

899 Derong Xu, Wei Chen, Wenjun Peng, Chao Zhang, Tong Xu, Xiangyu Zhao, Xian Wu, Yefeng
 900 Zheng, Yang Wang, and Enhong Chen. Large language models for generative information extrac-
 901 tion: A survey, 2024a. URL <https://arxiv.org/abs/2312.17617>.
 902

903 Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. Benchmarking benchmark leakage in
 904 large language models. *arXiv preprint arXiv:2404.18824*, 2024b. URL <https://arxiv.org/abs/2404.18824>.
 905

906 Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune,
 907 and David Ha. The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree
 908 search. *arXiv preprint arXiv:2504.08066*, 2025.
 909

910 John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan,
 911 and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering.
 912 In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL
 913 <https://openreview.net/forum?id=mXpq6ut8J3>.
 914

915 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
 916 ReAct: Synergizing reasoning and acting in language models. In *International Conference on*
 917 *Learning Representations (ICLR)*, 2023.

918 Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, William Song, Tiffany Zhao,
919 Pranav Vishnu Raja, Charlotte Zhuang, Dylan Z Slack, Qin Lyu, Sean M. Hendryx, Russell Ka-
920 plan, Michele Lunati, and Summer Yue. A careful examination of large language model per-
921 formance on grade school arithmetic. In *The Thirty-eight Conference on Neural Information*
922 *Processing Systems Datasets and Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=RJZRhMzZzH>.

924 Nan Zhang, Prafulla Kumar Choubey, Alexander Fabbri, Gabriel Bernadett-Shapiro, Rui Zhang,
925 Prasenjit Mitra, Caiming Xiong, and Chien-Sheng Wu. SireRAG: Indexing similar and related
926 information for multihop reasoning. In *The Thirteenth International Conference on Learning*
927 *Representations*, 2025. URL <https://openreview.net/forum?id=yp95goUAT1>.

928 Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin, Ji-
929 Rong Wen, and Jiawei Han. Don't make your llm an evaluation benchmark cheater, 2023. URL
930 <https://arxiv.org/abs/2311.01964>.

932 Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. ToolQA: A dataset for LLM
933 question answering with external tools. In *Thirty-seventh Conference on Neural Information*
934 *Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=pV1xV2RK6I>.

936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

972	APPENDIX TABLE OF CONTENTS	
973		
974		
975	A SYNTHWORLDS Framework	20
976		
977	B SYNTHWORLD-RM/SM Dataset Construction Details	21
978	B.1 Universe Construction	21
979	B.2 Surface-Form Perturbations	22
980	B.3 Parallel Document Generation	22
981	B.4 Multi-hop QA Construction	22
982	B.5 Human Validation	23
983	B.6 Discussion on Dataset Construction	23
984	B.7 Additional Figures and Tables	24
985	B.8 Qualitative Corpora Examples	28
986	B.9 Prompts for Corpora Construction	29
987	B.10 Prompts for Multi-hop QA Construction from Facts	36
988		
989		
990	C Experiment Details	43
991		
992		
993	C.1 Multihop QA Prompts	43
994	C.2 Page Navigation Prompts	46
995		
996	D Additional Experiment Tables	48
997		
998		
999		
1000		
1001		
1002		
1003		
1004		
1005		
1006		
1007		
1008		
1009		
1010		
1011		
1012		
1013		
1014		
1015		
1016		
1017		
1018		
1019		
1020		
1021		
1022		
1023		
1024		
1025		

1026 A SYNTHWORLDS FRAMEWORK

1028 In this section, we discuss the core formalization of the SYNTHWORLDS framework. Concrete de-
 1029 tails actualizing this framework in our SYNTHWORLD-RM/SM datasets are included in Appendix B.
 1030

1031 **World Knowledge Preliminaries.** Formally, our dataset generation takes as input a knowledge
 1032 base KG consisting of a collection of entities \mathcal{E} and a collection of relations \mathcal{R} . We define the set
 1033 of facts as $\mathcal{F} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$, and represent the corresponding graph as $G = (\mathcal{E}, \mathcal{F})$. Each entity
 1034 $e \in \mathcal{E}$ has an associated label $\ell(e) \in \mathcal{L}$, where \mathcal{L} denotes the space of surface-form names (e.g.,
 1035 textual strings such as “Albert Einstein”). In addition, each entity includes a relation of the form
 1036 $(e, \text{ent_type}, \tau(e))$, where $\tau(e) \in \mathcal{T}$ specifies the entity’s ontological type (e.g., person, house,
 1037 plane). $\tau(e)$ is intended to denote a general category, without mention of specific named entities.

1038 A universe of triplet facts is therefore defined by $U = (G, \ell)$.

1039 **Coherent Universe Construction.** To construct a coherent and connected universe we leverage
 1040 the facts from G . At a desired tractable size and complexity, we first sample a connected subgraph
 1041 $G' \subseteq G$ (Fig. 2a). G' is constructed by iteratively expanding the frontier from a seed set $\mathcal{Q}_0 \subseteq \mathcal{E}$.
 1042 At iteration t , given the current frontier $\mathcal{Q}_t \subseteq \mathcal{E}$, we sample neighbors $\mathcal{N}(v)$ for each $v \in \mathcal{Q}_t$ and
 1043 add them to the subgraph. Here, $\mathcal{N}(v)$ includes all entities u such that $(v, r, u) \in \mathcal{F}$ or $(u, r, v) \in \mathcal{F}$
 1044 for some $r \in \mathcal{R}$.

1045 After T expansion steps we obtain a sampled subgraph $G_T \subseteq G$. To ensure sufficient connectivity,
 1046 we extract the k -core subgraph (i.e., the maximal subgraph in which every node has degree at least
 1047 k) and then take its largest connected component, denoted $G_{T,k} \subseteq G_T$. For notational simplicity, in
 1048 the following we use G to refer to $G_{T,k}$.

1049 **Surface-Form Perturbations.** To obscure factual knowledge, we perturb surface forms, i.e., entity
 1050 names and timestamps tied to entities (Fig. 2b).⁴

1051 Simple renaming risks (a) **factual leakage**, where replacements still reveal real-world associations
 1052 (e.g., *Tokyo* → *Torioka*, which continues to suggest Japanese origins), or (b) **incoherence**, where
 1053 substitutions violate type or consistency constraints (e.g., *Ryan Reynolds* was born in *Vancouver*
 1054 → *Silvercrest Collegiate* was born in *Sarah Thompson*), thereby failing to preserve domain-general
 1055 knowledge. To prevent these issues, we systematically perturb all named entities and temporal labels
 1056 through controlled renaming that obscures underlying facts while preserving coherence.

1057 In particular, this entails: (i) **type-consistent naming**, where synthetic names respect the entity’s
 1058 ontological type (e.g., *Nile River* → *Lora River*, not *Lora Pavilion*), and (ii) **name-derivation con-
 1059 sistency**, where renames propagate to related surface forms (e.g., if *Vancouver* → *Metronis*, then
 1060 *George Vancouver*, after whom the city is named, → *Altheon Metronis*). These constraints pre-
 1061 serve semantic coherence and affiliation cues, preventing surface-level artifacts from confounding
 1062 evaluation.

1063 Let $\mathcal{E}_{\text{proper}} \subseteq \mathcal{E}$ denote the set of named-entity nodes subject to renaming and $\mathcal{L}_{\text{real}}$ denote the set of original real-mapped labels. We say that node u is *name-related* to node v if and only if
 1064 $u, v \in \mathcal{E}_{\text{proper}}$ and (i) $\ell(v)$ is a substring of $\ell(u)$ and (ii) $\exists r \in \mathcal{R} : (u, r, v) \in \mathcal{F}$ or $(v, r, u) \in \mathcal{F}$.
 1065 That is, name-relation requires both a lexical substring relationship and an explicit relation in the
 1066 knowledge graph. For instance, *Vancouver* is name-related to *Vancouver Canucks*.

1067 This induces a directed acyclic *name-related dependency graph* $G_{\text{dep}} = (\mathcal{E}_{\text{proper}}, E_{\text{dep}})$ where
 1068 $(u, r, v) \in E_{\text{dep}}$ if and only if u is name-related to v with relation r . We rename entities according
 1069 to a level-order (breadth-first) traversal of G_{dep} , processing all nodes at each level before moving to
 1070 the next level. This ensures that all entities at depth d are newly labeled before any entity at depth
 1071 $d + 1$, maintaining consistency across substring relationships.

1072 We define the updated labeling function $\ell' : \mathcal{E} \rightarrow \mathcal{L}$ through the following process. For each
 1073 $v \in \mathcal{E}_{\text{proper}}$ processed in level-order, we query a LM with input $(\tau(v), \{(\ell'(u), \tau(u), r) : (u, r, v) \in
 1074 E_{\text{dep}}\})$ to generate $\ell'(v)$. In other words, we rename entities by providing the LM with the target

1075 ⁴Other literals, e.g., population counts and physical measurements, are excluded because they could easily
 1076 (a) reveal real-world facts (e.g., “*Mount FakeMountain* is 8848m tall” still points to Mount Everest) or (b)
 1077 distort domain-general reasoning when perturbed.

1080 entity’s type and the new names of all related entities it depends on. For entities not being renamed,
 1081 we set $\ell'(e) = \ell(e)$. We include prompts for renaming in Appendix B.9.
 1082

1083 For timestamps, we apply a fixed offset δ per universe: for any timestamp x , we replace it with $x + \delta$,
 1084 preserving ordering and interval relations (e.g., a parent’s birth precedes a child’s), while removing
 1085 the potential for parametric knowledge to be leaked.

1086 After these perturbations, we produce a synth-mapped universe $U' = (G, \ell')$ where entities retain
 1087 their structure and types but receive new synthetic labels.

1088 **Parallel Corpora Generation.** For corpora generation (Fig. 2c), we first generate documents from
 1089 the synth-mapped universe U' such that the facts are faithful to G , then add symbolic references
 1090 to entity IDs in the text, before using these IDs references mapped to real-mapped labels to convert
 1091 each synthetic document into a real-mapped version. The output is two parallel corpora: one synth-
 1092 mapped and one real-mapped with identical sentence structures and world-consistent facts, differing
 1093 only in their surface-form labels.

1094 By generating documents from synth-mapped (as opposed to real-mapped) entities first, we exploit
 1095 the asymmetry that synthetic entity names $\ell'(e)$ have no connections to the LM’s parametric knowl-
 1096 edge. This prevents the LM from introducing auxiliary facts and makes it easier to stay faithful to
 1097 the provided triplets. For example, when writing about the synthetically named entity for *Austria*,
 1098 the LM cannot mention facts about *Vienna* based on external knowledge and must rely solely on the
 1099 provided facts.

1100 Concretely, for each entity $v \in \mathcal{E}$, we collect all incident edges $\{(u, r, v) \mid (u, r, v) \in \mathcal{F}\} \cup$
 1101 $\{(v, r, u) \mid (v, r, u) \in \mathcal{F}\}$ and retain only the majority orientation (i.e., whichever set is larger) to
 1102 define $N(v)$. We then query an LM to generate a document describing the facts in $N(v)$.⁵
 1103

1104 Next, following Hennigen et al. (2024), we instruct an LM to add symbolic references $\{e_1, e_2, \dots\}$
 1105 to the synth-mapped documents, adding to each mention of $\ell'(e)$ a symbolic identifier. This provides
 1106 both hyperlinks for document navigation (§4.1) and facilitates the conversion process described.

1107 Given a synthetic document with symbolic references and the entity mapping $\{(e, \ell(e), \ell'(e)) : e \in$
 1108 $\mathcal{E}\}$, we query an LM to generate an equivalent real-mapped document by replacing each symbolic
 1109 reference e_i with the original label $\ell(e_i)$. The symbolic references ensure that the correct entity
 1110 mapping is preserved during conversion. During this process, we apply programmatic and LM-based
 1111 checks to ensure document parallelism, factual consistency, and effective knowledge obfuscation.

1112

1113 B SYNTHWORLD-RM/SM DATASET CONSTRUCTION DETAILS

1114

1115 Our dataset construction pipeline follows the framework in Appendix A (overview in Fig. 2). All
 1116 prompts for dataset construction are in Appendix B.9- B.10. Table 2 summarizes the LM used and
 1117 LM API costs for each step of the pipeline including multi-hop QA task construction.

1118

1119 B.1 UNIVERSE CONSTRUCTION

1120

1121 For our specific SYNTHWORLDS corpora we start with the Wikidata KG (Vrandečić, 2012)
 1122 (01/20/2025 dump).

1123 Knowledge graphs such as Wikidata are heavily skewed toward a small set of high-frequency rela-
 1124 tions (e.g., instance of, subclass of, located in). If we sample subgraphs in strict proportion to this
 1125 distribution, the resulting universe is both narrow in structure and closely aligned with the original
 1126 world knowledge. This limits its usefulness for tasks where we want to probe reasoning in settings
 1127 that are not simply memorization of facts. To control edge-type diversity, we introduce a *uniformity*
 1128 *factor*. For $v \in \mathcal{E}$ at iteration t , let $\Gamma_t(r; v)$ denote the set of candidate triplets involving v with
 1129 relation r . We define

$$1130 P_t(r \mid v) = \frac{|\Gamma_t(r; v)|^\alpha}{\sum_k |\Gamma_t(k; v)|^\alpha}, \quad \alpha = 1 - \text{uniformity}.$$

1132
 1133 ⁵In initial experiments, including both orientations often led the LM to generate inconsistent documents,
 e.g., an entity described as both the son and the father of another.

1134 High uniformity yields diverse edge types ($\alpha = 0$: uniform), while low uniformity favors frequent
 1135 relations ($\alpha = 1$: frequency-proportional).

1136
 1137 To encourage diversity of entities, we initialize \mathcal{Q}_0 as the set of Wikidata entities across all categories
 1138 defined in Wikipedia’s popular pages (contributors, 2025) To ensure high-quality entities, we discard
 1139 Wikidata nodes that are time terms, Wikimedia-bookkeeping entities, unlabeled entries, or entities
 1140 whose names include numbers. We run the iterative sampling for $T = 11$ steps with uniformity
 1141 = 0.6, and take the 19-core subgraph $G' = G_{11,19}$.

1142 B.2 SURFACE-FORM PERTURBATIONS

1143
 1144 We rename entities identified via Wikidata’s entity naming rules.⁶

1145 Given all proper-name entities \mathcal{E}'_{proper} in G' that share a type description, we prompt a LM to
 1146 propose new names for that entity type following. In Wikidata, entity type is inferred through the
 1147 instance of relationship (P31). However, certain instances of continue to contain named entities.
 1148 For these cases we recursively apply the instance of until no named entities exist in the label. For
 1149 example, say Vancouver only has an instance of label “city in British Columbia” in this case we take
 1150 the instance of label for British Columbia which is “province of Canada”, finally we take the label
 1151 for Canada which is country so then the label becomes “city in province of country”.

1152 In addition, we incorporate Wikidata time qualifiers⁷ (e.g., Barack Obama → president → USA;
 1153 start time → 20 January 2009), which attach additional temporal information to fact triplets. To
 1154 prevent timestamps from trivially revealing real-world identities, we apply a $\delta = 39$.

1155 B.3 PARALLEL DOCUMENT GENERATION

1156 We prompt a LM to generate a factually consistent document from fact triplets (prompts in Ap-
 1157 pendix B.9). To ensure quality, we add the Wikidata entity id (prefixed with Q, e.g., Q15) when
 1158 generating symbolic references. These are unique identifiers for the underlying entity that we can
 1159 then use to check the correct label is used in the corresponding real-mapped and synth-mapped doc-
 1160 uments. We implement programmatic checks to guarantee that (1) only entities present in the facts
 1161 are included in the page, and (2) the display text for each entity matches the underlying link. When
 1162 converting from synth-mapped to real-mapped text, we additionally require that both documents
 1163 share the same set of symbolic references (thus inducing the same graph structure) and that no men-
 1164 tion of any synth-mapped entity remains. Finally, we enforce strict quality thresholds: we only keep
 1165 pages when (a) the similarity (measured using the Damerau-Levenshtein edit distance (Damerau,
 1166 1964)) between the initial generation and the symbolic-reference version exceeds 0.95, and (b) the
 1167 similarity between the synth-mapped and real-mapped versions exceeds with symbolic references
 1168 exceeds 0.85. Practically, this filtering ensures that only parallel documents with highly consistent
 1169 structure and minimal unintended variation are retained.

1170 To ensure that the generated pages are truly novel, we prompt the same LM to guess the underlying
 1171 entity from a synth-mapped document, providing it with the (unrealistic) clue that the page corre-
 1172 sponds to a real-mapped entity whose names have been perturbed. This constitutes a deliberately
 1173 strict check: in actual task settings, the LM would never be told that the page is based on a real-
 1174 world entity. Any page the LM gets correct we remove from our corpus. After each filtering step,
 1175 we retain only the largest connected component of the hyperlink graph, ensuring that the resulting
 1176 corpus remains navigable for downstream page-navigation tasks.

1177 B.4 MULTI-HOP QA CONSTRUCTION

1178
 1179 **Validating Facts for QA Construction.** Prior to the steps described in Section 4.1, we also first
 1180 validated what facts were actually in the generated corpora. This step accounts for cases where some
 1181 facts may have been omitted during generation. Given a document generated by a LM and the set of
 1182 source facts the generation based on, we use another LM to identify which of those facts are actually
 1183 present in the document. The prompt for this step is included in Appendix B.10.

1184
 1185
 1186
 1187⁶<https://www.wikidata.org/wiki/Help:Label>

⁷<https://www.wikidata.org/wiki/Help:Qualifiers>

1188 This step enables us to construct the directed fact graph $G_{\text{fact}} = (\mathcal{E}, \mathcal{F})$. Each fact is a directed triple
 1189

$$(e_i, r, e_j) \in \mathcal{F}, \quad e_i, e_j \in \mathcal{E},$$

1190 where r is a relation annotated with a property name, and the source page in our corpora from which
 1191 the fact was extracted. By construction, each edge originates from a distinct source page, ensuring
 1192 that multi-edge subgraphs aggregate knowledge across independent contexts.
 1193

1194 **Ensuring Diversity of Generated Questions.** Given the fact graph we sample graph motifs (i.e.,
 1195 the motifs in Table 3). A *motif* is a relational subgraph of G_{fact} , defined as
 1196

$$\mathcal{M} = (\mathcal{V}_M, \mathcal{F}_M), \quad \mathcal{V}_M \subseteq \mathcal{E}, \quad \mathcal{F}_M \subseteq \mathcal{F}.$$

1197 To ensure diversity and quality of questions generated, we sample graphs subject to the following
 1198 constraints:
 1199

- 1202 1. All entities in a motif must be distinct: $e_i \neq e_j \quad \forall i \neq j, e_i, e_j \in \mathcal{V}_M$.
 1203
- 1204 2. All facts in \mathcal{F}_M must come from different pages.
 1205
- 1206 3. For a given anchor configuration and relation sequence, at most one instantiation of
 1207 the motif is retained. For example, for motif A, we keep at most one subgraph
 1208 $\{(e_1, r_1, e_2), (e_2, r_2, e_3)\}$, for each tuple (e_1, r_1, r_2) . For motif E, we keep at most
 1209 one subgraph $\{(e_1, r_1, e_2), (e_3, r_2, e_4), (e_2, r_3, e_5), (e_4, r_4, e_5), (e_5, r_5, e_6)\}$, for each tu-
 1210 ple $(e_1, e_3, r_1, r_2, r_3, r_4, r_5)$. In other words, we ensure there is only one unique **reasoning**
 1211 **chain** for a given motif.
 1212
- 1213 4. Following Trivedi et al. (2022), we remove any n-hop question that is a sub-graph of any
 1214 m-hop question ($m > n > 1$).
 1215
- 1216 5. To prevent over-representation of any particular edge or intermediate node, we limit reuse
 1217 of facts and bridge entities within motifs. Concretely, each fact $(e_i, r, e_j) \in \mathcal{F}$ and each
 1218 bridge entity (i.e., entities that are neither roots nor terminal nodes of a motif) is sampled
 1219 at most five times per motif.
 1220

B.5 HUMAN VALIDATION

1221 To assess corpora quality, two researchers labeled each candidate fact as (i) *expressed in the document*,
 1222 (ii) *not expressed*, or (iii) *inconsistent with the document*. Across 28 unique pages ($n = 798$
 1223 facts), no inconsistencies were observed, giving a 95% upper bound of 0.4% on the true inconsis-
 1224 tency rate. On 7 double-annotated pages, agreement was 99.5% with Cohen’s $\kappa = 0.85$, indicating
 1225 almost perfect reliability. Corpus-level factual recall was 98.8% (95% CI [98.0, 99.7]), with mean
 1226 page recall 98.9%. These results demonstrate that the dataset is clean, reliable, and faithfully repre-
 1227 sents the intended facts.
 1228

1229 To validate question quality, another researcher inspected a sample of 30 parallel questions, covering
 1230 5 examples for each reasoning motif. For each question, the researcher verified three criteria: (i)
 1231 the questions were parallel (ii) the question led to a correct and unambiguous answer, and (iii) the
 1232 resulting question was coherent and natural. All questions were found satisfactory.
 1233

B.6 DISCUSSION ON DATASET CONSTRUCTION

1234 **Choice of Entities to Rename.** During corpus construction, we restrict renaming to Wikidata enti-
 1235 ties whose labels begin with a capital letter (e.g., *Geoffrey Hinton*, Q92894), which typically indi-
 1236 cates named entities. Entities whose labels begin with lowercase letters (e.g., *dog*, Q144; *oxygen*,
 1237 Q629) are not renamed. An edge case arises for entities such as *einsteinium* (Q1103), the element
 1238 named after Albert Einstein. Since *einsteinium* does not begin with a capital letter, it would not be
 1239 renamed, creating a potential factual knowledge leak (e.g., “*einsteinium* is named after [Renamed
 1240 Scientist]” implicitly revealing *Albert Einstein*). To mitigate this, we remove all synth-mapped pages
 1241 where such leakage could occur, ensuring that models cannot trivially recover world knowledge after
 1242 being told that entities have been renamed. Obfuscating *einsteinium*-style knowledge more broadly
 1243 and directly remains an avenue for future work.
 1244

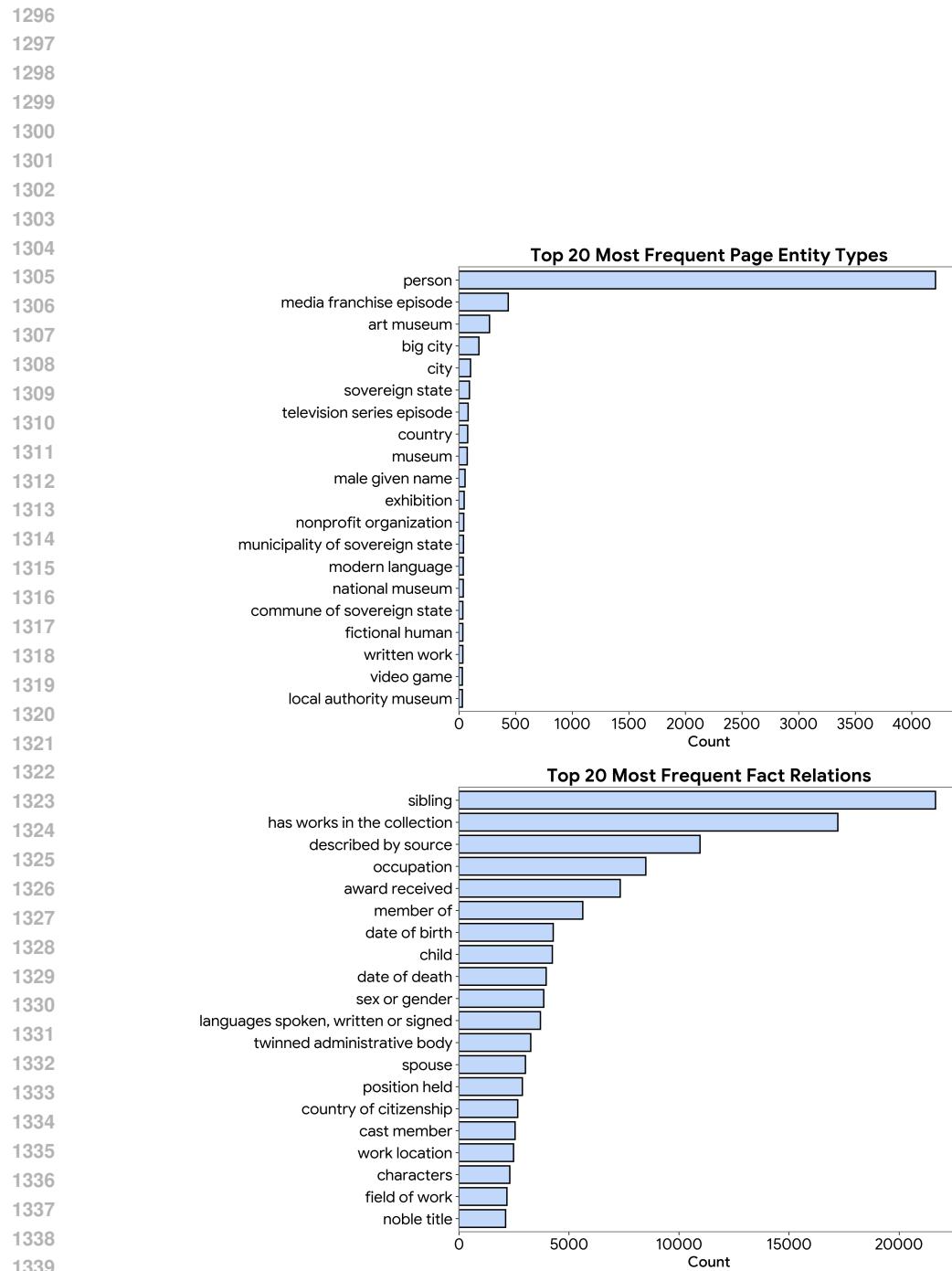
1242 **Controllability and Stochasticity in Data Generation.** To generate new instances of SYNTH-
 1243 WORLDs, we expose several controllable knobs. Different seed nodes (e.g., starting with AI re-
 1244 searchers) can be sampled to produce distinct yet structurally valid corpora. The uniformity factor
 1245 can be varied to influence graph connectivity. Subgraph sampling can also be restricted to entities of
 1246 specific types (e.g., researchers, institutions, students), or emphasize/de-emphasize particular edge
 1247 relations. Renaming strategies further contribute variability: alternative LMs, different temperature
 1248 settings, or varied timestamp perturbations can all yield distinct datasets. Finally, document gen-
 1249 eration may use different LMs to produce stylistic variation, while remaining consistent with the
 1250 underlying facts. Together, these controls balance the need for world consistency with stochastic
 1251 diversity across dataset instantiations.

1252 B.7 ADDITIONAL FIGURES AND TABLES

1253 Figure 6 shows the distributions of page entity types (based on Wikidata’s instance of prop-
 1254 erty) and relation types (across all facts) in the generated corpora. Figure 7 shows the in-degree and
 1255 out-degree distributions of the page graph in SYNTHWORLDs. Figure 8 visualizes the constructed
 1256 hyperlink graph used for Page Navigation. Table 2 provides the LM API cost of constructing SYN-
 1257 THWORLD-RM/SM. Table 3 includes all graph motifs and examples of constructed questions.

1260 Dataset Construction Step	1261 LM Used	1262 API Calls	1263 Inp Tok	1264 Out Tok	1265 Cost (\$)
1266 Surface Form Renaming	GPT-4o-mini	0.3K	237.9K	38.8K	\$0.06
1267 Corpora Generation	GPT-5-mini	35.3K	110.1M	74.3M	\$176
1268 Novelty Validation	GPT-5-mini	15.5K	6.7M	93.4M	\$188
1269 Multihop-QA Question Gen.	GPT-5-mini	4.8K	6.9M	3.0M	\$7.82
Total		55.9K	123.9M	170.7M	\$372

1270 Table 2: **Token usage and LM API costs for constructing SYNTHWORLD-RM/SM.** Totals are
 1271 shown in the last row. During the project period new LMs were released and we sought to use the
 1272 best models available to generate a public datasets. This means that GPT-4o-mini was used during
 1273 surface form renaming (a much simpler task) while all other steps used GPT-5-mini. The number
 1274 of API calls includes follow-up prompts when the initial LM output does not pass programmatic
 1275 validation checks. GPT-5-mini was used with the default reasoning effort set to `medium`. For
 1276 novelty validation, we enforced a very strict notion of novelty and explicitly instructed the model
 1277 to “think”, which inflated reasoning token usage (details in §B.3; prompt in §B.9). In practice,
 1278 one could reduce reasoning effort to `low`, since faithful evaluation on synth-mapped tasks would
 1279 not prompt LMs with the information that entities have been renamed. Such adjustments would
 1280 substantially lower costs, bringing the total closer to \$200.



1340 **Figure 6: Entity Type and Relation Type Distribution of SYNTHWORLD-RM/SM.** Documents
 1341 cover a broad range of entity types and relation types.
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

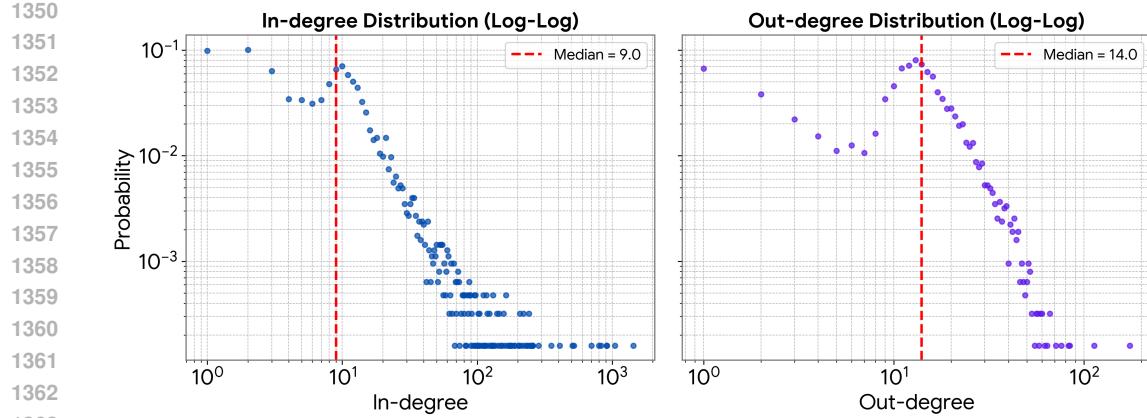


Figure 7: **Degree Distribution of SYNTHWORLD-RM/SM.** Our corpora preserve the interconnected and structured nature of knowledge networks (i.e., power-law degree distribution), matching the complexity of real-world information ecosystems.

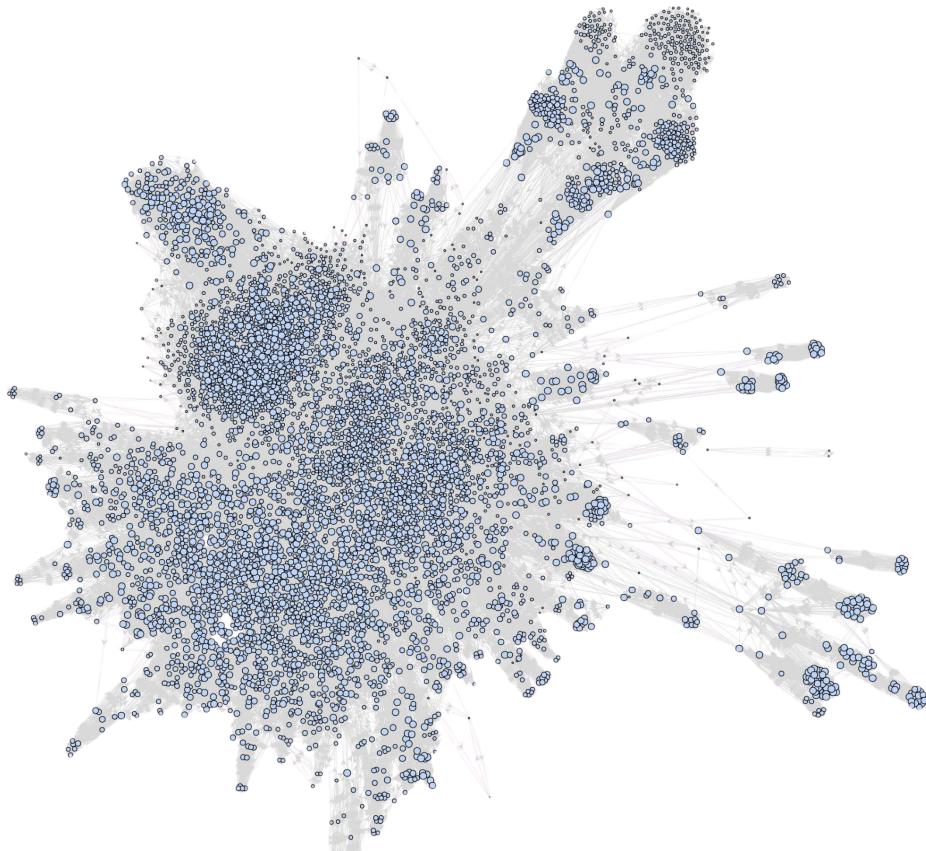


Figure 8: **SYNTHWORLD-RM/SM Hyperlink Graph** illustrating a scale-free topology, where a few highly connected hubs dominate while most nodes have relatively few links. Node size is determined by $\max(1, \min(4, \frac{\deg(v)}{8}))$.

Graph	Motif	Decomposition	Question
1404			
1405			
1406			
1407			
1408	A		1. Who was the screenwriter of The City on the Edge of Forever ? <i>Harlan Ellison</i> 2. In what year was Harlan Ellison nominated for Hugo Award for Best Short Story ? 1971
1409			In what year was the screenwriter of The City on the Edge of Forever nominated for Hugo Award for Best Short Story ? 1971
1410			
1411			
1412	B		1. Which family does Sirindhorn, Princess Royal belong to? <i>House of Mahidol</i> 2. Who is the chairperson of House of Mahidol ? Vajiralongkorn 3. Where does Vajiralongkorn live? <i>Grand Palace</i>
1413			Where does the chairperson of Sirindhorn, Princess Royal 's family live? <i>Grand Palace</i>
1414			
1415			
1416	C		1. Who is Johann Bernoulli 's doctoral student? <i>Daniel Bernoulli</i> 2. Who was Alexander R. Todd, Baron Todd's doctoral advisor? <i>Robert Robinson</i> 3. Which organization employs Daniel Bernoulli and has Robert Robinson as a member? <i>Russian Academy of Sciences</i>
1417			Which organization employs Johann Bernoulli 's doctoral student and has Alexander R. Todd, Baron Todd 's doctoral advisor as a member? <i>Russian Academy of Sciences</i>
1418			
1419			
1420	D		1. Who is the head of state of Kingdom of Bulgaria ? <i>Ferdinand I of Bulgaria</i> 2. Who is the mother of Ferdinand I of Bulgaria ? <i>Princess Clémentine, Princess of Koháry</i> 3. Who taught Princess Clémentine, Princess of Koháry ? <i>Jules Michelet</i> 4. When did Jules Michelet begin residing in Arathon ? June 1852
1421			When did the person who taught the mother of the head of state of Kingdom of Bulgaria begin residing in Arathon ? June 1852
1422			
1423			
1424	E		1. What country is Franz Xaver Winterhalter a citizen of? <i>German Empire</i> 2. Who is a relative of Princess Louise of Saxe-Gotha-Altenburg ? <i>Princess Margaret of Connaught</i> 3. Who is the head of state of the German Empire whose godparent is Princess Margaret of Connaught ? <i>William I, German Emperor</i> 4. Which conflict did William I, German Emperor participate in? <i>Napoleonic Wars</i>
1425			Which conflict did the head of state of the country Franz Xaver Winterhalter is a citizen of, whose godparent is a relative of Princess Louise of Saxe-Gotha-Altenburg , participate in? <i>Napoleonic Wars</i>
1426			
1427			
1428			
1429			
1430			
1431			
1432			
1433	F		1. Who won Matteucci Medal ? <i>Philipp Lenard</i> 2. Who was Philipp Lenard 's doctoral advisor? <i>Robert Bunsen</i> 3. Who is Henry Edward Armstrong 's employer? <i>University of London</i> 4. Who is both a student of Robert Bunsen and a director or manager at University of London ? Henry Enfield Roscoe
1434			Who is both a student of the doctoral advisor of the winner of Matteucci Medal and a director or manager at Henry Edward Armstrong 's employer? Henry Enfield Roscoe
1435			
1436			
1437			
1438			
1439			
1440			
1441			
1442			
1443			
1444			
1445			
1446			
1447			

Table 3: **Multi-hop Question Reasoning Graphs and Example Questions from SYNTHWORLDS**. Motifs in our fact triplet graph represent recurring subgraph patterns of triplet facts that form single-hop questions, which can be composed into multi-hop questions. SYNTHWORLDS follows the same multi-hop reasoning structures as the MuSiQue dataset Trivedi et al. (2022).

1452
1453
1454
1455
1456
1457

1458
1459

B.8 QUALITATIVE CORPORA EXAMPLES

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

Robert Silverberg (Q314553)

Robert Silverberg (born 15 January 1935) is an author, novelist, science fiction writer, screenwriter and writer whose work is primarily in the science fiction genre. His given name is **Robert** and he began his professional career in 1955.

Silverberg was born in **Brooklyn** and continues to reside there. He speaks **English**, which is his native language and the language in which he writes. His religion is **Judaism**. He has cited **Jack Vance** and **Roger Zelazny** as influences on his work.

Over the course of his career **Silverberg** has received several awards. He was awarded **Hugo Award for Best Novella** in 1969, **Locus Award for Best Fantasy Novel** in 1981 and **Locus Award for Best Novella** in 1988; he received **Science Fiction and Fantasy Hall of Fame** on 1 January 1999 and **Damon Knight Memorial Grand Master Award** in 2004.

He has also been nominated for numerous literary honors, including **Hugo Award for Best Novel** and **Hugo Award for Best Short Story** in 1970, **Locus Award for Best Short Story** in 1972, the **Locus Award for Best Novel** in 1973, **Hugo Award for Best Novella** in 1975, **Locus Award for Best Fantasy Novel** in 1985, **Locus Award for Best Science Fiction Novel** in 1987, **Locus Award for Best Novella** in 1999 and **Locus Award for Best Novelette** in 1990.

Silverberg is described by **Obálky knih**.

Yardley Raleth Quor

Yardley Raleth Quor (born 15 January 1974) is an author, novelist, science fiction writer, screenwriter and writer whose work is primarily in the science fiction genre. His given name is **Yardley** and he began his professional career in 1994.

Quor was born in **Myrthwood** and continues to reside there. He speaks **Velthar**, which is his native language and the language in which he writes. His religion is **Veltharion**. He has cited **Caelian Casado** and **Fythar Rees** as influences on his work.

Over the course of his career **Quor** has received several awards. He was awarded **The Storyteller's Legacy** in 2008, **The Literary Lantern** in 2020 and **The Storyteller's Connection** in 2027; he received **Exceptional Merit Recognition** on 1 January 2038 and **The Page Pen Award** in 2043.

He has also been nominated for numerous literary honors, including **The Prose Pursuit** and **The Wordsmith's Triumph** in 2009, **Echoes of Words** in 2011, the **Paper Pathway Award** in 2012, **The Storyteller's Legacy** in 2014, **The Literary Lantern** in 2024, **The Narrative Jewel** in 2026, **The Storyteller's Connection** in 2038 and **The Inked Imagination** in 2029.

Quor is described by **DataGalaxy**.

Mumbai (Q1156)

Mumbai is a large urban centre on the continent **Asia**. It functions as the state capital and is classified as a city, a metropolis and a megacity; it is also recognized as a locality and as a business cluster, reflecting a geographic concentration of interconnected businesses in a particular field.

The settlement began in 1507. Over its history **Mumbai** has been within different sovereign states: it lay in **Kingdom of England** from 11 May 1661 until 27 March 1668, and later lay in **British Raj** from 28 June 1858 until 14 August 1947.

Mumbai maintains formal twinning arrangements with several other administrative bodies. It is twinned with **London**, **Yokohama**, **Jakarta** and **Busan**, and with **Honolulu**—the partnership with **Honolulu** began on 20 January 1970.

Since 2019, **Mumbai** has been a member of the network **Creative Cities Network**.

The city appears in a range of published sources. It is described in the **Brockhaus and Efron Encyclopedic Dictionary** (a version, edition or translation), in the **Sytin Military Encyclopedia** (an encyclopedic dictionary), in **Jewish Encyclopedia of Brockhaus and Efron** (present in ethnoreligious group, nation and people encyclopedias), and in **The Nuttall Encyclopædia** (a literary work).

Crescendo

Crescendo is a large urban centre on the continent **Nystoria**. It functions as the state capital and is classified as a city, a metropolis and a megacity; it is also recognized as a locality and as a business cluster, reflecting a geographic concentration of interconnected businesses in a particular field.

The settlement began in 1546. Over its history **Crescendo** has been within different sovereign states: it lay in **Kytarathia** from 11 May 1700 until 27 March 1707, and later lay in **Lumeria** from 28 June 1897 until 14 August 1986.

Crescendo maintains formal twinning arrangements with several other administrative bodies. It is twinned with **Calidore**, **Celestport**, **Eldoria** and **Horizon Bay**, and with **Jaspis**—the partnership with **Jaspis** began on 20 January 2009.

Since 2058, **Crescendo** has been a member of the network **SyncSphere**.

The city appears in a range of published sources. It is described in the **Dreamt Compilation** (a version, edition or translation), in the **Factoid Fount** (an encyclopedic dictionary), in **Qylarans** (present in ethnoreligious group, nation and people encyclopedias), and in **The Midnight Library** (a literary work).

1512 B.9 PROMPTS FOR CORPORA CONSTRUCTION
15131514 **Generate Synthetic Names**
15151516 Give me {{ num_names }} fictional names for an entity X that is an instance of the following
1517 wikidata entity(ies):
15181519 EXAMPLE INPUT INSTANCE OF INFORMATION
1520 - business cluster (geographic concentration of interconnected businesses in a particular field)
1521 - city (large human settlement)1522 Your response should be a list of comma separated values, eg: 'foo, bar, baz' or 'foo,bar,baz' DO
1523 NOT include any other text in your response. DO NOT reference anything that already exists in the
1524 real world.1525 **Generate Synthetic Name with Substring Relation**
15261527 Given that the following facts related to the entity X are true:
15281529 EXAMPLE INPUT RELATED FACTS
1530 Cycle Ridge is a:
1531 - big city (city with a population of at least 100,000)
1532 - city (large human settlement)
1533 - cycling city (city designed for bicycle traffic)
15341535 Fact: Cycle Ridge → location → entity X
15361537 Give me a fictional name for the entity X that is an example of a:
15381539 EXAMPLE INPUT INSTANCE OF INFORMATION
1540 - public research university (type of higher learning institution; research university predominantly
funded by public means)1541 Entity X's name is likely to consist of the names of entities that it is connected to.
1542 Your response should be a single name for entity X on one line. DO NOT include any other text in
1543 your response. DO NOT reference anything that already exists in the real world.1544 **Generate Synth-Mapped and Real-Mapped Pages**
15451546 **System prompt:**
1547 You are a clear, neutral, and professional writer at the level expected for Wikipedia articles: precise,
1548 informative, and fluent, without unnecessary complexity.1549 **User:**1550 Given a page title, the page entity information (Wikidata instance of), and a set of facts between
1551 the page title and other entities, write a high-quality Wikipedia-style article.1552 You will be given the following information:
15531554 - Page title
1555 - Definitions for relation labels in page facts
1556 - Definitions for instance-of information about related entities in page facts
1557 - Page facts
1558 - Instance-of information about related entities in page facts
15591560 Your task is to produce an article that uses facts faithfully, organizes them into clear prose, and
1561 avoids contradictions.
15621563 **REQUIREMENTS:**1564 - Mention every fact faithfully; do not add or invent information.
1565 - Only use proper nouns that appear as entity names in the Page facts section.
1566 - Organize into thematic paragraphs.

1566

1567

1568

1569

1570

1571

1572

1573

1574

1575

1576

1577

1578

1579

1580

1581

1582

1583

1584

1585

1586

1587

1588

1589

1590

1591

1592

1593

1594

1595

1596

1597

1598

1599

1600

1601

1602

1603

1604

1605

1606

1607

1608

1609

1610

1611

1612

1613

1614

1615

1616

1617

1618

1619

RULES FOR INTERPRETING FACTS:

- Facts are always written as **Subject → Relation → Object**.
- Interpret the relation **relative to the Subject**.
- Examples:
 - *Albert Einstein → student → Nathan Rosen* → Nathan Rosen was a student of Albert Einstein.
 - *Albert Einstein → student of → Alfred Kleiner* → Albert Einstein was a student of Alfred Kleiner.

FACT RELATIONS AND DIRECTIONS:

- ALWAYS follow the exact meanings of the relation labels in "Definitions for relation labels in page facts."
- NEVER invert the direction for asymmetric relations (e.g., student/student of, parent/child, advisor/doctoral student, employer/employee).
- It can be easy to get this wrong—so check the label carefully and preserve its direction exactly.
- Do not generate contradictory statements.
- Normalize symmetric relations (e.g., sibling, spouse, collaborator) into one set, de-duplicate, and group entities naturally in one or more sentences for readability.

READABILITY/WRITING STYLE:

- Do not introduce speculative context, dates, regions, or concepts.
- Do not repeat the same facts in the writing.
- NEVER write "instance of" in the writing.
- If gender is not provided, always use they/them/their.
- If gender is provided, reflect it through pronouns (he/she/they) and NOT an explicit fact (to keep the writing natural and fluent).
- Vary sentence structure; avoid presenting every fact as an isolated clause or sentence.
- Group related facts into paragraphs rather than listing them line by line.
- Use connective phrasing for smoother flow (e.g., "Alongside his architectural work, he also painted...").
- When presenting multiple things, use natural connectors such as "among them," "including," or "as well as" instead of flat lists.
- Break long enumerations across sentences for readability.
- Make grammatical adjustments (articles, capitalization, punctuation) for natural flow.
- Use light connective narration ("They were part of a large family...") for readability.

OUTPUT: Return only the plain text article string (no Markdown).

Begin!

Page Title: Yardley Raleth Quor

Yardley Raleth Quor is an instance of the following entities:

EXAMPLE INPUT ON YARDLEY RALETH QUOR

- Person

Definitions for relation labels in page facts

IMPORTANT: Use the definitions below to correctly understand the page facts.

1620

1621

1622

1623

1624

1625

1626

1627

1628

1629

1630

1631

1632

1633

1634

1635

1636

1637

1638

1639

1640

1641

1642

1643

1644

1645

1646

1647

1648

1649

1650

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

EXAMPLE INPUT ON YARDLEY RALETH QUOR

- "award received": award or recognition received by a person, organization or creative work
- "date of birth": date on which the subject was born
- "described by source": work where this item is described
- "genre": creative work's genre or an artist's field of work (P101). Use main subject (P921) to relate creative works to their topic
- "given name": first name or another given name of this person; values used with the property should not link disambiguations nor family names
- "influenced by": this person, idea, etc. is informed by that other person, idea, etc., e.g. "Heidegger was influenced by Aristotle"
- "languages spoken, written or signed": language(s) that a person or a people speaks, writes or signs, including the native language(s)
- "native language": language or languages a person has learned from early childhood
- "nominated for": award nomination received by a person, organisation or creative work (inspired from "award received" (Property:P166))
- "occupation": occupation of a person; see also "field of work" (Property:P101), "position held" (Property:P39)
- "place of birth": most specific known birth location of a person, animal or fictional character
- "religion or worldview": religion of a person, organization or religious building, or associated with this subject
- "residence": the place where the person is or has been, resident
- "sex or gender": sex or gender identity of human or animal. For human: male, female, non-binary, intersex, transgender female, transgender male, agender, etc. For animal: male organism, female organism. Groups of same gender use subclass of (P279)
- "work period (start)": start of period during which a person or group flourished (fl. = "floruit") in their professional activity
- "writing language": language in which the writer has written their work

Definitions for instance-of information

EXAMPLE INPUT ON YARDLEY RALETH QUOR

- "award": something given to a person or a group of people to recognize their merit or excellence
- "ethnic religion": religion defined by the ethnicity of its adherents
- "language": particular system of communication, often named for the region or peoples that use it
- "lifestyle": interests, opinions, behaviours, and behavioural orientations of an individual, group, or culture
- "literary award": award for authors and literary associations
- "male given name": given name usually meant for boys and men
- "modern language": language in current use
- "natural language": language naturally spoken by humans, as opposed to "constructed" and "formal" languages
- "religion": social-cultural system
- "web portal": website that integrates applications, processes and services

Page facts (subject → relation property → object)

IMPORTANT: Facts are always written in the form Subject → Relation → Object. The relation definition is expressed relative to the Subject (the entity on the left). Always resolve the meaning by starting from the subject.

1674

1675

1676

1677

1678

1679

1680

1681

1682

1683

1684

1685

1686

1687

1688

1689

1690

1691

1692

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

1723

1724

1725

1726

1727

EXAMPLE INPUT ON YARDLEY RALETH QUOR

Yardley Raleth Quor → award received → Exceptional Merit Recognition
 - point in time → 2038-01-01
 Yardley Raleth Quor → award received → The Literary Lantern
 - point in time → 2020 (year)
 Yardley Raleth Quor → award received → The Page Pen Award
 - point in time → 2043 (year)
 Yardley Raleth Quor → award received → The Storyteller's Connection
 - point in time → 2027 (year)
 Yardley Raleth Quor → award received → The Storyteller's Legacy
 - point in time → 2008 (year)
 Yardley Raleth Quor → date of birth → 1974-01-15
 Yardley Raleth Quor → described by source → DataGalaxy
 Yardley Raleth Quor → genre → science fiction
 Yardley Raleth Quor → given name → Yardley
 Yardley Raleth Quor → influenced by → Caelian Casado
 Yardley Raleth Quor → influenced by → Fythar Rees
 Yardley Raleth Quor → languages spoken, written or signed → Velthar
 Yardley Raleth Quor → native language → Velthar
 Yardley Raleth Quor → nominated for → Echoes of Words
 - point in time → 2011 (year)
 Yardley Raleth Quor → nominated for → Paper Pathway Award
 - point in time → 2012 (year)
 Yardley Raleth Quor → nominated for → The Inked Imagination
 - point in time → 2029 (year)
 Yardley Raleth Quor → nominated for → The Literary Lantern
 - point in time → 2024 (year)
 Yardley Raleth Quor → nominated for → The Narrative Jewel
 - point in time → 2026 (year)
 Yardley Raleth Quor → nominated for → The Prose Pursuit
 - point in time → 2009 (year)
 Yardley Raleth Quor → nominated for → The Storyteller's Connection
 - point in time → 2038 (year)
 Yardley Raleth Quor → nominated for → The Storyteller's Legacy
 - point in time → 2014 (year)
 Yardley Raleth Quor → nominated for → The Wordsmith's Triumph
 - point in time → 2009 (year)
 Yardley Raleth Quor → occupation → author
 Yardley Raleth Quor → occupation → novelist
 Yardley Raleth Quor → occupation → science fiction writer
 Yardley Raleth Quor → occupation → screenwriter
 Yardley Raleth Quor → occupation → writer
 Yardley Raleth Quor → place of birth → Myrthwood
 Yardley Raleth Quor → religion or worldview → Veltharion
 Yardley Raleth Quor → residence → Myrthwood
 Yardley Raleth Quor → sex or gender → male
 Yardley Raleth Quor → work period (start) → 1994 (year)
 Yardley Raleth Quor → writing language → Velthar

Instance-of information about related entities in page facts (entity: wikidata entities which the entity is an instance of)

EXAMPLE INPUT ON YARDLEY RALETH QUOR

- "Caelian Casado": person
 - "DataGalaxy": multimedia database and web portal
 - "Echoes of Words": literary award
 - "Exceptional Merit Recognition": award
 - "Fythar Rees": person
 - "Myrthwood": borough of global city, borough of port settlement, borough of city, borough of big city, borough of largest city, borough of metropolis and borough of megacity, city in the sovereign state, city in the constitutional republic, city in the federal republic, city in the superpower, city in the country and city in the historical unrecognized state and consolidated city-county
 - "Paper Pathway Award": award for best book (by genre) and literary award
 - "The Inked Imagination": literary award
 - "The Literary Lantern": award for best book (by genre)
 - "The Narrative Jewel": literary award
 - "The Page Pen Award": literary award
 - "The Prose Pursuit": literary award
 - "The Storyteller's Connection": literary award
 - "The Storyteller's Legacy": literary award
 - "The Wordsmith's Triumph": literary award
 - "Velthar": language, modern language and natural language
 - "Veltharion": ethnic religion, lifestyle and religion
 - "Yardley": male given name

Page Content

1728

1729

**Assistant:**

1730

1731

1732

1733

1734

1735

1736

1737

1738

1739

1740

1741

1742

1743

1744

**User:**

1745

1746

1747

1748

1749

1750

1751

1752

1753

1754

1755

1756

1757

1758

1759

1760

1761

1762

1763

1764

1765

1766

1767

1768

1769

1770

1771

1772

1773

1774

1775

1776

1777

1778

1779

1780

1781

EXAMPLE OUTPUT GPT-5-MINI

Yardley Raleth Quor (born 15 January 1974) is an author, novelist, science fiction writer, screenwriter and writer whose work is primarily in the science fiction genre. His given name is Yardley and he began his professional career in 1994.

Quor was born in Myrthwood and continues to reside there. He speaks Velthar, which is his native language and the language in which he writes. His religion is Velharion. He has cited Caelian Casado and Fythar Rees as influences on his work.

Over the course of his career Quor has received several awards. He was awarded The Storyteller's Legacy in 2008, The Literary Lantern in 2020 and The Storyteller's Connection in 2027; he received Exceptional Merit Recognition on 1 January 2038 and The Page Pen Award in 2043.

He has also been nominated for numerous literary honors, including The Prose Pursuit and The Wordsmith's Triumph in 2009, Echoes of Words in 2011, the Paper Pathway Award in 2012, The Storyteller's Legacy in 2014, The Literary Lantern in 2024, The Narrative Jewel in 2026, The Storyteller's Connection in 2038 and The Inked Imagination in 2029.

Quor is described by DataGalaxy.

User:

Your job is to now rewrite the answer you provided above, but instead of responding directly with the text, transform the text by replacing related entity mentions with linked references using markdown and Jinja-style expressions: '[Entity Text]({{ Entity_ID }})'

You will be given the following:

1. Entity Reference JSON: A mapping of entity IDs to their name labels for each related entity.

For example, given:

```
```json
{
 "Q1397": "Ohio (U.S. state)",
 "Q30": "United States (sovereign state)"
}
```
```

You would transform:

"She was born in Ohio, USA."

To:

"She was born in [Ohio]({{ Q1397 }}), [USA]({{ Q30 }})."

IMPORTANT:

- Preserve all grammar, punctuation, and readability from the original text.
- The text and the spacing outside of the links should stay the same same.
- The text without the links should still be fluent and have proper grammar and punctuation.
- Only link proper nouns (capitalized entities like names, places, organizations)
- Only create links for entities that exist in the Entity Reference JSON. Never invent IDs or assume availability.
- As a general rule of thumb, link only the first occurrence of an entity in the text of the article.
- Links should not contain leading or trailing spaces within the square brackets, e.g., use '[North America]({{ Q49 }})', not '[North America]({{ Q49 }})'.
- DO NOT add square brackets to terms that are not in the Entity Reference JSON, e.g., "Mercedes-Benz is a [car manufacturer] founded in 1926" is incorrect.

BAD EXAMPLES:

- "They travelled to [South America]({{ Q30 }})" when Q30 refers to "United States"
- "They visited [Paris]({{ Q90 }})" when Q90 is not in the provided JSON.

Begin!

Entity Reference JSON:

1782

1783

1784

1785

1786

1787

1788

1789

1790

1791

1792

1793

1794

1795

1796

1797

1798

1799

1800

1801 **Assistant:**

1802

1803

1804

1805

1806

1807

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

1819

1820

1821

1822

1823

1824

1825

1826

1827

1828

1829

1830

1831

1832

1833

1834

1835

EXAMPLE INPUT ON YARDLEY RALETH QUOR

```
{
  "Q1860": "Velthar",
  "Q9268": "Veltharion",
  "Q18419": "Myrthwood",
  "Q255032": "The Prose Pursuit",
  "Q295406": "Fythar Rees",
  "Q312632": "Caelian Casado",
  "Q549884": "The Storyteller's Legacy",
  "Q595998": "Paper Pathway Award",
  "Q607354": "The Literary Lantern",
  "Q901462": "The Page Pen Award",
  "Q1056251": "The Wordsmith's Triumph",
  "Q2576795": "The Narrative Jewel",
  "Q4927937": "Yardley",
  "Q17008964": "The Inked Imagination",
  "Q17008972": "The Storyteller's Connection",
  "Q17008985": "Echoes of Words",
  "Q48961500": "Exceptional Merit Recognition",
  "Q67311526": "DataGalaxy"
}
```

New Output:**Assistant:**

EXAMPLE OUTPUT GPT-5-MINI

Yardley Raleth Quor (born 15 January 1974) is an author, novelist, science fiction writer, screenwriter and writer whose work is primarily in the science fiction genre. His given name is [Yardley]({{ Q4927937 }}) and he began his professional career in 1994.

Quor was born in [Myrthwood]({{ Q18419 }}) and continues to reside there. He speaks [Velthar]({{ Q1860 }}), which is his native language and the language in which he writes. His religion is [Veltharion]({{ Q9268 }}). He has cited [Caelian Casado]({{ Q312632 }}) and [Fythar Rees]({{ Q295406 }}) as influences on his work.

Over the course of his career Quor has received several awards. He was awarded [The Storyteller's Legacy]({{ Q549884 }}) in 2008, [The Literary Lantern]({{ Q607354 }}) in 2020 and [The Storyteller's Connection]({{ Q17008972 }}) in 2027; he received [Exceptional Merit Recognition]({{ Q48961500 }}) on 1 January 2038 and [The Page Pen Award]({{ Q901462 }}) in 2043.

He has also been nominated for numerous literary honors, including [The Prose Pursuit]({{ Q255032 }}) and [The Wordsmith's Triumph]({{ Q1056251 }}) in 2009, [Echoes of Words]({{ Q17008985 }}) in 2011, the [Paper Pathway Award]({{ Q595998 }}) in 2012, The Storyteller's Legacy in 2014, The Literary Lantern in 2024, [The Narrative Jewel]({{ Q2576795 }}) in 2026, The Storyteller's Connection in 2038 and [The Inked Imagination]({{ Q17008964 }}) in 2029.

Quor is described by [DataGalaxy]({{ Q67311526 }}).

**User:**

The previous response was perturbed, where names of entities and timestamp values were intentionally changed.

The actual page is about Robert Silverberg but in the above text, it was renamed to

Yardley Raleth Quor .

Your task is to now restore the above text by replacing all modified names and dates with their original values, using the provided JSON mapping.

The rest of the text must remain exactly as it is — only change the names and dates that appear in the mapping.

For example given:

```
```json
{
 "John Smith": "Elon Musk",
 "2023-07-14": "2020-07-14",
 "Acme Corp": "SpaceX"
}
```

1836

1837

1838

1839

1840

1841

1842

1843

1844

1845

1846

1847

1848

1849

1850

1851

1852

1853

1854

1855

1856

1857

1858

1859

1860

1861

1862

1863

1864

1865

1866

1867

1868

1869

1870

1871

1872

1873

1874

1875

1876

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

} ..

and text that said:

"On July 27, 2023, John Smith announced that Acme Corp would begin developing a reusable rocket. Smith was a leader in the company."

You would rewrite this to:

"On July 27, 2020, Elon Musk announced that SpaceX would begin developing a reusable rocket. Musk was a leader in the company."

**IMPORTANT.**

The final output must not contain any perturbed entity names and perturbed timestamp values. Outside of the renaming, the rest of the text should stay the exact same with the links are still included.

Do not rephrase or modify the surrounding text in any way.

Begin!

EXAMPLE INPUT ON YARDLEY RALETH QUOR

```
{
 "Q1860": "Velthar",
 "Q9268": "Veltharion",
 "Q18419": "Myrthwood",
 "Q255032": "The Prose Pursuit",
 "Q295406": "Fythar Rees",
 "Q312632": "Caelian Casado",
 "Q549884": "The Storyteller's Legacy",
 "Q595998": "Paper Pathway Award",
 "Q607354": "The Literary Lantern",
 "Q901462": "The Page Pen Award",
 "Q1056251": "The Wordsmith's Triumph",
 "Q2576795": "The Narrative Jewel",
 "Q4927937": "Yardley",
 "Q17008964": "The Inked Imagination",
 "Q17008972": "The Storyteller's Connection",
 "Q17008985": "Echoes of Words",
 "Q48961500": "Exceptional Merit Recognition",
 "Q67311526": "DataGalaxy"
}
```

**## Rewritten Output:****Assistant:**

EXAMPLE OUTPUT GPT-5-MINI

Robert Silverberg (born 15 January 1935) is an author, novelist, science fiction writer, screenwriter and writer whose work is primarily in the science fiction genre. His given name is [Robert]({{ Q4927937 }}) and he began his professional career in 1955.

Silverberg was born in [Brooklyn]({{ Q18419 }}) and continues to reside there. He speaks [English]({{ Q1860 }}), which is his native language and the language in which he writes. His religion is [Judaism]({{ Q9268 }}). He has cited [Jack Vance]({{ Q312632 }}) and [Roger Zelazny]({{ Q295406 }}) as influences on his work.

Over the course of his career Silverberg has received several awards. He was awarded [Hugo Award for Best Novella]({{ Q549884 }}) in 1969, [Locus Award for Best Fantasy Novel]({{ Q607354 }}) in 1981 and [Locus Award for Best Novella]({{ Q17008972 }}) in 1988; he received [Science Fiction and Fantasy Hall of Fame]({{ Q48961500 }}) on 1 January 1999 and [Damon Knight Memorial Grand Master Award]({{ Q901462 }}) in 2004.

He has also been nominated for numerous literary honors, including [Hugo Award for Best Novel]({{ Q255032 }}) and [Hugo Award for Best Short Story]({{ Q1056251 }}) in 1970, [Locus Award for Best Short Story]({{ Q17008985 }}) in 1972, the [Locus Award for Best Novel]({{ Q595998 }}) in 1973, Hugo Award for Best Novella in 1975, Locus Award for Best Fantasy Novel in 1985, [Locus Award for Best Science Fiction Novel]({{ Q2576795 }}) in 1987, Locus Award for Best Novella in 1999 and [Locus Award for Best Novelette]({{ Q17008964 }}) in 1990.

Silverberg is described by [Obálky knih]({{ Q67311526 }}).

1890

1891

1892

1893

1894

1895

1896

1897

1898

1899

1900

1901

1902

1903

1904

1905

1906

1907

1908

1909

1910

1911

1912

## B.10 PROMPTS FOR MULTI-HOP QA CONSTRUCTION FROM FACTS

1914

1915

1916

1917

1918

1919

1920

1921

1922

1923

1924

1925

1926

1927

1928

1929

1930

1931

1932

1933

1934

1935

1936

1937

1938

1939

1940

1941

1942

1943

### Validate Page Novelty

**System prompt:**

You are a master detective that identifies the main subject of a given page. Entity names in pages have been replaced with fictional ones to create a parallel fictional world. Your job is to identify the original real-world entity by analyzing patterns in dates, events, relationships, achievements, and biographical details.

**User:**

The following page is about a real-world wikidata entity, but some or all names have been changed to fictional ones. Identify the original main subject.

Hint: The entity currently named in the page is fictional - you need to determine which real person/organization/thing this fictional entity represents.

Your response starts with "Thought: ", where you will methodically break down the reasoning process, illustrating how you arrive at conclusions. Then, based on your reasoning, conclude with the prefix "Answer: " to present a concise, definitive response, devoid of additional elaborations (e.g., "Answer: Barack Obama").

**Page Content:**

INPUT SYNTH-MAPPED PAGE

Thought:

### Get Page Facts

**System prompt:**

You are an advanced reading comprehension assistant. Your task is to analyze a text passage and extract specific information to fill in triplet templates with placeholders marked as <ANSWER>.

**User:**

Given a page and a JSON mapping of partial facts (indicated by <ANS> placeholders), use the page content to extract the missing information.

Return your output as a list of answers for all triplets. For nested triplets (containing multiple <ANS> placeholders), ensure all parts are supported by the text.

**Guidelines:**

- Return empty array

if no relevant information is found

- For nested triplets, only include complete matches where all placeholder values are found

- If multiple valid answers exist, include all of them

- Use the relation descriptions to help you understand the meaning of the triplet relations. Triplets are always in the form of subject → relation → object.

**Example:****Page content:**

Pavel Cherenkov held roles as a nuclear physicist and a general physicist, and over the course of his career he received several honors, including the Nobel Prize in Physics, the Order of Lenin, the Order of the Red Banner of Labour and the Hero of Socialist Labour. He was nominated for the Nobel Prize in Physics in 1955 and later received that prize in 1958.

**Relation descriptions:**

- "award received": award or recognition received by a person, organization or creative work

- "nominated for": award nomination received by a person, organisation or creative work

(inspired from "award received" (Property:P166))

- "occupation": occupation of a person; see also "field of work" (Property:P101), "position held" (Property:P39)

1944

1945

1946

1947

1948

1949

1950

1951

1952

1953

1954

1955

1956

1957

1958

1959

1960

1961

1962

1963

1964

1965

1966

1967

1968

1969

1970

1971

1972

1973

1974

1975

1976

1977

1978

1979

1980

1981

1982

1983

1984

1985

1986

1987

1988

1989

1990

1991

1992

1993

1994

1995

1996

1997

1998

Partial fact templates:

```
```json
{
  "T1": "Pavel Cherenkov -> occupation -> <ANS>
  "T2": "Pavel Cherenkov -> award received -> <ANS>
  "T3": "Pavel Cherenkov -> award received -> <ANS1> AND <ANS1> -> point in time -> <ANS2>
  "T4": "Pavel Cherenkov -> nominated for -> <ANS1> AND <ANS1> -> point in time -> <ANS2>
}
```

```

Output:

```
```json
{
  "T1": ["nuclear physicist", "physicist"]
  "T2": ["Nobel Prize in Physics", "Order of Lenin", "Order of the Red Banner of Labour", "Hero of Socialist Labour"]
  "T3": [[{"Nobel Prize in Physics", "1955"}]]
  "T4": [[{"Nobel Prize in Physics", "1958"}]]
}
```

```

Output format:

Return only the JSON object with extracted answers. Use empty arrays for triplets where no information is found in the text.

Begin!

Page content:

## EXAMPLE INPUT ON YARDLEY RALETH QUOR

Yardley Raleth Quor (born 15 January 1974) is an author, novelist, science fiction writer, screenwriter and writer whose work is primarily in the science fiction genre. His given name is Yardley and he began his professional career in 1994.

Quor was born in Myrthwood and continues to reside there. He speaks Velthar, which is his native language and the language in which he writes. His religion is Veltharion. He has cited Caelian Casado and Fythar Rees as influences on his work.

Over the course of his career Quor has received several awards. He was awarded The Storyteller's Legacy in 2008, The Literary Lantern in 2020 and The Storyteller's Connection in 2027; he received Exceptional Merit Recognition on 1 January 2038 and The Page Pen Award in 2043.

He has also been nominated for numerous literary honors, including The Prose Pursuit and The Wordsmith's Triumph in 2009, Echoes of Words in 2011, the Paper Pathway Award in 2012, The Storyteller's Legacy in 2014, The Literary Lantern in 2024, The Narrative Jewel in 2026, The Storyteller's Connection in 2038 and The Inked Imagination in 2029.

Quor is described by DataGalaxy.

Relation descriptions:

## EXAMPLE INPUT ON YARDLEY RALETH QUOR

- award received: award or recognition received by a person, organization or creative work
- date of birth: date on which the subject was born
- described by source: work where this item is described
- genre: creative work's genre or an artist's field of work (P101). Use main subject (P921) to relate creative works to their topic
- given name: first name or another given name of this person; values used with the property should not link disambiguations nor family names
- influenced by: this person, idea, etc. is informed by that other person, idea, etc., e.g. "Heidegger was influenced by Aristotle"
- languages spoken, written or signed: language(s) that a person or a people speaks, writes or signs, including the native language(s)
- native language: language or languages a person has learned from early childhood
- nominated for: award nomination received by a person, organisation or creative work (inspired from "award received" (Property:P166))
- occupation: occupation of a person; see also "field of work" (Property:P101), "position held" (Property:P39)
- place of birth: most specific known birth location of a person, animal or fictional character
- religion or worldview: religion of a person, organization or religious building, or associated with this subject
- residence: the place where the person is or has been, resident
- sex or gender: sex or gender identity of human or animal. For human: male, female, non-binary, intersex, transgender female, transgender male, agender, etc. For animal: male organism, female organism. Groups of same gender use subclass of (P279)
- work period (start): start of period during which a person or group flourished (fl. = "floruit") in their professional activity
- writing language: language in which the writer has written their work

1998

1999

2000

2001

2002

2003

2004

2005

2006

2007

2008

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

Partial fact templates:

EXAMPLE INPUT ON YARDLEY RALETH QUOR

```
{
 "T0": "Yardley Raleth Quor -> award received -> <ANS1> AND <ANS1> -> point in
 time -> <ANS2>",
 "T1": "Yardley Raleth Quor -> date of birth -> <ANS>",
 "T2": "Yardley Raleth Quor -> described by source -> <ANS>",
 "T3": "Yardley Raleth Quor -> genre -> <ANS>",
 "T4": "Yardley Raleth Quor -> given name -> <ANS>",
 "T5": "Yardley Raleth Quor -> influenced by -> <ANS>",
 "T6": "Yardley Raleth Quor -> languages spoken, written or signed -> <ANS>",
 "T7": "Yardley Raleth Quor -> native language -> <ANS>",
 "T8": "Yardley Raleth Quor -> nominated for -> <ANS1> AND <ANS1> -> point in
 time -> <ANS2>",
 "T9": "Yardley Raleth Quor -> occupation -> <ANS>",
 "T10": "Yardley Raleth Quor -> place of birth -> <ANS>",
 "T11": "Yardley Raleth Quor -> religion or worldview -> <ANS>",
 "T12": "Yardley Raleth Quor -> residence -> <ANS>",
 "T13": "Yardley Raleth Quor -> sex or gender -> <ANS>",
 "T14": "Yardley Raleth Quor -> work period (start) -> <ANS>",
 "T15": "Yardley Raleth Quor -> writing language -> <ANS>"
}
```

Output:

**Assistant:**

EXAMPLE OUTPUT ON GPT-5-MINI

```
{
 "T0": [
 ["The Storyteller's Legacy", "2008"],
 ["The Literary Lantern", "2020"],
 ["The Storyteller's Connection", "2027"],
 ["Exceptional Merit Recognition", "1 January 2038"],
 ["The Page Pen Award", "2043"]
],
 "T1": ["15 January 1974"],
 "T2": ["DataGalaxy"],
 "T3": ["science fiction"],
 "T4": ["Yardley"],
 "T5": ["Caelian Casado", "Fythar Rees"],
 "T6": ["Velthar"],
 "T7": ["Velthar"],
 "T8": [
 ["The Prose Pursuit", "2009"],
 ["The Wordsmith's Triumph", "2009"],
 ["Echoes of Words", "2011"],
 ["Paper Pathway Award", "2012"],
 ["The Storyteller's Legacy", "2014"],
 ["The Literary Lantern", "2024"],
 ["The Narrative Jewel", "2026"],
 ["The Storyteller's Connection", "2038"],
 ["The Inked Imagination", "2029"]
],
 "T9": ["author", "novelist", "science fiction writer", "screenwriter", "writer"],
 "T10": ["Myrthwood"],
 "T11": ["Veltharion"],
 "T12": ["Myrthwood"],
 "T13": [],
 "T14": ["1994"],
 "T15": ["Velthar"]
}
```

Generate Single Hop Questions from Facts

System prompt:

2052

2053 You are an expert annotator for a question answering dataset. Your task is to generate high-quality,  
 2054 fact-based questions from the provided facts and relation descriptions.

2055

2056  **User:**

2057 You will be given a specific graph structure defining how questions are generated from facts (i.e.,  
 2058 entities/values (E) and relations (R)).  
 2059 Each E can represent either an entity (person, place, thing) or a value (date, number, text).

2060

2061 **### Structure:**

2062

2063 STRUCTURE FOR GRAPH A (WOULD CHANGE FOR OTHER GRAPHS)

2064

2065 Q1: E1 → R1 → E2 <ANS1>

2066

2067 Q2: E2 → R2 → E3 <ANS2>

2068

2069 Where E1, E2, E3 are entities/values and R1, R2 are relations.

2070

2071 <ANS1> and <ANS2> are different answers to Q1 and Q2 respectively.

2072

2073 **### Relation Types:**

2074

2075 - Simple Relations: A direct relationship between two entities/values (entity → relation → entity)

2076

2077 - Qualified Relations: When a relation needs additional context (time, location, role, etc.):

2078

2079 - 'Entity → [BaseRelation → Qualifier → Attribute] → Value'

2080

2081 - Interpretation: The Attribute of Qualifier's BaseRelation to Entity is Value

2082

2083 - Example:

2084

2085 - 'Paris → [mayor → Anne Hidalgo → start time] → 2014-04-05'

2086

2087 - Means: "The start time of Anne Hidalgo's role as mayor of Paris is 2014-04-05"

2088

2089 - Question: "When did Anne Hidalgo become mayor of Paris?"

2090

2091 - Facts are written in the form Subject → Relation → Object. The relation definition is expressed  
 2092 relative to the Subject (the entity on the left). Always resolve the meaning by starting from the  
 2093 subject.

2094

2095 **### Requirements:**

2096

2097 - Each question must be natural, fluent English and have a single, unambiguous correct answer.

2098

2099 - The answer to each question is exactly the entity/value tagged with <ANS>.

2100

2101 - The subject entity (the entity before → R...) must appear explicitly in the question text to ensure  
 2102 clarity.

2103

2104 - Phrase time-based relations naturally ("When did...?", "On what date...?", "In what year...?")  
 2105 matching the granularity of the <ANS> (date/year/etc.).

2106

2107 - Do not copy awkward relation phrasing verbatim if a more natural form exists ("Where was X  
 2108 born?", not "What is the place of birth of X?").

2109

2110 - Do not include <ANS> verbatim in the question text — the question must point to <ANS>  
 2111 naturally without revealing it.

2112

2113 **#### Relations**

2114

2115 - Use the relation wording from the Question Facts as the basis for your question, rephrasing only  
 2116 if needed for natural English.

2117

2118 - If the entity type is unclear (e.g., the relation description lists multiple possible types such as  
 2119 country or region), avoid inventing context (e.g., ask "What shares a border with X?" instead of  
 2120 "What country borders X?").

2121

2122 - Avoid using the word "entity" in the question text — questions should always sound natural.

2123

2124 **### Output Format:**

2125

2126 Respond only with questions in this JSON format:

2127

2128 

```
```json
```

2129

```
{
  "Q1": "Question 1",
  "Q2": "Question 2",
  "QN": "Question N"
}
```

2130

```
```
```

2131

2132 Do not include explanations, comments, or text outside the JSON object.

2133

2134 **### Example:**

2135

2136

2137

2106

2107

2108

2109

2110

2111

2112

2113

2114

2115

2116

2117

2118

2119

2120

2121

2122

2123

2124

2125

2126

2127

2128

2129

2130

2131

2132

2133

2134

2135

2136

2137

2138

2139

2140

2141

2142

2143

2144

2145

2146

2147

2148

2149

2150

2151

2152

2153

2154

2155

2156

2157

2158

2159

DEMONSTRATION FOR GRAPH A (WOULD CHANGE FOR OTHER GRAPHS)

Question Facts:

Q1: Stephen Hawking → place of birth → United Kingdom &lt;ANS1&gt;

Q2: United Kingdom → capital → London &lt;ANS2&gt;

Relation Descriptions:

- place of birth: most specific known birth location of a person, animal or fictional character
- capital: seat of government of a country, province, state or other type of administrative territorial entity

Output:

```
```json
{
  "Q1": "Where was Stephen Hawking born?",
  "Q2": "What is the capital of the United Kingdom?"
}
```
```

Begin!

Question Facts:

EXAMPLE INPUT FOR GRAPH A

Question Facts:

Q1: Jorith Luque → educated at → The Artistic Exchange <ANS1>  
 Q2: The Artistic Exchange → [founded by → Merith Watts → point in time] → 1864 (year) <ANS2>

Relation Descriptions:

EXAMPLE INPUT FOR GRAPH A

- "educated at": educational institution attended by subject
- "founded by": founder or co-founder of this organization, religion, place or entity
- "point in time": date something took place, existed or a statement was true; for providing time use the "refine date" property (P4241)

Output:



Assistant:

EXAMPLE OUTPUT GPT-5-MINI

```
```json
{
  "Q1": "Where was Jorith Luque educated?",
  "Q2": "In what year did Merith Watts found The Artistic Exchange?"
}
```
```

## Generate Multi-hop Questions from Single-hop Questions

### System prompt:

You are an expert annotator for a question answering dataset. Your task is to compose a coherent question from a list of decomposed questions.

Each decomposed question represents one atomic fact or relationship.

For example, given the following decomposed questions:

- Q1: Which university was Facebook launched in? → Harvard
- Q2: What city is <bridge> Harvard</bridge> located in? → Cambridge

They can be composed together into:

"Which city was Facebook launched in?"

Bridge entities are marked with <bridge> tags, each of which should be the answer of a decomposed question.

Characteristics of Good Questions:

2160  
 2161 - Fact-seeking: Questions that can be answered with a specific entity or concise explanation  
 2162 - Unambiguous: Has a single, clear correct answer  
 2163 - Requires comprehension: Demonstrates understanding beyond surface-level pattern matching  
 2164 - Natural language: Uses conversational phrasing that sounds like something a person would ask  
 Characteristics of Bad Questions:  
 2165 - Poorly formulated: Unclear or grammatically incorrect questions  
 2166 - False presuppositions: Questions based on incorrect assumptions  
 2167 - Opinion-based: Questions seeking subjective judgments rather than factual information  
 2168 - Not fact-seeking: Questions that don't clearly request factual information

2169  **User:**

2170 Your will be given a list of decomposed questions with marked bridge entities. Your task is to  
 2171 compose a coherent question from them.

2173 IMPORTANT The composed question SHOULD NOT include any bridge entities (i.e., those  
 2174 wrapped in <bridge> tags). If composed correctly, the bridge entities should not occur in the  
 2175 composed question.

2176 Requirements for composing questions:

- 2177 1. Use all decomposed questions: Incorporate information from all decomposed questions.
- 2178 2. Preserve meaning and answer: Retain the meaning and ensure the composed question's answer  
 2179 is the same as the last decomposed question's answer. Do not change the answer. Rephrasing is  
 2180 encouraged for fluency and clarity. (Incorrect example: "Which country was Facebook launched  
 2181 in?" — this changes both meaning and answer.)
- 2182 3. Keep it concise: Compress as much as possible without losing meaning. Prefer: "Which city  
 2183 was Facebook launched in?" over "Which city has a university, which Facebook was launched  
 2184 in?"
- 2185 4. Two-sentence fallback: If the composed question becomes too long to be coherent, you may  
 2186 split it into two sentences (1 assertion + 1 question), connected by coreference. Use only as a last  
 2187 resort.
- 2188 5. Answer alignment: The composed question must always have the same answer as the last  
 2189 decomposed question's answer.
- 2190 6. Do not remove necessary details that would make the question ambiguous. This means that  
 2191 non-bridge entities should be included
- 2192 7. Phrase time-based relations naturally ("When did...?", "On what date...?", "In what  
 2193 year...?") matching the granularity of the answer (date/year/etc.).

2194 **FAQ:**

- 2195 1. Should I paraphrase the question for clarity?

2196 You're encouraged to paraphrase the question to make it simple and coherent as long as the  
 2197 rephrased question leads to the associated answer. You do not need to use the exact same phrasing  
 2198 as the decomposed question, because sometimes they are awkward. E.g. replacing "terrain  
 2199 feature" → "mountain range", "administrative territorial entity" → "state/city/etc", "parental  
 2200 progenitor" → mother/father depending on the question context are all great

- 2201 2. Given the hard choice, do you prefer a shorter or more coherent question?

2202 Being able to parse and understand the question is more important to us than its length. So if you  
 2203 can't retain coherency of the question while keeping it short, write a longer but coherent  
 2204 composed question.

- 2205 3. What if the question is composable even if entities aren't exactly the same?

2206 There are some rare cases in which the marked bridge entity doesn't mean exactly the same, but  
 2207 yet are talking about the same entity. E.g.

- 2208 - Q1: Who was in charge of the US? → George Washington
- 2209 - Q2: Who was the creator of George Washington? → Donald De Lue

2210 Here, Q1's "George Washington" is a person, while Q2's refers to a monument of him. You can  
 2211 compose: "Who is the creator of the monument of the person in charge of the US?"

2212 But avoid nonsensical versions like: "Who is the creator of the person in charge of the US?"

2213 If the composition would be too awkward or confusing, respond with "No composition".

2214 **Output format:**

2215 Start your answer with "Thought: ", where you reason through your decision step-by-step.

2214

2215

2216

2217

2218

2219

2220

2221

2222

2223

2224

2225

2226

2227

2228

2229

2230

2231

2232

2233

2234

2235

2236

2237

2238

2239

2240

2241

2242

2243

2244

2245

2246

2247

2248

2249

2250

2251

2252

2253

2254

2255

2256

2257

2258

2259

2260

2261

2262

2263

2264

2265

2266

2267

Conclude with "Question: ", followed by the composed question (or "No composition") without any modification (i.e., no formatting, no bolding, and no markup) or further explanation.

Examples:

DEMONSTRATIONS BASED ON MUSIQUE

Decomposed questions:

- Q1: Who was the first President of Namibia? → Sam Nujoma
- Q2: Who succeeded <bridge>Sam Nujoma</bridge>? → Hifikepunye Pohamba

Thought:

- Q1 tells us that Sam Nujoma was the first President of Namibia.
- Q2 asks who succeeded <bridge>Sam Nujoma</bridge>, referring to the person identified in Q1.
- Since Sam Nujoma = first President of Namibia, we can substitute that description into Q2.

Question: Who succeeded the first President of Namibia?

Decomposed questions:

- Q1: At what location did Billy Giles die? → Belfast
- Q2: What part of the UK is <bridge>Belfast</bridge>located in? → Northern Ireland
- Q3: What is the unit of currency in <bridge>Northern Ireland</bridge>? → Pound sterling

Thought:

- Q1 says Billy Giles died in Belfast.
- Q2 tells us Belfast is in Northern Ireland.
- Q3 says Northern Ireland uses Pound sterling.
- So we can describe the place where Billy Giles died as "Northern Ireland."

Question: What currency is used where Billy Giles died?

Decomposed questions:

- Q1: What is McDonaldization named after? → McDonald's
- Q2: Which state is Horndean located in? → England
- Q3: When did the first <bridge>McDonald's</bridge>open in <bridge>England</bridge>? → 1974

Thought:

- Q1 says McDonaldization is named after McDonald's.
- Q2 says Horndean is located in England.
- Q3 asks when McDonald's first opened in England.
- Since Horndean is in England, we can use that substitution to generalize the location.

Question: When was the first establishment that McDonaldization is named after opened in the country Horndean is located?

Decomposed questions:

- Q1: Who brought Louis XVI style to the court? → Marie Antoinette
- Q2: Who is the mother of <bridge>Marie Antoinette</bridge>? → Maria Theresa
- Q3: In what city did <bridge>Maria Theresa</bridge>die? → Vienna
- Q4: When did Napoleon occupy <bridge>Vienna</bridge>? → 1805

Thought:

- Q1 tells us Marie Antoinette brought Louis XVI style to the court.
- Q2 tells us her mother is Maria Theresa.
- Q3 says Maria Theresa died in Vienna.
- Q4 tells us Napoleon occupied Vienna in 1805.
- So we can ask when Napoleon occupied the city where the mother of the woman who brought Louis XVI style to the court died.

Question: When did Napoleon occupy the city where the mother of the woman who brought Louis XVI style to the court died?

Decomposed questions:

- Q1: What continent is Aruba in? → South America
- Q2: What country is Prazeres from? → Portugal
- Q3: What was the colonial holding in <bridge>South America</bridge>governed by <bridge>Portugal</bridge>? → Brazil
- Q4: How many Germans live in <bridge>Brazil</bridge>? → 5 million

Thought:

- Q1 tells us Aruba is in South America.
- Q2 tells us Prazeres is from Portugal.
- Q3 tells us Brazil was the colonial holding in South America governed by Portugal.
- Q4 says 5 million Germans live in Brazil.
- So we can ask how many Germans live in the colonial holding in Aruba's continent that was governed by Prazeres's country.

Question: How many Germans live in the colonial holding in Aruba's continent that was governed by Prazeres's country?

Begin!

Decomposed questions:

EXAMPLE INPUT SINGLE HOP QUESTIONS FOR GRAPH A

Q1: Where was Jorith Luque educated? → The Artistic Exchange

Q2: In what year did Merith Watts found <bridge> The Artistic Exchange</bridge> ? → 1864  
(year)

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279

 **Assistant:**

EXAMPLE OUTPUT GPT-5-MINI

Thought: Q1 says Jorith Luque was educated at The Artistic Exchange. Q2 asks when Merith Watts founded that same institution and gives the year 1864. Combine by asking when Merith Watts founded the institution where Jorith Luque was educated, avoiding the bridge name.  
Question: In what year did Merith Watts found the institution where Jorith Luque was educated?

2280 **C EXPERIMENT DETAILS**

2281 For all experiments with Gemini-2.0-Flash with temperature=0. For all experiments with GPT-5-  
2282 mini, the reasoning effort is set to default medium.

2283 For the HippoRAG 2 baselines (Gutiérrez et al., 2025), we used the same model for NER and  
2284 OpenIE (i.e., GPT-5-mini for GPT-5-mini experiments and Gemini-2.0-Flash and Gemini-2.0-Flash  
2285 experiments). We follow (Gutiérrez et al., 2025) and use nvidia/NVEmbed-v2 (Lee et al., 2025) as  
2286 the retriever. For IRCoT, we run for a maximum of 10 steps.

2287 F1 error bars (Fig. 4) are calculated with 95% bootstrap confidence intervals. For each diffi-  
2288 culty bucket, we resample with replacement 5,000 times, compute the mean F1 (or F1-gap) for  
2289 each resample, and take the 2.5th and 97.5th percentiles as confidence bounds. Success-rate  
2290 error bars (Fig. 5) are calculated using 95% Wilson score confidence intervals, computed via  
2291 statsmodels.stats.proportion.proportion\_confint with method='wilson'  
2292 in Python.

2293 **C.1 MULTIHOP QA PROMPTS**

2294 For the QA baselines, we follow prior work whenever possible (Trivedi et al., 2022; Gutierrez et al.,  
2295 2024), including the use of prompt demonstrations.

2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307

**Multihop Question Answering — No Retrieval**

**System prompt:**

As an advanced question answering assistant, your task is to answer the question. Your response starts  
2304 after "Thought: ", where you will methodically break down the reasoning process, illustrating how  
2305 you arrive at conclusions. Conclude with "Answer: " to present a concise, definitive response, devoid  
2306 of additional elaborations. Your answer should be a single entity or timestamp.

2308  
2309  
2310  
2311  
2312  
2313  
2314

 **User:**

Question: {{ query }}

Thought:

2315  
2316  
2317  
2318  
2319  
2320  
2321

**Multihop Question Answering — Reading Comprehension (includes one demonstration)**

**System prompt:**

As an advanced question answering assistant, your task is to analyze text passages and corresponding  
2315 questions meticulously. Your response starts after "Thought: ", where you will methodically break  
2316 down the reasoning process, illustrating how you arrive at conclusions. Conclude with "Answer: " to  
2317 present a concise, definitive response, devoid of additional elaborations. Your answer should be a  
2318 single entity or timestamp.

 **User:**

2322

2323

2324

2325

2326

2327

2328

2329

2330

2331

2332

2333

2334

2335

2336

2337

2338

2339

2340

2341

2342

2343

2344

2345

2346

2347

2348

2349

2350

2351

2352

2353

2354

2355

2356

2357

2358

2359

2360

2361

2362

2363

2364

2365

2366

2367

2368

2369

2370

2371

2372

2373

2374

2375

## EXAMPLE DEMONSTRATION

The Last Horse (Spanish: El Último caballo) is a 1950 Spanish comedy film directed by Edgar Neville starring Fernando Fernán Gómez.

The University of Southampton, which was founded in 1862 and received its Royal Charter as a university in 1952, has over 22,000 students. The university is ranked in the top 100 research universities in the world in the Academic Ranking of World Universities 2010. In 2010, the THES - QS World University Rankings positioned the University of Southampton in the top 80 universities in the world. The university considers itself one of the top 5 research universities in the UK. The university has a global reputation for research into engineering sciences, oceanography, chemistry, cancer sciences, sound and vibration research, computer science and electronics, optoelectronics and textile conservation at the Textile Conservation Centre (which is due to close in October 2009.) It is also home to the National Oceanography Centre, Southampton (NOCS), the focus of Natural Environment Research Council-funded marine research.

Stanton Township is a township in Champaign County, Illinois, USA. As of the 2010 census, its population was 505 and it contained 202 housing units.

Neville A. Stanton is a British Professor of Human Factors and Ergonomics at the University of Southampton. Prof Stanton is a Chartered Engineer (C.Eng), Chartered Psychologist (C.Psychol) and Chartered Ergonomist (C.ErgHF). He has written and edited over forty books and over three hundred peer-reviewed journal papers on applications of the subject. Stanton is a Fellow of the British Psychological Society, a Fellow of The Institute of Ergonomics and Human Factors and a member of the Institution of Engineering and Technology. He has been published in academic journals including "Nature". He has also helped organisations design new human-machine interfaces, such as the Adaptive Cruise Control system for Jaguar Cars.

Finding Nemo

Theatrical release poster

Directed by Andrew Stanton

Produced by Graham Walters

Screenplay by Andrew Stanton Bob Peterson David Reynolds

Story by Andrew Stanton

Starring Albert Brooks Ellen DeGeneres Alexander Gould Willem Dafoe

Music by Thomas Newman

Cinematography Sharon Calahan Jeremy Lasky

Edited by David Ian Salter

Production company Walt Disney Pictures Pixar Animation Studios

Distributed by Buena Vista Pictures Distribution

Release date May 30, 2003 (2003-05-30)

Running time 100 minutes

Country United States

Language English

Budget \$94 million

Box office \$940.3 million

Question: When was Neville A. Stanton's employer founded?

Thought:

 **Assistant:**

The employer of Neville A. Stanton is University of Southampton. The University of Southampton was founded in 1862.

Answer: 1862

 **User:**

```
 {{ gold_and_distractor_passages }}
```

Question: {{ query }}

Thought:

## Multihop Question Answering — One-Step RAG (includes one demonstration)

**System prompt:**

As an advanced reading comprehension assistant, your task is to answer the question given the passages. Your response starts after "Thought: ", where you will methodically break down the reasoning process, illustrating how you arrive at conclusions. Conclude with "Answer: " to present a concise, definitive response, devoid of additional elaborations, explanations or extra information. Your answer should be a single entity or timestamp.

 **User:**

2376

2377

2378

2379

2380

2381

2382

2383

2384

2385

2386

2387

2388

2389

2390

2391

2392

2393

2394

2395

2396

2397

2398

2399

2400

2401

2402

2403

2404

2405

2406

2407

2408

2409

2410

2411

2412

2413

2414

2415

2416

2417

2418

2419

2420

2421

2422

2423

2424

2425

2426

2427

2428

2429

SAME AS READING COMPREHENSION

**Assistant:**

The employer of Neville A. Stanton is University of Southampton. The University of Southampton was founded in 1862.

Answer: 1862

**User:**

```
 {{ retrieved_passages }}
```

Question: {{ query }}

Thought:

## Multihop Question Answering — IRCoT + RAG (includes demonstration)

### System prompt:

You serve as an intelligent assistant, adept at facilitating users through complex, multi-hop reasoning across multiple documents. This task is illustrated through demonstrations, each consisting of a document set paired with a relevant question and its multi-hop reasoning thoughts. Your task is to generate one thought for the current step, DON'T generate the whole thoughts at once! If you reach what you believe to be the final step, start with "So the answer is:" Your answer should be a single entity or timestamp.

#### EXAMPLE DEMONSTRATION

The Last Horse (Spanish: El Último caballo) is a 1950 Spanish comedy film directed by Edgar Neville starring Fernando Fernán Gómez.

The University of Southampton, which was founded in 1862 and received its Royal Charter as a university in 1952, has over 22,000 students. The university is ranked in the top 100 research universities in the world in the Academic Ranking of World Universities 2010. In 2010, the THES - QS World University Rankings positioned the University of Southampton in the top 80 universities in the world. The university considers itself one of the top 5 research universities in the UK. The university has a global reputation for research into engineering sciences, oceanography, chemistry, cancer sciences, sound and vibration research, computer science and electronics, optoelectronics and textile conservation at the Textile Conservation Centre (which is due to close in October 2009.) It is also home to the National Oceanography Centre, Southampton (NOCS), the focus of Natural Environment Research Council-funded marine research.

Stanton Township is a township in Champaign County, Illinois, USA. As of the 2010 census, its population was 505 and it contained 202 housing units.

Neville A. Stanton is a British Professor of Human Factors and Ergonomics at the University of Southampton. Prof Stanton is a Chartered Engineer (C.Eng), Chartered Psychologist (C.Psychol) and Chartered Ergonomist (C.ErgHF). He has written and edited over forty books and over three hundred peer-reviewed journal papers on applications of the subject. Stanton is a Fellow of the British Psychological Society, a Fellow of The Institute of Ergonomics and Human Factors and a member of the Institution of Engineering and Technology. He has been published in academic journals including "Nature". He has also helped organisations design new human-machine interfaces, such as the Adaptive Cruise Control system for Jaguar Cars.

Finding Nemo

Theatrical release poster

Directed by Andrew Stanton

Produced by Graham Walters

Screenplay by Andrew Stanton Bob Peterson David Reynolds

Story by Andrew Stanton

Starring Albert Brooks Ellen DeGeneres Alexander Gould Willem Dafoe

Music by Thomas Newman

Cinematography Sharon Calahan Jeremy Lasky

Edited by David Ian Salter

Production company Walt Disney Pictures Pixar Animation Studios

Distributed by Buena Vista Pictures Distribution

Release date May 30, 2003 (2003-05-30)

Running time 100 minutes

Country United States

Language English

Budget \$94 million

Box office \$940.3 million

Question: When was Neville A. Stanton's employer founded?

Thought: The employer of Neville A. Stanton is University of Southampton. The University of Southampton was founded in 1862. So the answer is: 1862.

2430  
 2431     >User:  
 2432       {{ retrieved\_passages }}  
 2433       Question: {{ query }}  
 2434       Thought:  
 2435  
 2436  
 2437  
 2438

## C.2 PAGE NAVIGATION PROMPTS

2439  
 2440  
 2441     **Page Navigation — Links Only**  
 2442  
 2443     **System prompt:**  
 2444     You are a helpful assistant who can interact with a custom interface to expertly navigate information  
 2445     networks. The interface consists of a page viewer that shows you the current page contents with  
 2446     clickable links to related pages.

2447  
 2448     >User:  
 2449       <start\_page>  
 2450       START\_PAGE\_TITLE: {{ start\_page\_title }}  
 2451       START\_PAGE\_ID: {{ start\_page\_link\_id }}  
 2452       </start\_page>  
 2453       <target\_page>  
 2454       TARGET\_PAGE\_TITLE: {{ target\_page\_title }}  
 2455       TARGET\_PAGE\_ID: {{ target\_page\_link\_id }}  
 2456       </target\_page>  
 2457  
 2458       <instructions>  
 2459       # Task Instructions  
 2460       You need to navigate from the given START page to the TARGET page using the provided  
 2461       commands in as few steps as possible.  
 2462       You may use any strategy, but your goal is to reach the exact page\_id of the TARGET page, not  
 2463       just a similar title.  
 2464       You have not reached the TARGET page unless the CURRENT\_PAGE\_ID matches the  
 2465       TARGET\_PAGE\_ID exactly even if the titles seem similar.  
 2466       You succeed only when CURRENT\_PAGE\_ID == TARGET\_PAGE\_ID.  
 2467  
 2468       For example:  
 2469       TARGET\_PAGE\_ID: Elon\_Musk  
 2470       CURRENT\_PAGE\_ID: Elon  
 2471  
 2472       These are different pages. Despite the similarity in names, you must land on the exact page ID to  
 2473       complete the task.

2473       # Navigation Tips  
 2474       - Use hub pages (countries, years, broad categories) to bridge between different topics  
 2475       - Go broader before going narrower - find shared categories or themes  
 2476       - Look for pages with many links when you need more options  
 2477       - Think about what connects your start and target (time period, location, field, etc.)  
 2478  
 2479       IMPORTANT: This is an interactive process where you will think and issue ONE command via  
 2480       function calling, see its result, then think and issue your next command.  
 2481       In each step, please output your thinking so that we can follow along.  
 2482       Your thinking should be thorough and so it's fine if it's very long.  
 2483       </instructions>  
 2484       Begin!  
 2485       {{ start\_observation }}

2484

2485

2486

2487

## Page Navigation — Content + Links

**System prompt:**

You are a helpful assistant who can interact with a custom interface to expertly navigate information networks. The interface consists of a page viewer that shows you the current page contents with clickable links to related pages. Links are displayed in markdown format as [entity](link).

2489

2490

2491

**User:**

```

2492 <start_page>
2493 START_PAGE_TITLE: {{ start_page_title }}
2494 START_PAGE_ID: {{ start_page_link_id }}
2495 START_PAGE_CONTENT:
2496 {{ start_page_content }}
2497 </start_page>
2498 <target_page>
2499 TARGET_PAGE_TITLE: {{ target_page_title }}
2500 TARGET_PAGE_ID: {{ target_page_link_id }}
2501 TARGET_PAGE_CONTENT:
2502 {{ target_page_content }}
2503 </target_page>
2504
2505 <instructions>
2506 # Task Instructions
2507 You need to navigate from the given START page to the TARGET page using the provided
2508 commands in as few steps as possible.
2509 You may use any strategy, but your goal is to reach the exact page_id of the TARGET page, not
2510 just a similar title.
2511 You have not reached the TARGET page unless the CURRENT_PAGE_ID matches the
2512 TARGET_PAGE_ID exactly even if the titles seem similar.
2513 You succeed only when CURRENT_PAGE_ID == TARGET_PAGE_ID.
2514
2515 For example:
2516 TARGET_PAGE_ID: Elon_Musk
2517 CURRENT_PAGE_ID: Elon
2518
2519 These are different pages. Despite the similarity in names, you must land on the exact page ID to
2520 complete the task.
2521
2522 # Navigation Tips
2523 - Use hub pages (countries, years, broad categories) to bridge between different topics
2524 - Go broader before going narrower - find shared categories or themes
2525 - Look for pages with many links when you need more options
2526 - Think about what connects your start and target (time period, location, field, etc.)
2527
2528 IMPORTANT: This is an interactive process where you will think and issue ONE command via
2529 function calling, see its result, then think and issue your next command.
2530 In each step, please output your thinking so that we can follow along.
2531 Your thinking should be thorough and so it's fine if it's very long.
2532 </instructions>
2533
2534 Begin!
2535 {{ start_observation }}

```

2530

2531

2532

## Click Page Tool Definition

2533

2534

2535

2536

2537

```

def click_link_to_page(
 self,
 page_id: Annotated[
 str,
 "The ID of the page to navigate to. The ID must be a link on the LATEST CURRENT
 PAGE."],
)

```

```

2538
2539) : """
2540 """
2541 From the links on the current page, click on a link to the next page.
2542 """
2543
2544 Backtrack to Page in History Tool Definition
2545
2546 def backtrack_to_page_in_history(
2547 self,
2548 page_id: Annotated[
2549 str,
2550 "The ID of the page to navigate to. You may only click on links on the in the
2551 NAVIGATION_HISTORY.",
2552],
2553) :
2554 """
2555 Backtrack to a specific page in the navigation history.
2556 """

```

## D ADDITIONAL EXPERIMENT TABLES

Table 4 shows the results aggregated across all task instances for multi-hop QA. Table 5 shows the results aggregated across all task instances for page navigation.

| Model            | Baseline     | RM                |          | SM                |          | KA    |
|------------------|--------------|-------------------|----------|-------------------|----------|-------|
|                  |              | RM (F1)           | RM (R@5) | SM (F1)           | SM (R@5) |       |
| GPT-5-mini       | Closed-book  | 21.6 [19.5, 23.7] | –        | 0.2 [0.1, 0.4]    | –        | 21.4  |
|                  | Reading Comp | 88.1 [86.4, 89.8] | –        | 90.1 [88.5, 91.6] | –        | -2.0  |
|                  | One-step RAG | 49.8 [47.1, 52.4] | 56.1     | 24.4 [22.0, 26.8] | 45.0     | 25.4  |
|                  | IRCoT + RAG  | 54.3 [51.7, 56.9] | 58.9     | 38.1 [35.4, 40.7] | 52.2     | 16.2  |
| Gemini-2.0-Flash | Closed-book  | 19.4 [17.3, 21.4] | –        | 0.6 [0.3, 0.9]    | –        | 18.8  |
|                  | Reading Comp | 75.4 [73.0, 77.8] | –        | 80.3 [78.1, 82.4] | –        | -4.9  |
|                  | One-step RAG | 37.3 [34.7, 39.9] | 56.1     | 17.2 [15.2, 19.3] | 45.1     | 20.1  |
|                  | IRCoT + RAG  | 46.8 [44.1, 49.4] | 60.6     | 38.3 [35.5, 40.9] | 57.5     | 8.5   |
| gpt-oss-20b      | Closed-book  | 17.1 [13.3, 20.2] | –        | 1.6 [0.0, 1.9]    | –        | 15.5  |
|                  | Reading Comp | 66.4 [62.4, 69.6] | –        | 76.7 [72.3, 80.0] | –        | -10.3 |
|                  | One-step RAG | 41.0 [38.0, 44.1] | 55.9     | 20.2 [17.6, 22.8] | 44.6     | 20.8  |
|                  | IRCoT + RAG  | 30.7 [27.5, 33.9] | 45.2     | 18.4 [15.7, 21.2] | 38.5     | 12.3  |
| gpt-oss-120b     | Closed-book  | 21.6 [19.6, 23.7] | –        | 1.8 [1.4, 2.3]    | –        | 19.8  |
|                  | Reading Comp | 73.6 [71.3, 75.8] | –        | 82.8 [80.8, 84.8] | –        | -9.2  |
|                  | One-step RAG | 43.5 [41.0, 46.1] | 55.9     | 22.1 [19.8, 24.4] | 44.6     | 21.5  |
|                  | IRCoT + RAG  | 38.3 [35.6, 41.0] | 43.5     | 23.6 [20.5, 26.8] | 38.2     | 14.7  |

Table 4: **Multi-hop QA Performance on SYNTHWORLD-RM/SM.** Metrics are F1 scores (with 95% confidence intervals) for answer correctness and Recall@5 (R@5) for retrieval. The rightmost column reports the knowledge advantage gap (KA).

2581  
2582  
2583  
2584  
2585  
2586  
2587  
2588  
2589  
2590  
2591

2592

2593

2594

2595

2596

2597

2598

2599

2600

2601

2602

2603

2604

2605

2606

2607

2608

2609

2610

2611

2612

2613

| Model            | Environment     | RM                | SM                | KA   |
|------------------|-----------------|-------------------|-------------------|------|
| GPT-5-mini       | Links Only      | 50.8 [47.7, 53.9] | 19.8 [17.4, 22.4] | 31.0 |
|                  | Content + Links | 52.3 [49.2, 55.4] | 30.6 [27.8, 33.5] | 21.7 |
| Gemini-2.0-Flash | Links Only      | 36.1 [33.2, 39.1] | 15.6 [13.5, 18.0] | 20.5 |
|                  | Content + Links | 41.5 [38.5, 44.6] | 28.0 [25.3, 30.9] | 13.5 |
| gpt-oss-20b      | Links Only      | 28.6 [25.9, 31.5] | 13.2 [11.2, 15.4] | 15.4 |
|                  | Content + Links | 31.5 [28.7, 34.4] | 22.4 [19.9, 25.1] | 9.1  |
| gpt-oss-120b     | Links Only      | 39.9 [36.9, 43.0] | 16.2 [14.0, 18.6] | 23.7 |
|                  | Content + Links | 45.6 [42.5, 48.7] | 30.3 [27.5, 33.2] | 15.3 |
| Kimi-K2-Instruct | Links Only      | 45.3 [42.2, 48.4] | 17.4 [15.2, 19.9] | 27.9 |
|                  | Content + Links | 49.4 [46.3, 52.5] | 31.6 [28.8, 34.5] | 17.8 |
| Kimi-K2-Thinking | Links Only      | 43.6 [40.6, 46.7] | 16.2 [14.0, 18.6] | 27.4 |
|                  | Content + Links | 45.4 [42.3, 48.5] | 28.5 [25.8, 31.4] | 16.9 |

Table 5: **Navigation Success Rate on SYNTHWORLD-RM/SM.** Success rates are reported with 95% confidence intervals. The rightmost column reports the knowledge advantage gap (KA).

2627

2628

2629

2630

2631

2632

2633

2634

2635

2636

2637

2638

2639

2640

2641

2642

2643

2644

2645