

# TRIDENT: Selection-Conditional Facet Tests and Episode-Frozen Min-Cost Cover for Budgeted RAG

Anonymous ACL submission

## Abstract

Budgeted RAG is a decision problem: under a hard evidence cap, which passages you keep determines both accuracy and what you can credibly claim about the evidence at query time. We introduce TRIDENT, a framework that mines auditable reasoning facets, tests facet support with a calibrated verifier, and selects evidence under an explicit token budget. In the Safe-Cover regime, we freeze the retrieval pipeline into a replayable episode, map verifier scores to selection-conditional conformal  $p$ -values under a logged contract, and apply per-query multiple-testing control to yield facet-support certificates—or return a machine-checkable abstention with a reason code. In Pareto-Knapsack, we drop per-query guarantees and optimize a quality–cost frontier for throughput. On HotpotQA at a 500-token evidence cap, TRIDENT Pareto-500 improves EM/F1 from 30.81/39.61 to 45.30/58.22 (+47% relative), while using 3% fewer evidence tokens and 5% lower latency than naive top- $k$  truncation. These results show that under tight budgets, selection rigor and query-time evidence accountability matter as much as retrieval strength.

## 1 Introduction

Retrieval-augmented generation can improve factuality by grounding a model in retrieved text; however, under strict context limits and latency targets, the bottleneck shifts: the question is no longer whether relevant information can be found, but which evidence is worth retaining under a hard cap. Retrieving more passages can raise accuracy, yet it directly increases token consumption and end-to-end latency while enlarging the surface area for redundant or spurious citations (Lewis et al., 2020; Karpukhin et al., 2020; Izacard and Grave, 2020). Under production constraints—such as hard context limits, SLOs, and GPU contention—the system must do more than rank documents; it must also decide what to keep and when to stop.

Despite progress in adaptive pipelines, such as multi-step retrieval and retrieve–reflect control (Guu et al., 2020; Asai et al., 2023), budgeted evidence selection is still commonly implemented as top- $k$  concatenation plus truncation. Under a hard cap, this policy is brittle for multi-hop questions: truncation can drop a required bridge passage while retaining redundant or merely related text. The opportunity cost is substantial. Even with the retrieval stack held fixed, selection alone can dramatically change outcomes: at an approximate 500-token evidence cap on HotpotQA, our Pareto mode improves EM/F1 from 30.81/39.61 to 45.30/58.22.

A second pain point is query-time accountability. Many RAG systems output citations, but they provide little quantitative statement that the selected passages satisfy the intermediate requirements that a query depends on. Under tight budgets, this matters operationally: a system should sometimes answer, but it should also sometimes refuse—and that refusal should be an auditable decision derived from a rigorous contract, not a post-hoc narrative.

We propose TRIDENT, a budgeted evidence-selection framework that makes the choice of evidence explicit under strict caps. TRIDENT utilizes typed reasoning facets to structure what must be supported, and a verifier signal to determine what evidence is worth paying for. It operates in two regimes: a *Pareto* mode that optimizes the quality–cost frontier under a hard evidence cap, and a *Safe-Cover* mode that emits machine-checkable evidence receipts, or cleanly abstains when the audit contract cannot be satisfied. The end-to-end design is summarized in Section 3 (Figure 1).

## Contributions.

- We formalize budgeted RAG as facet-level evidence selection under a hard context budget, with an explicit separation between certifying support and producing an answer, and with

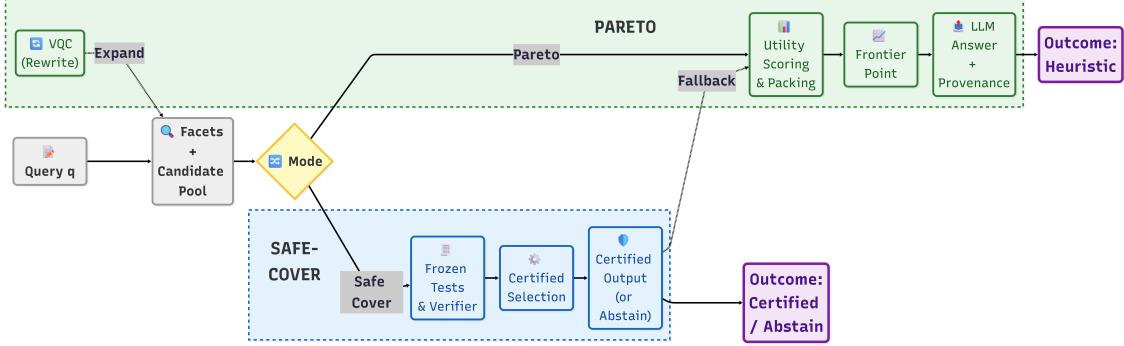


Figure 1: TRIDENT overview. Facets define verifiable requirements, and selection packs evidence under a hard token budget. Safe-Cover certifies facet support or abstains while Pareto-Knapsack optimizes a quality–cost frontier without per-query validity claims.

standardized abstention outcomes for evaluation integrity.

- We introduce Safe-Cover, a certified regime that maps verifier scores to selection-conditional conformal p-values, enabling auditable minimum-cost cover or a replay-valid infeasibility certificate under a query-level error budget.
- We evaluate on HotpotQA, 2WikiMulti-HopQA, and MuSiQue, reporting EM/F1 together with evidence tokens and latency, with diagnostics that attribute performance to coverage, calibration, and verifier behavior rather than generic failure labels.

## 2 Related Work

**Structured retrieval and reasoning granularity.** Multi-hop QA has pushed the field beyond simple independent top- $k$  passage selection. This evolution began with stronger dense retrievers (DPR; Karpukhin et al., 2020) and multi-hop retrieval variants (e.g., MDR; Xiong et al., 2021; BeamDR; Zhao et al., 2021), and has recently matured into structured approaches that organize evidence as graphs, propositions, or memory-like indices (e.g., HippoRAG; Gutiérrez et al., 2025; PropRAG; Wang and Han, 2025). These methods enhance candidate discovery and chaining by modifying indexing granularity or retrieval structure, primarily addressing the recall problem. They do not by themselves address the accounting problem: defining a query-time, machine-checkable statement of which intermediate reasoning requirements are satisfied under a strict evidence budget.

**Verifiable selection and efficient RAG.** A broad line of reliability work focuses on verifying or repairing model outputs after generation, using self-critique loops (e.g., SELF-RAG; Asai et al., 2023; Corrective RAG; Yan et al., 2024) and attribution frameworks (e.g., FActScore; Min et al., 2023; Li et al., 2024). In parallel, systems-oriented work targets deployment constraints directly, reducing time-to-first-token through caching (Lu et al., 2025), shrinking context via selective augmentation (Mao et al., 2025), or filtering evidence based on estimated utility (Wang et al., 2025). Recent advances in context selection reconceptualize quality assessment as a data valuation problem, measuring each context’s marginal contribution through influence-based metrics that capture query-aware relevance, list-aware uniqueness, and generator-aware alignment (Deng et al., 2025). A remaining gap is ex ante evidence verification: validating intermediate support before generation to avoid spending tokens on low-value or non-entailing evidence, while making the selection decision auditable at query time rather than defended retroactively.

**Statistical guarantees in dynamic retrieval.** Conformal prediction offers rigorous tools for calibration and uncertainty quantification (Angelopoulos and Bates, 2022), but applying such guarantees in retrieval settings is a subtle task. The policy itself induces the hypotheses being tested: retrieval and shortlisting determine which (passage, requirement) pairs become candidates. Agentic loops that rewrite queries or expand pools can shift the data distribution mid-flight, complicating any per-query validity claim unless the selection mechanism is explicitly controlled and logged. This motivates the “frozen episode” approach, which involves certificate-style guarantees that are con-

ditional on a locked selection contract (retriever snapshot, shortlist policy, binning), thereby explicitly separating certified regimes from adaptive optimization regimes where no per-query statistical validity is claimed.

### 3 Framework

We frame retrieval-augmented answering as budgeted, auditable evidence selection with an explicit separation between (i) certifying that evidence supports query requirements and (ii) producing an answer. Given a query  $q$ , the system extracts a finite set of intermediate requirements (facets)  $F(q)$  and selects a small set of passages  $S$  under a hard evidence budget. The framework supports two serving regimes: a certified regime that emits machine-checkable facet-support certificates under a replayable selection contract, and a frontier regime that optimizes quality–cost operating points without per-query statistical claims. In both regimes, abstention is a first-class outcome: rather than hallucinating an answer, the system returns a fixed output token accompanied by an auditable reason code. Figure 1 summarizes the end-to-end pipeline. Both regimes share the same verification backbone: facets define what must be supported, the verifier supplies the support signal, and selection decides what evidence is worth paying for under  $B_{ctx}$ . Figure 2 shows Safe-Cover’s certified trace on a two-hop query.

#### 3.1 Facets: Auditable Reasoning Requirements

A facet miner maps  $q$  to  $F(q) = \{f_1, \dots, f_m\}$ , where each facet is intended to be checkable against a single passage. Facets are typed (e.g., ENTITY, RELATION, BRIDGE-HOP1, BRIDGE-HOP2, TEMPORAL, NUMERIC) to support stratified calibration and targeted diagnostics. Each facet includes (i) a hypothesis/claim template and (ii) a deterministic shortlisting key (anchors, triggers, or entity bindings) that defines exactly which passages will be tested by the verifier, ensuring the process remains auditable. Multi-hop queries may introduce placeholder facets whose hypotheses are not meaningful until a binding value is available. To avoid scoring non-instantiated hypotheses, placeholder facets are skipped during the initial scoring pass and evaluated only after instantiation (e.g., after Hop-1 binds an intermediate entity that populates a Hop-2 template). This implementation

detail prevents spurious failures caused by testing ill-formed hypotheses. Facet mining is the functional bottleneck for certified behavior: the system can only certify what it can express as facets. Accordingly, all guarantee statements are conditional on the mined set  $F(q)$  and the logged contract used to test it.

#### 3.2 Candidates and Cost Model

A retriever produces a candidate pool  $P = \{p_1, \dots, p_n\}$ . Each passage  $p$  has a nonnegative cost  $c(p)$  representing the token cost of serializing that evidence into the prompt. We impose a hard evidence cap  $B_{ctx}$  and require  $\sum_{p \in S} c(p) \leq B_{ctx}$  for any selected set  $S$ . Costs count evidence tokens only; generation tokens are tracked separately. The cost model and candidate pool form part of the immutable selection contract for the certified regime.

#### 3.3 Verification as fixed tests and selection-conditional p-values

**Deterministic shortlisting defines the tested set.** Verifying every pair  $(p, f)$  is computationally wasteful and statistically unstable. For each facet  $f$ , a deterministic shortlister selects up to  $T_f$  passages from  $P$  using fixed rules and fixed tie-breaks. Only these shortlisted pairs are scored by the verifier. This turns "what was tested" into a declared, replayable object—a prerequisite for both valid multiple-testing control and external auditing.

**Event of interest and verifier scores.** For each shortlisted pair  $(p, f)$ , a verifier produces a score  $s(p, f) \in R$  intended to correlate with the event  $\Sigma(p, f)$ : "passage  $p$  is sufficient evidence for facet  $f$ ". In practice, the verifier may be two-stage (lexical gate followed by an NLI score); the score  $s(p, f)$  is the final verification statistic passed to calibration. This is the key semantic point: the p-value measures sufficiency for the facet, not merely whether a string appears.

**Selection-conditional conformal p-values.** Crucially, verifier score distributions are artifacts of the retrieval policy. We therefore calibrate under the same policy used at test time. Each shortlisted  $(p, f)$  is assigned to a bin  $b = b(p, f)$  and mapped to a conformal p-value using a negative calibration pool  $\mathcal{N}_b$  generated by replaying the identical selection logic:

$$\pi(p, f) = \frac{1 + \sum_{u \in \mathcal{N}_b} I[s(u) \geq s(p, f)]}{|\mathcal{N}_b| + 1}, \quad (1)$$

249 with randomized tie-handling when scores are discrete. These p-values are valid only under the re-  
 250 played contract.  
 251

252 **Mondrian binning and feasibility.** To reduce  
 253 heterogeneity, we use Mondrian calibration with  
 254 a bin key that includes at least facet type and a  
 255 passage length bucket:

$$256 \quad b(p, f) = (\text{canonical\_type}(f), \text{len\_bucket}(p)). \quad (2)$$

257 Deterministic conformal p-values have a floor  
 258  $1/(|\mathcal{N}_b| + 1)$ . If the target threshold falls below  
 259 this floor, certification is mathematically impossi-  
 260 ble. We apply a fixed bin-merging rule to reach a  
 261 minimum effective pool size, logging merge depth  
 262 to the replay record. If feasibility still cannot be  
 263 met, the certified regime does not guess—it returns  
 264 an explicit reason for infeasibility.

### 265 3.4 Certified regime: Safe-Cover as 266 risk-controlled min-cost facet cover

267 The certified regime converts calibrated facet tests  
 268 into a fixed set system and solves a budgeted  
 269 min-cost cover. It guarantees one of two out-  
 270 comes: a verified answer backed by per-facet cer-  
 271 tificates (with a bounded query-level error rate), or  
 272 a machine-checkable abstention. Figure 2 provides  
 273 a concrete trace of this process.

274 **Episode contract (frozen knobs).** Within an  
 275 episode (one query execution), all parameters defin-  
 276 ing the tested set and thresholds are frozen and  
 277 logged: retriever snapshot and candidate cap; deter-  
 278 ministic shortlisting rules and  $T_f$ ; verifier version;  
 279 binning scheme; calibration pools; evidence ser-  
 280 alization rules; and the evidence budget  $B_{\text{ctx}}$ . No  
 281 mid-episode adaptation is permitted. This is what  
 282 makes the guarantee checkable: if any contract  
 283 element changes, certificate validity is void.

284 **Bonferroni allocation and fixed coverage sets.**  
 285 We allocate a query-level error budget  $\alpha_{\text{query}}$  across  
 286 facets and tests:

$$287 \quad \alpha_f = \frac{\alpha_{\text{query}}}{|F(q)|}, \quad \bar{\alpha}_f = \frac{\alpha_f}{T_f}. \quad (3)$$

288 A facet is covered by passage  $p$  when  $\pi(p, f) \leq \bar{\alpha}_f$ .  
 289 This induces a fixed coverage set per passage:

$$290 \quad C(p) = \{f \in F(q) : \pi(p, f) \leq \bar{\alpha}_f\}. \quad (4)$$

291 Because shortlisting and thresholds are frozen,  
 292  $C(p)$  is fixed within the episode, restoring classical  
 293 set-cover semantics.

294 **Budgeted min-cost cover and reproducible  
 295 greedy selection.** Safe-Cover solves a min-cost  
 296 cover under a hard evidence budget:

$$297 \begin{aligned} \min_{S \subseteq P} \quad & \sum_{p \in S} c(p) \\ \text{s.t.} \quad & \bigcup_{p \in S} C(p) \supseteq F(q), \\ & \sum_{p \in S} c(p) \leq B_{\text{ctx}}. \end{aligned} \quad (5)$$

298 We use a deterministic greedy algorithm that maxi-  
 299 mizes newly covered facets per unit cost, ensuring  
 300 exact reproducibility.

301 **Winning-passage alignment and answer extrac-  
 302 tion.** Certification is useless if the answer is gen-  
 303 erated from an unverified context. For each facet,  
 304 the system records a winning passage (typically  
 305 the passage with the smallest p-value among can-  
 306 didates that passed the threshold). Answer extraction  
 307 is aligned to these winning passages: typed extrac-  
 308 tors run on the certifying evidence, and the system  
 309 does not silently substitute unrelated context if ex-  
 310 traction fails. This resolves a standard failure mode  
 311 where a certificate is computed on one piece of  
 312 evidence, but the answer is generated from another.

313 **Abstention certificates and reason codes.** If  
 314 any required facet has no covering passage, Safe-  
 315 Cover returns a NO\_COVER outcome. If cov-  
 316 erage exists in principle but cannot be achieved  
 317 under the remaining budget, Safe-Cover returns  
 318 INFEASIBLE\_BUDGET (optionally supported by a  
 319 dual lower bound on remaining cover cost). If  
 320 threshold feasibility fails due to bin floors after  
 321 merging, Safe-Cover returns INFEASIBLE\_PVALUE.  
 322 All abstentions are normalized to a fixed output  
 323 token (e.g., ABSTAIN) for evaluation integrity.

324 **Certificate payload and query-level statement.**  
 325 A facet-support certificate includes the facet ID,  
 326 winning passage ID, p-value, threshold, bin key,  
 327 shortlist metadata, and version hashes. Under the  
 328 frozen contract and Bonferroni allocation, the prob-  
 329 ability of emitting any false facet-support certificate  
 330 within a query is bounded by  $\alpha_{\text{query}}$ .

### 331 3.5 Heuristic tier: provenance-checked 332 answering

333 Not all queries are naturally expressible as a small  
 334 set of atomic, single-passage facets. When Safe-  
 335 Cover abstains, or the query falls outside the certifi-

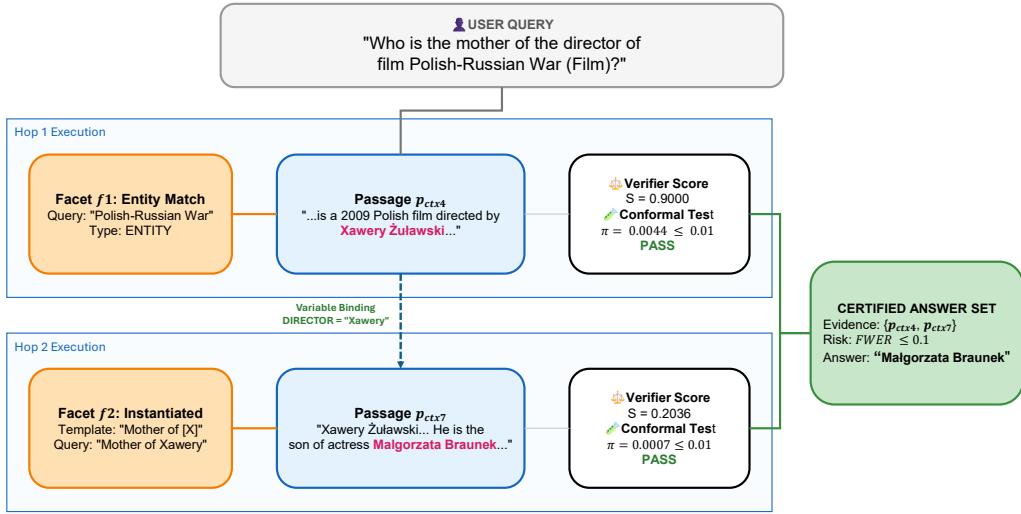


Figure 2: Safe-Cover trace on a two-hop query: Hop-1 binds an intermediate entity, Hop-2 instantiates a dependent facet, and passing tests yield certificates and a certified evidence set.

able regime, the system invokes an LLM-based answerer over the complete context pool to maximize recall. This tier enforces a strict provenance check: the final answer must be extractable from the supplied passages via Unicode-normalized matching. It does not claim the certified regime’s query-level error bound.

### 3.6 Frontier regime: Pareto-Knapsack

The frontier regime drops per-query certification and instead optimizes a quality–cost objective under the same evidence cap. It treats verifier scores as continuous utility signals and greedily packs the context window to maximize evidence strength per token. Because this mode may employ adaptive mechanisms (e.g., query rewrites), it reports empirical quality–cost trade-offs without per-query validity claims. Figure 3 details the complete Pareto pipeline: facet mining and retrieval produce a candidate pool, batch scoring assigns verifier-derived utilities, and lazy greedy packing selects passages under the evidence budget  $B_{ctx}$  by maximizing marginal utility per cost. The system enforces strict provenance checking—answers must be extractable from the selected evidence—and abstains when provenance fails.

### 3.7 Operational safeguards and diagnostics

The system logs certificate payloads, shortlist sizes, bin depths, near-threshold counts, abstention codes, and version metadata sufficient for replay. These signals support debugging (coverage sparsity, veri-

fier discrimination) and operational monitoring. If drift alarms trigger—indicating the test-time distribution has diverged from the calibration pool—the certified regime fails closed until re-calibration is performed.

## 4 Experiments

TRIDENT targets a regime where multi-hop questions require evidence, but evidence is expensive: both context length and latency must be actively controlled. We therefore evaluate TRIDENT as a quality–cost system rather than a pure-accuracy model. Across all experiments, we hold the generator and upstream retrieval/verifier stack constant; we vary only the serving regime (Pareto-Knapsack vs. Safe-Cover), the evidence budget, and (for Safe-Cover) the certification parameters. This isolates the impact of evidence selection from retrieval quality and generation capacity.

### 4.1 Experimental setup

**Datasets.** We evaluate on three multi-hop QA benchmarks: HotpotQA (Yang et al., 2018), 2Wiki-MultiHopQA (Ho et al., 2020), and MuSiQue (Trivedi et al., 2022). Each requires aggregating evidence across multiple passages—exactly the setting where top-k truncation is brittle under a hard cap. HotpotQA and 2Wiki provide annotated supporting facts for selection analysis under budget pressure. MuSiQue is a harder stress-test: longer reasoning chains and less redundant evidence expose when selection or certification be-

Method / Config	Model	2Wiki				HotpotQA				MuSiQue			
		EM (%)	F1 (%)	Evidence tokens (EvTok)	Lat <sub>50</sub> (ms)	EM (%)	F1 (%)	Evidence tokens (EvTok)	Lat <sub>50</sub> (ms)	EM (%)	F1 (%)	Evidence tokens (EvTok)	Lat <sub>50</sub> (ms)
<i>TRIDENT (Pareto-Knapsack; evidence-budget sweep)</i>													
Pareto-400	Llama-3-8B-Instruct	25.95	33.83	251.39	1044.99	43.80	55.80	348.01	446.12	15.02	26.93	344.75	1029.04
Pareto-500	Llama-3-8B-Instruct	28.16	36.25	331.48	1603.85	45.30	58.22	446.22	467.80	19.21	30.74	446.30	1642.87
Pareto-1000	Llama-3-8B-Instruct	33.81	42.45	631.08	1967.42	50.23	63.68	874.63	548.55	20.16	33.54	847.73	514.44
<i>TRIDENT (Safe-Cover; certified mode)</i>													
Safe-2000 (equal)	Llama-3-8B-Instruct	34.70	43.61	869.07	1431.62	49.24	62.55	1251.46	1631.55	3.82	8.74	1899.39	15310.80
Safe-4000 (loose)	Llama-3-8B-Instruct	34.52	43.42	871.20	2016.33	48.89	62.27	1251.33	505.17	4.53	10.09	1898.32	15555.27

Table 1: TRIDENT main results. EvTok: evidence tokens passed to the generator. Lat<sub>50</sub>: per-query end-to-end latency median (ms) over all queries, including abstentions. Latency and token accounting follow App. §A–C. Abstention, confidence intervals and additional protocol details are in the appendix.

comes feasibility-limited.

**Models and decoding.** All budget-matched comparisons use the same model (Llama-3-8B-Instruct (AI@Meta, 2024) or Qwen3-8B (Team, 2025)) with identical decoding parameters (temperature 0.0, consistent stopping criteria, same max\_new\_tokens). We apply uniform answer normalization and extraction rules across systems (App. §A) and report both generators throughout but interpret within-generator comparisons most heavily, since absolute latency and output behavior can differ across model families under identical caps in our serving setup.

**Retrieval and budget enforcement.** All TRIDENT variants share the same retriever, re-ranker, and verifier; differences arise solely from selection logic. For baselines, we log retrieval conditions and—when dataset-provided contexts are used—the exact ordering rule applied before truncation. This matters on 2Wiki especially: passage ordering can determine whether a bridge passage survives a hard cap or gets truncated away. We enforce deterministic ordering and report it verbatim (App. §B), and every run logs the final evidence list passed to the generator. We distinguish evidence tokens (EvTok) from total input tokens (TotTok) and enforce a hard evidence-token cap by truncating after selection and before generation using a deterministic policy (App. §B). We report EvTok and latency alongside EM/F1 because cost is central to the objective—accuracy gains that require 3× the tokens are real but less interesting for budget-constrained deployment.

**Serving regimes.** Pareto-Knapsack optimizes quality against cost under a hard evidence cap. We sweep budgets (400/500/1000 tokens) to trace the quality–cost frontier; main results use a relaxed ac-

ceptance threshold, with Section 4.4 showing that tightening  $\alpha$  primarily raises abstention rather than improving accuracy proportionally.

Safe-Cover is the certified regime, reported at two configurations (Safe-2000 and Safe-4000) corresponding to different evidence caps and certification stringency. All episode-frozen knobs (shortlisting policy, tests-per-facet  $T_f$ , binning key, threshold allocation, tie-breaks, and serialization format) are locked before selection and logged as part of the audit record. This is what makes the certificate replayable: if any contract element changes, the validity claim does not transfer.

Latency is measured end-to-end (retrieve → re-rank → verify → select → answer/abstain), including abstentions, under fixed batching and cache policies (App. §A). When comparing latency in prose, we specify whether we reference mean or percentile values and cite the corresponding table.

## 4.2 Diagnostics for certified behavior

Safe-Cover’s core claim is a query-time certificate under a replayable contract. That claim is only credible if failures are auditable rather than hidden behind aggregate metrics. The system logs three classes of diagnostic signals: (1) per-facet coverage status with abstention reason codes (NO\_COVER, INFEASIBLE\_BUDGET, INFEASIBLE\_PVALUE), (2) calibration feasibility indicators (bin sizes, merge depth, threshold floors), and (3) verifier discrimination metrics. Figure 2 shows a representative execution trace: tests run on fixed  $(p, f)$  pairs, passing facets record their winning passages, and answers are extracted from those winners. If extraction fails or certified coverage is infeasible, the system returns ABSTAIN with a machine-readable reason code that enables post-hoc diagnosis of whether failures stem from evidence scarcity, calibration power limits, verifier discrimination, or budget con-

Method / Config	Model	2Wiki					HotpotQA					MuSiQue				
		EM	F1	EvTok	TotTok	Lat	EM	F1	EvTok	TotTok	Lat	EM	F1	EvTok	TotTok	Lat
<i>Llama-3-8B-Instruct Comparisons</i>																
TRIDENT Pareto-500	Llama-3-8B	28.16	36.25	<b>331</b>	<b>878</b>	5,650	45.30	58.22	<b>446</b>	<b>1,006</b>	<b>2,541</b>	19.21	30.74	446	<b>920</b>	<b>2,829</b>
VanillaRAG-500	Llama-3-8B	26.13	29.26	<u>361</u>	<u>1,002</u>	2,559	30.81	39.61	<u>461</u>	<u>1,193</u>	2,686	5.13	9.16	<b>445</b>	<u>1,155</u>	2,912
TRIDENT Pareto-1000	Llama-3-8B	<b>33.81</b>	<b>42.45</b>	631	1,179	5,476	<b>50.23</b>	<b>63.68</b>	875	1,438	<u>2,570</u>	<b>20.16</b>	<b>33.54</b>	848	1,310	<b>793</b>
VanillaRAG-1000	Llama-3-8B	27.05	30.28	441	1,163	<b>2,531</b>	31.57	40.55	544	1,360	2,673	5.83	9.62	565	1,396	2,834
HippoRAG2	Llama-3-8B	27.42	31.32	569	1,472	3,509	37.79	47.82	593	1,488	4,185	15.35	23.95	663	1,639	7,835
<i>Qwen3-8B Comparisons</i>																
TRIDENT Pareto-500	Qwen3-8B	13.15	17.41	<b>332</b>	<u>1,332</u>	25,118	28.11	36.13	<b>451</b>	<u>1,468</u>	25,728	3.64	6.40	446	<b>1,427</b>	21,814
VanillaRAG-500	Qwen3-8B	20.27	21.98	<u>361</u>	<b>1,241</b>	8,549	27.08	33.77	<u>461</u>	<b>1,465</b>	<b>8,677</b>	3.14	5.21	<b>445</b>	<u>1,463</u>	<b>9,857</b>
TRIDENT Pareto-1000	Qwen3-8B	<b>36.93</b>	<b>43.06</b>	632	1,661	26,273	<b>44.14</b>	<b>56.12</b>	873	1,913	29,963	<u>7.29</u>	<u>10.94</u>	847	1,847	42,463
VanillaRAG-1000	Qwen3-8B	<u>21.32</u>	23.23	441	1,410	8,793	27.87	34.70	544	1,640	8,718	3.14	5.35	565	1,716	9,980
HippoRAG2	Qwen3-8B	<u>21.14</u>	<u>23.75</u>	518	1,716	<b>4,783</b>	34.77	44.01	573	1,896	22,621	<b>14.77</b>	<b>20.47</b>	625	2,221	53,329
<i>Reference only</i>																
Self-RAG	Self-RAG	2.74	17.44	877	935	886	14.56	30.68	1,261	1,319	959	1.65	9.16	1,168	1,235	1,053

Table 2: Baseline comparison across varying evidence budgets (approx. 500 and 1000 tokens). EM/F1 are in %. Lat is average latency in ms. EvTok denotes evidence tokens passed to the generator; TotTok denotes total input tokens. The best and second-best results in each column are highlighted in bold and underlined, respectively.

471 straints.

472 **Auditability artifacts.** Safe-Cover emits a re-  
473 playable certificate payload for each passed facet  
474 (facet id, winning passage id, p-value/threshold, bin  
475 key, version hashes) plus machine-checkable out-  
476 come codes under the episode-frozen contract. If  
477 drift alarms trigger, the certified regime fails closed  
478 until re-calibration.

### 479 4.3 Main results

480 Tables 1 and 2 show that TRIDENT’s gains come  
481 from treating budgeted RAG as a selection problem  
482 with explicit intermediate requirements, not from  
483 retrieving more text or relying on longer contexts.  
484 Under strict caps, the dominant failure mode of top-  
485  $k$  truncation is not that relevant evidence is absent,  
486 but that the critical bridge passage gets dropped or  
487 diluted within the packed context.

488 **Selection quality under a hard cap.** At matched  
489 evidence budgets, verification-aware selection  
490 yields large gains over naive concatenation. On  
491 HotpotQA at a 500-token cap, Pareto-500 improves  
492 EM/F1 from 30.81/39.61 (VanillaRAG-500) to  
493 45.30/58.22—a +14.49/+18.61 absolute gain, 47%  
494 relative improvement on both metrics—while us-  
495 ing slightly fewer evidence tokens (446 vs. 461).  
496 Using the mean end-to-end latency reported in Ta-  
497 ble 2, Pareto-500 is also modestly faster (2,541ms  
498 vs. 2,686ms), consistent with selecting less redun-  
499 dant context.

500 On MuSiQue, at essentially the same evidence  
501 budget (446 vs. 445 EvTok), Pareto-500 improves

502 EM/F1 from 5.13/9.16 to 19.21/30.74 ( $3.7 \times$  EM;  
503 3.4× F1). On 2WikiMultiHopQA, gains are  
504 smaller but consistent: 28.16/36.25 vs. 26.13/29.26  
505 EM/F1 while using fewer evidence tokens (331 vs.  
506 361). The pattern throughout: Pareto improves ac-  
507 curacy by packing better evidence, not by spending  
508 more tokens.

509 **The budget knob behaves predictably.** Sweep-  
510 ing the evidence budget yields a monotone, inter-  
511 pretable operating curve (Table 1): increasing the  
512 allowance from 400 to 500 tokens improves EM/F1,  
513 and increasing to 1000 improves further, with pre-  
514 dictable cost increases and modest latency changes.  
515 This is exactly what a budgeted framework should  
516 provide—an operator can choose a target quality  
517 level and read off the corresponding token cost  
518 without redesigning the system at each operating  
519 point.

520 Abstention at Pareto-500 remains low but non-  
521 zero (HotpotQA 1.54%, 2Wiki 2.12%, MuSiQue  
522 6.04%), consistent with fail-closed behavior: when  
523 evidence cannot support a reliable output under  
524 the cap, the system abstains rather than fabricating  
525 plausible citations.

526 **Selection complements stronger retrieval.** TRI-  
527 DENT is not a retrieval substitute; it improves how  
528 a fixed candidate pool gets distilled under bud-  
529 get. Compared to HippoRAG (Table 2), Pareto-500  
530 achieves higher accuracy while using fewer tokens  
531 and lower mean latency: on HotpotQA, +7.5 EM  
532 / +10.4 F1 with 25% fewer EvTok (446 vs. 593)  
533 and 39% lower latency (2,541ms vs. 4,185ms); on

534 MuSiQue, +3.9 EM / +6.8 F1 with 33% fewer Ev-  
 535 Tok (446 vs. 663) and 64% lower latency (2,829ms  
 536 vs. 7,835ms). Better retrieval and better selection  
 537 compound rather than compete—expanding the  
 538 candidate pool helps, but packing the right subset  
 539 into a tight context window remains decisive.

540 **Safe-Cover: trading cost for auditability.** Safe-  
 541 Cover targets a different axis: it answers only  
 542 when it can certify facet coverage under an episode-  
 543 frozen contract, otherwise abstaining with an ex-  
 544 plicit reason code. Where certification is fea-  
 545 sible, Safe-Cover remains accurate (HotpotQA  
 546 49.24/62.55 EM/F1; 2Wiki 34.70/43.61 in Table 1),  
 547 but uses substantially more evidence tokens—the  
 548 intended cost of certifying coverage rather than  
 549 merely optimizing average accuracy. For applica-  
 550 tions where a wrong answer is worse than no  
 551 answer, this trade-off makes sense.

552 **The MuSiQue collapse.** On MuSiQue, Safe-  
 553 Cover achieves very low EM/F1 despite large  
 554 evidence budgets. This does not invalidate the  
 555 framework; it exposes where certification becomes  
 556 power-limited. Safe-Cover can only certify what  
 557 the verifier can separate and what calibration  
 558 can support under bin floors and multiple-testing  
 559 thresholds. Safe-Cover’s failure mode is itself au-  
 560 ditable—it fails closed with explicit reason codes.  
 561 A system that confidently hallucinated answers  
 562 would score higher on EM but be far less trust-  
 563 worthy in deployment.

Config		EM / F1	Abstain %	EvTok	TotTok	Lat p50
$\alpha=0.01$ (Strict)	HotpotQA	43.67% / 56.37%	10.93%	454.25	1017.32	514.95
$\alpha=0.02$	HotpotQA	44.39% / 57.57%	8.17%	454.40	1018.23	506.62
$\alpha=0.05$	HotpotQA	44.63% / 56.53%	16.21%	454.20	1015.11	516.44
$\alpha=0.1$	HotpotQA	44.77% / 57.37%	6.89%	454.05	1012.57	512.94
$\alpha=0.2$	HotpotQA	44.94% / 57.39%	4.17%	452.65	1016.12	1030.03
$\alpha=0.5$	HotpotQA	45.56% / 58.45%	1.54%	451.02	1011.33	1166.79
$\alpha=0.6$ (Loose)	HotpotQA	45.30% / 58.22%	1.54%	446.22	1006.11	467.80
No-Rerank ( $\alpha=0.6$ )	HotpotQA	45.94% / 58.08%	1.54%	450.51	1011.62	2275.23
Soft Sigmoid ( $\alpha=0.6$ )	HotpotQA	46.14% / 59.20%	1.54%	466.90	1025.37	531.21
Soft Sigmoid ( $\alpha=0.6$ )	2Wiki	28.90% / 36.81%	2.12%	334.58	879.78	1148.53
Soft Sigmoid ( $\alpha=0.6$ )	Musique	17.28% / 28.62%	6.04%	460.32	934.35	1504.08

Table 3: Pareto ablations under a fixed evidence budget. Different  $\alpha$  configurations sweeps the relaxed threshold  $\alpha$  (HotpotQA), config No-Rerank removes reranking, and config soft sigmoid replaces binary marginal gain with a soft-sigmoid utility. We report EM/F1, abstention, EvTok, and Lat<sub>50</sub>.

#### 4.4 Ablations and analysis

Our ablations test whether Pareto’s gains rely on brittle implementation choices—a fragile reranker ordering, a knife-edge utility at  $p = \alpha$ , or a narrow threshold setting—or whether they reflect a stable interaction between verification signals and

budgeted selection. The distinction matters: brittle gains evaporate when conditions shift; stable gains transfer.

**Threshold sensitivity.** Sweeping the relaxed acceptance threshold  $\alpha$  (detailed plots in Appendix A) shows that tightening  $\alpha$  sharply increases abstention—the system cannot assemble a budget-feasible cover—while EM/F1 changes more modestly. The evidence budget, not the threshold, is the binding constraint on token consumption. We use a relaxed setting in the main results to prioritize answering when evidence is plausibly sufficient; operators with stricter requirements can tighten  $\alpha$  and accept higher abstention.

**Reranking and utility design.** Removing reranking leaves HotpotQA EM/F1 and abstention essentially unchanged (Table 3), while increasing latency—Pareto’s 500-token gains are not artifacts of fragile reranker-induced ordering, as reranking primarily improves efficiency by reducing noise in the candidate list. This suggests the selection logic is robust to upstream perturbations. Similarly, replacing the binary utility with a soft-sigmoid score (utility\_tau of 0.05) yields comparable and slightly improved HotpotQA EM/F1 (46.14/59.20 vs. 45.30/58.22) at the same abstention rate (1.54%), with a modest EvTok increase. Pareto does not depend on a knife-edge decision boundary at  $p = \alpha$ ; smoothing near-threshold support can help the greedy solver in borderline cases without changing the semantic definition of coverage.

## 5 Conclusion

TRIDENT reframes budgeted RAG as facet-level evidence selection under an explicit token cap, treating cost as a first-class constraint rather than an afterthought. Empirically, Pareto-Knapsack establishes a stable quality-cost frontier and substantially outperforms naive top- $k$  truncation—a 47% relative gain on HotpotQA at 500 tokens demonstrates that under strict budgets, what you keep matters more than how much you retrieve. For deployments requiring auditability, Safe-Cover converts this selection logic into a rigorous contract—either a replayable facet-support certificate under an episode-frozen protocol, or an explicit abstention—ensuring statistical guarantees are claimed only when the selection contract is satisfied.

## 618 Limitations

619 **Contract-scoped validity.** Safe-Cover certificates  
620 are claimed only under the logged selection  
621 contract: the same retriever snapshot, deterministic  
622 shortlisting policy (including tie-breaks and  $T_f$ ),  
623 verifier, calibration bins, and evidence serialization  
624 rules used during calibration must be replayed at  
625 test time. If any element changes, certificate validity  
626 is not claimed.

627 **Sensitivity to verifier and facet design.** Certification  
628 depends on the verifier’s ability to separate sufficient  
629 from insufficient support at the facet granularity.  
630 Poorly specified facets, weak verifier discrimination,  
631 or sparse calibration bins can make certification  
632 infeasible even when a correct answer exists in the  
633 candidate pool.

634 **Conservatism under diffuse evidence.** Safe-  
635 Cover can be conservative on tasks where support  
636 is distributed across passages or does not align  
637 cleanly with single-passage facet checks, as seen  
638 on MuSiQue. This is less a failure of the abstraction  
639 than a reminder that auditable certification  
640 fundamentally trades coverage for strictness.

641 **Scope of evaluation.** We focus on multi-hop QA  
642 benchmarks and report accuracy, evidence tokens,  
643 and latency. Broader tasks (summarization, long-  
644 form synthesis) may require different facet designs  
645 and verification signals.

## 646 References

647 AI@Meta. 2024. [Llama 3 model card](#).

648 Anastasios N. Angelopoulos and Stephen Bates. 2022.  
649 [A Gentle Introduction to Conformal Prediction and](#)  
650 [Distribution-Free Uncertainty Quantification](#). *arXiv*  
651 *preprint*. ArXiv:2107.07511 [cs].

652 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and  
653 Hannaneh Hajishirzi. 2023. [Self-RAG: Learning](#)  
654 [to Retrieve, Generate, and Critique through Self-](#)  
655 [Reflection](#). *arXiv preprint*. ArXiv:2310.11511 [cs].

656 Jiale Deng, Yanyan Shen, Ziyuan Pei, Youmin Chen,  
657 and Linpeng Huang. 2025. [Influence Guided Context](#)  
658 [Selection for Effective Retrieval-Augmented Genera-](#)  
659 [tion](#). *arXiv preprint*. ArXiv:2509.21359 [cs].

660 Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michi-  
661 hiro Yasunaga, and Yu Su. 2025. [HippoRAG:](#)  
662 [Neurobiologically Inspired Long-Term Memory](#)  
663 [for Large Language Models](#). *arXiv preprint*.  
664 ArXiv:2405.14831 [cs] TLDR: HippoRAG is intro-  
665 duced, a novel retrieval framework inspired by the  
666 hippocampal indexing theory of human long-term

667 memory to enable deeper and more efficient knowl-  
668 edge integration over new experiences and can tackle  
669 new types of scenarios that are out of reach of exist-  
670 ing methods.

671 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat,  
672 and Ming-Wei Chang. 2020. [REALM: Retrieval-](#)  
673 [Augmented Language Model Pre-Training](#). *arXiv*  
674 *preprint*. ArXiv:2002.08909 [cs].

675 Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara,  
676 and Akiko Aizawa. 2020. [Constructing a multi-](#)  
677 [hop QA dataset for comprehensive evaluation of](#)  
678 [reasoning steps](#). In *Proceedings of the 28th Interna-*  
679 *tional Conference on Computational Linguistics*,  
680 pages 6609–6625, Barcelona, Spain (Online). Interna-  
681 tional Committee on Computational Linguistics.

682 Gautier Izacard and Edouard Grave. 2020. [Leveraging](#)  
683 [Passage Retrieval with Generative Models for Open](#)  
684 [Domain Question Answering](#). *arXiv preprint*.

685 Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick  
686 Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and  
687 Wen-tau Yih. 2020. [Dense Passage Retrieval for](#)  
688 [Open-Domain Question Answering](#). *arXiv preprint*.

689 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio  
690 Petroni, Vladimir Karpukhin, Naman Goyal, Hein-  
691 rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-  
692 täschel, Sebastian Riedel, and Douwe Kiela. 2020.  
693 [Retrieval-Augmented Generation for Knowledge-](#)  
694 [Intensive NLP Tasks](#). *arXiv preprint*.

695 Yifei Li, Xiang Yue, Zeyi Liao, and Huan Sun. 2024.  
696 [AttributionBench: How Hard is Automatic Attribu-](#)  
697 [tion Evaluation?](#) *arXiv preprint*. ArXiv:2402.15089  
698 [cs].

699 Songshuo Lu, Hua Wang, Yutian Rong, Zhi Chen,  
700 and Yaohua Tang. 2025. [TurboRAG: Accelerating](#)  
701 [Retrieval-Augmented Generation with Precomputed](#)  
702 [KV Caches for Chunked Text](#). In *Proceedings of*  
703 *the 2025 Conference on Empirical Methods in Natu-*  
704 *ral Language Processing*, pages 6599–6612, Suzhou,  
705 China. Association for Computational Linguistics.

706 Yuren Mao, Xuemei Dong, Wenyi Xu, Yunjun Gao, Bin  
707 Wei, and Ying Zhang. 2025. [FIT-RAG: Black-Box](#)  
708 [RAG with Factual Information and Token Reduc-](#)  
709 [tion](#). *ACM Transactions on Information Systems*,  
710 43(2):1–27. TLDR: A novel black-box RAG frame-  
711 work which utilizes the factual information in the  
712 retrieval and reduces the number of tokens for aug-  
713 mentation, dubbed FIT-RAG, which achieves both  
714 superior effectiveness and efficiency.

715 Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike  
716 Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer,  
717 Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023.  
718 [FActScore: Fine-grained Atomic Evaluation of Fac-](#)  
719 [tual Precision in Long Form Text Generation](#). *arXiv*  
720 *preprint*. ArXiv:2305.14251 [cs].

721 Qwen Team. 2025. [Qwen3 technical report](#). *Preprint*,  
722 arXiv:2505.09388.

723	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	777
724	and Ashish Sabharwal. 2022. MuSiQue: Multi-	778
725	hop questions via single-hop question composition.	779
726	<i>Transactions of the Association for Computational</i>	780
727	<i>Linguistics</i> .	781
728	Jingjin Wang and Jiawei Han. 2025. PropRAG: Guiding	782
729	Retrieval with Beam Search over Proposition Paths.	783
730	In <i>Proceedings of the 2025 Conference on Empirical</i>	784
731	<i>Methods in Natural Language Processing</i> , pages	785
732	6223–6238, Suzhou, China. Association for Compu-	
733	tational Linguistics.	
734	Zihan Wang, Zihan Liang, Zhou Shao, Yufei Ma,	786
735	Huangyu Dai, Ben Chen, Lingtao Mao, Chenyi Lei,	787
736	Yuqing Ding, and Han Li. 2025. InfoGain-RAG:	788
737	Boosting Retrieval-Augmented Generation through	789
738	Document Information Gain-based Reranking and	790
739	Filtering. In <i>Proceedings of the 2025 Conference on</i>	791
740	<i>Empirical Methods in Natural Language Processing</i> ,	792
741	pages 7201–7215, Suzhou, China. Association for	793
742	Computational Linguistics.	
743	Wenhan Xiong, Xiang Lorraine Li, Srinivas Iyer, Jingfei	794
744	Du, Patrick Lewis, William Yang Wang, Yashar	
745	Mehdad, Wen-tau Yih, Sebastian Riedel, Douwe	
746	Kiela, and Barlas Oğuz. 2021. Answering Complex	
747	Open-Domain Questions with Multi-Hop Dense Re-	
748	trieval. <i>arXiv preprint</i> . ArXiv:2009.12756 [cs].	
749	Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling.	795
750	2024. Corrective Retrieval Augmented Generation.	796
751	<i>arXiv preprint</i> . ArXiv:2401.15884 [cs].	797
752	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-	798
753	gio, William W. Cohen, Ruslan Salakhutdinov, and	799
754	Christopher D. Manning. 2018. HotpotQA: A dataset	800
755	for diverse, explainable multi-hop question answer-	
756	ing. In <i>Conference on Empirical Methods in Natural</i>	
757	<i>Language Processing (EMNLP)</i> .	
758	Chen Zhao, Chenyan Xiong, Jordan Boyd-Graber, and	801
759	Hal Daumé. 2021. Multi-Step Reasoning Over Un-	802
760	structured Text with Beam Dense Retrieval. <i>arXiv</i>	803
761	<i>preprint</i> . ArXiv:2104.05883 [cs].	804
762	<b>A Appendix: Deferred System Details</b>	805
763	<b>A. Evaluation protocol, statistical treatment,</b>	806
764	<b>and latency definition</b>	807
765	<b>Confidence intervals.</b> The main paper reports	808
766	point estimates for readability. Here we report 95%	809
767	confidence intervals (CIs) for EM and token-level	810
768	F1. For EM (a Bernoulli metric), we use Wilson	811
769	score intervals. For F1, we use a stratified bootstrap	812
770	over queries ( $B=1000$ resamples) with a fixed ran-	813
771	dom seed. When comparing two methods on the	814
772	same dataset, we additionally report paired boot-	815
773	strap deltas (method A minus method B) with 95%	816
774	intervals. We include these intervals for auditabil-	817
775	ity rather than as a substitute for careful protocol	818
776	design.	819
777	<b>Latency definition and inclusion policy.</b> Lat-	
778	ency is measured end-to-end per query over the	
779	full pipeline: retrieve → rerank → verify/score	
780	→ select → generate (or abstain). All reported	
781	latency percentiles include abstentions. All meth-	
782	ods use the same batching configuration for re-	
783	trieval, reranking, and verification. We keep the	
784	caching policy fixed within a run and compute	
785	$\text{Lat}_{50}/\text{Lat}_{90}/\text{Lat}_{95}$ under that policy.	
786	<b>Token accounting.</b> We distinguish (i) evidence	
787	tokens (EvTok), the number of tokens contributed	
788	by the evidence passages passed to the generator,	
789	and (ii) total input tokens (TotTok), the full genera-	
790	tor input including the fixed prompt template and	
791	evidence. Both are computed using the generator	
792	tokenizer. Completion lengths are standardized via	
793	a shared decoding configuration across all methods.	
794	<b>B Baseline fairness and budget enforcement</b>	
795	<b>Evidence-token budget.</b> A budget of $B$ evidence	
796	tokens denotes a hard cap on the concatenated evi-	
797	dence text passed to the generator after a method’s	
798	final selection step. Any method that constructs an	
799	explicit evidence context is forced to respect this	
800	cap.	
801	<b>Deterministic truncation policy.</b> If the selected	
802	evidence exceeds $B$ , we truncate deterministically:	
803	passages are concatenated in the method’s final	
804	priority order until the cap is reached. The final	
805	passage is truncated at the exact token boundary.	
806	This enforces budget compliance without altering	
807	the method’s internal ranking logic.	
808	<b>Generator and candidate parity.</b> Unless oth-	
809	erwise stated, head-to-head comparisons use the	
810	same generator model and decoding parameters	
811	(temperature, top- $p$ , max_new_tokens) to isolate	
812	the impact of selection from generation capacity.	
813	When a baseline uses the shared retrieval stack, it	
814	receives the same candidate pool and reranker or-	
815	dering as the frontier regime. When a baseline (e.g.,	
816	HippoRAG) uses its own retrieval strategy, we treat	
817	it as a retrieval-strength reference but still enforce	
818	the same evidence-token cap. All runs log the final	
819	evidence list actually passed to the generator.	
820	<b>C Abstention definition and reporting</b>	
821	<b>Definition.</b> An abstention occurs when the sys-	
822	tem intentionally returns a dedicated token be-	
823	cause it cannot produce an evidence-conditioned	

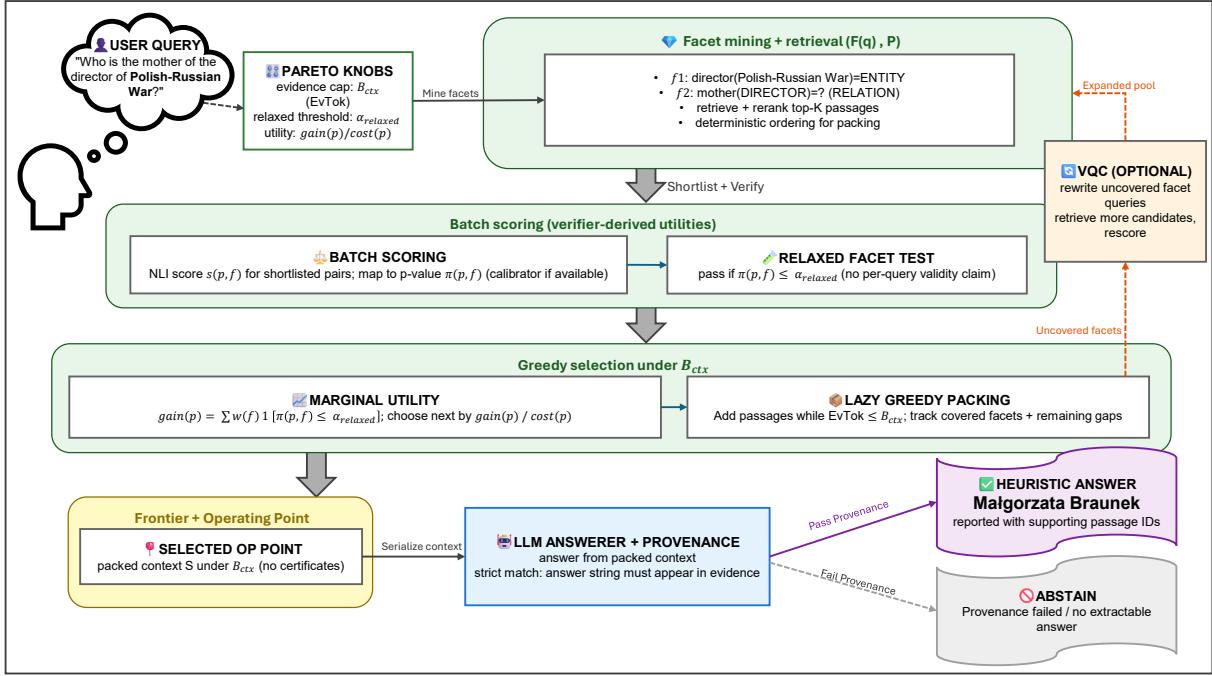


Figure 3: Pareto-Knapsack selection pipeline. Facets and candidates flow through batch scoring with relaxed thresholds ( $\alpha_{relaxed}$ ), greedy selection maximizes marginal utility per cost under the evidence budget  $B_{ctx}$ , and the LLM answerer enforces strict provenance checking. Optional VQC (Variable Query Completion) can retrieve additional candidates for uncovered facets.

824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835

answer under the configured constraints. Abstention rates are reported as a percentage of evaluated queries. In the certified regime, abstentions carry specific reason codes: NO\_COVER (some required facet has no passing passage), INFEASIBLE\_BUDGET (a valid cover exists but exceeds the cap), or INFEASIBLE\_PVALUE (calibration bins too sparse to support the required threshold).

**Counts at the Pareto-500 operating point.** Table 4 details the raw abstention counts for the primary operating point reported in the main paper.

Dataset	Method	Abstained / N	Rate (%)
HotpotQA	Pareto-500	114 / 7405	1.54
2WikiMultiHopQA	Pareto-500	266 / 12576	2.12
MuSiQue	Pareto-500	146 / 2417	6.04

Table 4: Abstention rate summary at the Pareto-500 operating point.

## D. Episode contract and replay requirements

836  
837  
838  
839  
840  
841

**Contract elements.** Certificate validity is claimed only under a replayed selection contract, as described in Section 3.4. The contract includes, at minimum: retriever snapshot and candidate cap; deterministic shortlisting rules and tie-breaks

(including  $T_f$ ); verifier version and scoring configuration; binning key and bin-merging policy; calibration pools (or their hashes); evidence serialization format and tokenization; and the evidence budget  $B_{ctx}$ . All contract elements are logged as version hashes plus runtime knobs.

**Failure behavior.** If any contract hash mismatches (e.g., index snapshot changes, verifier weights change), certificates are declared invalid and the certified regime is disabled. The system fails closed—falling back to abstention or non-certified operation—rather than emitting certificates under unknown conditions. This keeps the guarantee honest.

## E. Shift monitoring and conservative fallback

856  
857  
858  
859  
860  
861

**Drift signals.** We log per-bin score summaries and near-threshold counts to detect distributional shift in verifier scores. This monitoring is treated as an operational safeguard rather than a formal proof obligation within the core guarantee.

862  
863  
864  
865  
866

**Fallback policy.** If drift alarms trigger, the certified regime suspends certificate emission until re-calibration is performed under the updated contract. During this window, the system may continue to answer in the non-certified regime but does not

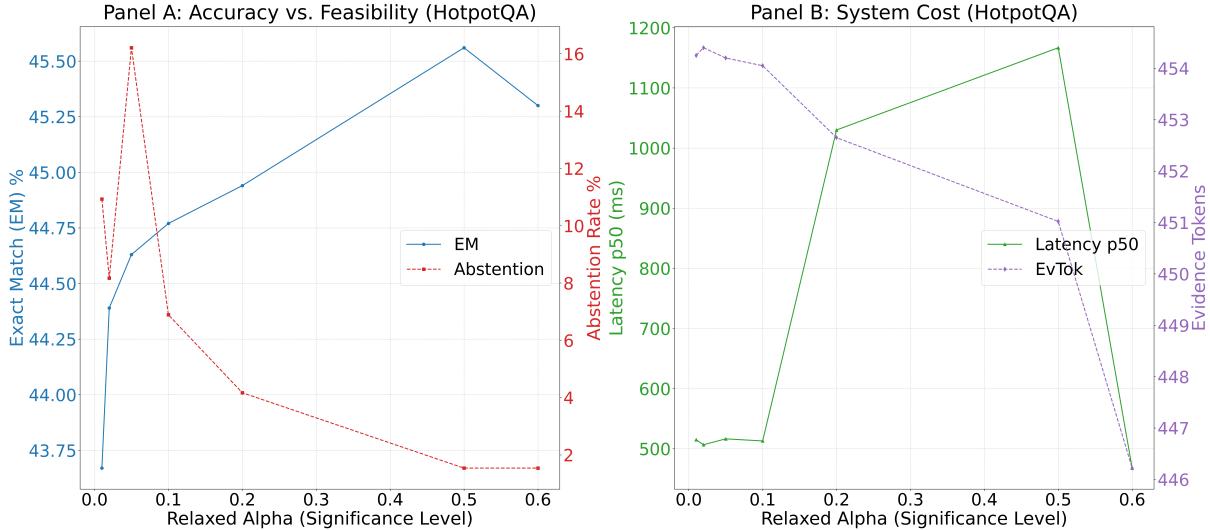


Figure 4: Pareto threshold sensitivity on HotpotQ.

claim the certified error bound.

## F. Deferred components

**Verifier-driven query rewriting.** As shown in Figure 3, typed rewrites can improve recall by retrieving additional candidates for uncovered facets. However, because rewriting changes the candidate distribution—and thus invalidates calibration conditions—we restrict rewriting to non-certified operation in the main protocol. In Pareto mode, VQC operates as an optional expansion layer that rewrites uncovered facet queries, retrieves more candidates, and rescores them before final selection.

**Long-run budget control.** A lightweight controller can manage token budgets across a stream of queries (e.g., amortization across sessions), but this is orthogonal to the per-query certification focus of this paper.

## G. Label noise sensitivity

If a fraction  $\epsilon$  of calibration negatives are mislabeled, conformal p-values become more conservative than intended. We optionally include a sensitivity analysis that inflates denominators to model this noise, reporting trends across  $\epsilon \in \{0, 0.01, 0.05, 0.10\}$ . This analysis aids deployment planning but is not required for the core validity claims.

## H. Relaxed- $\alpha$ sweep: feasibility vs. accuracy

Figure 4 analyzes the trade-off between strictness and feasibility on HotpotQA by varying the relaxed acceptance threshold  $\alpha$ .

Panel A shows that the dominant effect of tightening  $\alpha$  is on feasibility: abstention rates rise sharply as the system struggles to assemble a budget-feasible cover, while EM changes more gradually.

Panel B confirms that evidence usage (EvTok) remains essentially flat across values of  $\alpha$ . The changes in performance are not driven by "spending" more tokens—the evidence cap is binding and stable. Latency variations instead reflect pipeline dynamics, specifically how easily the greedy solver finds passing support under stricter thresholds.

867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896