

Q2EI: Query-to-Entity Inference for Semantic Condensation in Domain-Specific Retrieval

Anonymous ACL submission

Abstract

Retrieval-Augmented Generation (RAG) remains unreliable in specialized domains due to semantic and lexical mismatch between lay queries and professional terminology, and existing generative expansion often introduces redundancy or hallucinations that cause semantic drift. We propose Generative Query Condensation (GQC), a query rewriting strategy that reframes rewriting as semantic condensation rather than expansion. To operationalize GQC, we introduce Query-to-Entity Inference (Q2EI), an entity-centric rewriting method that realizes semantic condensation through explicit inference of the underlying target entity. By moving semantic alignment from retrieval-time vector matching to the rewriting stage, Q2EI produces information-dense query representations. Experimental results on medical and legal benchmarks show that Q2EI consistently outperforms strong baselines across retrievers, improving retrieval effectiveness while substantially reducing rewriting token consumption compared to generative expansion methods. Further analysis confirms that these gains primarily arise from accurate entity inference, and that Q2EI’s semantic condensation design limits error amplification when inference is imperfect, leading to more stable and interpretable retrieval behavior¹.

1 Introduction

Retrieval-Augmented Generation (RAG) alleviates knowledge hallucination and temporal obsolescence in Large Language Models (LLMs) by grounding generation in external non-parametric knowledge, and has achieved strong performance on general tasks such as open-domain question answering (Lewis et al., 2020; Guu et al., 2020).

¹We release code and data in an anonymized GitHub repository at <https://anonymous.4open.science/r/Q2EI-EFE0>

However, retrieval still faces significant bottlenecks in specialized domains such as medicine, law, and IT operations. A key challenge is that user queries often reside in a non-professional semantic space, whereas documents are written in a professional one. This semantic misalignment is difficult to resolve by vector retrieval alone: retrievers primarily capture surface-level semantic similarity, but cannot reliably infer the underlying domain concept implied by a lay description. We refer to such non-expert user queries as *lay queries*. This representation mismatch leads to lexical mismatch and semantic gaps, thereby compromising retrieval effectiveness (see Figure 1 (a)).

Specifically, due to a lack of domain knowledge, user queries are often phenomenon-level descriptions. However, the documents use precise, domain-specific-level terminology. For instance, a patient might describe symptoms as “*Ate home-canned food, now blurred vision, droopy eyelids, hard to swallow—could this be deadly food poisoning?*” whereas the relevant document is indexed under the term “*Botulism.*” On the one hand, this representation misalignment makes it difficult for sparse retrievers (e.g., TF-IDF (Salton and Buckley, 1988) and BM25 (Robertson et al., 2009)) to retrieve target documents due to their reliance on lexical overlap. On the other hand, dense retrievers (such as Multilingual E5 (mE5) (Wang et al., 2024), BGE-M3 (Multi-Granularity, 2024), or Contriever (Izacard et al., 2021)) also struggle to establish a semantic mapping from phenomenon-level descriptions to professional terminology.

To alleviate these issues, existing works predominantly adopt Generative Query Expansion (GQE). Generation-Augmented Retrieval (GAR) (Mao et al., 2021) and its successors, such as HyDE (Gao et al., 2023) and Query2Doc (Q2D) (Wang et al., 2023), expand query semantics by generating pseudo-documents to improve

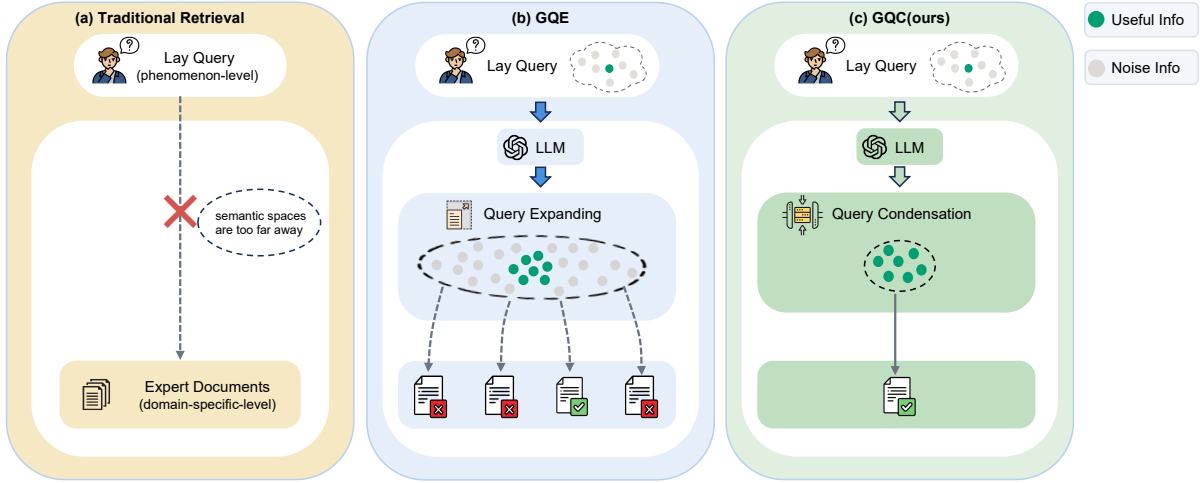


Figure 1: Schematic comparison of different retrieval strategies. (a) **Traditional Retrieval** suffers from significant lexical mismatch between lay queries and expert documents. (b) **GQE** bridges the gap but often introduces redundancy and noise, which may mislead the retriever toward irrelevant documents. (c) The proposed **GQC** (ours) employs an inference engine to condense semantics into high-density representations.

semantic coverage and lexical overlap (see Figure 1 (b)). However, our motivation experiments in section 5.1 indicate that in specialized domains, such verbose generated content tends to dilute information density and induce semantic drift. The result is also consistent with recent studies (Weller et al., 2024). Another category of research employs Knowledge Graph (KG) methods, utilizing explicit relational constraints to align semantics (Luo et al., 2023; Zhu et al., 2025; Mavromatis and Karypis, 2025). However, their high graph construction costs and dependence on graph coverage (Pan et al., 2024) limit their application.

In this paper, we propose the Generative Query Condensation (GQC) strategy (see Figure 1 (c)). Unlike conventional dense retrieval methods that rely on implicit vector matching during retrieval, this strategy leverages the parametric knowledge of LLMs to perform explicit semantic extraction and compression at the query rewriting stage. GQC compresses the semantic information within the query into a high-density representation, thereby effectively improving the signal-to-noise ratio. We introduce Query-to-Entity Inference (Q2EI) as a concrete implementation of GQC. Unlike GQE methods, Q2EI utilizes LLMs to condense ambiguous user descriptions into a core domain-specific entity and reconstructs an entity-centric query representation. This effectively mitigates both lexical and semantic mismatch.

Experimental results on two specialized-

domain benchmarks (MedQuAD and COLIEE) show that Q2EI yields substantial and consistent improvements across retrievers. On COLIEE, mE5 achieves an nDCG@10 improvement of 3.35 over the strongest baseline. On MedQuAD, BM25 yields a Recall@10 improvement of 13.97 compared with the best-performing baseline, while maintaining significantly lower computational overhead.

The main contributions of this paper are summarized as follows:

- We propose GQC as a query rewriting strategy for specialized domain retrieval. By prioritizing semantic condensation over text expansion, it alleviates lexical mismatch and semantic gaps. We further analyze the critical role of high-information-density representations in specialized domain retrieval.
- We introduce Q2EI, an entity inference method that requires no fine-tuning and does not rely on manually constructed knowledge bases. By guiding LLMs with specific inference instructions to infer potential entities from lay queries, we generate high signal-to-noise ratio entity-centric queries.
- Extensive experiments on multiple specialized domain datasets and various retriever architectures demonstrate that Q2EI significantly outperforms strong baseline methods. Further analysis attributes the gains

$$\max_f \text{Sim}(\text{Enc}(f(q_{\text{lay}})), \text{Enc}(d_{\text{target}})), \quad (1)$$

$$d_{\text{target}} \in \mathcal{D}_{\text{expert}}.$$

where $\text{Enc}(\cdot)$ denotes the retriever encoder, and $\text{Sim}(\cdot, \cdot)$ represents the similarity function.

3.2 Strategy Comparison: GQE vs. GQC

To implement the mapping f , existing works predominantly adopt Generative Query Expansion (GQE), whereas this paper employs the Generative Query Condensation (GQC) strategy.

GQE. Methods represented by HyDE (Gao et al., 2023) introduce an intermediate variable c_{gen} to approximate the semantic distribution of the target document by supplementing context. The mapping process can be expressed as:

$$f_{\text{GQE}}(q_{\text{lay}}) = \mathcal{T}(q_{\text{lay}}, c_{\text{gen}}), \quad (2)$$

$$c_{\text{gen}} \sim P_{\text{LLM}}(\cdot | q_{\text{lay}}; \mathcal{E}).$$

where \mathcal{T} denotes a text concatenation or fusion operation, and \mathcal{E} represents the set of expansion instructions used to guide the LLM in generating supplementary context c_{gen} to improve the semantic coverage of the query.

GQC (Semantic Condensation). GQC does not approximate the document distribution by supplementing redundant context; instead, it leverages the parametric knowledge of LLMs to perform semantic condensation on the query, outputting a rewritten query with high information density:

$$f_{\text{GQC}}(q_{\text{lay}}) = q_{\text{GQC}}, \quad (3)$$

$$q_{\text{GQC}} \sim P_{\text{LLM}}(\cdot | q_{\text{lay}}; \mathcal{C}).$$

where \mathcal{C} is the instruction set for semantic condensation, used to constrain the output to focus on retrievable core semantics. This process establishes a semantic bridge between ambiguous phenomenon-level descriptions and precise professional knowledge. Compared to GQE, GQC significantly reduces redundant information and improves the signal-to-noise ratio of the query.

This distinction is illustrated by the motivating example in Appendix A (Table 2). It can be observed that while GQE methods cover key semantics, the introduced redundancy and hallucinations can lead the retriever to incorrect documents. In contrast, the proposed condensation strategy (implemented as Q2EI, detailed in section 3.3) retains

only core entity semantics, thereby avoiding such retrieval errors caused by noise.

3.3 Q2EI: Entity-Centric Query Rewriting

Based on the GQC strategy, we propose Query-to-Entity Inference (Q2EI), an entity-centric query rewriting method, as shown in Figure 2. This method comprises three steps:

Step 1: Instruction Construction and Constraint Setting. We construct instructions $\mathcal{I}_{\text{prompt}}$ to guide the model’s inference and rewriting process. As shown in Figure 9 in Appendix B, the instructions include three key elements:

(1) *Persona-based Prompting:* Guides the LLM to act as a domain expert via specific role-play instructions (Xu et al., 2023; Kong et al., 2024), thereby enhancing the model’s ability to understand and parse domain semantics;

(2) *Entity-first inference with reformulation constraint:* Requires the model to first infer the normalized entity e_{core} and generate an entity-centric query q_{ent} based on this entity. This ensures that the rewritten query focuses on core entity semantics and aligns with domain-specific expression conventions;

(3) *Adaptive Demonstration:* Supports both zero-shot and few-shot settings, relying on the model’s internal knowledge and a small number of in-context mapping examples (Brown et al., 2020), respectively.

Step 2: Model Inference and Entity-Centric Rewriting. Under the instruction constraints constructed in Step 1, given q_{lay} and $\mathcal{I}_{\text{prompt}}$, the LLM infers the core entity within the query and normalizes it into an entity-centric query q_{ent} . Through entity rewriting, this process ensures the query closely adheres to the core semantics of q_{lay} .

Step 3: Entity-Centric Retrieval. We utilize q_{ent} for retrieval, transforming the match from “Phenomenon Description \rightarrow Professional Document” to “Entity-Centric Query \rightarrow Professional Document.” This effectively alleviates the semantic gap and improves the signal-to-noise ratio of the query.

4 Experiments

4.1 Experimental Setup

Datasets. We evaluate Q2EI in two specialized domains: medicine and law, utilizing the MedQuAD (Ben Abacha and Demner-Fushman, 2019) and COLIEE (Kim et al., 2022) datasets,

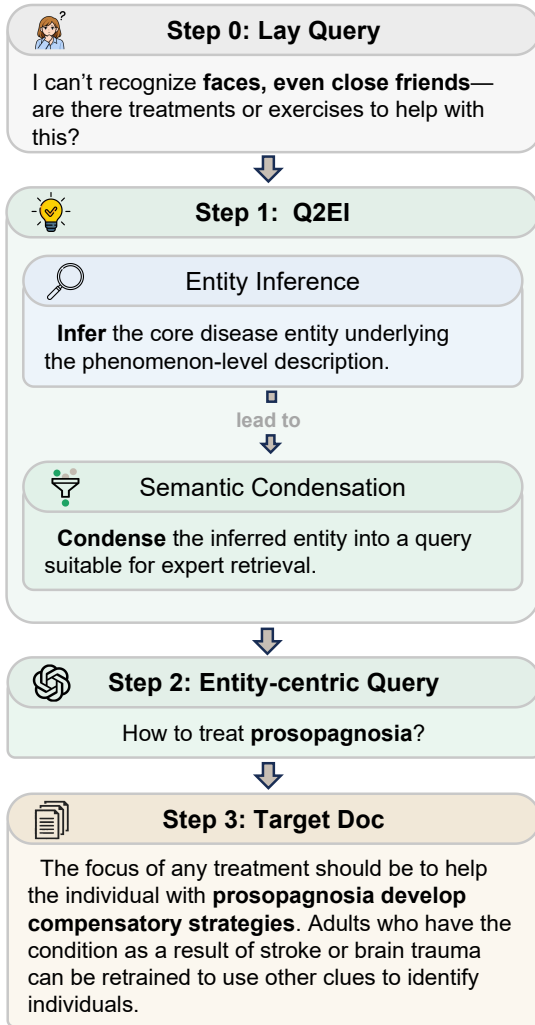


Figure 2: Overview of the Q2EI method.

327 respectively. Existing benchmarks typically use
 328 expert-written normalized questions as queries. To
 329 simulate the questioning style of lay users in
 330 specialized domains, we constructed lay queries:
 331 we utilized LLMs to rewrite original professional
 332 questions into lay descriptions of surface-level de-
 333 scriptions (for MedQuAD) or real-life dilemmas
 334 (for COLIEE). All rewritten queries underwent
 335 manual review by multiple reviewers to ensure the
 336 original semantic intent was maintained and pro-
 337 fessional terminology was removed. The specific
 338 prompt templates are provided in Appendix B.

339 **Baselines.** We compare Q2EI with the four
 340 methods: (1) **Lay Query:** Directly uses the lay
 341 query for retrieval. (2) **Keyword Extraction:** Uti-
 342 lizes an LLM to extract keywords from the query
 343 as retrieval input. This baseline is included to
 344 demonstrate that the performance gains of Q2EI
 345 stem from deep entity inference rather than shal-
 346 low lexical extraction. (3) **HyDE** (Gao et al.,

2023): A generative retrieval method that expands
 347 queries by generating hypothetical documents. (4)
 348 **Query2Doc (Q2D)** (Wang et al., 2023): Gen-
 349 erates pseudo-documents using few-shot prompt-
 350 ing to introduce contextual information. The spe-
 351 cific prompts employed for all the LLM-based
 352 baselines mentioned above are provided in Ap-
 353 pendix B. 354

355 **Implementation Details.** We primarily em-
 356 ploy GPT-5 (OpenAI, 2025) as the generative
 357 language model for the experiments reported in
 358 this section. To further verify the robustness of
 359 our method across different models, we conduct
 360 supplementary experiments using Claude Sonnet
 361 4.5 (Anthropic, 2025); these results are discussed
 362 in Appendix C. We evaluate performance across
 363 three retrievers, including the sparse retriever
 364 BM25 (Robertson et al., 2009) and the dense re-
 365 trievers Multilingual E5 (mE5) (Wang et al., 2024)
 366 and BGE-M3 (Multi-Granularity, 2024). Q2EI is
 367 tested under both zero-shot and few-shot settings.
 368 In the few-shot setting, we use the same examples
 369 as Q2D to ensure a fair comparison. All meth-
 370 ods are evaluated using Recall@1, Recall@10,
 371 and nDCG@10 metrics. All experiments are im-
 372 plemented by the open-source LlamaIndex frame-
 373 work (Liu, 2022).

4.2 Main Results 374

4.2.1 Performance on MedQuAD 375

376 Table 1 compares the retrieval accuracy of Q2EI
 377 and various baselines on the medical (MedQuAD)
 378 and legal (COLIEE) datasets. Overall, Q2EI
 379 demonstrates clear and consistent gains across do-
 380 mains and retriever architectures, with the largest
 381 improvements on dense retrievers—suggesting that
 382 entity-level semantic condensation is particularly
 383 effective for alignment in embedding-based re-
 384 trieval spaces.

385 Taking mE5 on MedQuAD as an example, the
 386 nDCG@10 scores for HyDE and Q2D are 18.29
 387 and 18.70, respectively, whereas Q2EI (zero-shot)
 388 improves this to 42.38, and Q2EI (few-shot) fur-
 389 ther achieves 43.13. These results indicate that,
 390 compared to query expansion methods relying on
 391 long text generation, entity-centric query rewriting
 392 aligns better with the semantic alignment require-
 393 ments in dense representation spaces. Compared
 394 with Keyword Extraction, Q2EI achieves consis-
 395 tently higher performance across retrievers. For
 396 example, on BGE-M3, Q2EI (few-shot) reaches an

Method	BM25			BGE-M3			mE5		
	R@1	R@10	N@10	R@1	R@10	N@10	R@1	R@10	N@10
MedQuAD									
Lay Query	7.35	21.88	13.94	7.17	31.25	17.94	5.70	24.82	14.71
HyDE	–	–	–	16.91	46.88	30.93	8.46	30.33	18.29
Q2D	13.97	38.05	24.77	18.01	49.08	32.35	7.90	31.25	18.70
Keyword Extraction	12.32	41.18	25.90	12.13	35.29	23.27	9.38	34.19	21.04
Q2EI (zero-shot)	14.34	49.63	30.90	22.06	55.51	38.39	24.82	60.85	42.38
Q2EI (few-shot)	13.79	52.02	31.31	23.16	57.17	39.71	25.92	61.40	43.13
COLIEE									
Lay Query	2.59	9.43	6.00	7.02	16.08	11.12	6.47	14.23	10.23
HyDE	–	–	–	6.28	22.55	15.62	6.65	15.71	11.26
Q2D	12.94	27.91	20.53	11.28	27.54	19.80	8.32	22.00	15.77
Keyword Extraction	5.55	13.68	9.51	5.91	14.60	9.73	6.84	15.16	11.11
Q2EI (zero-shot)	8.13	22.74	15.44	10.91	23.66	17.58	9.24	19.59	14.36
Q2EI (few-shot)	13.49	28.10	21.01	13.49	28.10	20.67	12.20	25.69	19.12

Table 1: Retrieval performance (%) comparison of Q2EI and baseline methods on medical (MedQuAD) and legal (COLIEE) domain datasets. Note: R@k stands for Recall@k, N@k stands for nDCG@k.

nDCG@10 of 39.71, substantially higher than the 23.27 achieved by Keyword Extraction, indicating that the gains cannot be explained by shallow lexical cues alone but instead stem from explicit entity inference.

4.2.2 Performance on COLIEE

In the legal domain, the few-shot setting proves particularly critical for performance improvement. On the COLIEE dataset, after introducing the same number of examples, Q2EI (few-shot) achieves the best results across all retrievers. For instance, on mE5, its nDCG@10 improves from 15.77 (Q2D) to 19.12. This indicates that, under equal utilization of in-context learning, semantic condensation-based query rewriting exhibits more stable retrieval advantages in cross-domain scenarios.

5 Analysis

5.1 Information Density Hypothesis: Redundancy Induces Semantic Drift

To verify the information density hypothesis, we designed a controlled experiment. Using 500 original queries from the MedQuAD test set as a baseline, we employed an LLM to generate query-relevant pathological descriptions of target lengths ($k \in \{100, 200, 300\}$ words) (approximately k words each). These descriptions are appended to the original queries (see Appendix D for prompts). All experiments are evaluated using the Multilingual E5 (mE5) retriever.

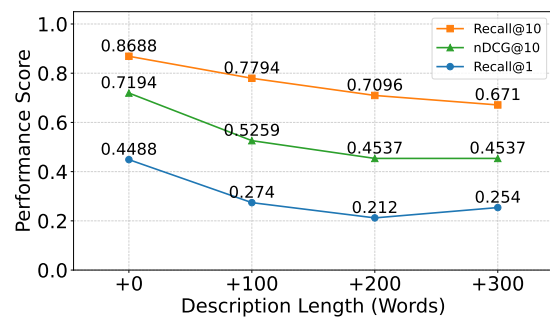


Figure 3: Validation of Information Density Hypothesis: Trend analysis of the impact of appending the generative content description length on retrieval performance.

As shown in Figure 3, retrieval performance exhibits a decline as the length of the appended description increases. Compared to the original query, appending just 100 words causes Recall@10 to drop from 86.88 to 77.94; when the length increases to 300 words, Recall@10 further degrades to 67.10. These results indicate that retrieval effectiveness can weaken even if the generated content is semantically relevant. Appending the generated content induces a shift in the embedding centroid. This shift ultimately leads to the degradation of retrieval performance.

5.2 Attribution Analysis: Impact of Entity Inference Accuracy

To analyze the relationship between retrieval performance gains and the accuracy of entity in-

ference, we divided the evaluation queries in the MedQuAD test set into two mutually exclusive subsets: *Entity-Aligned Group*, where the inferred entity holds an equivalence (same, hyponym, or hypernym) relationship with the target entity of the original question; and *Entity-Misaligned Group*, where this relationship is not met. In the zero-shot setting, *Entity-Aligned Group* accounts for 83.4 of queries, while in the few-shot setting, it accounts for 84.0. Notably, failed inference cases constitute only a small portion of the evaluation set; for completeness, we report the overall retrieval performance (Recall@1 and Recall@10 over all queries) in Appendix F.

We evaluated the retrieval performance of both groups on BM25, BGE-M3, and mE5. As shown in Figure 4 (a), the *Entity-Aligned Group* consistently outperforms the *Entity-Misaligned Group* by large absolute margins across retriever architectures. Taking mE5 in the zero-shot setting as an illustrative example, the nDCG@10 of the *Entity-Aligned Group* reaches 52.12, whereas the *Entity-Misaligned Group* achieves 2.72, yielding an absolute gap of 49.40 points. Comparable absolute gaps are observed on BM25 (38.42 vs. 5.60) and BGE-M3 (48.24 vs. 9.16). These results demonstrate that Q2EI’s retrieval effectiveness is strongly correlated with entity inference accuracy. When entity inference is misaligned with the target concept, retrieval performance drops sharply, indicating that accurate entity inference constitutes a critical prerequisite for effective retrieval within the Q2EI framework.

5.3 Misalignment Analysis: Entity Misalignment Reflects Query-Intrinsic Hardness

To characterize retrieval behavior under entity misalignment, we compare all evaluated retrieval methods (Lay Query, HyDE, and Q2D) using the same group partition. We focus on nDCG@10 in this section; additional metrics and full experimental results are reported in Appendix E.

As shown in Figure 4 (b), all methods yield consistently low nDCG@10 across BM25, BGE-M3, and mE5. Specifically, in the Entity-Misaligned Group, nDCG@10 ranges from 2.72 to 9.10 across all methods and retrievers, whereas in the Entity-Aligned Group it spans 15.96 to 38.92. This uniform degradation across methods suggests that performance bottlenecks under entity misalignment are largely method-agnostic and instead

stem from query-intrinsic issues: the phenomenon-level descriptions are often vague, incomplete, or even incorrect, making it difficult-and in some cases highly unreliable-to infer the correct underlying condition from the described manifestations alone.

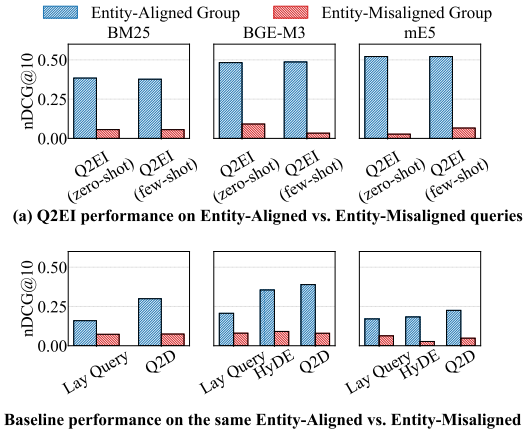


Figure 4: Retrieval performance (nDCG@10) under an entity-alignment split. (a) Performance comparison of Q2EI (zero-shot and few-shot) in Entity-Aligned vs. Entity-Misaligned Groups. (b) Performance of baseline methods (Lay Query, HyDE, Q2D) under the same grouping.

In the Entity-Aligned Group, Q2EI consistently achieves the best (or tied-best) retrieval performance across retriever architectures. For example, on mE5, Q2EI (few-shot) reaches an nDCG@10 of 52.11, substantially outperforming HyDE (18.34) and Q2D (22.53) under the same retriever. Comparable margins are observed on BM25 and BGE-M3. This pattern suggests that in the Entity-Aligned Group, Q2EI infers an entity aligned with the target and condenses the query into its core domain semantics, producing a high-information-density representation. By contrast, HyDE and Q2D generate substantially longer expansions; while these expansions may include relevant content, they often contain redundant information that lowers the effective signal-to-noise ratio for retrieval.

Notably, Q2EI remains competitive in the Entity-Misaligned Group, which represents a particularly challenging regime where the correct condition is often hard to disambiguate from the query alone. In this setting, performance differences mainly reflect how much a rewriting method drifts when its hypothesis is off-target. Even when Q2EI infers an incorrect entity, the predicted condition may still share overlapping symptom clusters or

526 manifestations with the true target, leading to a
527 comparatively mild semantic mismatch. In con-
528 trast, HyDE and Q2D tend to generate verbose
529 expansions under misalignment, which can intro-
530 duce additional off-target details and increase the
531 risk of semantic drift. As a result, Q2EI does not
532 exhibit disproportionate degradation and is often
533 among the best-performing methods in this group;
534 for instance, on mE5, Q2EI (few-shot) achieves
535 an nDCG@10 of 6.61, compared to 4.82 for Q2D.
536 Overall, these findings suggest that semantic con-
537 densation is relatively failure-tolerant: when infer-
538 ence is imperfect, compact, high-density rewrites
539 help limit error amplification and yield more sta-
540 ble retrieval behavior.

541 **5.4 Efficiency and Deployment Cost: Token** 542 **Consumption Comparison**

543 To compare the computational overhead of differ-
544 ent methods during the query rewriting stage, we
545 measure the average token consumption per query
546 on the MedQuAD dataset. This metric accounts
547 for the entire rewriting process, including the in-
548 put prompt, intermediate reasoning tokens gener-
549 ated by the LLM, and the final rewritten output.

550 The results reveal substantial differences in gener-
551 ation cost across methods. HyDE incurs the
552 highest overhead, consuming 13,408.55 tokens per
553 query on average due to the need to generate mul-
554 tiple hypothetical documents. Q2D reduces this
555 cost to 2,767.93 tokens per query, but still requires
556 generating relatively long pseudo-contexts. In con-
557 trast, Q2EI dramatically lowers generation cost:
558 the zero-shot variant consumes only 905.21 tokens
559 per query, corresponding to 6.8% of HyDE’s to-
560 ken usage, while the few-shot variant consumes
561 1,718.34 tokens per query, or 12.8% of HyDE’s
562 cost.

563 These reductions directly reflect the design prin-
564 ciple of the GQC strategy. Rather than produc-
565 ing verbose textual expansions, Q2EI infers a com-
566 pact, entity-centric representation that preserves
567 core semantics while minimizing redundant gener-
568 ation. As a result, Q2EI substantially reduces
569 deployment cost at the query rewriting stage, offer-
570 ing significantly improved efficiency without sac-
571 rificing retrieval effectiveness.

572 **6 Conclusion**

573 To address lexical mismatch and semantic gaps in
574 specialized-domain retrieval, this paper proposes

575 the GQC strategy and introduces Q2EI as its con-
576 crete instantiation. Unlike the additive paradigm
577 of GQE, GQC adopts a subtractive approach
578 via semantic condensation. Q2EI operationalizes
579 GQC through explicit entity inference. It maps
580 users’ phenomenon-level descriptions to a core do-
581 main entity and rewrites them into entity-centric
582 queries. This reasoning-driven condensation con-
583 structs a semantic bridge between user intent and
584 domain knowledge, improving the query’s signal-
585 to-noise ratio.

586 Experimental results demonstrate that Q2EI sig-
587 nificantly outperforms existing baselines across
588 multiple datasets and diverse retriever architec-
589 tures. Our in-depth analysis confirms four key ob-
590 servations: (1) redundant generated content tends
591 to induce semantic drift; (2) the performance gains
592 of Q2EI are attributed to accurate entity infer-
593 ence; (3) Q2EI exhibits failure-tolerant behavior-
594 when entity inference is imperfect, its semantic
595 condensation design limits error amplification and
596 leads to more stable retrieval behavior compared
597 to generative expansion methods; and (4) Q2EI
598 achieves these improvements with significantly re-
599 duced computational overhead.

600 In summary, Q2EI offers an efficient and cost-
601 effective query rewriting strategy for specialized
602 domain retrieval. This approach demonstrates ex-
603 ceptional robustness and exhibits significant poten-
604 tial practical value in professional scenarios such
605 as medicine and law.

606 **Limitations**

607 While Q2EI demonstrates superior performance in
608 specialized domain retrieval tasks, we acknowl-
609 edge the following limitations, which also illumi-
610 nate directions for future research.

611 **Inference Latency & Cost:** Although the com-
612 putational overhead of Q2EI is significantly lower
613 than that of GQE, on-the-fly inference inevitably
614 introduces higher latency compared to traditional
615 sparse or dense retrieval methods. For industrial-
616 grade deployment scenarios with strict real-time
617 requirements, the current inference cost remains
618 a primary bottleneck. Fortunately, this bottle-
619 neck is likely to ease over time, as hardware cost-
620 efficiency continues to improve and model com-
621 pression/acceleration techniques increasingly en-
622 able faster and cheaper LLM inference in produc-
623 tion.

624 **Dependency on LLM Capabilities:** The per-

625 formance of Q2EI is constrained by the paramet-
 626 ric domain knowledge of the underlying LLM.
 627 When handling extremely ambiguous queries or
 628 those involving highly long-tail knowledge, incor-
 629 rect entity inference can trigger a cascading effect,
 630 directly propagating errors to the retrieval stage
 631 and degrading retrieval effectiveness. Conversely,
 632 Q2EI is expected to benefit directly from ongo-
 633 ing advances in foundation models: as LLMs be-
 634 come more capable and better grounded in domain
 635 knowledge, the quality of entity inference and thus
 636 retrieval performance should improve accordingly.

637 **Domain & Task Applicability:** The core
 638 value of Q2EI lies in realizing the inference from
 639 “phenomenon-level descriptions” to “core entity.”
 640 In simple matching tasks that do not require com-
 641 plex reasoning or in non-entity-centric retrieval
 642 scenarios, the marginal utility of this method may
 643 diminish.

644 References

645 Anthropic. 2025. [Claude sonnet 4.5 system card](#). API
 646 version: cclaude-sonnet-4-5-20250929.

647 Asma Ben Abacha and Dina Demner-Fushman. 2019.
 648 A question-entailment approach to question answer-
 649 ing. *BMC bioinformatics*, 20(1):511.

650 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
 651 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
 652 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
 653 Askell, and 1 others. 2020. Language models are
 654 few-shot learners. *Advances in neural information
 655 processing systems*, 33:1877–1901.

656 Ilias Chalkidis, Manos Fergadiotis, Prodromos
 657 Malakasiotis, Nikolaos Aletras, and Ion Androut-
 658 sopoulos. 2020. Legal-bert: The muppets straight
 659 out of law school. *arXiv preprint arXiv:2010.02559*.

660 Linyi Ding, Sizhe Zhou, Jinfeng Xiao, and Ji-
 661 awei Han. 2024. Automated construction of
 662 theme-specific knowledge graphs. *arXiv preprint
 663 arXiv:2404.19146*.

664 Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie
 665 Callan. 2023. Precise zero-shot dense retrieval with-
 666 out relevance labels. In *Proceedings of the 61st An-
 667 nual Meeting of the Association for Computational
 668 Linguistics (Volume 1: Long Papers)*, pages 1762–
 669 1777.

670 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pa-
 671 supat, and Mingwei Chang. 2020. Retrieval aug-
 672 mented language model pre-training. In *Internat-
 673 ional conference on machine learning*, pages 3929–
 674 3938. PMLR.

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Se-
 675 bastian Riedel, Piotr Bojanowski, Armand Joulin,
 676 and Edouard Grave. 2021. Unsupervised dense in-
 677 formation retrieval with contrastive learning. *arXiv
 678 preprint arXiv:2112.09118*. 679

Vitor Jeronimo, Luiz Bonifacio, Hugo Abonizio,
 680 Marzieh Fadaee, Roberto Lotufo, Jakub Zavrel, and
 681 Rodrigo Nogueira. 2023. Inpars-v2: Large language
 682 models as efficient dataset generators for informa-
 683 tion retrieval. *arXiv preprint arXiv:2301.01820*. 684

Vladimir Karpukhin, Barlas Oguz, Sewon Min,
 685 Patrick SH Lewis, Leda Wu, Sergey Edunov, Danqi
 686 Chen, and Wen-tau Yih. 2020. Dense passage re-
 687 trieval for open-domain question answering. In
 688 *EMNLP (1)*, pages 6769–6781. 689

Mi-Young Kim, Juliano Rabelo, Randy Goebel, Masa-
 690 haru Yoshioka, Yoshinobu Kano, and Ken Satoh.
 691 2022. Coliee 2022 summary: Methods for legal
 692 document retrieval and entailment. In *Jsai inter-
 693 national symposium on artificial intelligence*, pages
 694 51–67. Springer. 695

Aobo Kong, Shiwan Zhao, Hao Chen, Qicheng Li,
 696 Yong Qin, Ruiqi Sun, Xin Zhou, Enzhi Wang, and
 697 Xiaohang Dong. 2024. Better zero-shot reasoning
 698 with role-play prompting. In *Proceedings of the
 699 2024 Conference of the North American Chapter of
 700 the Association for Computational Linguistics: Hu-
 701 man Language Technologies (Volume 1: Long Pa-
 702 pers)*, pages 4099–4113. 703

Victor Lavrenko and W Bruce Croft. 2003. Rele-
 704 vance models in information retrieval. In *Language
 705 modeling for information retrieval*, pages 11–56.
 706 Springer. 707

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio
 708 Petroni, Vladimir Karpukhin, Naman Goyal, Hein-
 709 rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-
 710 täschel, and 1 others. 2020. Retrieval-augmented
 711 generation for knowledge-intensive nlp tasks. *Ad-
 712 vances in neural information processing systems*,
 713 33:9459–9474. 714

Jerry Liu. 2022. [Llamaindex](#). 715

Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and
 716 Shirui Pan. 2023. Reasoning on graphs: Faithful and
 717 interpretable large language model reasoning. *arXiv
 718 preprint arXiv:2310.01061*. 719

Yuning Mao, Pengcheng He, Xiaodong Liu, Yelong
 720 Shen, Jianfeng Gao, Jiawei Han, and Weizhu Chen.
 721 2021. Generation-augmented retrieval for open-
 722 domain question answering. In *Proceedings of the
 723 59th Annual Meeting of the Association for Compu-
 724 tational Linguistics and the 11th International Joint
 725 Conference on Natural Language Processing (Vol-
 726 ume 1: Long Papers)*, pages 4089–4100. 727

Gary Marcus. 2018. Deep learning: A critical ap-
 728 praisal. *arXiv preprint arXiv:1801.00631*. 729

Target Passage	Laboratory tests that are used to diagnose CCHF include antigen-capture enzyme-linked immunosorbent assay (ELISA), real time polymerase chain reaction (RT-PCR), virus isolation attempts... Laboratory diagnosis... by using the combination of detection of the viral antigen... Later in the course... antibodies can be found in the blood...
Query	After forest hike, tick bite, now high fever, bad headache, body aches, vomiting, nosebleeds. Should I worry?
Q2D	Tick-borne infections can cause abrupt high fever... including Rocky Mountain spotted fever (RMSF), ehrlichiosis/anaplasmosis, Colorado tick fever... Nosebleeds (epistaxis), easy bruising... suggest thrombocytopenia... seen in severe RMSF... Doxycycline is first-line...
Q2EI	What is Crimean-Congo hemorrhagic fever (CCHF) in a patient with a recent tick bite, sudden high fever, severe myalgias, and epistaxis/easy bruising?
Wrong passage	The first symptoms of Rocky Mountain spotted fever (RMSF)... spotted (petechial) rash of RMSF is usually not seen... Rickettsia rickettsii infects the endothelial cells... can be performed on a skin biopsy... Recommended Dosage Doxycycline is the first line treatment... such as ehrlichiosis and anaplasmosis...

Table 2: Case study from the MedQuAD dataset. Comparisons between Q2D and our Q2EI. Yellow indicates parts similar to the target passage, Pink marks "distractors" that induced retrieval errors; Blue represents redundant context.

Method	BM25			BGE-M3			mE5		
	R@1	R@10	N@10	R@1	R@10	N@10	R@1	R@10	N@10
Lay query	7.35	21.88	13.94	7.17	31.25	17.94	5.70	24.82	14.71
Q2D	12.50	47.92	29.72	25.00	52.08	37.26	12.50	37.50	25.06
Q2EI (zero-shot)	14.58	62.50	35.89	45.83	68.70	57.82	39.58	72.92	55.93
Q2EI (few-shot)	12.50	56.25	31.75	29.17	64.58	45.04	29.17	68.75	48.39

Table 3: Retrieval performance (%) comparison. Note: All values are presented as percentages.

B Prompt Design

We detail the specific prompts used in our experiments.

Lay Query Generation. Figure 5 shows the prompt used to rewrite professional questions into lay queries. The rewritten queries simulate the questioning style of non-expert users in specialized domains. This process removes domain-specific terminology while preserving the original intent.

Keyword Extraction. Figure 6 presents the prompt used for the Keyword Extraction baseline. The prompt instructs the model to extract salient keywords from the query. This baseline demonstrates that the performance gains of Q2EI stem from deep entity inference rather than shallow lexical extraction.

HyDE. Figure 7 shows the prompt used by the HyDE method. For cost efficiency, we generate four hypothetical documents per query. We average their embeddings with the lay query embedding for retrieval.

Q2D. Figure 8 shows the prompt used for Q2D. Following the original implementation, we use few-shot prompting to generate pseudo-passages. For each query, we randomly select three pairs of lay queries and target passages from the dataset as in-context demonstrations.

Q2EI. Figure 9 presents the instruction prompt template used in Q2EI. The prompt include three key elements (i) a persona to place the model in a domain-expert role, (ii) an entity-first constraint that requires inferring and normalizing the most likely core entity from the lay query, and (iii) a format constraint that enforces a standard professional-question form and asks the model to output only the rewritten question. Optionally, a small set of in-context examples can be appended for the few-shot setting.

C Robustness to the Generative Backbone

We examine the robustness of Q2EI to the choice of the generative backbone. We repeat the query

Lay Query Generation Prompt

You are a simulation of an anxious patient with NO medical background. Your task is to generate a CASUAL, LAYMAN search query based ****ONLY**** on the medical condition mentioned in the provided Professional Question.

INSTRUCTIONS:

1. IDENTIFY THE DISEASE: Look at the Professional Question and identify the specific disease or condition.

2. USE INTERNAL KNOWLEDGE: Use your own internal knowledge to imagine how a regular person would describe the symptoms, causes, or transmission ****WITHOUT** knowing the medical name******.

3. SIMULATE THE SCENARIO: Write a query as if you are experiencing the issue, but don't know what it is.

4. STRICTLY NO ENTITIES: Do NOT use the disease name. Use vague descriptions (e.g., 'weird virus', 'stomach bug').

5. LENGTH: Keep it under 20 words.

Query:

Professional Question: {query}

Rewritten Question :

Figure 5: Prompt for rewriting professional questions into lay queries.

Keyword Extraction Prompt

You are a keyword extractor. Extract the most important keywords from the query.

Query:

Query: {query}

Rewritten Question :

Figure 6: Prompt used for Keyword Extraction from a lay query.

HyDE Prompt

Please write a passage to answer the question.

Query:

Query: {query}

Rewritten Question :

Figure 7: Prompt used by HyDE to generate a hypothetical passage for retrieval.

Q2D Prompt

Write a passage that answers the given query.

Examples (optional):

Query : {Example Query 1}

Output : {Example Output 1}

Query : {Example Query 2}

Output : {Example Output 2}

Query : {Example Query 3}

Output : {Example Output 3}

Query : {Example Query 4}

Output : {Example Output 4}

Query:

Query: {query}

Rewritten Question :

Figure 8: Few-shot prompt used by Q2D to generate a pseudo-document.

Q2EI Prompt

You are a Medical Search Specialist optimizing queries for a Medical Database.

1: INFER THE SPECIFIC ENTITY

You MUST infer the most likely Specific Disease Name or Parasite Name based on the transmission method or unique symptoms described.

2: STANDARDIZE FORMAT

Use standard question formats. Such as 'What is [Disease Name]?' or 'How to diagnose [Disease Name]?' or 'How to prevent [Disease Name]?'.

Note:

Output ONLY the rewritten professional question.

Examples (optional):

Query : {Example Query 1}

Output : {Example Output 1}

Query : {Example Query 2}

Output : {Example Output 2}

Query : {Example Query 3}

Output : {Example Output 3}

Query : {Example Query 4}

Output : {Example Output 4}

Query:

Query : {query}

Rewritten Question :

Figure 9: Prompt used by Q2EI to infer the core entity and rewrite a lay query into a professional, entity-centric query (MedQuAD example).

rewriting step using Claude Sonnet 4.5 (API: claude-sonnet-4-5-20250929).

Setup. We compare Q2EI (zero-shot and few-shot) with Lay Query and Q2D. For cost efficiency, we evaluate a random subset of 50 queries from MedQuAD. We exclude HyDE, which requires multiple hypothetical documents. This substantially increases token usage (Section 5.4). We treat this experiment as a small-scale robustness check.

Results. Results are reported in Table 3. The overall trend is consistent with results obtained using GPT-based backbones. Q2EI consistently outperforms Lay Query and Q2D across retrievers. For instance, on the Multilingual E5 (mE5) retriever, Q2EI (zero-shot) achieves a Recall@10 of 72.92, compared to 37.50 for Q2D and 24.82 for the lay query. Interestingly, in this setting, the zero-shot performance of Q2EI occasionally surpasses the few-shot setting. One possible explanation is the strong inference capability of this backbone. This may reduce the reliance on in-context examples. This observation aligns with the objective of semantic condensation.

D Controlled Redundancy Prompts

Short Description

```
Identify the main medical entity in the following query and write a 100-word definition of it. Do not answer the query itself.
```

```
Query:  
### Query: {query}  
### Rewritten Question :
```

Figure 10: Prompt for generating a short description to control redundancy.

Medium Description

```
Identify the main medical entity in the following query and write a 100-word definition of it. Do not answer the query itself.
```

```
Query:  
### Query: {query}  
### Rewritten Question :
```

Figure 11: Prompt for generating a medium-length description to control redundancy.

Long Description

```
Identify the main medical entity in the following query and write a 100-word definition of it. Do not answer the query itself.
```

```
Query:  
### Query: {query}  
### Rewritten Question :
```

Figure 12: Prompt for generating a long description to control redundancy.

Figures 10, 11, and 12 show the prompts used in our motivation analysis (Section 5.1). We control the description length to vary the amount of redundant information.

E Additional Results for Attribution Analysis

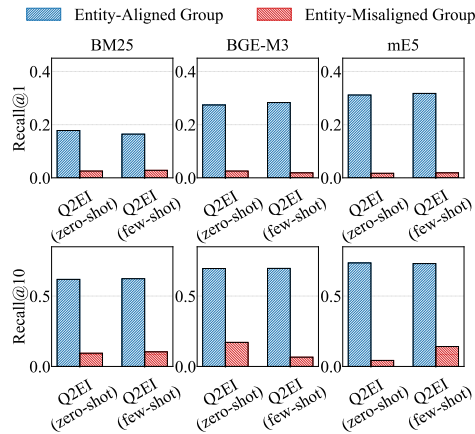
We provide additional metrics to support the attribution analysis. Figure 13 reports Recall@1 and Recall@10 under different retrievers and rewriting methods. Results are shown separately for the Entity-Aligned Group and the Entity-Misaligned Group. The observed trends are consistent with the nDCG@10 results reported in Section 5.2 and Section 5.3. Specifically, Q2EI achieves substantially higher recall when entity inference is correct. When inference fails, all methods exhibit uniformly low performance. These results further support the conclusion that retrieval gains mainly stem from accurate entity inference. They also indicate that Q2EI does not amplify failure cases.

F Entity Candidate Ablation

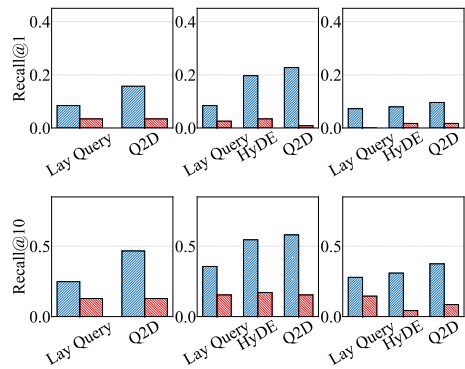
Number of Candidates (N)	Accuracy (%)
$N = 1$	83.4
$N = 2$	87.4
$N = 3$	89.0
$N = 4$	90.2

Table 4: Entity inference accuracy (%) vs. number of candidates N for Q2EI (zero-shot).

We analyze the impact of inferring multiple potential entities. Figures 14, 15, and 16 illustrate the prompts used to infer varying numbers of entities. Table 4 reports the inference accuracy. Specifically, it shows the probability that the set of inferred entities includes one that holds an equiva-



(a) Q2EI performance on Entity-Aligned vs. Entity-Misaligned queries



(b) Baseline performance on the same Entity-Aligned vs. Entity-Misaligned split

Figure 13: Retrieval performance (Recall@1 and Recall@10) under an entity-alignment split. (a) Performance comparison of Q2EI (zero-shot and few-shot) in Entity-Aligned vs. Entity-Misaligned Groups. (b) Performance of baseline methods (Lay Query, HyDE, Q2D) under the same grouping.

Two Entities Prompt

```

You are a Medical Search Specialist optimizing queries for a Medical Database.
### 1: INFER THE SPECIFIC ENTITY
You MUST infer the most likely **Specific Disease Name** or **Parasite Name** based on the transmission method or unique symptoms described. If you can't make sure, you can guess up to two possible diseases and use 'OR' to connect them.
### 2: STANDARDIZE FORMAT
Use standard question formats. Such as 'What is [Disease Name]?' or 'How to diagnose [Disease Name]?' or 'How to prevent [Disease Name]?'.
### Note:
Output ONLY the rewritten professional question.
Query:
### Query : {query}
### Rewritten Question :

```

Figure 14: Prompt for Q2EI with two potential entities candidates on MedQuAD.

Three Entities Prompt

```

You are a Medical Search Specialist optimizing queries for a Medical Database.
### 1: INFER THE SPECIFIC ENTITY
You MUST infer the most likely **Specific Disease Name** or **Parasite Name** based on the transmission method or unique symptoms described. If you can't make sure, you can guess up to three possible diseases and use 'OR' to connect them.
### 2: STANDARDIZE FORMAT
Use standard question formats. Such as 'What is [Disease Name]?' or 'How to diagnose [Disease Name]?' or 'How to prevent [Disease Name]?'.
### Note:
Output ONLY the rewritten professional question.
Query:
### Query : {query}
### Rewritten Question :

```

Figure 15: Prompt for Q2EI with three potential entities candidates on MedQuAD.

Four Entities Prompt

```

You are a Medical Search Specialist optimizing queries for a Medical Database.
### 1: INFER THE SPECIFIC ENTITY
You MUST infer the most likely **Specific Disease Name** or **Parasite Name** based on the transmission method or unique symptoms described. If you can't make sure, you can guess up to four possible diseases and use 'OR' to connect them.
### 2: STANDARDIZE FORMAT
Use standard question formats. Such as 'What is [Disease Name]?' or 'How to diagnose [Disease Name]?' or 'How to prevent [Disease Name]?'.
### Note:
Output ONLY the rewritten professional question.
Query:
### Query : {query}
### Rewritten Question :

```

Figure 16: Prompt for Q2EI with four potential entities candidates on MedQuAD.

931 lence (same, hyponym, or hypernym) relationship
932 with the target entity. The results show that in-
933 creasing the number of inferred entities improves
934 coverage. However, this also introduces noise. Fu-
935 ture work will explore confidence scores to bal-
936 ance coverage gains against noise.