ZEROGR: A GENERALIZABLE AND SCALABLE FRAME-WORK FOR ZERO-SHOT GENERATIVE RETRIEVAL

Weiwei Sun^{1,*} Keyi Kong^{2,*} Xinyu Ma³ Shuaiqiang Wang³ Dawei Yin³ Maarten de Rijke⁴ Zhaochun Ren^{5,†} Yiming Yang¹ Carnegie Mellon University ²Shandong University ³Baidu Inc ⁴University of Amsterdam ⁵Leiden University sunnweiwei@gmail.com,luxinyayaya012@gmail.com z.ren@liacs.leidenuniv.nl

ABSTRACT

Generative retrieval (GR) reformulates information retrieval (IR) by framing it as the generation of document identifiers (docids), thereby enabling an end-to-end optimization and seamless integration with generative language models (LMs). Despite notable progress under supervised training, GR still struggles to generalize to zero-shot IR scenarios, which are prevalent in real-world applications. To tackle this challenge, we propose ZEROGR, a zero-shot generative retrieval framework that leverages natural language instructions to extend GR across a wide range of IR tasks. Specifically, ZEROGR is composed of three key components: (i) an LM-based docid generator that unifies heterogeneous documents (e.g., text, tables, code) into semantically meaningful docids; (ii) an instruction-tuned query generator that generates diverse types of queries from natural language task descriptions to enhance corpus indexing; and (iii) a reverse annealing decoding strategy to balance precision and recall during docid generation. We investigate the impact of instruction fine-tuning scale and find that performance consistently improves as the number of IR tasks encountered during training increases. Empirical results on the BEIR and MAIR benchmarks demonstrate that ZEROGR outperforms strong dense retrieval and generative baselines in zero-shot settings, establishing a new state-of-the-art for instruction-driven GR.

1 Introduction

Dense retrieval (DR) (Karpukhin et al., 2020; Izacard et al., 2021), which encodes documents and queries as embedding vectors, is arguably the most effective and widely adopted paradigm (Thakur et al., 2021; Muennighoff et al., 2022) in information retrieval (IR). Despite its success, DR's expressivity is fundamentally limited by the embedding dimensionality (Cao et al., 2020) and does not fully leverage the capabilities of generative language models (LMs) (Tay et al., 2022). As an alternative, generative retrieval (GR) (Metzler et al., 2021) introduces a paradigm shift that encodes corpus information into the model parameters, enabling document retrieval by generating (relevant) document identifiers (docids). GR has demonstrated competitive performance on various IR tasks when large-scale supervised data is available (Tay et al., 2022; Sun et al., 2023; Chen et al., 2022), spanning both traditional web search (Campos et al., 2016) and knowledge-intensive retrieval applications (Petroni et al., 2020).

Despite its promising performance on in-domain tasks, GR exhibits limited generalization to out-of-distribution IR tasks. Existing GR models are typically trained on specific corpora and queries, and prior studies has shown that such training leads to poor performance on unseen tasks (Zhang et al., 2025; Liu et al., 2023b). In contrast, real-world IR models are typically evaluated in a broader setting, characterized by substantial diversity and heterogeneity. These often involve heterogeneous corpora and queries (Thakur et al., 2021), task-specific relevance criteria (Su et al., 2022; Asai et al., 2022), and predominantly zero-shot scenarios where no supervised data is available (Thakur et al., 2021;

^{*}Equal contribution

[†]Corresponding author

Muennighoff et al., 2022). Consequently, GR approaches designed for supervised conditions struggle to generalize to such heterogeneous and data-scarce retrieval scenarios.

To address the limitations of GR in zero-shot and heterogeneous IR scenarios, we draw inspiration from recent advancements in instructed DR methods (Su et al., 2022; Asai et al., 2022) and propose ZEROGR, a generalizable framework for **ZERO**-shot Generative information Retrieval. ZEROGR is a simple yet effective way to adapt GR to diverse IR tasks in a zero-shot setting by leveraging natural language task instructions. Specifically, we advance GR along three dimensions: (i) for *docid design*, we propose a docid generator to efficiently convert a document of any format (e.g., paragraph, table, code) into a unified text-based docid representation; (ii) for *corpus indexing*, we propose an instructed query generator to generate diverse types of queries based on different task instructions; (iii) for *docid decoding*, we propose a reverse annealing strategy that more effectively trades off precision and recall of docid decoding than prior work.

Building on ZEROGR, we investigate *instruction fine-tuning scaling* (Chung et al., 2022) in the context of GR along two key axes: the size of instruction tuning data and the size of the underlying model. We find that increasing both the diversity and quantity of training tasks yields substantial improvements in zero-shot retrieval performance on unseen tasks. Beyond training data scaling, we also examine model size scaling and inference-time scaling for corpus indexing, observing consistently promising scaling trends in both cases..

Our best model, based on Llama-3B LM, surpasses strong dense retrieval and generative retrieval baselines across heterogeneous IR benchmarks, including BEIR (Thakur et al., 2021) and MAIR (Sun et al., 2024). Notably, ZEROGR outperforms OpenAI Embed-v3 on zero-shot MAIR tasks, highlighting its strong generalization to unseen retrieval tasks.

In summary, our contributions are as follows: (i) We propose ZEROGR, a zero-shot GR framework that can construct task-specific GR search indices based on natural language instructions. (ii) Within ZEROGR, we enhance GR by introducing three key components: a unified text-based docid generator, an instruction-conditioned pseudo-query generator, and a reverse annealing decoding strategy. And (iii) ZEROGR achieves competitive performance on heterogeneous IR benchmarks, establishing it as the first GR approach capable of generalizing to diverse tasks in a zero-shot setting.

2 RELATED WORK

2.1 GENERATIVE RETRIEVAL

Unlike traditional dense retrieval methods (Karpukhin et al., 2020; Xiong et al., 2020), GR formulates information retrieval as a docid generation task, enabling end-to-end optimization of the inference-time search index (Tay et al., 2022; Metzler et al., 2021). Previous research on GR has largely focused on three key aspects: (i) Docid design: Early approaches employed rule-based formats such as titles (Cao et al., 2020; Chen et al., 2022), URLs (Zhou et al., 2022), or text spans/summaries (Bevilacqua et al., 2022; Li et al., 2023a). More recent work has shifted toward learning-based docid designs that capture corpus semantics more effectively, including embedding clustering (Tay et al., 2022) and RQ-VAE-based approaches (Sun et al., 2023; Zeng et al., 2023; Wang et al., 2023b). (ii) Corpus indexing: Several strategies have been explored to enrich corpus representations, such as document chunking (Tay et al., 2022), pseudo-query generation (Zhuang et al., 2022), rehearsal-based augmentation (Tang et al., 2023), multi-granular indexing (Wen et al., 2025), and continual training for dynamic corpora (Mehta et al., 2022; Chen et al., 2023; Zhang et al., 2025). (iii) Docid decoding: The dominant approach has been constrained beam search (Cao et al., 2020; Tay et al., 2022). More advanced strategies include multi-stage decoding (Ren et al., 2023), multi-docid decoding (Li et al., 2023b), and simultaneous decoding (Zeng et al., 2024). Despite steady progress, existing work primarily remains confined to supervised fine-tuning, relying heavily on training data and failing to generalize to zero-shot retrieval tasks.

2.2 Instruction Fine-tuning in IR

Inspired by the studies in LLM instruction tuning (Chung et al., 2022; Wang et al., 2022b), instruction fine-tuning for retrieval has gained increased attention to improve zero-shot IR performance (Su et al., 2022; Asai et al., 2022). Instruction-tuned models are able to adapt to various tasks based on natural

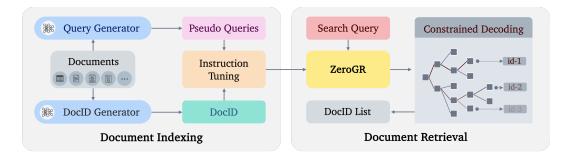


Figure 1: **An overview of ZEROGR.** Given a document collection, ZEROGR converts them into unified DocID representations, generates diverse pseudo-queries, and builds a generative retrieval index. During online retrieval, ZEROGR decodes docids with reverse-annealed temperature scheduling to balance precision and recall.

language instructions that specify the relevance criteria. Recent studies in this direction include multi-task fine-tuning (Lee et al., 2024a), LLM-generated instruction data (Wang et al., 2023a; Lee et al., 2024b; Oh et al., 2024), and instruction-negatives (Weller et al., 2024). These efforts have primarily focused on dense retrieval or cross-encoder rerankers (Sun et al., 2024). To the best of our knowledge, we are the first to investigate instruction fine-tuning for GR and to conduct a systematic study of the factors that influence instruction fine-tuning in IR models.

3 PRELIMINARIES

Zero-shot document retrieval. We formulate the task of zero-shot document retrieval as follows. Given a corpus $\mathcal{D}=(d_1,\ldots,d_n)$ containing n documents, a *corpus indexing* function \mathcal{I} takes \mathcal{D} as input and constructs a search index $m=\mathcal{I}(\mathcal{D})$. Then, a *retrieval* function \mathcal{F} takes the index m and a query q as input, and returns a list of relevant documents: $(d_i,\ldots)=\mathcal{F}(m,q)$. Note that in a typical zero-shot document retrieval setting, no training data is available. However, a natural language task instruction $instr_t$ specifying the retrieval task is generally assumed to be available, as it is usually easier to obtain (Muennighoff et al., 2022).

In dense retrieval (Karpukhin et al., 2020), the indexing function can be defined as using a document encoder to encode the corpus as an embedding matrix such as $\mathbf{E} \in \mathbb{R}^{n \times k}$ and the index is defined as the matrix $m := \mathbf{E}$. For the retrieval function, a query encoder encodes the query q as $\mathbf{q} \in \mathbb{R}^{1 \times k}$ and then perform maximum inner-product search (MIPS) over index m to find the closest document in the embedding space.

Generative retrieval. GR aims to retrieve the document d_i by generating the corresponding document identifier (docid) given the query q. To this end, GR assigns an identifier (docid) to each document in the corpus, e.g. (z_1, \ldots, z_n) , where each z_i is a sequence of tokens $z_i = \{z_i^{(1)}, \ldots, z_i^{(T)}\}$ with a maximum length of T. Based on this, the indexing function $\mathcal{I}(\mathcal{D})$ of GR is to train a language model (LM) \mathcal{LLM} on the corpus \mathcal{D} , encoding the corpus information and also document-docid mapping. The retrieval function F is instantiated by the same \mathcal{LLM} , and it generates the relevant document identifiers (docids) (z_1, \ldots, z_n) given the query q: $(z_i, \ldots) = \mathcal{LLM}(q)$.

4 ZeroGR

We propose ZEROGR, a zero-shot GR framework that can adapt LMs into task-specific generative search indexes based on task instructions. As shown in Figure 1, the proposed ZEROGR framework consists of three key components: (i) a docid generator G_{ψ} , which takes a document d_i as input and outputs its docid z_i ; (ii) an instructed query generator, which takes a task instruction *instr* and a document d_i as input and outputs multiple pseudo-queries; (iii) a generative retriever \mathcal{LLM} , which takes the instruction and a query as input and generates a list of docids.

The ZEROGR pipeline proceeds as follows: (i) given a new corpus \mathcal{D} and its associated task instruction *instr*, the docid generator assigns each document d_i a docid z_i ; (ii) the instructed query generator G_{θ}

samples B queries $\{q_{i,1},\ldots,q_{i,B}\}$ for each document $d_i\in\mathcal{D}$, thereby creating $\langle q_{i,j},z_i\rangle$ pairs; and (iii) the generative retriever is trained to predict the corresponding docid z_i given the concatenation of *instr* and a sampled query $q_{i,j}$. After training, the generative retriever $\mathcal{LLM}(z\mid q,instr)$ serves as the search index m. For a given query q, a newly proposed reverse annealing decoding strategy is employed to generate a ranked list of docids as retrieval results.

4.1 Unified Docid Representation

Documents in downstream IR tasks can be heterogeneous, e.g., financial tables (Zhu et al., 2022), code files (Liu et al., 2023a), meeting transcripts (Golany et al., 2024), or legal cases (Bhattacharya et al., 2019). Existing simple docid strategies, such as using document titles, URLs, or spans (Cao et al., 2020; Bevilacqua et al., 2022), often fail to generalize to user-customized data. ZEROGR therefore introduces a model-based **docid generator** G_{ψ} that maps any document to a short, keyword-rich sentence (typically 6–8 words) ranked by coverage. Formally, for a document d_i we define

$$z_i = G_{\psi}(d_i) = \underset{t \in \mathcal{V}^{\leq L}}{\arg \max} G_{\psi}(t \mid d_i), \tag{1}$$

where t is a token sequence of length $\leq L$ (with L=8) drawn from the vocabulary $\mathcal V$. To instantiate G_ψ , we first prompt a powerful LM (e.g., GPT-40) with detailed instructions and in-context examples to create a training set of $\langle d_i, z_i \rangle$ pairs. A smaller model (Llama-3.2-1B) is then fine-tuned on this data, enabling fast, scalable generation of unified docids across diverse IR tasks. See Section 5.1 for details of training data.

4.2 Instructed Corpus Indexing

Corpus indexing in GR encodes each document $d_i \in \mathcal{D}$ into the model's parameters so that, at inference time, the model can recover d_i by *generating* its document identifier z_i . DSI-QG (Zhuang et al., 2022) accomplishes this by pairing every document with a set of pseudo-queries, but its effectiveness diminishes when the pseudo-query distribution diverges from real user queries (Pradeep et al., 2023; Dai et al., 2022). This gap is especially large in heterogeneous IR scenarios, such as conversational, code, or multimodal search.

We mitigate the distribution gap with an **instructed query generator** G_{θ} , obtained by instructiontuning a 1B-parameter Llama model on diverse IR datasets verbalized through task-specific instructions. Given a document d_i and a task instruction *instr*, the generator produces a pseudo-query $q_{i,j}$ from the conditional distribution

$$q_{i,j} \sim G_{\theta}(\cdot \mid d, instr).$$
 (2)

For each document we draw B queries with temperature of 1:

$$Q_i = \{ q_{i,1}, \dots, q_{i,B} \}. \tag{3}$$

These $\langle d_i, z_i \rangle$ pairs are used to train the generative retriever \mathcal{LLM} by minimizing the cross-entropy loss

$$\mathcal{L}(\phi) = -\sum_{q_{i,j} \in \mathcal{Q}_{i}} \sum_{q_{i,j} \in \mathcal{Q}_{i}} \log \mathcal{LLM}(z_{i} \mid q_{i,j}, instr), \tag{4}$$

thereby embedding the corpus into the model's parameters. Appendix E summarizes the instruction-tuning datasets.

4.3 REVERSE-ANNEALED DOCID GENERATION

During inference, a GR model must decode each docid z_i as a sequence of tokens. Standard beam search often collapses to a few high-probability sequences, hurting recall. We therefore propose **reverse-annealed sampling**: each z_i is generated token-by-token, while the sampling temperature is gradually increased to encourage diversity. Let $f(\cdot)$ denote the trained decoder after corpus indexing, and let T be a prefix tree whose leaves correspond to valid docids. For the i-th docid we decode a token sequence $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,L_i})$ using temperature $t_i = g(i)$. At position j we sample

$$x_{i,j} \sim \operatorname{Softmax}\left(\frac{\ell_{i,j}}{t_i}\right)\Big|_{T_{i,j}},$$
 (5)

where $\ell_{i,j}$ are the logits conditioned on the current prefix $(x_{i,1:j-1})$, and the subscript $T_{i,j}$ masks probabilities to tokens that keep the prefix inside the tree. After the complete sequence \mathbf{x}_i is produced, its leaf is removed from T so no subsequent iteration can repeat the same docid. The per-iteration temperature t_i follows a *normalized sigmoid*:

$$t_i = g(i) = T_{\text{max}} \cdot \frac{\sigma(k(\frac{i}{K} - m)) - \sigma(-km)}{\sigma(k(1 - m)) - \sigma(-km)}, \quad \sigma(z) = \frac{1}{1 + e^{-z}},\tag{6}$$

where K is the total number of docids to generate, k > 0 controls the slope, and $m \in (0,1)$ sets the midpoint. Starting from a low temperature yields high-precision early selections; increasing t_i over iterations boosts exploration, thereby balancing precision and recall across the final ranked list.

5 EVALUATION SETUP

5.1 Training Datasets

To support the development of ZEROGR, we collect training data covering a diverse range of IR tasks. Specifically, we use MAIR (Sun et al., 2024), a multi-task IR evaluation benchmark comprising 126 tasks, and extract the training splits of these tasks when available. As shown in Figure 6 in Appendix E, **ZeroGR-Train** is a dataset spanning 69 IR tasks across 6 domains, containing 41 million query-document pairs. ZeroGR-Train is the largest open-source IR training corpus to date. It offers greater domain and task diversity,

Domain	#Tasks	#Samples
Medical	5	421, 430
Financial	8	31, 315
Academic	18	744, 160
Coding	13	1, 969, 586
Legal	7	23, 086, 948
Web-based	18	15, 319, 445

Table 1: Statistics of ZeroGR-Train

includes detailed instructional annotations, and provides reliable relevance labels. See Table 6 for details.

5.2 EVALUATION DATASETS

To evaluate zero-shot GR on diverse downstream tasks, we use the BEIR and MAIR benchmarks: (i) **BEIR** (Thakur et al., 2021). We evaluate models on all 12 tasks from BEIR collections. (ii) **MAIR** (Sun et al., 2024). As we collect training data from a subset of MAIR tasks, we divide MAIR into seen and unseen subsets, where the unseen subset contains tasks not present in the ZeroGR-Training data, to validate the zero-shot generalization of models. In constructing this benchmark, we curated a diverse set of long-tail tasks across 6 domains, and intentionally omitted redundant tasks (e.g., different years of the same competition) and structurally complex ones (e.g., IFEval) that would introduce evaluation overhead. Given the large size of the MAIR dataset, we also develop a Dev subset of MAIR for model ablation.

5.3 EVALUATION METRIC

We evaluate models using the following metrics: (i) *Top-1 accuracy*, which measures retrieval precision by checking whether the top-ranked document is relevant to the query; (ii) *nDCG@10*, a popular metric that evaluates the quality of the top-10 ranked results by considering both the relevance and position of retrieved documents; and (iii) *Recall@100*, which assesses recall by calculating the percentage of relevant documents retrieved within the top-100 ranked list.

5.4 IMPLEMENTATION DETAILS

We implement the three components of ZEROGR, i.e., query generator, docid generator, and final generative retriever, all with Llama-based LMs. For the docid generator, a Llama-1B-Instruct model is trained on our curated document-docid pairs for 5 epochs with a constant learning rate of 5e-5. Similarly, for the query generator, a Llama-1B-Instruct model is trained on the ZeroGR-Training set for 5 epochs with a constant learning rate of 5e-5. For the generative retriever, the model is trained for each evaluated task on data generated by the query generator and docid generator, based on our "Document Indexing" workflow described in Fig 1.

		BEII	R (11 7	Tasks)		MAIR (38 Tasks)							
Model	Avg	Web.	Aca.	Med.	Fin.	Avg	Web.	Aca.	Legal	Med.	Fin.	Cod.	
BM25	42.4	45.4	38.8	32.7	41.6	36.1	34.3	39.2	34.5	42.4	40.0	17.3	
Contriever GTR-T5-base	47.6 45.3 48.0	51.5 50.7 53.3	43.0 35.3 37.4	33.9 32.7 33.4	45.3	33.6 32.5 35.4	39.8 36.0 39.8	33.4 33.6 39.6	26.8 25.3 27.8	30.8 31.9 31.8	37.3 37.4	17.7 18.7 24.0	
GTR-T5-large E5-Base	48.9	51.8	46.1	35.4		37.2	36.2	48.6	28.5	35.3	44.9		
E5-Large BGE-Base	49.2 50.5	51.7 52.5	47.9 47.1	37.4 36.0		37.0	38.6 38.6	51.0 40.2	25.0 25.8	35.6 37.6	46.6 42.2		
BGE-Large	51.8	53.8	47.9	38.1		39.4	39.4	46.2	36.0	37.2	45.1	29.0	
OpenAI-Embed E5-mistral-7B GritLM-7B	54.2 55.7 45.0	56.3 56.4 47.7	47.2 48.6 48.2	37.6 39.6 36.9		46.8	40.6 45.4 44.1	48.2 55.4 58.2	31.0 42.3 43.3	39.7 43.1 42.6	49.4 55.3 57.6	28.7 40.0 40.0	
ZeroGR-3B	48.1	49.2	45.8	34.7	53.8	41.1	42.7	47.4	40.0	38.3	39.2	36.3	

Table 2: Combined Domain-wise Results on BEIR (nDCG@10) and MAIR (Hit@1). Performance of different retrieval models across various domains.

5.5 Baselines

We evaluate ZEROGR against several representative IR baselines, spanning different retrieval paradigms to provide a comprehensive comparison. (i) For sparse retrieval, we adopt the classical term-based model BM25, implemented using the BM25S package (Lù, 2024), which remains a strong baseline in many IR tasks due to its simplicity and effectiveness. (ii) For traditional dense retrieval models trained on a single task, we include Contriever-MARCO, GTR-base, and GTR-Large, all of which are pretrained or fine-tuned on the MS MARCO dataset (Ni et al., 2021; Izacard et al., 2021), representing a common practice in dense retrieval pipelines. (iii) For multi-task-trained dense retrievers, we incorporate E5-Base and E5-Large (Wang et al., 2022a), BGE-base and BGE-Large (Xiao et al., 2023), as well as OpenAI-Embedding-v3-Small, all of which use supervision from multiple tasks to enhance generalization across diverse domains. (iv) For instruction-tuned dense retrieval models, which aim to align the retriever with human instructions, we include E5-Mistral-7B-instruct (Wang et al., 2023a), and GritLM-7B (Muennighoff et al., 2024), which are trained on large-scale, diverse instruction datasets to follow task-specific intents effectively.

6 EXPERIMENTS

Our experiments aim to address the following research questions:

- 1. How does ZEROGR compare with state-of-the-art retrieval methods? We evaluate ZEROGR against leading models on the MAIR benchmark (Section 6.1) and conduct additional analysis on the BEIR datasets (Section 6.2).
- 2. How do model design and training strategies influence the performance of ZEROGR on unseen IR tasks?

To answer this, we conduct a systematic study on the development set, investigating key factors in generative retrieval. Specifically, we analyze how instruction tuning task diversity (Section 6.3), docid design (Section 6.4), corpus indexing strategy, model size (Section 6.5), and decoding strategy (Section 6.6) affect performance.

6.1 EVALUATION RESULTS ON MAIR

As shown in Table 2 (MAIR), our proposed ZEROGR framework demonstrates strong performance across a wide range of retrieval tasks. It achieves an average score of 41.1 (Acc@1), substantially outperforming traditional sparse retrieval methods like BM25 and widely adopted dense retrieval

Method	Avg	Argu.	SciFact	NFC.	FiQA	SciDocs	Covid
GENRE (Cao et al., 2020)	23.0	42.5	42.3	20.0	11.6	6.8	14.7
GENRET (Sun et al., 2023)	41.1	34.3	63.9	31.6	30.2	14.9	71.8
GLEN (Lee et al., 2023)	-	17.6	_	15.9	_	_	_
TIGER (Rajput et al., 2023)	31.0	14.0	37.0	39.5	16.0	14.0	65.7
ZEROGR	44.9	35.4	72.8	34.7	34.1	18.7	73.5

Table 3: Performance of different generative retrieval models across various datasets on BEIR.

models such as Contriever, GTR, E5, BGE, and even the strong instruction-tuned OpenAI-Embedding-v3-Small. These results highlight the effectiveness of our instruction-based generative retrieval approach in capturing deeper semantic relevance.

The performance gains of ZEROGR are not limited to familiar tasks but also generalize well to unseen domains. Notably, the model achieves state-of-the-art results on several previously unseen datasets, including Apple, MB, PM.A, DD, and NCL (cf Table 4). This demonstrates the robustness and transferability of the approach, as it adapts effectively to new retrieval settings without requiring additional task-specific supervised data. See Figure 5 for a comparison on the MAIR unseen subset, where ZEROGR achieves competitive performance against recent dense retrieval methods.

Despite having a relatively modest model size of 3B parameters, ZEROGR delivers competitive or even superior results compared to larger instruction-tuned dense retrievers with 7B parameters. This indicates that our design is highly parameter-efficient, achieving strong performance across diverse tasks without relying on massive model scaling.

6.2 EVALUATION RESULTS ON BEIR

As shown in Table 2 (BEIR), ZEROGR consistently outperforms several strong retrieval baselines across a diverse set of IR benchmarks. Compared to the sparse retrieval method BM25, it achieves higher performance on six datasets, highlighting its strength in capturing deep semantic relevance beyond surface-level token overlap. When evaluated against Contriever—a dense retriever trained on MS MARCO—ZEROGR leads on four datasets, suggesting its superior generalization and contextual understanding. Table 3 compares ZEROGR with previous SOTA generative retrieval baselines on BEIR, which we can see our method achieve best performance among most datasets.

6.3 SCALING INSTRUCTION FINE-TUNING

A key factor in enhancing the performance of LM-based tasks is scaling, i.e., increasing model size or data volume. The effectiveness of ZEROGR stems from instruction fine-tuning on multi-task IR datasets, which improves the instruction-following abilities of both the query generator and the title generator models. To investigate the impact of multi-task training, we curate training data with varying numbers of tasks: (a) *MS MARCO*, which contains a single task (i.e., MS MARCO (Campos et al., 2016)) and is commonly used in previous GR work; (b) + *OpenQA*, which adds popular open-domain question answering datasets, including NQ (Kwiatkowski et al., 2019) and HotpotQA; (c) + *BEIR-Train*, which incorporates the training splits of BEIR (Thakur et al., 2021), such as NFCorpus and Quora; (d) + *MTEB-Train*, which includes additional tasks from MTEB (Muennighoff et al., 2022) that are not covered in BEIR, such as NLI (we use the public BGE training split to collect these data); and (e) + *ZeroGR-Train*, which includes the data we collected from the training split of the MAIR (Sun et al., 2024) task collection, comprising 69 tasks from 6 domains, as shown in Figure 2.

Figure 2 shows the evaluation results of models (both query generator and docid generator) trained with different levels of task diversity, evaluated on the unseen task subset (i.e., tasks not included in any training set) of MAIR. The left plot in Figure 2 shows the distribution of average query length across tasks. We observe that models trained on more IR tasks generate queries with greater length diversity, indicating task-aware query generation strategies. In contrast, the baseline model trained only on MS MARCO produces short queries, averaging 8 words. The middle plot shows the docid conflict rate, i.e., the percentage of documents in the corpus assigned the same docid by the docid generator. Models trained on diverse tasks exhibit lower conflict rates, suggesting a stronger ability to

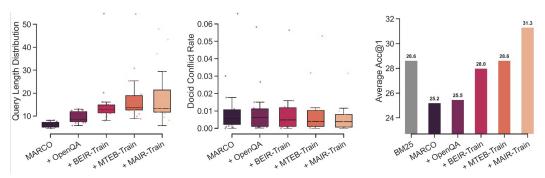


Figure 2: Model performance on unseen-dev tasks as a function of the number of instruction-tuning tasks. We gradually increase the number of instruction-tuning tasks, starting from MS MARCO, and incrementally add open-domain QA datasets (e.g., NQ), BEIR-Train sets (e.g., NFC), MTEB-Train data (e.g., NLI), and finally the ZeroGR-Train collection, which includes 60 tasks across 6 domains. Left: More instruction-tuning tasks lead to more diverse queries. Middle: More instruction-tuning tasks reduce docid conflicts. Right: More instruction-tuning tasks improve the Acc@1 score.

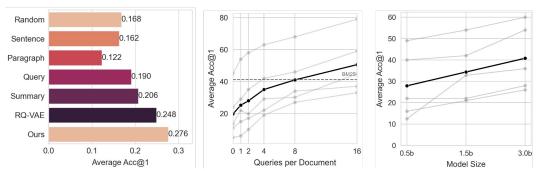


Figure 3: **Left:** Comparison of different docid designs. **Middle:** Acc@1 vs. generated queries per document. **Right:** Acc@1 vs. model size.

process heterogeneous corpora. The MS MARCO baseline shows higher conflict on several diverse tasks. Finally, the right plot reports retrieval performance (top-1 accuracy) for different models. We observe consistent performance improvements on unseen tasks as training data diversity increases.

6.4 Comparisons of Different Docid Designs

Figure 3 compares our proposed unified docid with previous GR docid designs, while keeping all other factors (e.g., query generator, model choice, optimization strategy) constant to ensure an apple-to-apple comparison of docid effectiveness. The compared docid designs include: (i) **Random** (Tay et al., 2022), a baseline that assigns each document a random string as its docid; (ii) **Sentence** (Bevilacqua et al., 2022), which uses all sentences of each document as its docid; (iii) **Paragraph** (Tay et al., 2022), which takes the first paragraph of each document as its docid; (iv) **Query** (Tang et al., 2023), which uses a query generator to produce a single query per document as its docid; (v) **Summary**, as introduced in (Li et al., 2024), which uses the output of a summarization model as the docid; (vi) **RQ-VAE** (Zeng et al., 2023), which trains a RQ-VAE model on document embeddings produced by the BGE-Large model, enabling quantization of document embeddings into a sequence of tokens. This is a widely adopted docid representation in competitive GR systems.

From the results, we observe that among the various docid designs, our proposed docid generator consistently achieves the best performance on unseen development tasks. In particular, it significantly outperforms other text-based approaches such as *Summary* and *Query*, highlighting its superior ability to encode meaningful and discriminative document representations. This suggests that our design not only captures richer document semantics but is also better aligned with the generative retrieval objective, enabling more accurate and robust document retrieval. We further find that the performance of the *RQ-VAE* method is relatively unstable across different tasks, often requiring longer training to

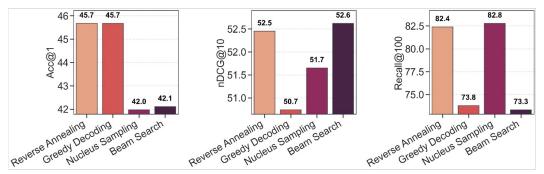


Figure 4: Ablation study of decoding algorithms across different metrics. Our proposed reverse annealing decoding achieves a good balance between precision and recall.

converge effectively. In contrast, our text-based docid benefits from the pretrained LM's inherent understanding of natural language, which facilitates more efficient learning and faster convergence. This synergy between instruction-driven docid generation and LM capabilities underpins the strong performance and generalization ability observed in our experiments.

6.5 SCALING QUERIES NUMBER AND MODEL SIZE

The middle section of Figure 3 illustrates the impact of the number of pseudo-queries generated per document on the average top-1 accuracy of ZEROGR. We observe a clear upward trend: as the number of pseudo-queries increases, the retrieval performance improves steadily. This highlights the importance of diverse query views for better semantic coverage during indexing. Notably, when using eight queries per document, ZEROGR already reaches performance on par with the strong sparse baseline BM25. Further increasing the query count to sixteen enables ZEROGR to surpass BM25, suggesting that high query diversity provides richer signals for matching user queries to relevant documents.

The right section of Figure 3 examines how the size of the backbone language model affects retrieval performance. For this analysis, we adopt a series of Qwen2.5 (Qwen et al., 2025) models with varying parameter scales. The results demonstrate a consistent gain in top-1 accuracy on unseen IR tasks as the model size grows, implying that larger models benefit from enhanced generalization and better understanding of the instruction-based retrieval formulation. This finding underscores the value of scaling up model capacity in generative retrieval frameworks, particularly in zero-shot settings.

6.6 Analysis of Decoding Strategies

In Figure 4, we compare our reverse annealing decoding with other popular decoding algorithms, including greedy decoding (i.e., greedily sampling from the GR model without replacement), nucleus sampling with a top-p of 0.9, and beam search. All methods decode the top-100 docids for evaluation. From the results, we observe that greedy decoding achieves the best performance in terms of Acc@1, but lacks diversity and yields low recall. Nucleus sampling performs poorly on Acc@1 but achieves high recall. In contrast, reverse annealing strikes a good balance between precision and recall, achieving competitive results across all metrics.

7 CONCLUSION

This work presents ZEROGR, an instruction-driven framework that extends generative retrieval to zero-shot scenarios. By unifying three key components—a model-based docid generator, an instruction-conditioned query generator, and a reverse-annealed decoding algorithm—ZEROGR transforms a corpus and a natural-language task description into a task-specific generative index without requiring supervision. The framework introduces natural-language instructions to align corpus indexing and docid decoding, enabling seamless adaptation to heterogeneous IR tasks. It incorporates a unified, keyword-centric docid representation, an instruction-tuned pseudo-query generator, and a reverse annealing decoding strategy that jointly balance precision and recall while

maintaining parameter efficiency. Systematic scaling studies along task diversity, query volume, and model size reveal consistent performance improvements. Empirical evaluations on 69 MAIR tasks and seven BEIR datasets demonstrate that ZEROGR consistently outperforms strong baselines, including BM25, traditional dense retrievers, and recent instruction-tuned models, while remaining competitive with much larger 7B-parameter systems.

REFERENCES

- Akari Asai, Timo Schick, Patrick Lewis, Xilun Chen, Gautier Izacard, Sebastian Riedel, Hannaneh Hajishirzi, and Wen tau Yih. Task-aware retrieval with instructions. *ArXiv*, abs/2211.09260, 2022.
- Michele Bevilacqua, Giuseppe Ottaviano, Patrick Lewis, Wen tau Yih, Sebastian Riedel, and Fabio Petroni. Autoregressive search engines: Generating substrings as document identifiers. *ArXiv*, abs/2204.10628, 2022.
- Paheli Bhattacharya, Kripabandhu Ghosh, Saptarshi Ghosh, Arindam Pal, Parth Mehta, Arnab Bhattacharya, and Prasenjit Majumder. Overview of the fire 2019 aila track: Artificial intelligence for legal assistance. In *Fire*, 2019.
- Daniel Fernando Campos, Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, Li Deng, and Bhaskar Mitra. Ms marco: A human generated machine reading comprehension dataset. *ArXiv*, abs/1611.09268, 2016.
- Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. Autoregressive entity retrieval. *ArXiv*, abs/2010.00904, 2020.
- Jiangui Chen, Ruqing Zhang, J. Guo, Y. Liu, Yixing Fan, and Xueqi Cheng. Corpusbrain: Pre-train a generative retrieval model for knowledge-intensive language tasks. *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, 2022.
- Jiangui Chen, Ruqing Zhang, J. Guo, M. de Rijke, Wei Chen, Yixing Fan, and Xueqi Cheng. Continual learning for generative retrieval over dynamic corpora. *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, 2023.
- Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *ArXiv*, abs/2210.11416, 2022.
- Zhuyun Dai, Vincent Zhao, Ji Ma, Yi Luan, Jianmo Ni, Jing Lu, Anton Bakalov, Kelvin Guu, Keith B. Hall, and Ming-Wei Chang. Promptagator: Few-shot dense retrieval from 8 examples. *ArXiv*, abs/2209.11755, 2022.
- Lotem Golany, Filippo Galgani, Maya Mamo, Nimrod Parasol, Omer Vandsburger, Nadav Bar, and Ido Dagan. Efficient data generation for source-grounded information-seeking dialogs: A use case for meeting transcripts. *ArXiv*, abs/2405.01121, 2024.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning. *Trans. Mach. Learn. Res.*, 2022, 2021.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Yu Wu, Sergey Edunov, Danqi Chen, and Wen tau Yih. Dense passage retrieval for open-domain question answering. *ArXiv*, abs/2004.04906, 2020.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc V. Le, and Slav Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.

- Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. Nv-embed: Improved techniques for training llms as generalist embedding models. *ArXiv*, abs/2405.17428, 2024a.
- Jinhyuk Lee, Zhuyun Dai, Xiaoqi Ren, Blair Chen, Daniel Cer, Jeremy R. Cole, Kai Hui, Michael Boratko, Rajvi Kapadia, Wen Ding, Yi Luan, Sai Meher Karthik Duddu, Gustavo Hernández Abrego, Weiqiang Shi, Nithi Gupta, Aditya Kusupati, Prateek Jain, Siddhartha R. Jonnalagadda, Ming-Wei Chang, and Iftekhar Naim. Gecko: Versatile text embeddings distilled from large language models. *ArXiv*, abs/2403.20327, 2024b.
- Sunkyung Lee, Minjin Choi, and Jongwuk Lee. Glen: Generative retrieval via lexical index learning. In *Conference on Empirical Methods in Natural Language Processing*, 2023.
- Alan Li, Daniel Cheng, Phillip Keung, Jungo Kasai, and Noah A. Smith. Summarization-based document IDs for generative retrieval with language models. In Lucie Lucie-Aimée, Angela Fan, Tajuddeen Gwadabe, Isaac Johnson, Fabio Petroni, and Daniel van Strien, editors, *Proceedings of the First Workshop on Advancing Natural Language Processing for Wikipedia*, pages 126–135, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.wikinlp-1.18.
- Haoxin Li, Phillip Keung, Daniel Cheng, Jungo Kasai, and Noah A. Smith. Summarization-based document ids for generative retrieval with language models. *Proceedings of the First Workshop on Advancing Natural Language Processing for Wikipedia*, 2023a.
- Yongqing Li, Nan Yang, Liang Wang, Furu Wei, and Wenjie Li. Multiview identifiers enhanced generative retrieval. *ArXiv*, abs/2305.16675, 2023b.
- Tianyang Liu, Canwen Xu, and Julian McAuley. Repobench: Benchmarking repository-level code auto-completion systems. *ArXiv*, abs/2306.03091, 2023a.
- Yu-An Liu, Ruqing Zhang, Jiafeng Guo, Wei Chen, and Xueqi Cheng. On the robustness of generative retrieval models: An out-of-distribution perspective. *arXiv* preprint arXiv:2306.12756, 2023b.
- Xing Han Lù. Bm25s: Orders of magnitude faster lexical search via eager sparse scoring. *ArXiv*, abs/2407.03618, 2024.
- Sanket Vaibhav Mehta, Jai Gupta, Yi Tay, Mostafa Dehghani, Vinh Q. Tran, Jinfeng Rao, Marc Najork, Emma Strubell, and Donald Metzler. Dsi++: Updating transformer memory with new documents. In *Conference on Empirical Methods in Natural Language Processing*, 2022.
- Donald Metzler, Yi Tay, Dara Bahri, and Marc Najork. Rethinking search: Making experts out of dilettantes. *ArXiv*, abs/2105.02274, 2021.
- Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. Mteb: Massive text embedding benchmark. In *Conference of the European Chapter of the Association for Computational Linguistics*, 2022.
- Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. Generative representational instruction tuning. *ArXiv*, abs/2402.09906, 2024.
- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernández Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei Yang. Large dual encoders are generalizable retrievers. *ArXiv*, abs/2112.07899, 2021.
- Hanseok Oh, Hyunji Lee, Seonghyeon Ye, Haebin Shin, Hansol Jang, Changwook Jun, and Minjoon Seo. Instructir: A benchmark for instruction following of information retrieval models. *ArXiv*, abs/2402.14334, 2024.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vassilis Plachouras, Tim Rocktaschel, and Sebastian Riedel. Kilt: a benchmark for knowledge intensive language tasks. In *North American Chapter of the Association for Computational Linguistics*, 2020.

- Ronak Pradeep, Kai Hui, Jai Gupta, Ádám Dániel Lelkes, Honglei Zhuang, Jimmy J. Lin, Donald Metzler, and Vinh Q. Tran. How does generative retrieval scale to millions of passages? *ArXiv*, abs/2305.11841, 2023.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025.
- Shashank Rajput, Nikhil Mehta, Anima Singh, Raghunandan H. Keshavan, Trung Hieu Vu, Lukasz Heldt, Lichan Hong, Yi Tay, Vinh Q. Tran, Jonah Samost, Maciej Kula, Ed H. Chi, and Maheswaran Sathiamoorthy. Recommender systems with generative retrieval. *ArXiv*, abs/2305.05065, 2023.
- Ruiyang Ren, Wayne Xin Zhao, J. Liu, Huaqin Wu, Ji rong Wen, and Haifeng Wang. Tome: A two-stage approach for model-based retrieval. *ArXiv*, abs/2305.11161, 2023.
- Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen tau Yih, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. One embedder, any task: Instruction-finetuned text embeddings. *ArXiv*, abs/2212.09741, 2022.
- Weiwei Sun, Lingyong Yan, Zheng Chen, Shuaiqiang Wang, Haichao Zhu, Pengjie Ren, Zhumin Chen, Dawei Yin, M. de Rijke, and Zhaochun Ren. Learning to tokenize for generative retrieval. *ArXiv*, abs/2304.04171, 2023.
- Weiwei Sun, Zhengliang Shi, Jiulong Wu, Lingyong Yan, Xinyu Ma, Yiding Liu, Min Cao, Dawei Yin, and Zhaochun Ren. Mair: A massive benchmark for evaluating instructed retrieval. In *Conference on Empirical Methods in Natural Language Processing*, 2024.
- Yubao Tang, Ruqing Zhang, J. Guo, Jiangui Chen, Zuowei Zhu, Shuaiqiang Wang, Dawei Yin, and Xueqi Cheng. Semantic-enhanced differentiable search index inspired by learning strategies. *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2023.
- Yi Tay, Vinh Q. Tran, Mostafa Dehghani, Jianmo Ni, Dara Bahri, Harsh Mehta, Zhen Qin, Kai Hui, Zhe Zhao, Jai Gupta, Tal Schuster, William W. Cohen, and Donald Metzler. Transformer memory as a differentiable search index. *ArXiv*, abs/2202.06991, 2022.
- Nandan Thakur, Nils Reimers, Andreas Ruckl'e, Abhishek Srivastava, and Iryna Gurevych. Beir: A heterogenous benchmark for zero-shot evaluation of information retrieval models. *ArXiv*, abs/2104.08663, 2021.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. *ArXiv*, abs/2212.03533, 2022a.
- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Improving text embeddings with large language models. *ArXiv*, abs/2401.00368, 2023a.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Maitreya Patel, Kuntal Kumar Pal, M. Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Shailaja Keyur Sampat, Savan Doshi, Siddhartha Mishra, Sujan Reddy, Sumanta Patro, Tanay Dixit, Xudong Shen, Chitta Baral, Yejin Choi, Noah A. Smith, Hannaneh Hajishirzi, and Daniel Khashabi. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Conference on Empirical Methods in Natural Language Processing*, 2022b.
- Zihan Wang, Yujia Zhou, Yiteng Tu, and Zhicheng Dou. Novo: Learnable and interpretable document identifiers for model-based ir. *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, 2023b.

- Orion Weller, Benjamin Van Durme, Dawn Lawrie, Ashwin Paranjape, Yuhao Zhang, and Jack Hessel. Promptriever: Instruction-trained retrievers can be prompted like language models. *ArXiv*, abs/2409.11136, 2024.
- Haoyang Wen, Jiang Guo, Yi Zhang, Jiarong Jiang, and Zhiguo Wang. On synthetic data strategies for domain-specific generative retrieval. *ArXiv*, abs/2502.17957, 2025.
- Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. C-pack: Packaged resources to advance general chinese embedding. *ArXiv*, abs/2309.07597, 2023.
- Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning for dense text retrieval. *ArXiv*, abs/2007.00808, 2020.
- Hansi Zeng, Chen Luo, Bowen Jin, Sheikh Muhammad Sarwar, Tianxin Wei, and Hamed Zamani. Scalable and effective generative information retrieval. *Proceedings of the ACM on Web Conference* 2024, 2023.
- Hansi Zeng, Chen Luo, and Hamed Zamani. Planning ahead in generative retrieval: Guiding autoregressive generation through simultaneous decoding. *ArXiv*, abs/2404.14600, 2024.
- Zhen Zhang, Xinyu Ma, Weiwei Sun, Pengjie Ren, Zhumin Chen, Shuaiqiang Wang, Dawei Yin, Maarten de Rijke, and Zhaochun Ren. Replication and exploration of generative retrieval over dynamic corpora. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3325–3334, 2025.
- Yujia Zhou, Jing Yao, Zhicheng Dou, Ledell Yu Wu, Peitian Zhang, and Ji rong Wen. Ultron: An ultimate retriever on corpus with a model-based indexer. *ArXiv*, abs/2208.09257, 2022.
- Fengbin Zhu, Wenqiang Lei, Fuli Feng, Chao Wang, Haozhou Zhang, and Tat seng Chua. Towards complex document understanding by discrete reasoning. *Proceedings of the 30th ACM International Conference on Multimedia*, 2022.
- Shengyao Zhuang, Houxing Ren, Linjun Shou, Jian Pei, Ming Gong, G. Zuccon, and Daxin Jiang. Bridging the gap between indexing and retrieval for differentiable search index with query generation. *ArXiv*, abs/2206.10128, 2022.

A Prompts

- 1. **Length**: Strictly 6-8 words (terms/words)
 2. **Term Inclusion**: Must include 3-5 core terms directly from
 the document
- 3. **Term Positioning**: Rank by relevance and importance (highest → lowest, general → specific)
- 4. **Formatting**:
 - Use lowercase letters, numbers, and spaces only
 - Preserve special terms/symbols (e.g., PD3.1)
- **No articles** (a, the), **linking verbs**, or auxiliary verbs
 - **No verbs** (use nouns/adjectives only)
- 5. **Requirements**:
 - Terms must be derivable from the document
 - Ensure uniqueness and precise core content representation

B Decoding

Algorithm 1 DocID Generation with Reverse Annealing

Require: T (total number of docids), model, query, max_temperature

Ensure: List of generated docids

- 1: **for** t = 1, 2, ..., T **do**
- 2: # Compute normalized decoding temperature (Eq. 6)
- 3: temperature, \leftarrow reverse-annealing $(t, T, \text{max_temperature})$
- 4: # Generate next tokens with temperature control
- 5: $\operatorname{docid}_t \leftarrow \operatorname{model}(\operatorname{query}, \operatorname{temperature}_t)$
- 6: end for
- 7: return List of generated docids

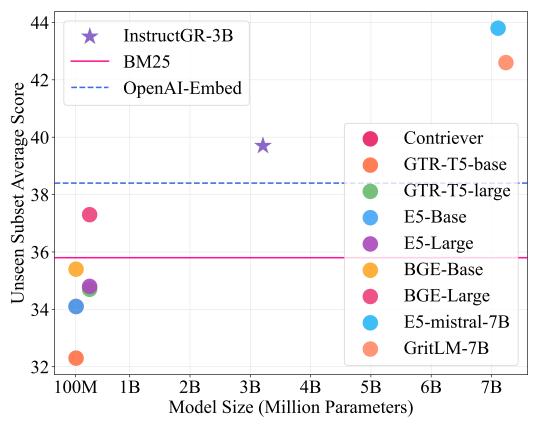


Figure 5

C USE OF LARGE LANGUAGE MODELS

In the preparation of this manuscript, we utilized Large Language Models (LLMs) to enhance the clarity and linguistic quality of our academic writing. Specifically, we employed Claude Sonnet 4 (Anthropic) and GPT-5 (OpenAI) for language polishing and refinement purposes.

D EXPERIMENTAL RESULTS ON MAIR AND BEIR

Our experimental results on MAIR and BEIR are shown in Table 4 and Table 5.

E THE ZEROGR-TRAIN DATASET

We show the statistics of ZeroGR-Train Dataset in the Table 6 and Figure 6.

Table 4: Model performance (top-1 retrieval accuracy) on seen and unseen subset of MAIR.

Dataset									Seen S	ubset								Unseen	Subset	
Model	Avg	FiQA	NFC.s	SciD.	SciF.	ГоQА.	TAT	CoF. I	LeetC.	LitSe. I	3iSum (CodeSe.	Math	ConvF. C	Conala	StMath.	Apple	FinBen .	AILAC .	AILAS
BM25	35.4	24.0	45.5	16.0	53.0	11.0	67.1	47.9	12.0	66.0	69.0	33.0	41.0	47.9	7.0	20.0	52.1	9.0	14.0	10.
Contriever	33.0	33.0	43.5	17.0	62.0	11.0	54.3	42.7	7.0	39.0	54.0	37.0	36.0	42.7	9.0	13.0	50.7	6.0	12.0	6.0
GTR-T5-base	31.8	33.0	43.0	12.0	50.0	16.0	58.6	37.5	6.0	41.0	47.0	41.0	49.0	37.5	9.0	16.0	47.9	10.0	4.0	10.0
GTR-T5-large	34.6	45.0	43.5	14.0	55.0	12.0	40.0	42.7	10.0	43.0	59.0	51.0	63.0	42.7	11.0	23.0	50.7	14.0	6.0	6.
E5-Base	36.6	41.0	40.0	17.0	63.0	16.0	61.4	52.1	11.0	49.0	61.0	59.0	78.0	52.1	10.0	36.0	52.1	18.0	6.0	12.
E5-Large	37.7	45.0	45.5	20.0	67.0	13.0	70.0	57.3	9.0	49.0	52.0	57.0	75.0	57.3	11.0	44.0	47.9	13.0	12.0	6.0
BGE-Base	36.1	43.0	43.5	20.0	62.0	13.0	58.6	41.7	10.0	44.0	67.0	64.0	50.0	41.7	13.0	25.0	47.9	20.0	8.0	8.
BGE-Large	38.5	51.0	46.5	22.0	65.0	13.0	60.0	44.8	13.0	56.0	68.0	66.0	66.0	44.8	8.0	22.0	46.6	23.0	8.0	8.0
OpenAI-Embed	39.8	51.0	51.0	22.0	60.0	19.0	62.9	51.0	6.0	53.0	59.0	67.0	73.0	51.0	13.0	33.0	52.1	30.0	10.0	10.0
GTE-Qwen2-1.5B	43.8	54.0	50.0	24.0	69.0	25.0	65.7	65.6	41.0	63.0	79.0	70.0	84.0	65.6	20.0	40.0	47.9	33.0	12.0	10.0
E5-mistral-7B	45.7	60.0	50.5	17.0	67.0	14.0	67.1	64.6	36.0	68.0	74.0	54.0	78.0	64.6	30.0	47.0	43.8	41.0	12.0	38.0
GritLM-7B	46.2	63.0	49.5	29.0	69.0	17.0	85.7	62.5	46.0	60.0	74.0	53.0	87.0	62.5	21.0	46.0	43.8	33.0	12.0	42.0
ZeroGR-3B	40.4	37.0	36.5	24.0	51.0	13.0	38.6	57.3	36.0	41.0	81.0	61.0	81.0	57.3	12.0	40.0	52.1	11.0	12.0	22.0
Dataset										Uns	seen Su	bset								
Model	ACOR.	CPCD	CORE	MB.	PM.	PM.A	CliDS	CliT23	DD	Table (QuanT	PoRec l	Monant	NCL. N	NCL.T	Legal	Geno.	Touche	CliT21 N	News2
BM25	32.8	1.0	37.5	83.8	53.9	6.5	28.3	51.4	15.6	10.0	86.9	24.5	67.4	50.7	22.2	45.0	52.8	59.2	33.3	10.9
Contriever	40.4	1.0	52.5	89.2	32.9	0.0	6.7	37.8	21.3	8.3	76.8	46.7	65.0	60.0	34.2	35.0	27.8	52.0	32.7	23.
GTR-T5-base	31.3	1.0	47.5	89.2	36.8	1.6	13.3	36.5	13.6	12.5	77.8	37.7	70.0	48.0	22.2	40.0	25.0	55.1	29.3	15.0
GTR-T5-large	34.8	3.0	60.0	91.9	31.6	1.6	8.3	39.2	16.2	10.0	78.8	48.7	68.0	53.3	23.9	40.0	25.0	65.3	37.3	19.
E5-Base	40.4	3.0	42.5	81.1	43.4	6.5	11.7	39.2	13.2	5.8	78.8	54.2	72.0	55.3	27.4	35.0	33.3	37.8	36.0	14.
E5-Large	38.4	3.0	45.0	86.5	36.8	4.8	15.0	40.5	12.3	7.5	80.8	52.5	71.0	55.3	47.9	30.0	38.9	41.8	32.0	17.
BGE-Base	39.4	0.0	45.0	91.9	48.7	0.0	30.0	31.1	16.4	8.3	81.8	44.4	70.0	57.3	42.7	20.0	36.1	41.8	41.3	19.9
BGE-Large	36.9	0.0	52.5	94.6	42.1	3.2	23.3	28.4	17.8	5.0	81.8	50.4	74.0	54.0	40.2	60.0	38.9	51.0	41.3	15.
OpenAI-Embed	32.8	1.0	55.0	86.5	46.1	3.2	25.0	44.6	13.5	12.5	86.9	47.3	76.0	57.3	49.6	45.0	33.3	52.0	38.0	13.
GTE-Qwen2-1.5B	37.9	5.1	70.0	81.1	14.5	8.1	20.0	17.6	14.8	14.2	85.9	58.6	74.2	61.3	51.3	40.0	61.1	65.3	31.3	23.
E5-mistral-7B	41.9	5.0	60.0	83.8	43.4	1.6	46.7	48.6	19.4	11.7	83.8	66.1	71.0	62.7	58.1	45.0	36.1	58.2	46.7	25.
GritLM-7B	35.4	7.0	65.0	70.3	59.2	0.0	28.3	45.9	14.8	10.8	86.9	72.2	77.0	63.3	50.4	45.0	33.3	57.1	47.3	22.

Table 5: nDCG@10 on BEIR benchmark datasets.

Category	Method	Avg.	ArguAna	SciFact	NFCorpus	FiQA	SciDocs	Touche
Sparse	BM25	40.36	32.70	65.10	32.70	24.80	12.40	59.00
DR	Contiever	45.88	32.10	70.30	33.90	35.50	15.70	42.50
DR	GTR	43.68	32.60	58.60	32.60	34.40	12.00	48.00
DR	GritLM-7B	45.75	41.70	76.70	36.90	44.10	19.50	21.50
DR	OpenAI Embed	52.36	37.10	73.10	37.60	48.50	21.20	47.50
GR	GENRE	22.98	42.50	42.30	20.00	11.60	6.80	_
GR	GENRET	41.12	34.30	63.90	31.60	30.20	14.90	_
GR	GLEN	16.75	17.60	-	15.90	_	_	_
GR	TIGER (Llama-3B)	33.61	14.00	37.00	39.50	16.00	14.00	58.10
GR	ZeroGR-3B	46.19	35.40	72.80	34.70	34.10	18.70	37.50
Category	Method	Avg.	TREC-News	Fever	Robust04	Quora	Covid	CQADupStack
Sparse	BM25	40.36	20.70	58.30	18.50	73.80	58.30	28.00
DR	Contiever	45.88	27.30	90.60	26.60	86.60	59.60	29.90
DR	GTR	43.68	22.40	83.20	26.80	88.60	56.04	28.92
DR	GritLM-7B	45.75	31.00	84.90	29.90	66.80	50.20	_
DR	OpenAI Embed	52.36	26.20	92.80	32.10	89.90	78.20	44.10
GR	GENRE	22.98	_	_	_	_	14.70	_
GR	GENRET	41.12	_	-	_	-	71.80	_
GR	GLEN	16.75	_	_	_	-	_	_
GR	TIGER (Llama-3B)	33.61	16.40	-	15.80	59.60	65.70	_
GR	ZeroGR-3B	46.19	23.50	86.67	25.50	76.74	73.45	35.20

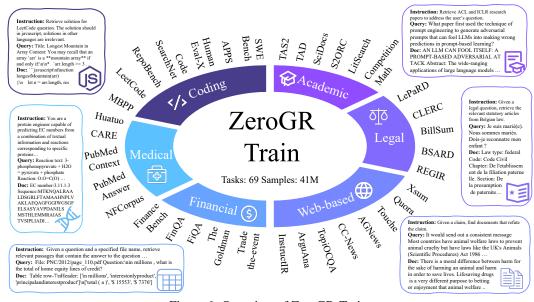


Figure 6: Overview of ZeroGR-Train.

Table 6: Dataset statistics grouped by domain and sorted by sample count

Dataset	Samples	Dataset	Samples					
	Acadeı	nic						
S2ORC-title-citation	100,000	TAD	208,255					
S2ORC-abstract-citation	100,000	TAS2	107,700					
S2ORC-title-abstract	100,000	StackMathQA	47,142					
ProofWiki-Proof	15,520	ProofWiki-Reference	2,098					
ProofWiki-Proof	15,520	ProofWiki-Reference	2,098					
Stacks-Proof	10,928	Stacks-Reference	9,022					
Stacks-Reference	9,022	Competition-Math	7,500					
Competition-Math	7,500	SciDocs	900					
SciFact	809	LitSearch	146					
	Code	ė						
CodeSearchNet	1,880,853	CodeEditSearch	21,395					
SWE-Bench	18,817	RepoBench	16,655					
HF-API	8,191	TLDR	6,414					
TensorAPI	6,190	APPS	5,000					
LeetCode	2,260	Conala	1,794					
PyTorchAPI	837	HumanEval-X	720					
MBPP	374							
	Finan	ce						
USnews	9,999	FinQA	6,251					
FiQA	5,500	HC3Finance	3,104					
ConvFinQA	3,037	TheGoldman	1,512					
TAT-DQA	1,012	Trade-the-event	900					
	Lega	l						
LePaRD	22,734,882	CLERC	327,414					
BillSum	18,949	REGIR-UK2EU	2,100					
REGIR-EU2UK	2,000	BSARD	886					
CUAD	717							
	Medic	eal						
PubMedQA-Context	196,696	PubMedQA-Answer	196,696					
Huatuo	25,371	NFCorpus	2,590					
CARE	77	1						
Web								
Reddit	12,704,958	AGNews	1,157,745					
CC-News	708,241	Xsum	204,045					
zsRE	147,909	ToT	109,454					
Fever	109,810	WoW	63,734					
TopiOCQA	45,450	AY2	18,395					
CQADupStack	13,045	InstructIR	9,806					
Quora	9,900	WnCw	5,499					
TREx	4,900	ExcluIR	3,352					
NevIR	1,896	ArguAna	1,306					