

Constrained Auto-Regressive Decoding Constrains Generative Retrieval

Shiguang Wu

Shandong University
Qingdao, China
shiguang.wu@mail.sdu.edu.cn

Zhaochun Ren*

Leiden University
Leiden, The Netherlands
z.ren@liacs.leidenuniv.nl

Xin Xin

Shandong University
Qingdao, China
xinxin@sdu.edu.cn

Jiyuan Yang

Shandong University
Qingdao, China
jiyuan.yang@mail.sdu.edu.cn

Mengqi Zhang

Shandong University
Qingdao, China
mengqi.zhang@sdu.edu.cn

Zhumin Chen

Shandong University
Qingdao, China
chenzhumin@sdu.edu.cn

Maarten de Rijke

University of Amsterdam
Amsterdam, The Netherlands
m.derijke@uva.nl

Pengjie Ren*

Shandong University
Qingdao, China
jay.ren@outlook.com

Abstract

Generative retrieval (GR) seeks to replace traditional search index data structures with a single large-scale neural network, offering the potential for improved efficiency and seamless integration with generative large language models. As an end-to-end paradigm, GR adopts a learned differentiable search index to conduct retrieval by directly generating document identifiers through corpus-specific constrained decoding. The generalization capabilities of generative retrieval on out-of-distribution corpora have gathered significant attention. Recent advances primarily focus on the problems arising from training strategies, and addressing them through various learning techniques. However, the fundamental challenges of generalization arising from constrained auto-regressive decoding still remain unexplored and systematically understudied.

In this paper, we examine the inherent limitations of constrained auto-regressive generation from two essential perspectives: *constraints* and *beam search*. We begin with the Bayes-optimal setting where the GR model exactly captures the underlying relevance distribution of all possible documents. Then we apply the model to specific corpora by simply adding corpus-specific constraints. Our main findings are two-fold: (i) For the effect of constraints, we derive a lower bound of the error, in terms of the KL divergence between the ground-truth and the model-predicted step-wise marginal distributions. This error arises due to the *unawareness* of future constraints during generation and is shown to depend on the average Simpson diversity index of the relevance distribution.

(ii) For the beam search algorithm used during generation, we reveal that the usage of marginal distributions may not be an ideal approach. Specifically, we prove that for sparse relevance distributions, beam search can achieve perfect top-1 precision but suffer from poor top- k recall performance. To support our theoretical findings, we conduct experiments on synthetic and real-world datasets, validating the existence of the error from adding constraints and the recall performance drop due to beam search. This paper aims to improve our theoretical understanding of the generalization capabilities of the auto-regressive decoding retrieval paradigm, laying a foundation for its limitations and inspiring future advancements toward more robust and generalizable generative retrieval.

CCS Concepts

• Information systems → Retrieval models and ranking.

Keywords

Generative retrieval; Constrained decoding; Beam search

ACM Reference Format:

Shiguang Wu, Zhaochun Ren, Xin Xin, Jiyuan Yang, Mengqi Zhang, Zhumin Chen, Maarten de Rijke, and Pengjie Ren. 2025. Constrained Auto-Regressive Decoding Constrains Generative Retrieval. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '25)*, July 13–18, 2025, Padua, Italy. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3726302.3729934>

1 Introduction

The advent of generative models has catalyzed the emergence of generative retrieval (GR) as a new paradigm in information retrieval. GR provides a potential way to replace the conventional index structure, such as inverted index and vector-based index, with a single large-scale neural network [34]. By integrating the retrieve-then-rank pipeline into an end-to-end framework, GR offers the promise of enhanced efficiency. Typically, GR adopts auto-regressive generative models, e.g., BART [23] and T5 [43], trained to generate document identifiers (docIDs) given a query.

*Corresponding authors.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '25, Padua, Italy

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-1592-1/2025/07
<https://doi.org/10.1145/3726302.3729934>

Generalization in neural information retrieval. Generalization is a key problem in neural information retrieval models [16, 18, 50, 51, 56, 56, 59]. Similar to the scaling principle in large language models (LLMs) [7], many dense retrieval approaches consider training on a high-resource dataset and then evaluated on different domains [18, 39, 44, 50, 51]. In particular, GTR [39] demonstrates significant improvements in out-of-domain performance by successfully scaling up the model size and the training corpus. Unlike dense retrieval, GR, as a retrieval paradigm itself, is much more concerned with the generalization abilities to unseen corpora [2, 12, 26, 46].

Bayes-optimal generative retrieval. Given the success of generative LLMs, GR is expected to capture the universal relevance distribution when trained on sufficiently large retrieval datasets. While extensive studies have already proposed effective training strategies at large scale [12, 26, 61, 62], the performance of an ideal, *viz.* a Bayes-optimal GR model, on unseen corpora has not been systematically studied. In this paper, we study an amortized Bayes-optimal auto-regressive GR model that fully encapsulates the underlying relevance distribution over the complete corpus containing all possible documents. On top of this, we apply the constrained auto-regressive generation process to provide a valid docID in any given downstream corpus which is a subset of the complete one.

Constrained beam search and generalization. The core components of GR are the differentiable index and the generation of docID [54]. In this paper, we focus on how the generation process affects the generalization of GR. Most existing GR models adopt constrained beam search to generate the top- k docIDs as the default retrieval strategy [46, 47, 49, 52]. However, Zeng et al. [62] point out the pitfalls of this strategy, *i.e.*, a greedy local search algorithm is likely to prune false negative docIDs and may thus not be sufficient for developing effective GR models. While these effects are typically entangled with model errors in relevance prediction, the impact of constrained beam search using a Bayes-optimal model with the correct relevance distribution still remains unclear. Within our setting, we address this gap by isolating and analyzing the theoretical root of limitations of constrained auto-regressive beam search.

Main findings. We study the inherent limitations from two essential perspectives: *constraints* and *beam search*. (i) **Constraints:** In Section 4.1, we derive a lower bound on the error, in terms of the KL divergence between the ground-truth and model-predicted step-wise marginal distributions. This error is influenced by the concentration of the relevance distribution and arises from the model’s lack of awareness of future constraints during generation. (ii) **Beam search:** In Section 4.2, we show that the usage of marginal distributions may not be suitable. We prove that for some sparse relevance distribution with a *thick* tail, beam search can achieve perfect top-1 precision, but poor top- k recall performance. We are aware of existing techniques on improving recall of GR model, and summarize two directions based on our analysis: (i) *aggregation*, *i.e.*, aggregating the relevant documents within a narrow branch, and (ii) *amplification*, *i.e.*, amplifying the scores of relevant documents to be significantly large. Both directions are essentially meant to provide a high concentration prior on the relevance distribution. We show that although they are empirically beneficial for small scale corpus, difficulty will arise if we aim to construct a Bayes-optimal

GR. Our results provide theoretical grounding of constrained beam search in GR, and shows the inherent limitations of this retrieval strategy towards a reliable Bayes-optimal GR model. Our results also imply the importance of balancing concentration during model training for mitigating this problem.

Contributions. Our main contributions are as follows: (i) We provide theoretical results concerning Bayes-optimal GR to determine how constrained beam search affects the generalization. (ii) We decompose the negative effect from two angles, constraints and beam search, and identify a trade-off factor, *i.e.*, the concentration of relevance distribution. (iii) Our theoretical results are verified by experiments on synthetic and real-world datasets.

2 Related work

Generative retrieval (GR) is an emerging direction in neural information retrieval, exploring the possibility of replacing traditional index structures in retrieval systems with a single large-scale neural networks [27, 54]. It leverages generative models to directly generate the relevant docIDs given a query. This paradigm originated with Cao et al. [9], Metzler et al. [34] and has garnered considerable attention [4, 30, 38, 46–49, 52, 55, 57, 61, 66, 67, 70] in the information retrieval community.

Generalization remains a challenge for GR, especially when applied to out-of-distribution corpora [2, 11, 22, 27, 31, 32, 46]. Previous research attributes this challenges to limited model capacity [22, 60], lack of learning in the docID construction [46, 57, 66], and difficulties in learning semantic representations [47, 53]. In contrast, our work focuses on the constrained auto-regressive decoding strategy widely applied in GR, which is crucial for adapting GR models to new corpora dynamically. Our setting aligns closely to the few-shot indexing approach [2], where a pre-trained LLM generates docIDs solely based on its pre-trained knowledge and generalization capabilities, without additional training. We treat their method as a conceptual blueprint for a fully generalizable GR system and aim to investigate the inference stage under this setting.

Updatable generative retrieval is another critical task on dynamic corpora. The primary challenges indicate the cost of updating the model with new documents and the catastrophic forgetting problem [11, 17, 21, 33]. Previous efforts have concentrated on developing efficient continual learning strategies by fixing the indexing construction procedure. We consider an idealized scenario where the model has full knowledge of all possible documents and focuses solely on generating relevant docIDs on dynamic corpus.

Constrained decoding. Constrained decoding has been widely studied for guiding machine learning models to produce outputs that satisfy specific conditions [1, 37]. Instead of learning to satisfy through training, constraints are more often preferable only in the inference time due to the flexibility and efficiency. Nishino et al. [40, 41] demonstrate the preservation of relative errors of certain loss functions in realizable setting. We instead provide a failure case via establishing the existence of a lower-bound error for auto-regressive models operating under step-wise inference-time constraints. Recent work on controllable text generation (CTG) in LLMs [see, *e.g.*, 65] also explores imposing constraints during inference without updating the underlying model [10, 19, 35, 36].

However, many of these approaches do not focus on strictly enforcing constraint satisfaction. A few studies [20, 63, 64] propose methods to produce outputs that strictly adhere to constraints, mainly hard keywords inclusion constraints, using tractable probabilistic models or policy gradient techniques. Our work differs by focusing on a specific corpus-level constraint, i.e., the set of valid docIDs is sampled from the complete corpus, a problem unique to retrieval tasks.

Beam search is a widely used heuristic algorithm for decoding structured predictors and has been applied as a non-maximum inner product search (MIPS) setup for large-scale retrieval systems with explicit tree structures [24, 68, 69, 71]. Beam search is known to have a performance deterioration, and only few works provided theoretical insights into this issue. As far as we know, only Zhuo et al. [71] demonstrate a training-test discrepancy in tree-structured models using binary cross-entropy loss. They showed that pseudo-labeling during training does not guarantee that beam search will get the most relevant targets. In our work, we analyze the marginal distribution of an auto-regressive distribution and provide a theoretical result on the top-1 and top- k performance under sparse relevance situations. Zhuo et al. [71] also provide a Bayes-optimal tree structure, which is often called max-heap assumption [24], and we will discuss the difficulty of enforcing this assumption in our setting in Section 5. In GR, some work have reached the same conclusion that beam search is not sufficient for retrieval as it is likely to prune the relevant docIDs and the model is not able to recover from this [26, 28, 62]. They propose to use a hybrid retrieval strategy to help bypassing this problem. We instead focus on understanding the root cause of this problem, i.e., the usage of marginal distribution.

3 Preliminaries

We formulate GR and introduce key notations in this section.

Generative retrieval. Following Tay et al. [49], we formulate GR where the mapping from documents and docIDs is one-to-one function. A corpus, denoted as \mathcal{D} , is a set of documents d , with each document represented as a sequence of tokens, i.e., $d = (d_1, \dots, d_m)$, where m is the length. In this paper, we assume all documents have the same length. A generative retrieval model f , typically implemented using a sequence-to-sequence architecture such as T5 [43] or BART [23], generates a ranked list of the most relevant docIDs in \mathcal{D} for the given query. The ranking list is computed through beam search during generation. To ensure the model reliably generates valid docIDs from the corpus, a constrained auto-regressive decoding process g together with the beam search is used.

Bayes-optimal generative retrieval. We first assume the complete corpus \mathcal{D} contains all possible documents of length m . Any downstream corpus is therefore a subset of \mathcal{D} . We denote f as the Bayes-optimal GR model on \mathcal{D} which has the ability to predict the ground-truth relevance distribution over \mathcal{D} given any query. The Bayes-optimal model f is considered as an ideal and oracle prototype model which helps us understand the behavior of the generation process. When f is applied to a downstream corpus $\mathcal{D}^c \subset \mathcal{D}$, it uses the corresponding constrained decoding process g^c to predict relevant documents in \mathcal{D}^c . This induced GR model on

Table 1: Glossary.

Symbol	Description
k, m	the vocabulary size and document length
\mathcal{D}	the complete corpus, i.e., $[k]^m$
\mathcal{D}^c, C	downstream corpus and constraints
f, f^c	Bayes-optimal and induced downstream GR
$\Pr(\cdot q)$	relevance distribution given query q
$d = (d_1, \dots, d_m)$	document in \mathcal{D}

\mathcal{D}^c is denoted as f^c . Note that a similar setting has recently been proposed as zero-shot indexing [2].

Notation. Table 1 lists the main notation used in the paper. For an integer n , we denote the set $\{1, \dots, n\}$ by $[n]$. We use bold face to denote random variables. Tokens in a document d are integers from $[k]$, making $d \in [k]^m$, and k is the vocabulary size. We set the complete corpus $\mathcal{D} = [k]^m$. We denote the underlying relevance distribution over \mathcal{D} given a query q as $\Pr(\cdot | q)$, where $\mathbf{d} \sim \Pr(\mathbf{d} | q) = \prod \Pr(d_i | \mathbf{d}_{<i}, q)$. For simplicity, we focus on a single query q and omit it in some context. We use subscripts to indicate a sliced subset or operations at specific steps, i.e., \square_i , and $\square_{\geq i}$, etc. Particularly, we refer to $\mathcal{D}_{\geq i}$ or $\mathcal{D}_{\geq i}^c$ a branch with root d_i , under some prefix $d_{<i}$.

Downstream corpus and constraints. We construct the downstream corpus \mathcal{D}^c by sampling. For simplicity, each document is sampled with an equal probability p . In practice, the sampling is agnostic to the future user queries, and thus independent with $\Pr(\cdot)$. We use C to be the event that \mathbf{d} is in \mathcal{D}^c , and $\Pr(\cdot | C)$ is the distribution under the corpus \mathcal{D}^c . C_i means the i -th token of document \mathbf{d} is valid with respect to the downstream corpus constraints, given some context-clear prefix tokens $d_{<i}$, i.e., d_i is valid if it appears in \mathcal{D}^c for the prefix $d_{<i}$. We use “constraints”, and “downstream corpus” to represent the result of sampling interchangeably. The constrained generation process g_i^c is applied at the i -th step of f . It first zeros out the invalid tokens and then re-normalizes the remaining probabilities.

4 Theoretical analysis

We investigate the inherent limitations of GR arising from constraints and beam search individually. Our analysis disentangles these factors to isolate their individual effects. For constraints, we analyze its effect on the marginal distribution at each generation step. For beam search, we analyze how it independently degrades recall performance on the complete corpus, disregarding constraints. These results are asymptotic with respect to the vocabulary size k . In Appendix A, we examine the required magnitude of k for a sufficiently expressive model and show that it needs to be exponentially large in terms of the ratio of raw document length and docID length.

4.1 Constraints cause marginal distribution mismatch

We begin by studying the effects of applying constraints to the model performance. We first identify the factor that causes the error during the generation. Then we quantitatively analyze the magnitude of this error in (i) uniform, and (ii) general relevance distribution given some query q . We consider the first generation

step without loss of generality. Other cases can be reduced to it by adjusting the document length m or vocabulary size k .

Unawareness of future constraints. Since f^c is *unaware* of the constraints in the future steps, there may exist biases between the distributions of complete corpus \mathcal{D} and downstream corpus \mathcal{D}^c . Specifically, after applying constrained decoding g^c , $f^c(q)$ returns

$$\Pr(\mathbf{d}_1 \mid q, C_1) = g^c[\Pr(\mathbf{d}_1 \mid q)] \quad (1)$$

$$\propto \mathbb{I}[d_1 \in \mathcal{D}_1^c] \sum_{\forall d_{>1} \in \mathcal{D}_{>1}} \Pr(d_1 d_{>1} \mid q), \quad (2)$$

where the C_1 constraint is satisfied through g^c , and $\mathbb{I}[d_1 \in \mathcal{D}_1^c]$ means d_1 is a valid first token in the downstream corpus. In the contrary, the ground-truth marginal distribution only sum over the documents in \mathcal{D}^c , so we have

$$\Pr(\mathbf{d}_1 \mid q, C_1, \mathcal{C}_{>1}) = g^c[\Pr(\mathbf{d}_1 \mid q, \mathcal{C}_{>1})] \quad (3)$$

$$\propto \mathbb{I}[d_1 \in \mathcal{D}_1^c] \sum_{\forall d_{>1} \in \mathcal{D}_{>1}^c} \Pr(\mathbf{d}_1 d_{>1} \mid q), \quad (4)$$

where $\Pr(\mathbf{d}_1 \mid q, C_1, \mathcal{C}_{>1}) = \Pr(\mathbf{d}_1 \mid q, C)$. Here we use the **red** mark to highlight the differences from Eq. 2. Note that this gap would not arise if the downstream corpus were preset, with both training and inference performed on it, as the model would learn $\Pr(\mathbf{d} \mid q, C)$ directly. We then analyze the Kullback–Leibler (KL) divergence as follows,¹

$$\text{KL}[\Pr(\cdot \mid C) \parallel \Pr(\cdot \mid C_1)] \quad (5)$$

$$= \mathbb{E}_{\mathbf{d}_1 \sim \Pr(\cdot \mid C)} \left[\log \frac{\Pr(\mathbf{d}_1 \mid C)}{\Pr(\mathbf{d}_1 \mid C_1)} \right] \quad (6)$$

$$= \mathbb{E}[\log \Pr(C_{>1} \mid \mathbf{d}_1)] - \log \Pr(C_{>1} \mid C_1). \quad (7)$$

In Eq. 7, we see that the KL divergence is the gap between the average proportion of constraints locally within each branch $\mathcal{D}_{\geq 1}^c$ and the global average constraints. This indicates that the variation across branches may contribute the mismatch. Recall that we construct our constraints via an i.i.d. sampling, and we have the expectation of the constrained marginal distribution $\Pr(d_1 \mid q, C_1)$ as follows:

$$\mathbb{E}_C[\Pr(d_1 \mid q, C)] \propto \mathbb{E}[\mathbb{I}[\mathbf{d} \in \mathcal{D}^c]] \Pr(d_1 \mid q) \propto \Pr(d_1 \mid q). \quad (8)$$

Therefore, the downstream corpus will follow the same distribution as the complete corpus on average. However, as we will see in the next, the gap is not concentrated at zero with high probability in terms of KL divergence.

Uniform relevance distribution. We first give an example case when the relevance distribution $\Pr(\mathbf{d} \mid q)$ is a uniform distribution. If the size of the downstream corpus is $k^r \ll k^m = |\mathcal{D}|$, we derive an asymptotic lower bound of the error for large k as follows

$$\text{KL}[\Pr(\cdot \mid C) \parallel \Pr(\cdot \mid C_1)] \gtrsim \frac{0.05}{k^{r-1}}. \quad (9)$$

In particular, we have a constant error 0.05 if the size is $O(k)$. For detailed illustration of the proof, see Theorem B.2. Here is the intuition: GR model will predict a uniform distribution over the

¹The reason we adopt the KL divergence is that it is not only part of the training loss, i.e., empirical cross-entropy, but also related to the ranking performance. Several publications have showed that the cross-entropy is a bound of several commonly used metrics, e.g., Normalized Discounted Cumulative Gain (nDCG) and Mean Reciprocal Rank (MRR) [8, 58] for binary relevance score. KL divergence and cross-entropy only differ by the entropy.

valid first tokens, but the ground-truth should be proportional to the number of valid documents in each valid branch. Due to the variance of the sampling, the ground-truth distribution will not be exactly uniform.

General relevance distribution. Following the same idea of the uniform case, we further give the result for general relevance distributions. The key is the Simpson diversity index [45], which is used for measuring the degree of concentration. We introduce the average Simpson diversity index. It is computed as the squared expectation of root sum of squared probabilities $\Pr(d \mid d_1)$, i.e.,

$$\mathbb{E}_{d_1}^2 \left[\sqrt{\sum_d \Pr(d \mid d_1)^2} \right]. \quad (10)$$

Recall that each document is selected with probability p , and let A be the average Simpson diversity index, we have an asymptotic lower bound of the KL divergence for large k in Eq. 11.

$$\text{KL}[\Pr(\cdot \mid C) \parallel \Pr(\cdot \mid C_i)] \gtrsim \frac{0.05A}{p}. \quad (11)$$

For detailed illustration of the proof, see Theorem B.3. We also show that the lower bound reaches its minimum in the uniform relevance distribution case. Note that we have $A \geq \frac{k}{|\mathcal{D}|}$, where the equality is obtained when $\Pr(\cdot)$ is uniform. Let the selected corpus size be $k^r \approx p|\mathcal{D}|$, we have the same error shown in Eq. 9.

$$\frac{0.05A}{p} \approx \frac{0.05k}{|\mathcal{D}|p} = \frac{0.05}{k^{r-1}}. \quad (12)$$

Note that for concentrated distribution, the Simpson diversity index considerably exceeds p , resulting in a corresponding larger error.

In conclusion, we infer a lower bound of the KL divergence between the predicted and ground-truth marginal distribution. The bound is proportional to the degree of concentration of the underlying relevance distribution.

4.2 The impact of beam search on recall

Beam search often fails to retain the correct top- k prefix candidates, resulting in the exclusion of relevant documents during generation [62]. We attribute this limitation to the usage of conditional decomposition of the joint distribution over the corpus.

To formalize this analysis, we model the decomposition as a tree where: (i) each node at layer i represents a token d_i generated under a specific prefix $d_{<i}$, (ii) each sub-tree represents possible continuations of a prefix, and (iii) the value of node d_i with prefix $d_{<i}$, denoted $V(d_i \mid d_{<i})$, equals $\sum_{d_{\geq i}} \Pr(d_{\geq i} \mid d_{<i})$, i.e., the marginal probability $\Pr(d_i \mid d_{<i})$. The property of this structure is that node values represent marginal probabilities aggregated over all possible future paths. However, the objective of the retrieval model requires identifying specific document with maximal joint probability $\Pr(\mathbf{d})$. This creates a fundamental mismatch: nodes with high values $V(d_i \mid d_{<i})$ may not contain the highest-joint-probability documents in their sub-trees.

Non-relevant branches overtaking relevant ones. We setup a scenario with a sparse relevance distribution to elucidate this issue. We present how branches containing relevant documents can be overtaken by non-relevant peer branches during generation. At the first generation step, the model have to choose within k nodes, each associated with a sub-tree of k^{m-1} documents (leaves). We

assign a logit uniformly sampled from $[-1, 1]$ to each document in the corpus. A subset of $\lambda k \ll |\mathcal{D}|$ documents is randomly selected as relevant, and each is assigned a logit within $O(\log k \pm \log \log k)$. The exponential of each logit is the final relevance score of the document, i.e., the score of the corresponding root-to-leaf path.

We prove that the recall of the top- λk valued branches is upper-bounded by $0.5 + o(1)$ with high probability. This indicates that many relevant branches are excluded from the highest-valued branches. However, the top-1 branch is highly likely to contain the most relevant documents. A detailed formal statement and proof of this result is provided in Theorem C.1 in Appendix C. The thickness of the tail distribution or the sharpness of relevant branches determines the probability of overtaking. If the scores for relevant branches are sufficiently high, this issue becomes less pronounced.

In summary, GR models relying on the sum of sub-tree values for ranking branches struggle to achieve high recall performance while maintaining top-1 precision.

4.3 Concentration as a trade-off factor

We find a common dependence on the concentration of relevance distribution from both components, i.e., the Simpson diversity index, and the thickness of the tail distribution. However, the concentration takes effect in opposite directions. Less concentrated distribution is more stable against sampling and thus have better alignment with the ground-truth constrained marginal distribution. In contrast, high concentration reduces the false positive branches and thus improve the recall performance using beam search, but it requires an accurate marginal distribution.

5 Difficulty of concentration in Bayes-optimal GR

While increasing the degree of concentration is a promising approach to improve recall performance on the complete corpus—if the effects of constraints are ignored—achieving this in the context of Bayes-optimal GR is still challenging. In this section, we show details about the challenge in the context of Bayes-optimal GR in two directions:

- (1) **Aggregation:** Aggregating relevant documents within a narrow branch to ensure that its value significantly exceeds that of others throughout the generation process.
- (2) **Amplification:** Amplifying the scores of relevant documents relative to non-relevant ones, so that the branch value serves as a reliable sketch for branch quality.

While these strategies have shown their effectiveness [24, 61] on small scale, in-domain datasets, they are much more challenging to implement in Bayes-optimal GR models. Specifically, as we will discuss, achieving this concentrated structure requires additional computational and data resources beyond those needed to obtain the Bayes-optimal model itself. Detailed discussions are shown in the following.

5.1 Aggregation introduces redundancy

Although the corpus size $|\mathcal{D}| = k^m$ is sufficiently large, the model does not necessarily fully utilize the entire code space. Here we

study the entropy of the corpus distribution $\Pr(\cdot)$, which is defined as $H(\mathbf{d}) = -\sum_{d \in \mathcal{D}} \Pr(d) \log \Pr(d)$. It can be decomposed into the entropy of the marginal distribution at each step, i.e., $H(\mathbf{d}) = \sum_{i=1}^m H(d_i | d_{<i})$. When the distribution is uniform, the entropy is maximized, i.e., $H(\mathbf{d}) = \log |\mathcal{D}| = m \log k$. If we introduce external prior knowledge to aggregate the relevant documents, the relevant branches will stand out at a very early stage of generation, and the entropy at that layer $H(d_i | d_{<i})$ will be low and even approaching zero. The more aggregated the relevant documents are, the lower the entropy becomes. In other words, for the same corpus size, the model will waste more code space on redundant structures. In practice, this approach for concentration is often realized by conducting hierarchical clustering on the corpus, which has been shown to be effective in small scale retrieval tasks. The trade-off between concentration and redundancy can be ignored to some extent when the corpus is small. However, if we would like to build a Bayes-optimal GR model, learning a sufficiently large code space is already expensive, see Appendix A. Effective concentration will introduce more redundancy, which is computationally inefficient.

5.2 Amplification requires high-quality data

We treat the amplification strategy as an approximation of the max-heap structure discussed in Li et al. [24]. In this structure, the value of each node is the maximum value of its children instead of the sum, i.e., $V(d_i | d_{<i}) = \max_{d_{>i}} \Pr(d_i, d_{>i} | d_{<i})$. One can prove that this structure can achieve perfect recall performance by preserving the relevant documents in the top- n branches [71]. Note that this structure is no longer a chain decomposition of the joint distribution, and the distribution at each step is different from the original distribution. The new distribution can still be learned through empirical risk minimization by carefully filtering the training data. Considering the i -th step, the negative log-likelihood loss is as follows:

$$\mathcal{L}^i(\theta) = - \sum_{d \in \tilde{\mathcal{D}}} \log \Pr(d_i | d_{<i}; \theta), \quad (13)$$

where $\tilde{\mathcal{D}}$ is the training set, and θ is the model parameter to be optimized. In order to learn the max-heap structure, the predicted distribution $\Pr(d_i | d_{<i}; \theta)$ should be proportional to $V(d_i | d_{<i})$. Therefore the loss function should filter out the non-maximum successors in the training set, i.e.,

$$\mathcal{L}^i(\theta) = - \sum_{d \in \tilde{\mathcal{D}}} \mathbb{I}[d \in \arg \max \Pr(d_{>i} | d_{\leq i})] \log \Pr(d_i | d_{<i}; \theta). \quad (14)$$

As we can see, only the most relevant successors are allowed to contribute to the loss function. It not only needs to throw away large amounts of data but also requires more high-quality data and careful filtering strategies. Most retrieval models can be seen to be trained on this loss [24], as they are only trained on the labeled relevant documents. However, the quality and quantity of retrieval data are often limited considering that the amount of data needed for a generative model is generally much larger than that for a discriminative model.

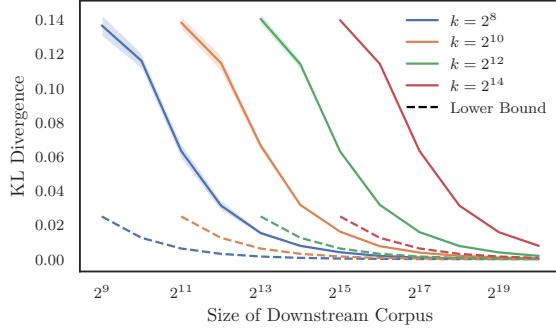


Figure 1: The KL divergence error in the first generation step on synthetic uniform relevance distribution data with uniformly sampled downstream corpus.

6 Synthetic experiments

In this section, we experimentally verify the theoretical results and investigate scenarios beyond the assumed data distributions presented in Section 4.1 and Section 4.2.

6.1 Effects of constraints

Uniform relevance distribution. We begin by simulating the case of a uniform relevance distribution, as discussed in Section 4.1. Since the lower bound is expressed in terms of the size of the downstream corpus, we vary its size to observe the behavior across different vocabulary sizes. We choose a sufficiently large m so that the complete corpus of size k^m is much larger than the downstream corpus. For a downstream corpus of size n , the sampling probability is given by $p = \frac{n}{k^m}$. The simulation results are shown in Figure 1.

As illustrated, the KL divergence for different k values exhibits a consistent decreasing trend. For a fixed downstream corpus size, a smaller k results in significantly lower error. This is expected, as a smaller k leads to each branch covering a larger number of selected documents, thereby reducing variance. In this case, the model-predicted marginal distribution closely follows the actual one.

General relevance distribution. Next, we simulate the case of general relevance distributions, as described in Section 4.1. For simplicity, we ensure that each branch has the same Simpson diversity index. We assign random weights to documents within the range $[1, e^{100}]$, where the assignment probability decreases exponentially with larger weights. The concentration of the distribution is controlled by varying the rate of decrease: slower rates result in less concentrated distributions, corresponding to smaller Simpson diversity indices. When the Simpson diversity index approaches 1, our lower bound no longer holds.

The simulation results are shown in Figure 2. Compared to the uniform distribution case, the KL divergence is significantly larger and decreases more slowly as the downstream corpus size increases. For data with a Simpson diversity index approaching 1, the KL divergence reaches approximately 6. In contrast, the uniform case consistently maintains a very low KL divergence. For example, the lowest Simpson diversity index in the figure is around $1e-06$, which matches the magnitude of the uniform distribution, $\frac{1}{2^{20}} \approx 1e-06$, and the corresponding KL divergence is approximately 0.19, as seen in Figure 1.

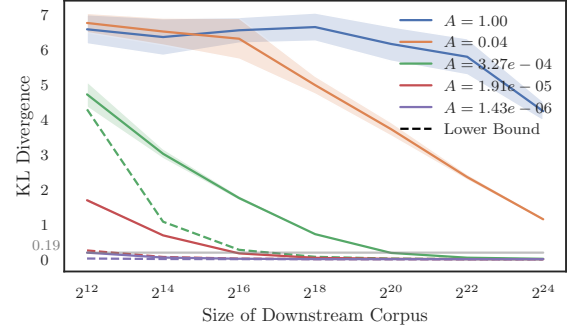


Figure 2: The KL divergence error in the first generation step on several synthetic relevance distribution data with different degrees of concentration. The vocabulary size is 2^{10} , the docID length is 3. A is the Simpson diversity index of the relevance distribution. As for highly concentrated distributions, e.g., $A = 1$, and $A = 0.04$, the Lyapunov's condition will no longer hold (see Theorem B.3 for more details), we do not draw their lower bounds.

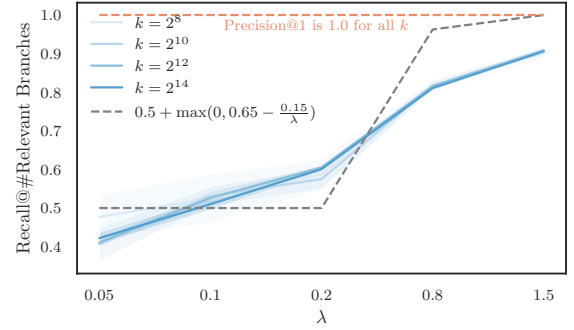


Figure 3: The recall of relevant branches cut off at the total number of relevant branches in the first generation step. The synthetic relevance distribution is constructed as Section 4.2. The total number of sampled relevant document is λk . The size of each branch is fixed as $n = 2^{25}$.

6.2 Effects of beam search

Verification of theoretical results. We first construct a sparse relevance distribution as outlined in Section 4.2. The results are presented in Figure 3. We evaluate the recall of relevant branches under varying values of λ , which controls the number of relevant documents sampled. Specifically, we compute how many relevant branches are preserved in the top positions during the first generation step. We also examine large λ values, which are beyond the scope of Theorem C.1.

For small λ , the recall performance is approximately 0.5. For larger λ values, the recall still aligns with the theoretical bounds. In all cases, the precision@1 remains consistently perfect. Regarding score magnitudes, as each branch has a fixed size of $n = 2^{25}$, relevant documents achieve scores around $e^{8.6} \approx 5400$, while non-relevant documents score around $e^1 \approx 2.7$.

Effects on different degrees of concentration. We also investigate the impact of varying the sharpness of the relevance distribution by introducing a temperature parameter T . For a document

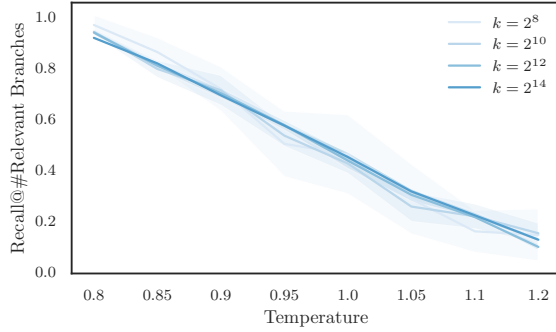


Figure 4: The recall of relevant branches cut off at the total number of relevant branches in the first generation step. The synthetic relevance distribution is constructed as Section 4.2. The total number of sampled relevant document is $0.05k$. The size of each branch is fixed as $n = 2^{25}$.

with logit s , the score is computed as $\exp(s/T)$. Lower temperatures increase the gap between scores for relevant and non-relevant documents. The results are shown in Figure 4.

Our findings reveal that recall performance is highly sensitive to temperature, achieving perfection within a narrow range. This confirms the advantage of constructing GR models that capture concentrated distributions effectively.

6.3 Summary

The results from the synthetic data distributions validate the theoretical findings presented in earlier sections. Although these synthetic settings are not practical for real-world scenarios, they provide a controlled environment to clearly demonstrate the negative effects of constrained decoding and beam search.

7 Experiments on real-world dataset

7.1 Experimental setups

Document identifier design. For the docID design, we adopt the codebook and semantic ID mapping from Zeng et al. [61], which introduces the first effective GR model that outperforms conventional retrieval models on the full MS MARCO passage corpus [3]. The codebook size is set to 256.

Datasets. We evaluate our approach using the MS MARCO V1 passage corpus [3], which contains 8.8 million passages, along with three evaluation datasets: (i) MS MARCO-dev, consisting of 7k queries; (ii) TREC DL 2019 [14], with 43 queries; and (iii) TREC DL 2020 [13], with 54 queries.

Relevance distribution. While recent advance improves GR performance [61], there remains a gap compared to state-of-the-art models. We use SLIM++ [25] to compute relevance scores due to its strong overall performance on the MS MARCO V1 passage re-ranking leaderboard.² We use the top-10,000 scores and extrapolate the remaining scores by taking 1% of the lowest score from the ranked list, adding small Gaussian perturbations.

7.2 Effect of constraints

We examine two sampling strategies: uniform sampling, as discussed in previous sections, and sampling directly based on the

²MS MARCO V1 Passage Regressions are available at [Pyserini](https://pyserini.org).

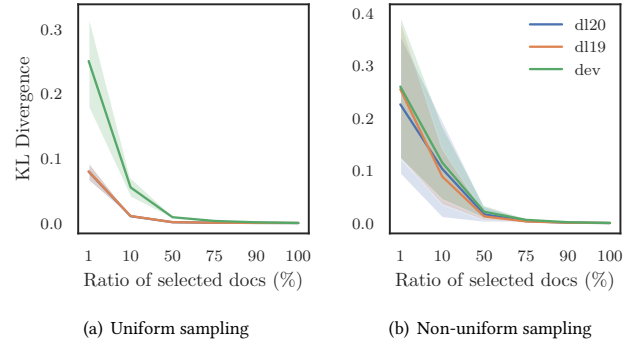


Figure 5: The KL divergence error in the first generation step on MS MARCO V1 passage corpus. A subset of top-10,000 rank list from SLIM++ [25] is treated as the complete corpus, and the downstream corpus is (a) uniformly sampled, or (b) sampled with the relevance distribution.

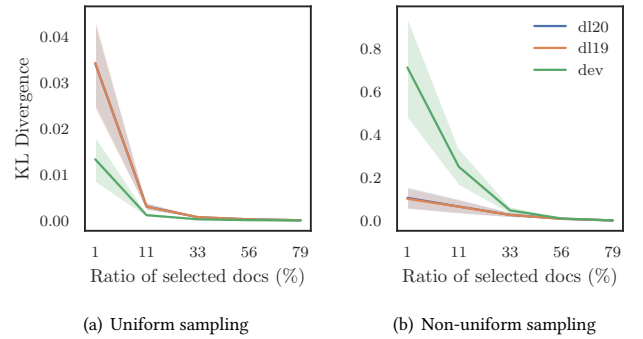


Figure 6: The KL divergence error in the first generation step on MS MARCO V1 passage corpus. The relevance distribution is constructed as Section 7.1. The downstream corpus is (a) uniformly sampled, or (b) sampled with the relevance distribution.

relevance distribution. In practice, sampling should remain agnostic to future relevance, but this analysis illustrates what occurs if the downstream corpus distribution aligns with the complete corpus.

For uniform sampling, Figures 5(a) and 6(a) depict the KL divergence across varying complete corpus sizes. We observe that the error decreases rapidly as the size of the downstream corpus increases. For small k , the Simpson diversity index in each branch is lower due to the higher number of elements, which helps reduce the KL divergence as defined in Eq. 11.

In contrast, for non-uniform sampling, Figures 5(b) and 6(b) show that the error is substantially higher than in the uniform case. Non-uniform sampling introduces greater variance and exacerbates the mismatch between distributions.

7.3 Effect of beam search

We evaluate recall performance by examining the retrieval of relevant branches. The cutoff is set to the number of relevant documents, which does not exceed the number of relevant branches. Results are presented in Table 2. The recall is relatively low, indicating that branches containing highly relevant documents may not always be

retrieved within the top- k branches. However, the top-1 precision is higher, particularly for the MS MARCO-dev evaluation set.

Table 2: Recall@50 and Precision@1 of relevant branches (256 in total) at the first generation step, which are denoted as “R@50” and “P@1”. The highest 50 documents are labeled relevant, and the branches contain these documents are labeled relevant. The relevance distribution is constructed in Section 7.1.

Model	R@50	P@1
TREC DL 19	53.7	69.8
TREC DL 20	63.3	75.9
MS MARCO-dev	67.5	90.5

7.4 Summary

Our experimental results demonstrate that the theoretical findings reveal the negative effects of constrained beam search under real-world dataset distributions to some extent. Since Zeng et al. [61] designed and trained the GR model specifically on the MS MARCO passage corpus, the docID structure is highly adapted to that query-document training distribution. Additionally, the size of the MS MARCO corpus is far from ideal for representing a complete corpus. As a result, the experimental findings only provide an approximate reflection of the theoretical results.³

8 Limitations

Concerning the theoretical aspects of our work, as we do not have an evidence of what an optimal GR model should look like, we fail to provide practical assumptions on the relevance distribution and the structure of docID. Our results are sensitive to the parameters and assumptions, e.g., the sharpness of the relevance distribution for Theorem C.1, and may not accurately reflect practical real-world situations. Besides, we have not studied how the two factors affect each other when using constrained beam search.

Concerning the experimental aspects of our work, we only use the docID from Zeng et al. [61] and the MS MARCO V1 passage corpus.

9 Conclusion

In this paper, we have provided theoretical results on the effect of constrained beam search for a Bayes optimal GR model. We have considered two separate aspects, constraints and beam search, and examine the root cause of the negative effect on generalization. Both aspects are intrinsically connected to the degree of concentration of the relevance distribution across the complete corpus. When it is more concentrated, the model achieves decent recall performance, provided the marginal distribution aligns closely with the actual one. However, applying downstream corpus constraints increases this marginal distribution gap at the same time. To validate our theoretical findings, we have conducted experiments on synthetic and real-world dataset and have shown the case beyond the assumed conditions in the theorem. Overall, we give a systematical investigation from both theory and experiments to the limitation of constrained decoding on retrieval generalization.

³The source code is available at <https://github.com/Furyton/constrained-generation-in-gr>.

Based on our findings, practitioners in the field may consider balancing the degree of concentration when designing and training GR model, and using post-calibration to fix the errors when using the model on a downstream corpus. Other forms of decoding strategies beyond constrained beam search are also suggested.

As to future work, we will continue to study how these results can be used for analyzing training properties of corpus-specific GR models. Incorporating learnable decoding strategies during the training of a differentiable search index may also be of interest in this field.

Acknowledgments

Shiguang Wu gratefully acknowledges Yixiao Yu for insightful discussions regarding Lemma B.1, which were pivotal in establishing the first main result of this work. This research was (partially) funded by the Natural Science Foundation of China (62472261, 62372275, 62272274, 62202271, T2293773), the National Key R&D Program of China with grant No. 2024YFC3307300 and No. 2022YFC-3303004, the Provincial Key R&D Program of Shandong Province with grant No. 2024CXGC010108, the Natural Science Foundation of Shandong Province with grant No. ZR2024QF203, the Technology Innovation Guidance Program of Shandong Province with grant No. YDZX2024088, the Hybrid Intelligence Center, a 10-year program funded by the Dutch Ministry of Education, Culture and Science through the Netherlands Organisation for Scientific Research, <https://hybrid-intelligence-centre.nl>, project LESSEN with project number NWA.1389.20.183 of the research program NWA ORC 2020/21, which is (partly) financed by the Dutch Research Council (NWO), project ROBUST with project number KICH3.LTP.-20.006, which is (partly) financed by the Dutch Research Council (NWO), DPG Media, RTL, and the Dutch Ministry of Economic Affairs and Climate Policy (EZK) under the program LTP KIC 2020-2023, and the FINDHR (Fairness and Intersectional Non-Discrimination in Human Recommendation) project that received funding from the European Union’s Horizon Europe research and innovation program under grant agreement No 101070212. All content represents the opinion of the authors, which is not necessarily shared or endorsed by their respective employers and/or sponsors.

Appendix

A On the magnitude of vocabulary size

In this work, we use a sequence of tokens as the docID to represent the document as a general case. In order for the GR model to generate the docID of relevant documents, the docID faithfully representing the semantic information is highly preferred. Therefore, the vocabulary size of the docID is a key factor to determine the capacity of the GR model. In this section, we will give a brief perspective of its magnitude by considering the average information content in real-world documents.

Bits per byte (bpb) is a widely used metric to measure the information content of a document. It is defined as the number of bits required to encode the content in a lossless way. Let α be the bpb of general English text, and n be the length (in bytes) of a document. The information content of the document is $n\alpha$ bits on average. Since the vocabulary size of the GR model is k , the maximum average information in a single token is $\log_2 k$ bits. Each docID has m

tokens, so there will be $m \log_2 k$ bits in total. One would expect that the information content of the docID should be larger than the one of the document, i.e., $m \log_2 k \geq n\alpha$, which implies $k \geq 2^{\frac{n\alpha}{m}}$. From OpenAI [42] and several publications [6, 29], the bpb is usually around 1. For a document of length 512 bytes, and our docID length m is 32, the vocabulary size k should be about $2^{16\alpha} \approx 65,536$. This is about the same size of the vocabulary in a language model. The size of the complete corpus is astronomically high, and for a regular size downstream corpus, the sampling probability is approaching zero, and will raise large KL divergence according to Section 4.1. This is similar to the case where we want to use the language model as GR model and some textual content as the docID.

B Lower bound of KL divergence

LEMMA B.1 (LOWER BOUND OF KL DIVERGENCE BETWEEN BINOMIAL AND UNIFORM DISTRIBUTION). *Let $S_i \sim \text{Binomial}(m, p)$, where $i \in [k]$, and $Z = \sum_{i=1}^k S_i$. We define a normalized distribution P as*

$$P_i := \frac{S_i}{Z}, i \in [k], \quad (15)$$

and a uniform distribution Q on $\text{supp}(P) := \{i | S_i > 0\}$ as

$$Q_i := \begin{cases} \frac{1}{|\text{supp}(P)|}, & \text{if } S_i > 0 \\ 0, & \text{otherwise} \end{cases} \quad (16)$$

Then, we have a lower bound of the KL divergence between P and Q for large k ,

$$\text{KL}[P \parallel Q] = \sum_i P_i \ln \frac{P_i}{Q_i} \gtrsim \frac{0.05}{mp}. \quad (17)$$

PROOF. By the De Moivre–Laplace theorem [5], as k grows large, for n in the neighborhood of mp , we have

$$S_i \sim \mathcal{N}(mp, mpq), \quad p + q = 1. \quad (18)$$

Therefore, let $n = mp + \sqrt{mpq}$, we have

$$\Pr[S_i \geq n] \simeq 1 - \Phi(1). \quad (19)$$

Next, we focus on S_i that deviates from the mean. Let $A_n := \{S_i \geq n\}$ be the deviated elements, and $Y_n = |A_n|$, we have

$$P(A_n) = \frac{\sum_i S_i \mathbb{I}[S_i \geq n]}{Z} \geq \frac{nY_n}{Z}, \text{ and} \quad (20)$$

$$Q(A_n) = \frac{\sum_i \mathbb{I}[S_i \geq n]}{\sum_i \mathbb{I}[S_i \geq 1]} = \frac{Y_n}{Y_1}. \quad (21)$$

Since

$$\mathbb{E}[Z] = kmp, \quad (22)$$

$$\mathbb{E}[Y_n] = k \Pr[S_i \geq n] = O(k), \quad (23)$$

$$\mathbb{E}[Y_1] = k(1 - \Pr[S_i = 0]) \simeq k(1 - e^{-k^s}), \quad (24)$$

by the multiplicative Chernoff bound, with probability at least $1 - \exp(-\delta^2 \Omega(k))$, $|Y_n - \mathbb{E}[Y_n]| \leq \delta \mathbb{E}[Y_n]$, $|Y_1 - \mathbb{E}[Y_1]| \leq \delta \mathbb{E}[Y_1]$,

and $|Z - mpq| \leq \delta mpq$, we set $\delta = o(1)$ and $\delta^2 k = k^{O(1)}$, and have

$$d_{\text{TV}}(P, Q) \geq |P(A_n) - Q(A_n)| \quad (25)$$

$$\geq \frac{nY_n}{Z} - \frac{Y_n}{Y_1} \quad (26)$$

$$\geq \frac{(1 - \delta)n\mathbb{E}[Y_n]}{(1 + \delta)kmp} - \frac{(1 + \delta)\mathbb{E}[Y_n]}{(1 - \delta)\mathbb{E}[Y_1]} \quad (27)$$

$$\simeq \frac{1 - \Phi(1)}{\sqrt{mp}}. \quad (28)$$

where $d_{\text{TV}}(P, Q) = \sup_{A \subseteq [k]} |P(A) - Q(A)|$ is the total variation distance of two distributions.

Lastly, we use Pinsker's inequality [15] to give the asymptotic lower bound in Eq. 17,

$$\text{KL}[P \parallel Q] \geq 2d_{\text{TV}}(P, Q)^2 \gtrsim \frac{0.05}{mp}. \quad \square$$

THEOREM B.2 (LOWER BOUND OF KL DIVERGENCE FOR UNIFORM RELEVANCE DISTRIBUTION). *Let $r = m - 1$ and the sampling probability $p = \frac{1}{k^{r-s}}$. $\Pr(\cdot)$ is a uniform distribution. We have, for the KL divergence in Eq. 7,*

$$\text{KL}[\Pr(\cdot|C) \parallel \Pr(\cdot|C_i)] \gtrsim \frac{0.05}{k^s}. \quad (29)$$

Here s is small and hence the right hand side converges slowly with respect to k . When $s = 0$, we have a constant lower bound 0.05.

PROOF. We consider each possible token $d_1 \in [k]$ at first position. Since the selection of each document follows $\text{Bernoulli}(p)$, the number of documents selected with first token being d_1 follows $\text{Binomial}(k^r, p)$. As the GR model can only consider the constraint in the current step, it will return a uniform distribution over the valid d_1 tokens. By revoking Lemma B.1, we have the lower bound. \square

THEOREM B.3 (LOWER BOUND OF KL DIVERGENCE FOR GENERAL RELEVANCE DISTRIBUTION). *Let S_{ij} be a weighted Bernoulli random variable, with parameter p and weight w_{ij} , where $i \in [k]$ and $j \in [n]$. Suppose for some $\delta > 0$, $\{S_{ij} \mid j \in [n]\}$ satisfy Lyapunov's condition [5], i.e.,*

$$\lim_{n \rightarrow \infty} \frac{pq^{2+2\delta} + p^{2+2\delta}q}{\left(\sum_j w_{ij}^2 pq\right)^{1+\delta}} \sum_j w_{ij}^{2+2\delta} = 0. \quad (30)$$

We define P and Q similar to Theorem B.2, as

$$P[i] = \frac{S_i}{Z}, \quad Q[i] = \frac{w_i}{W}, \quad (31)$$

where $Z = \sum_{i=1}^k S_i$ and $W = \sum_{i=1}^k w_i$. We have a lower bound of the KL divergence between P and Q for large k ,

$$\text{KL}[P \parallel Q] \gtrsim \frac{0.05\mathbb{E}^2[A_i]}{p}, \quad (32)$$

where $A_i^2 = \sum_j w_{ij}^2 / w_i^2$.

PROOF. Here we use the Lyapunov central limit theorem [5], to approximate the distribution of W'_i . We have

$$S_i \sim \mathcal{N}(w_i p, pqw_i^2 A_i^2) = \mathcal{N}(\mu_i, \sigma_i^2), \quad p + q = 1. \quad (33)$$

As we have done in Theorem B.2, we choose a subset of $[k]$ to compute a lower bound of total variation distance. Let $I = \{i \mid S_i \geq$

$\mu_i + \sigma_i$, $\delta = o(1)$ and $k\delta^2 = k^{O(1)}$, we have, with probability at least $1 - 3 \exp(-\frac{2\delta^2 k^2}{B^2})$, where $B = \sum_{i \in [k]} w_i^2$,

$$d_{TV}(P, Q) \geq |P(I) - Q(I)| \quad (34)$$

$$\geq \frac{\sum_{i \in I} S_i}{Z} - \frac{\sum_{i \in I} w_i}{W} \quad (35)$$

$$\geq \frac{0.16}{\sqrt{p}} \frac{\sum_{i \in [k]} w_i A_i}{\sum_{i \in [k]} w_i}. \quad (36)$$

If we set $W = 1$, it can be simplified as $\frac{0.16 \mathbb{E}[A_i]}{\sqrt{p}}$. \square

C Recall performance using beam search

THEOREM C.1 (TOP- λk RECALL AND TOP-1 PRECISION OF DATA MODEL IN SECTION 4.2). *Suppose we have k branches, each with $n = k^{m-1}$ documents (leaves). We randomly select λk docs from all branches as the relevant ones, where $\lambda \ll 1$. We assign each non-relevant doc a score uniformly at random from $[-1, 1]$, and each relevant doc from $[\delta - \Delta, \delta + \Delta]$, where $\delta = 0.5 \log(0.8n)$, and $\Delta = O(0.5 \log \log k)$. We then take the exponential of the score for each doc and use the sum of the scores as the ranking score for each branch. Then, we have the following results with high probability:*

- (1) *The top- λk recall is lower than $0.5 + \max\{0.65 - 0.15/\lambda, 0\}$.*
- (2) *The top-1 precision is 1.*

PROOF. First, we show the distribution of values for non-relevant branches. For each non-relevant branch, the score S is the sum of n i.i.d. random variables from $Y_i = e^{X_i}$, where $X_i \sim \text{Uniform}(-1, 1)$. As $\mathbb{E}[Y] \approx 1.175$, and $\mathbb{V}[Y] \approx 0.4$, we have $\mathbb{E}[S] = 1.175n$ and $\mathbb{V}[S] = 0.4n$. According to the central limit theorem, the distribution of S is approximately $\mathcal{N}(1.175n, 0.4n)$. Therefore, for two non-relevant branches, the difference of their scores is approximately $\mathcal{N}(0, 0.8n)$, and $\Pr[S_1 - S_2 \geq \sqrt{0.8n}] \approx 0.15$.

Next, we estimate the number of relevant branches. As the proportion $p = \frac{\lambda k}{kn}$ is small, and $np = \lambda$, we can approximate the number of relevant docs, $\#R$, in each branch as a Poisson distribution with parameter λ . Then we have $\Pr[\#R = 0] = e^{-\lambda} \approx 1 - \lambda$, and $\Pr[\#R = 1] = \lambda e^{-\lambda} \approx \lambda$. Therefore, there are approximately λk relevant branches and $k - \lambda k$ non-relevant branches.

We consider a barrier $B = 1.175n + \sqrt{0.8n}$ and the event that some non-relevant branches are above the barrier and half of relevant branches are below the barrier. According to the data model, with high probability, half of the relevant documents have a score below $\sqrt{0.8n}$. Since each non-relevant branch has a probability of 0.15 to exceed another by $\sqrt{0.8n}$, we have that, with high probability, at least $0.15(k - \lambda k)$ non-relevant branches will exceed one of the low-score relevant branches. Thus at most $0.5\lambda k + \max\{0.5\lambda k - 0.15(k - \lambda k), 0\}$ relevant branches will be in the top- λk . Then the recall is at most $0.5 + \max\{0.65 - 0.15/\lambda, 0\}$.

For the top-1 precision, the maximum score of the relevant documents can approach $\delta + \Delta$ w.h.p. And for a normal distribution D with mean 0 and variance $\delta^2 = 0.8n$, $\Pr\left[\frac{D}{\delta} \geq \epsilon\right] \leq O(\frac{1}{\epsilon} e^{-\epsilon^2/2})$. If we let $\epsilon = e^\Delta$, we have that the probability there is a non-relevant branch exceeding by $e^\Delta \sqrt{0.8n}$ is $o(1)$. Then w.h.p. the top-1 precision is 1. \square

REMARK C.2. *In fact, the recall will be lower than 0.5 if λ is small enough because more non-relevant branches will be much higher*

than the barrier. The result mainly comes from the carefully designed score of relevant documents which is linear in n . It may not hold for extremely skewed distribution of scores, e.g., the relevant score is exponentially large, which actually corresponds to the “amplification” discussed in Section 5.

References

- [1] Kareem Ahmed, Kai-Wei Chang, and Guy Van den Broeck. 2023. Semantic Strengthening of Neuro-Symbolic Learning. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics (Proceedings of Machine Learning Research, Vol. 206)*, Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent (Eds.). PMLR, 10252–10261. <https://proceedings.mlr.press/v206/ahmed23a.html>
- [2] Arian Askari, Chuan Meng, Mohammad Aliannejadi, Zhaochun Ren, Evangelos Kanoulas, and Suzan Verberne. 2024. Generative Retrieval with Few-shot Indexing. *CoRR abs/2408.02152* (2024). doi:10.48550/ARXIV.2408.02152 arXiv:2408.02152
- [3] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew Bm McNamara, Bhaskar Mitra, Tri Nguyen, et al. 2016. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268* (2016).
- [4] Michele Bevilacqua, Giuseppe Ottaviano, Patrick Lewis, Scott Yih, Sebastian Riedel, and Fabio Petroni. 2022. Autoregressive Search Engines: Generating Substrings as Document Identifiers. *Advances in Neural Information Processing Systems* 35 (2022), 31668–31683.
- [5] Patrick Billingsley. 2017. *Probability and measure*. John Wiley & Sons.
- [6] Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego De Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassier, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack Rae, Erich Elsen, and Laurent Sifre. 2022. Improving Language Models by Retrieving from Trillions of Tokens. In *Proceedings of the 39th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (Eds.). PMLR, 2206–2240. <https://proceedings.mlr.press/v162/borgeaud22a.html>
- [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 1877–1901. https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf
- [8] Sebastian Bruch, Xuanhui Wang, Michael Bendersky, and Marc Najork. 2019. An Analysis of the Softmax Cross Entropy Loss for Learning-to-Rank with Binary Relevance. In *Proceedings of the 2019 ACM SIGIR International Conference on Theory of Information Retrieval (Santa Clara, CA, USA) (ICTIR '19)*. Association for Computing Machinery, New York, NY, USA, 75–78. doi:10.1145/3341981.3344221
- [9] Nicola De Cao, Gautier Izcard, Sebastian Riedel, and Fabio Petroni. 2021. Autoregressive Entity Retrieval. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=5k8f6U039V>
- [10] Souradip Chakraborty, Soumya Suvra Ghosal, Ming Yin, Dinesh Manocha, Mengdi Wang, Amrit Singh Bedi, and Furong Huang. 2024. Transfer Q Star: Principled Decoding for LLM Alignment. *arXiv preprint arXiv:2405.20495* (2024).
- [11] Jiangui Chen, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Wei Chen, Yixing Fan, and Xueqi Cheng. 2023. Continual Learning for Generative Retrieval over Dynamic Corpora. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM '23)*. Association for Computing Machinery, New York, NY, USA, 306–315. doi:10.1145/3583780.3614821
- [12] Jiangui Chen, Ruqing Zhang, Jiafeng Guo, Yiqun Liu, Yixing Fan, and Xueqi Cheng. 2022. CorpusBrain: Pre-train a Generative Retrieval Model for Knowledge-Intensive Language Tasks. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17–21, 2022*, Mohammad Al Hasan and Li Xiong (Eds.). ACM, 191–200. doi:10.1145/3511808.3557271
- [13] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, and Daniel Campos. 2021. Overview of the TREC 2020 deep learning track. *arXiv:2102.07662* [cs.IR] <https://arxiv.org/abs/2102.07662>
- [14] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. Overview of the TREC 2019 deep learning track. *arXiv:2003.07820* [cs.IR] <https://arxiv.org/abs/2003.07820>
- [15] Imre Csiszár and János Körner. 2011. *Information theory: coding theorems for discrete memoryless systems*. Cambridge University Press.
- [16] Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023. Precise Zero-Shot Dense Retrieval without Relevance Labels. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 1762–1777. doi:10.18653/v1/2023.acl-long.99
- [17] Jiafeng Guo, Changjiang Zhou, Ruqing Zhang, Jiangui Chen, Maarten de Rijke, Yixing Fan, and Xueqi Cheng. 2024. CorpusBrain++: A Continual Generative Pre-Training Framework for Knowledge-Intensive Language Tasks. *arXiv:2402.16767* [cs.IR] <https://arxiv.org/abs/2402.16767>
- [18] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *EMNLP. Association for Computational Linguistics*, 6769–6781.
- [19] Minbeom Kim, Hwanhee Lee, Kang Min Yoo, Joonsuk Park, Hwaran Lee, and Kyomin Jung. 2023. Critic-Guided Decoding for Controlled Text Generation. In *Findings of the Association for Computational Linguistics: ACL 2023*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 4598–4612. doi:10.18653/v1/2023.findings-acl.281
- [20] Minbeom Kim, Thibaut Thonet, Jos Rozen, Hwaran Lee, Kyomin Jung, and Marc Dymetman. 2024. Guaranteed Generation from Large Language Models. *arXiv:2410.06716* [cs.CL] <https://arxiv.org/abs/2410.06716>
- [21] Varsha Kishore, Chao Wan, Justin Lovelace, Yoav Artzi, and Kilian Q. Weinberger. 2023. IncDSI: incrementally updatable document retrieval. In *Proceedings of the 40th International Conference on Machine Learning (Honolulu, Hawaii, USA) (ICML '23)*. JMLR.org, Article 704, 13 pages.
- [22] Hyunji Lee, JaeYoung Kim, Hoyeon Chang, Hanseok Oh, Sohee Yang, Vladimir Karpukhin, Yi Lu, and Minjoon Seo. 2023. Nonparametric Decoding for Generative Retrieval. In *ACL (Findings)*. Association for Computational Linguistics, 12642–12661.
- [23] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5–10, 2020*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault (Eds.). Association for Computational Linguistics, 7871–7880. doi:10.18653/V1/2020.ACL-MAIN.703
- [24] Haitao Li, Qingyao Ai, Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Zheng Liu, and Zhao Cao. 2023. Constructing Tree-based Index for Efficient and Effective Dense Retrieval. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (Taipei, Taiwan) (SIGIR '23)*. Association for Computing Machinery, New York, NY, USA, 131–140. doi:10.1145/3539618.3591651
- [25] Minghan Li, Sheng-Chieh Lin, Xueguang Ma, and Jimmy Lin. 2023. SLIM: Sparsified Late Interaction for Multi-Vector Retrieval with Inverted Indexes. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (Taipei, Taiwan) (SIGIR '23)*. Association for Computing Machinery, New York, NY, USA, 1954–1959. doi:10.1145/3539618.3591977
- [26] Xiaoxi Li, Zhicheng Dou, Yujia Zhou, and Fangchao Liu. 2024. CorpusLM: Towards a Unified Language Model on Corpus for Knowledge-Intensive Tasks. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July 14–18, 2024*, Grace Hui Yang, Hongning Wang, Sam Han, Claudia Hauff, Guido Zuccon, and Yi Zhang (Eds.). ACM, 26–37. doi:10.1145/3626772.3657778
- [27] Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yuyao Zhang, Peitian Zhang, Yutao Zhu, and Zhicheng Dou. 2024. From Matching to Generation: A Survey on Generative Information Retrieval. *CoRR abs/2404.14851* (2024). doi:10.48550/ARXIV.2404.14851 arXiv:2404.14851
- [28] Xiaoxi Li, Yujia Zhou, and Zhicheng Dou. 2024. UniGen: A Unified Generative Framework for Retrieval and Question Answering with Large Language Models. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20–27, 2024, Vancouver, Canada*, Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (Eds.). AAAI Press, 8688–8696. doi:10.1609/AAAI.V38I8.28714
- [29] Yucheng Li, Yunhao Guo, Frank Guerin, and Chenchua Lin. 2024. Evaluating Large Language Models for Generalization and Robustness via Data Compression. *arXiv:2402.00861* [cs.CL] <https://arxiv.org/abs/2402.00861>
- [30] Yongqi Li, Nan Yang, Liang Wang, Furu Wei, and Wenjie Li. 2024. Learning to Rank in Generative Retrieval. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20–27, 2024, Vancouver, Canada*, Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (Eds.). AAAI Press, 8716–8723. doi:10.1609/AAAI.V38I8.28717
- [31] Yu-An Liu, Ruqing Zhang, Jiafeng Guo, Wei Chen, and Xueqi Cheng. 2023. On the Robustness of Generative Retrieval Models: An Out-of-Distribution Perspective. *CoRR abs/2306.12756* (2023). doi:10.48550/ARXIV.2306.12756 arXiv:2306.12756
- [32] Yu-An Liu, Ruqing Zhang, Jiafeng Guo, Changjiang Zhou, Maarten de Rijke, and Xueqi Cheng. 2024. On the Robustness of Generative Information Retrieval Models. *arXiv:2412.18768* [cs.IR] <https://arxiv.org/abs/2412.18768>

- [33] Sanket Vaibhav Mehta, Jai Gupta, Yi Tay, Mostafa Dehghani, Vinh Q. Tran, Jinpeng Rao, Marc Najork, Emma Strubell, and Donald Metzler. 2023. DSI++: Updating Transformer Memory with New Documents. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 8198–8213. doi:10.18653/v1/2023.emnlp-main.510
- [34] Donald Metzler, Yi Tay, Dara Bahri, and Marc Najork. 2021. Rethinking search: making domain experts out of dilettantes. *SIGIR Forum* 55, 1, Article 13 (July 2021), 27 pages. doi:10.1145/3476415.3476428
- [35] Fatemehsadat Mireshghallah, Kartik Goyal, and Taylor Berg-Kirkpatrick. 2022. Mix and Match: Learning-free Controllable Text Generation using Energy Language Models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, Dublin, Ireland, 401–415. doi:10.18653/v1/2022.acl-long.31
- [36] Sidharth Mudgal, Jong Lee, Harish Ganapathy, YaGuang Li, Tao Wang, Yanping Huang, Zhifeng Chen, Heng-Tze Cheng, Michael Collins, Trevor Strohman, Jilin Chen, Alex Beutel, and Ahmad Beirami. 2025. Controlled decoding from language models. In *Proceedings of the 41st International Conference on Machine Learning (Vienna, Austria) (ICML '24)*. JMLR.org, Article 1484, 18 pages.
- [37] Waleed Mustafa, Yunwen Lei, Antoine Ledent, and Marius Kloft. 2021. Fine-grained Generalization Analysis of Structured Output Prediction. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21, Zhi-Hua Zhou (Ed.)*. International Joint Conferences on Artificial Intelligence Organization, 2841–2847. doi:10.24963/ijcai.2021/391 Main Track.
- [38] Thong Nguyen and Andrew Yates. 2023. Generative Retrieval as Dense Retrieval. *CoRR abs/2306.11397* (2023). doi:10.48550/ARXIV.2306.11397 arXiv:2306.11397
- [39] Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith Hall, Ming-Wei Chang, and Yinfei Yang. 2022. Large Dual Encoders Are Generalizable Retrievers. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 9844–9855. doi:10.18653/v1/2022.emnlp-main.669
- [40] Masaaki Nishino, Kengo Nakamura, and Norihito Yasuda. 2024. Generalization analysis on learning with a concurrent verifier. In *Proceedings of the 36th International Conference on Neural Information Processing Systems (New Orleans, LA, USA) (NIPS '22)*. Curran Associates Inc., Red Hook, NY, USA, Article 302, 12 pages.
- [41] Masaaki Nishino, Kengo Nakamura, and Norihito Yasuda. 2025. Understanding the impact of introducing constraints at inference time on generalization error. In *Proceedings of the 41st International Conference on Machine Learning (Vienna, Austria) (ICML '24)*. JMLR.org, Article 1551, 11 pages.
- [42] OpenAI. 2024. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL] <https://arxiv.org/abs/2303.08774>
- [43] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research* 21, 140 (2020), 1–67. <http://jmlr.org/papers/v21/20-074.html>
- [44] Ruiyang Ren, Yingqi Qu, Jing Liu, Xin Zhao, Qifei Wu, Yuchen Ding, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2023. A Thorough Examination on Zero-shot Dense Retrieval. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 15783–15796. doi:10.18653/v1/2023.findings-emnlp.1057
- [45] E. H. SIMPSON. 1949. Measurement of Diversity. *Nature* 163, 4148 (April 1949), 688–688. doi:10.1038/163688a0
- [46] Weiwei Sun, Lingyong Yan, Zheng Chen, Shuaiqiang Wang, Haichao Zhu, Pengjie Ren, Zhumin Chen, Dawei Yin, Maarten Rijke, and Zhaochun Ren. 2023. Learning to Tokenize for Generative Retrieval. In *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (Eds.), Vol. 36. Curran Associates, Inc., 46345–46361. https://proceedings.neurips.cc/paper_files/paper/2023/file/91228b942a4528cdac031c1b68b127e8-Paper-Conference.pdf
- [47] Yubao Tang, Ruqing Zhang, Jiafeng Guo, Jiangui Chen, Zuowei Zhu, Shuaiqiang Wang, Dawei Yin, and Xueqi Cheng. 2023. Semantic-Enhanced Differentiable Search Index Inspired by Learning Strategies. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (Long Beach, CA, USA) (KDD '23)*. Association for Computing Machinery, New York, NY, USA, 4904–4913. doi:10.1145/3580305.3599903
- [48] Yubao Tang, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Wei Chen, and Xueqi Cheng. 2024. Generative Retrieval Meets Multi-Graded Relevance. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. <https://openreview.net/forum?id=2xTkeyJfJB>
- [49] Yi Tay, Vinh Tran, Mostafa Dehghani, Jianmo Ni, Dara Bahri, Harsh Mehta, Zhen Qin, Kai Hui, Zhe Zhao, Jai Prakash Gupta, Tal Schuster, William W. Cohen, and Donald Metzler. 2022. Transformer Memory as a Differentiable Search Index. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (Eds.). http://papers.nips.cc/paper_files/paper/2022/hash/892840a6123b5ec99ebaab8be1530fba-Abstract-Conference.html
- [50] Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*. <https://openreview.net/forum?id=wCu6T5xFje>
- [51] Kexin Wang, Nandan Thakur, Nils Reimers, and Iryna Gurevych. 2022. GPL: Generative Pseudo Labeling for Unsupervised Domain Adaptation of Dense Retrieval. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). Association for Computational Linguistics, Seattle, United States, 2345–2360. doi:10.18653/v1/2022.naacl-main.168
- [52] Yujing Wang, Yingyan Hou, Haonan Wang, Ziming Miao, Shibin Wu, Qi Chen, Yuqing Xia, Chengmin Chi, Guoshuai Zhao, Zheng Liu, Xing Xie, Hao Sun, Weiwei Deng, Qi Zhang, and Mao Yang. 2022. A Neural Corpus Indexer for Document Retrieval. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (Eds.). http://papers.nips.cc/paper_files/paper/2022/hash/a46156bd3579c3b268108ea6aca71d13-Abstract-Conference.html
- [53] Zihan Wang, Yujia Zhou, Yiteng Tu, and Zhicheng Dou. 2023. NOVO: Learnable and Interpretable Document Identifiers for Model-Based IR. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (Birmingham, United Kingdom) (CIKM '23)*. Association for Computing Machinery, New York, NY, USA, 2656–2665. doi:10.1145/3583780.3614993
- [54] Ryan W. White and Chirag Shah. 2025. *Information Access in the Era of Generative AI*. Springer Cham. <https://doi.org/10.1007/978-3-031-73147-1>
- [55] Shiguang Wu, Wenda Wei, Mengqi Zhang, Zhumin Chen, Jun Ma, Zhaochun Ren, Maarten de Rijke, and Pengjie Ren. 2024. Generative Retrieval as Multi-Vector Dense Retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July 14-18, 2024*, Grace Hui Yang, Hongning Wang, Sam Han, Claudia Hauff, Guido Zuccon, and Yi Zhang (Eds.). ACM, 1828–1838. doi:10.1145/3626772.3657697
- [56] Shicheng Xu, Liang Pang, Huawei Shen, and Xueqi Cheng. 2023. BERM: Training the Balanced and Extractable Representation for Matching to Improve Generalization Ability of Dense Retrieval. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 6620–6635. doi:10.18653/v1/2023.acl-long.365
- [57] Tianchi Yang, Minghui Song, Zihan Zhang, Haizhen Huang, Weiwei Deng, Feng Sun, and Qi Zhang. 2023. Auto Search Indexer for End-to-End Document Retrieval. In *EMNLP (Findings)*. Association for Computational Linguistics, 6955–6970.
- [58] Weiqin Yang, Jiawei Chen, Xin Xin, Sheng Zhou, Binbin Hu, Yan Feng, Chun Chen, and Can Wang. 2024. PSL: Rethinking and Improving Softmax Loss from Pairwise Perspective for Recommendation. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. <https://openreview.net/forum?id=PhjnK9KWox>
- [59] Yue Yu, Chenyan Xiong, Si Sun, Chao Zhang, and Arnold Overwijk. 2022. COCODR: Combating the Distribution Shift in Zero-Shot Dense Retrieval with Contrastive and Distributionally Robust Learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 1462–1479. doi:10.18653/v1/2022.emnlp-main.95
- [60] Peiwen Yuan, Xinglin Wang, Shaoxiong Feng, Boyuan Pan, Yiwei Li, Heda Wang, Xupeng Miao, and Kan Li. 2024. Generative Dense Retrieval: Memory Can Be a Burden. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024 - Volume 1: Long Papers, St. Julian's, Malta, March 17-22, 2024*, Yvette Graham and Matthew Purver (Eds.). Association for Computational Linguistics, 2835–2845. <https://aclanthology.org/2024.eacl-long.173>
- [61] Hansi Zeng, Chen Luo, Bowen Jin, Sheikh Muhammad Sarwar, Tianxin Wei, and Hamed Zamani. 2024. Scalable and Effective Generative Information Retrieval. In *Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, May 13-17, 2024*, Tat-Seng Chua, Chong-Wah Ngo, Ravi Kumar, Hady W. Lauw, and Roy Ka-Wei Lee (Eds.). ACM, 1441–1452. doi:10.1145/3589334.3645477
- [62] Hansi Zeng, Chen Luo, and Hamed Zamani. 2024. Planning Ahead in Generative Retrieval: Guiding Autoregressive Generation through Simultaneous Decoding. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July 14-18, 2024*, Grace Hui Yang, Hongning Wang, Sam Han, Claudia Hauff, Guido Zuccon, and Yi Zhang (Eds.). ACM, 469–480. doi:10.1145/3626772.3657746

- [63] Honghua Zhang, Meihua Dang, Nanyun Peng, and Guy Van Den Broeck. 2023. Tractable control for autoregressive language generation. In *Proceedings of the 40th International Conference on Machine Learning* (Honolulu, Hawaii, USA) (ICML '23). JMLR.org, Article 1716, 14 pages.
- [64] Honghua Zhang, Po-Nien Kung, Masahiro Yoshida, Guy Van den Broeck, and Nanyun Peng. 2024. Adaptable Logical Control for Large Language Models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. <https://openreview.net/forum?id=58X9v92zRd>
- [65] Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. 2023. A Survey of Controllable Text Generation Using Transformer-based Pre-trained Language Models. *ACM Comput. Surv.* 56, 3, Article 64 (Oct. 2023), 37 pages. doi:10.1145/3617680
- [66] Peitian Zhang, Zheng Liu, Yujia Zhou, Zhicheng Dou, and Zhao Cao. 2023. Term-Sets Can Be Strong Document Identifiers For Auto-Regressive Search Engines. *ArXiv abs/2305.13859* (2023). <https://api.semanticscholar.org/CorpusID:258841428>
- [67] Yujia Zhou, Jing Yao, Zhicheng Dou, Ledell Wu, and Ji-Rong Wen. 2023. DynamicRetriever: A Pre-trained Model-based IR System Without an Explicit Index. *Mach. Intell. Res.* 20, 2 (April 2023), 276–288. <https://doi.org/10.1007/s11633-022-1373-9>
- [68] Han Zhu, Daqing Chang, Ziru Xu, Pengye Zhang, Xiang Li, Jie He, Han Li, Jian Xu, and Kun Gai. 2019. *Joint optimization of tree-based index and deep model for recommender systems*. Curran Associates Inc., Red Hook, NY, USA.
- [69] Han Zhu, Xiang Li, Pengye Zhang, Guozheng Li, Jie He, Han Li, and Kun Gai. 2018. Learning Tree-based Deep Model for Recommender Systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (London, United Kingdom) (KDD '18). Association for Computing Machinery, New York, NY, USA, 1079–1088. doi:10.1145/3219819.3219826
- [70] Shengyao Zhuang, Houxing Ren, Linjun Shou, Jian Pei, Ming Gong, Guido Zuccon, and Daxin Jiang. 2022. Bridging the Gap Between Indexing and Retrieval for Differentiable Search Index with Query Generation. *ArXiv abs/2206.10128* (2022). <https://api.semanticscholar.org/CorpusID:249890267>
- [71] Jingwei Zhuo, Ziru Xu, Wei Dai, Han Zhu, Han Li, Jian Xu, and Kun Gai. 2020. Learning optimal tree models under beam search. In *Proceedings of the 37th International Conference on Machine Learning (ICML '20)*. JMLR.org, Article 1080, 10 pages.