# Retrieval-Augmented Generation for AI-Generated Content: A Survey

Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng,
Fangcheng Fu, Ling Yang, Wentao Zhang, Jie Jiang, Bin Cui

*Abstract*—Advancements in model algorithms, the growth of foundational models, and access to high-quality datasets have propelled the evolution of Artificial Intelligence Generated Content (AIGC). Despite its notable successes, AIGC still faces hurdles such as updating knowledge, handling long-tail data, mitigating data leakage, and managing high training and inference costs. Retrieval-Augmented Generation (RAG) has recently emerged as a paradigm to address such challenges. In particular, RAG introduces the information retrieval process, which enhances the generation process by retrieving relevant objects from available data stores, leading to higher accuracy and better robustness. In this paper, we comprehensively review existing efforts that integrate RAG technique into AIGC scenarios. We first classify RAG foundations according to how the retriever augments the generator, distilling the fundamental abstractions of the augmentation methodologies for various retrievers and generators. This unified perspective encompasses all RAG scenarios, illuminating advancements and pivotal technologies that help with potential future progress. We also summarize additional enhancements methods for RAG, facilitating effective engineering and implementation of RAG systems. Then from another view, we survey on practical applications of RAG across different modalities and tasks, offering valuable references for researchers and practitioners. Furthermore, we introduce the benchmarks for RAG, discuss the limitations of current RAG systems, and suggest potential directions for future research. Github: https://github.com/PKU-DAIR/RAG-Survey.

*Index Terms*—Retrieval-augmented generation, AI-generated content, generative models, information retrieval.

## I. INTRODUCTION

### A. Background

RECENT years have witnessed the surge in interests surrounding Artificial Intelligence Generated Content (AIGC). Various content generation tools have been meticulously crafted to produce diverse outputs across various modalities, such as Large Language Models (LLMs) including the GPT series [1]–[3] and the LLAMA series [4]–[6] for texts and codes, DALL-E [7]–[9] and Stable Diffusion [10] for images, and Sora [11] for videos. The word "AIGC" emphasizes that the contents are produced by advanced generative models other than human beings or rule-based approaches. These generative models have achieved remarkable performance due to the utilization of novel model algorithms,

● Penghao Zhao and Hailin Zhang contributed equally to this paper.
● Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng, Fangcheng Fu, Ling Yang, Wentao Zhang and Bin Cui are with Peking University (e-mail: penghao.zhao@stu.pku.edu.cn, z.hl@pku.edu.cn, yuqinhan@stu.pku.edu.cn, wzr@stu.pku.edu.cn, 1800012997@pku.edu.cn, ccchengff@pku.edu.cn, yangling0818@163.com, wentao.zhang@pku.edu.cn, bin.cui@pku.edu.cn).
● Jie Jiang is with Tencent Inc. (email: zeus@tencent.com)
● Bin Cui is Corresponding Author.

explosive scale of foundation models, and massive high-quality datasets. Specifically, sequence-to-sequence tasks have transitioned from utilizing Long Short-Term Memory (LSTM) networks [12] to Transformer-based models [13], and image-generation tasks have shifted from Generative Adversarial Networks (GANs) [14] to Latent Diffusion Models (LDMs) [10] as well. Notably, the architecture of foundation models, initially constituted by millions of parameters [15], [16], has now grown to billions or even trillions of parameters [1], [4], [17]. These advancements are further bolstered by the availability of rich, high-quality datasets [1], [18], which provide ample training samples to fully optimize model parameters.

Information retrieval is another pivotal application within the field of computer science. Different from generation, retrieval aims to locate relevant existing objects from a vast pool of resources. The most prevalent application of retrieval lies in web search engines, which primarily focus on the task of document retrieval [19], [20]. In the present era, efficient information retrieval systems can handle document collections on the order of billions [21], [22]. Besides documents, retrieval has also been applied for many other modalities [23]–[26].

Despite the remarkable progress made by advanced generative models, AIGC continues to face a number of well-known challenges, including the struggle to maintain up-to-date knowledge, the inability to incorporate long-tail knowledge [27], and the risk of leaking private training data [28]. Retrieval-Augmented Generation (RAG) is proposed to alleviate, if not completely address, the aforementioned challenges through its adaptable data repository [29]. The knowledge stored for retrieval can be conceptualized as non-parametric memory, which is easily modifiable, capable of accommodating broad long-tail knowledge, and also able to encode confidential data. In addition, retrieval can also be employed to reduce the generation costs. For example, RAG can reduce the size of large generative models [30], provide support for long contexts [31], and eliminate certain generation steps [32].

A typical RAG process is depicted in Fig. 1. Given an input query, the retriever locates and looks up relevant data sources, then the retrieved results interact with the generator to enhance the overall generation process. There are several *foundational paradigms* (*foundations* in short) according to how the retrieved results augment the generation: they can serve as augmented input to the generator [33], [34]; they can join at the middle stage of generation as latent representations [35], [36]; they can contribute to the final generation results in the form of logits [37], [38]; they can even influence or omit certain generation steps [32], [39].
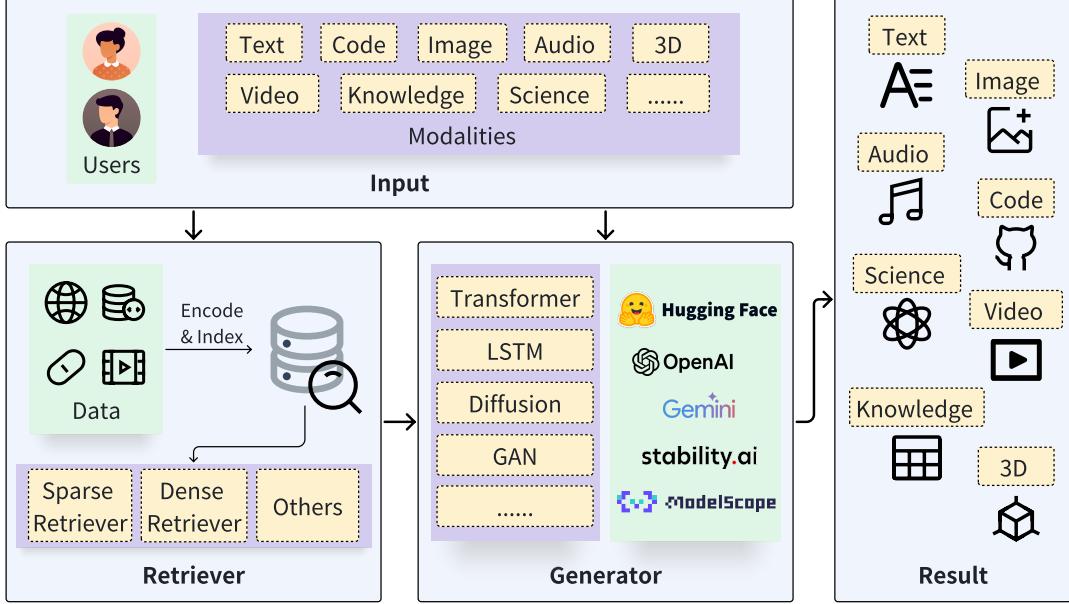
Fig. 1: A generic RAG architecture. The user queries, spanning different modalities, serve as input to both the retriever and the generator. The retriever extracts relevant information from data sources. The generator interacts with the retrieval results and ultimately produces outcomes of various modalities.

Moreover, beyond the foundational RAG process, researchers have proposed numerous *enhancements* to elevate the overall quality. These methods encompass specific optimizations for individual components as well as holistic enhancements aimed at the entire pipeline.

In addition, while the concept of RAG initially emerged in text-to-text generation [34], this technique has also found *applications* across various domains, including codes [40]–[42], audios [43], [44], images [45]–[47], videos [48], [49], 3D [50], [51], knowledge [52]–[54], and AI for science [55], [56]. In particular, the essential idea and process of RAG are largely consistent across modalities. However, it necessitates minor adjustments in augmentation techniques, and the selection of retrievers and generators varies depending on the specific modalities and applications.

Despite the rapid growth in recent research on RAG and the booming applications, a systematic review encompassing all foundations, enhancements, and applications is notably absent, hindering the development of this field. For one thing, the absence of discussion on RAG foundations significantly undermines the practical value of the research in this domain, leaving the potential of RAG not fully explored. While the majority of research interest, particularly among LLM researchers, centers on query-based RAG in text-generation tasks, it is essential to acknowledge that other RAG foundations are also effective and with significant potential for usage and further development. For another, the lack of an overview on RAG applications causes researchers and practitioners to overlook RAG's progress across multiple modalities and remain unaware of how RAG can be effectively applied. Although text generation is typically considered as the main application of RAG, we emphasize that the development of

RAG in other modalities has also begun to catch on and has yielded promising advancements. Certain modalities have a rich historical connection to retrieval techniques, infusing RAG with distinctive characteristics. Inspired by this, in this paper, our objective is to present a comprehensive survey to provide a systematic overview of RAG.

### B. Contribution

This survey offers a comprehensive overview of RAG, covering foundations, enhancements, applications, benchmarks, limitations, and potential future directions. While retrievers and generators exhibit variations across modalities and tasks, we distill the fundamental abstractions of RAG foundations, considering applications as adaptations stemming from these abstractions. We aim to offer references and guidelines to researchers and practitioners, providing valuable insights for advancing RAG methodologies and related applications. In summary, we list our contributions as follows:

- We conduct a comprehensive review of RAG, and distill the abstractions of RAG foundations for various retrievers and generators.
- We investigate the enhancements in the literature of RAG, elaborating the techniques leveraged to enable more effective RAG systems.
- For various modalities and tasks, we survey existing AIGC methods that incorporate RAG techniques, exhibiting how RAG contributes to current generative models.
- We discuss the limitations and promising research directions of RAG, shedding light on its potential future development.

### C. Related Work

As the field of RAG advances, several surveys have emerged; yet they address only specific facets of the area. In

particular, they either exclusively focus on a single RAG foundation or provide only a brief overview of RAG augmentation methodologies for limited scenarios.

Most of the existing works focus on text-related RAG tasks that are facilitated by LLMs, without in-depth investigation in other modalities. The survey by Li et al. [57] offers a basic overview of RAG and discusses specific applications within the scope of text generation tasks. In a similar vein, the tutorial crafted by Asai et al. [58] centers on retrieval-based language models, detailing their structures and training strategies. Meanwhile, a recent survey by Gao et al. [59] explores RAG in the context of LLMs, with a particular emphasis on enhancement approaches for query-based RAG. Recognizing that RAG has extended beyond the text domain, our work broadens its reach to the entire AIGC landscape, facilitating a more comprehensive coverage of RAG research.

In addition, another survey proposed by Zhao et al. [60] introduces RAG applications across multiple modalities, but ignoring the discussion on RAG foundations. While existing research has explored various aspects of RAG, there remains a need for a comprehensive overview that covers RAG foundations, enhancements, and its applicability across different domains. In this paper, we aim to address the gap by presenting a systematic survey of RAG.

### D. Roadmap

The rest of the paper is organized as follows. Section II elaborates on the preliminary of RAG, introducing retrievers and generators. Section III presents RAG foundations and further enhancements on RAG. Section IV reviews existing research on RAG across various applications. Section V investigates the benchmark frameworks for RAG. Section VI discusses current limitations of RAG and potential future directions. Finally, Section VII concludes this paper.

## II. PRELIMINARY

In this section, we provide an overview of the general RAG architecture and explore the generators and the retrievers in today's RAG-based AIGC.

### A. Overview

As shown in Fig. 1, the entire RAG system consists of two core modules: the retriever and the generator, where the retriever searches for relevant information from the data store and the generator produces the required contents. The RAG process unfolds as follows: (i) the retriever initially receives the input query and searches for relevant information; (ii) then, the original query and the retrieval results are fed into the generator through a specific augmentation methodology; (iii) finally, the generator produces the desired outcomes.

### B. Generator

The remarkable performance of generative AI across diverse tasks has ushered in the era of AIGC. The generation module plays a crucial role within the RAG system. Different generative models are applied for different scenarios, such as transformer models for text-to-text tasks, VisualGPT [61] for image-to-text tasks, Stable Diffusion [10] for text-to-image tasks, Codex [2] for text-to-code tasks, etc. Here we introduce 4 typical generators that are frequently used in RAG: transformer model, LSTM, diffusion model, and GAN.

*1) Transformer Model:* Transformer models are one of the best performing models in the field of Natural Language Processing (NLP), consisting of self-attention mechanisms, feed-forward networks, layer normalization modules, and residual networks [62]. As shown in Fig. 2, the input of the transformer is mapped to a tensor $x_{in}$ with a shape of ($b$, $s$, $h$) after the tokenization process and embedding model, where $b$ represents batch size, $s$ represents sequence length and $h$ represents hidden dimension. Next, the position encoding will be sent to the self attention layer along with this tensor. The input $x_{in}$ and the output $x_{out}$ of the self-attention module will be connected by the residual network and the layer normalization module. Finally, the output of the "Add & Norm" module $x_{out}$ will be sent to the feed forward network.

The entire process can be defined as follows:

$$Q = x_{in} * W_q + b_q$$

$$K = x_{in} * W_k + b_k$$

$$V = x_{in} * W_v + b_v$$

$$x_{out} = LayerNorm1(Softmax(\frac{Q * K^T}{\sqrt{h}}) * V * W_o + b_o) + x_{in}$$

$$y = LayerNorm2((x_{out} * W_1 + b_1) * W_2 + b_2) + x_{out}$$

It should be noted that $w_q, w_k, w_v, w_o$ are learnable tensors with shape ($h$, $h$); $b_q, b_k, b_v, b_o$ are a learnable tensors with shape ($h$,).

*2) LSTM:* Long Short-Term Memory (LSTM) [63] is a special Recurrent Neural Network (RNN) model that overcomes the exploding/vanishing gradient problems of RNN in processing long-term dependency information by introducing cell state and gate mechanisms. The LSTM model consists of three gates: Input Gate, Forget Gate, and Output Gate. These gates update the cell state by controlling the information flow, enabling the model to remember long-term dependent information. Cell State is the core module of the LSTM model which can memorize and maintain information. The Input Gate decides which input data should be retained in the cell state. Forget Gate determines which cell state information should be discarded to avoid excessive memory. Output Gate determines how the information in the cell state affects the current output. The flow of data and the collaborative work process between components are shown in the Fig. 2.

The entire process can be defined as follows:

$$f = sigmoid(W^f * x^t + U^f * y^{t-1} + b^f)$$

$$z = tanh(W^z * x^t + U^z * y^{t-1} + b^z)$$

$$i = sigmoid(W^i * x^t + U^i * y^{t-1} + b^i)$$

$$o = sigmoid(W^o * x^t + U^o * y^{t-1} + b^o)$$

$$c^t = z \odot i + f \odot c^{t-1}$$

$$y^t = o \odot tanh(c^t)$$

(a) General transformer model architecture.



(b) General LSTM block architecture.



(c) General latent diffusion model architecture.
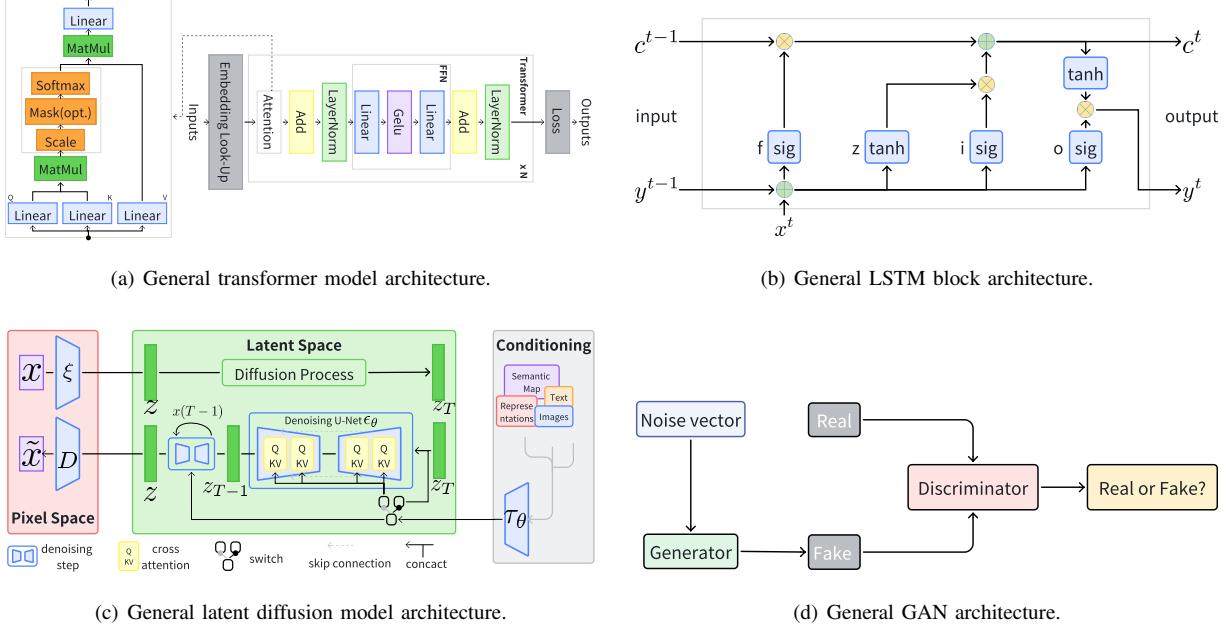


(d) General GAN architecture.

Fig. 2: General architectures of several generators.

*3) Diffusion Model:* Diffusion models [64] are a family of deep generative models that can create realistic and diverse samples of data [65]–[72], such as images [73]–[79], texts [80]–[83], videos [84]–[88], and molecules [89]–[93]. As shown in Fig. 2, diffusion models work by gradually adding noise to data until it becomes random, then reversing the process to generate new data from noise. This process is based on probabilistic modeling and neural networks. Diffusion models mainly have three equivalent formulations: denoising diffusion probabilistic models [65]–[67], score-based generative models [68], [69], and stochastic differential equations [70], [71], with following improvements like DDIM [94], Rectified Flow [95], Consistency Model [96] and RPG-DiffusionMaster [79].

Especially, let $x_0$ be a random variable that follows the data distribution $q(x_0)$, and let $x_t$ be a random variable that follows the distribution $q(x_t|x_0)$ after adding noise at time step $t$. Then, DDPM can be formulated as follows:

- **Forward Process** The forward process perturbs data with a sequence of Gaussian noise injections, transforming the data distribution $q(x_0)$ into a simple prior distribution $q(x_T) \approx N(0, I)$. The transition kernel at each time step is given by

$$q(x_t|x_{t-1}) = N(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I),$$

where $\beta_t \in (0,1)$ is a hyperparameter. The marginal distribution of $x_t$ conditioned on $x_0$ is

$$q(x_t|x_0) = N(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I),$$

where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=0}^{t} \alpha_s$.

- **Reverse Process** The reverse process generates new data samples by reversing the forward process with a

learnable Markov chain. The prior distribution is $p(x_T) = N(x_T; 0, I)$ and the transition kernel is

$$p_\theta(x_{t-1}|x_t) = N(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)),$$

where $\theta$ denotes model parameters, and $\mu_\theta(x_t, t)$ and $\Sigma_\theta(x_t, t)$ are parameterized by deep neural networks. The reverse process starts from sampling $x_T \sim p(x_T)$ and iteratively samples $x_{t-1} \sim p_\theta(x_{t-1}|x_t)$ until $t = 0$.

- **Model Training** For each sample $x_0 \sim q(x_0)$, the model training objective is to maximizing the variational lower bound (VLB) of the log-likelihood of the data $x_0$. The simplified form $\mathcal{L}_{\text{VLB}}(x_0)$ is given by

$$\mathbb{E}_{q(x_{1:T}|x_0)}\left[ -\log p(x_T) - \sum_{t=1}^{T} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1})} \right]$$

With simplication and reparameterization trick, the overall objective $\mathbb{E}_{q(x_0)}\left[\mathcal{L}_{\text{VLB}(x_0)}\right]$ can be simplified into the final form

$$\mathbb{E}_{t \sim U[1,T], x_0 \sim q(x_0), \epsilon \sim \mathcal{N}(0,I)} \left[\lambda(t)\|\epsilon - \epsilon_\theta(x_t, t)\|\right]$$

where $\lambda(t)$ is a positive weighting function, $U[1,T]$ is a uniform distribution over the set $\{1, 2, \ldots, T\}$, and $\epsilon_\theta$ is a deep neural network with parameter $\theta$ that predicts the noise vector $\epsilon$ given $x_t$ and $t$. Note that, the overall objective is also equivalent to matching the joint distribution of the reverse process $p_\theta(x_0, x_1, \ldots, x_T)$ to that of the forward process $q(x_0, x_1, \ldots, x_T)$ by minimizing the KL divergence between them.

*4) GAN:* Generative Adversarial Networks (GANs) [14] are highly anticipated deep learning models with amazing capabilities which can simulate and generate realistic images, audio, and other data. Due to its outstanding performance, GANs have achieved significant achievements in various

fields [97]. The design inspiration of GANs comes from the zero-sum game in game theory.

As shown in Fig. 2, a typical GAN consists of two main components: a generator and a discriminator. These two parts compete with each other through adversarial learning, allowing the generator to continuously improve its ability to generate realistic samples, while the discriminator continuously improves its ability to distinguish between true and false samples.

## C. Retriever

Retrieval is to identify and obtain relevant information given an information need. Specifically, let's consider information resources that can be conceptualized as a key-value store $\{(k_i, v_i)\}_{i=1}^N$, where each key $k_i$ corresponds to a value $v_i$ ($k_i$ and $v_i$ can be identical). Given a query $q$, the objective is to search for the top-$k$ most similar keys using a similarity function $s$, and obtain the paired values. Based on different similarity functions, existing retrieval methods can be categorized into sparse retrieval, dense retrieval, and others. For widely used sparse and dense retrieval, the whole process can be divided into two distinct phases: in the first phase, each object is encoded into a specific representation; and in the second phase, an index is constructed to organize the data source for efficient search.

*1) Sparse Retriever:* Sparse retrieval methods are widely used in document retrieval, where the keys are actually documents to be searched (values are the same documents in this scenario). These methods leverage term matching metrics such as TF-IDF [98], query likelihood [99], and BM25 [19], which analyze word statistics from texts and construct inverted indices for efficient searching. Among them, BM25 is a hard-to-beat baseline in industrial-scale web search. For a query $q$ containing keywords $\{q_i\}_{i=1}^n$, the BM25 score of a document $D$ is:

$$s(q, D) = \sum_{i=1}^n IDF(q_i) \cdot \frac{f(q_i, D) \cdot (a+1)}{f(q_i, D) + a \cdot (1 - b + b \cdot \frac{|D|}{avgdl})}$$

where $IDF$ is the inverse document frequency weight, $f(q_i, D)$ is the number of times that $q_i$ occurs in the document $D$, $|D|$ is the length of $D$, $avgdl$ is the average document length in the corpus collection, $a$ and $b$ are tunable parameters.

$IDF$ is computed as:

$$IDF(q_i) = \ln(\frac{N - n(q_i) + 0.5}{n(q_i) + 0.5} + 1)$$

where $N$ is the number of documents, and $n(q_i)$ is the number of documents containing $q_i$. $IDF$ score is also used in TF-IDF.

To enable efficient search, sparse retrieval typically leverages an inverted index to organize documents. Concretely, each term from the query performs a lookup to obtain a list of candidate documents, which are subsequently ranked based on their statistical scores.

*2) Dense Retriever:* Unlike sparse retrieval, dense retrieval methods represent queries and keys using dense embedding vectors, and build approximate nearest neighbor (ANN) index to speed up the search. This can be applied to all modalities. For text data, recent advancements in pre-trained models,

including BERT [15] and RoBERTa [100], have been employed to encode queries and keys individually [20], [101]–[104]. Similar to text, models have been proposed to encode code data [25], [105], [106], audio data [26], [107], image data [24], [108], video data [109], [110]. etc. The similarity score between dense representations are usually computed with metrics such as cosine, inner product, L2-distance.

During training, dense retrieval usually follows a contrastive learning paradigm, making positive samples more similar and negative samples less similar. Several hard negative techniques [101], [111] have been proposed to further improve the model quality. During inference, approximate nearest neighbor (ANN) methods are applied for efficient searching. Various indices are developed to serve ANN search, such as tree [112], [113], locality sensitive hashing [114], neighbor graph index (e.g., HNSW [115], DiskANN [116], HMANN [117]), and the combination of graph index and inverted index (e.g., SPANN [22]).

*3) Others:* In addition to sparse retrieval and dense retrieval, there are alternative methods for retrieving relevant objects [118], [119]. Instead of calculating representations, some research works directly use the edit distance between natural language texts [120] or abstract syntax trees (AST) of code snippets [121], [122]. For knowledge graph, entities are linked with relations, which can be regarded as a pre-built index for retrieval searching. Therefore, RAG methods which involve knowledge graph can use $k$-hop neighbor search as retrieval process [123], [124]. Named entity recognition (NER) [125] is another way of retrieval, where the input is the query and the entites are the keys.

## III. METHODS

In this section, we introduce RAG foundations and outline enhancement methods that further improve the effectiveness.

### A. RAG Foundations

Based on how the retriever augments the generator, we categorize RAG foundations into 4 classes, as shown in Fig. 3.

*1) Query-based RAG:* Query-based RAG, originated from the idea of prompt augmentation, integrates the user's query with insights from information fetched during the retrieval process, directly into the initial stage of the language model's input. This paradigm stands as a widely adopted approach within the applications of RAG. After retrieval, the retrieved content will be merged with the original user query to create a composite input sequence which is subsequently fed into the generator for response generation. Query-based RAG has been widely applied across various modalities.

For Text Generation, REALM [33] employs a dual-BERT framework to streamline knowledge retrieval and integration, marrying pre-trained models with knowledge extractors. The initial BERT module processes the input question alongside documents to facilitate retrieval, utilizing MIPS for selecting the top-k documents with the highest probability and periodically updating the index. The document snippets obtained are then integrated with the query, feeding into the second BERT module to produce multiple outputs that are aggregated into a singular, comprehensive response. Lewis et al. [34] synergized pre-trained language models with knowledge retrieval
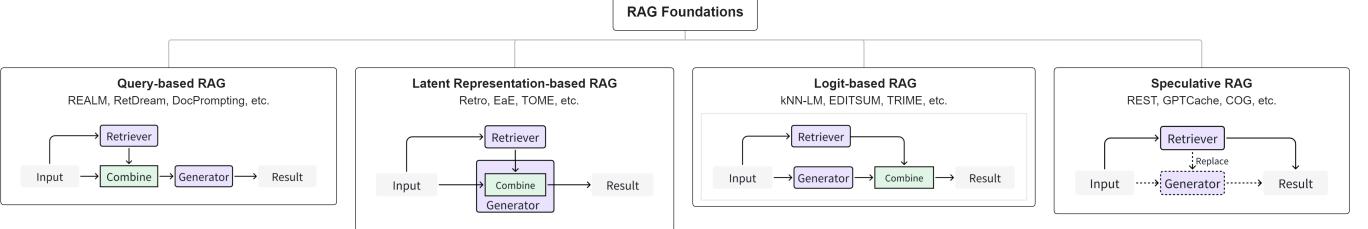
Fig. 3: Taxonomy of RAG foundations.

mechanisms, leveraging DPR and BART structures to accomplish retrieval-augmented generation tasks. DPR serves as the retrieval component, sourcing pertinent information from vast document databases, while BART uses this information for text generation. RAG-Token and RAG-Sequence differ in their retrieval timings, with the former retrieving information at each token generation and the latter conducting a single retrieval for the entire sequence. SELF-RAG [126] enhances the accuracy and relevance of responses by integrating a retrieval and critique strategy. Initially, the model employs a retriever to search for information paragraphs closely related to the input question. Subsequently, the critique model evaluates these paragraphs to determine their relevance and level of support of the retrieved text, assessing their impact on the generation of responses. Finally, the generator model constructs responses based on this information and evaluates the quality of these responses through critique marks. In addition to being compatible with local generators, Query-based RAG is also applicable to scenarios that use LLM through API calls. REPLUG [127] illustrates this methodology by treating the language model as a "black box", utilizing Contriever to seamlessly incorporate relevant external documents into the query. REPLUG LSR, a variant with LM-Supervised Retrieval, further refines this process by optimizing retrieval through language model-supervised insights, aiming to reduce perplexity scores and improve model performance by enriching its contextual understanding. In-Context RALM [128] uses BM25 for document retrieval and trains a predictive reranker to reorder and integrate the top-ranked documents.

In contemporary research on other modalities, augmenting inputs with retrieved contents (which are not limited to texts) has proven highly effective in enhancing the performance of various tasks. This strategy is applicable across several critical domains, including code generation, audio generation, and Knowledge Base Question Answering (KBQA).

For Text-to-Code task, APICoder [129] and DocPrompting [42] demonstrate how effectively integrating retrieved information into language models can improve the accuracy and relevance of generated code. In Automatic Program Repair task, CEDAR [130] and InferFix [131] utilize retrieved code snippets to aid the repair process, enhancing the model's understanding and application of repair strategies by combining them with the original input. For Code Completion task, ReACC [132] employs a prompting mechanism, leveraging retrieved code snippets as part of the new input to increase the accuracy and efficiency of code completion.

Recent researches in Knowledge Base Question Answering (KBQA) has also shown significant effects of combining retrieval and language models. For instance, Uni-Parser [133], RNG-KBQA [123], and ECBRF [134] effectively improve the performance and accuracy of QA systems by merging queries and retrieved information into prompts. BLLM augmentation [135] represents an innovative attempt at zero-shot KBQA using black-box large language models. This method, by directly integrating retrieved information into the model input without the need for additional sample training, demonstrates the great potential of combining retrieval and language models to enhance the model's generalization ability in understanding and answering unseen questions.

In the AI-for-Science domain, Chat-Orthopedist [136] provides support for shared decision-making among adolescents with idiopathic scoliosis. It enhances the application effectiveness and information accuracy of LLMs by integrating retrieved information into the prompts of the model.

In the task of Image Generation, RetrieveGAN [45] enhances the relevance and accuracy of generated images by integrating retrieved information, including selected image patches and their corresponding bounding boxes, into the input stage of the generator. IC-GAN [137] modulates the specific conditions and details of the generated images by concatenating noise vectors with instance features.

For 3D Generation, RetDream [50] initially utilizes CLIP [24] to retrieve relevant 3D assets, then merges the retrieved contents with the user input during the input phase.

Query-based RAG, often paired with LLM generators, offers modular flexibility, allowing swift integration of pre-trained components for quick deployment. Prompt design is crucial for utilizing retrieved data within this setup.

*2) Latent Representation-based RAG:* In Latent Representation-based RAG framework, retrieved objects are incorporated into generative models as latent representations. This enhances the model's comprehension abilities and improves the quality of the generated content.

The Fusion-in-Decoder (FiD) [35] technique leverages both BM25 and DPR for sourcing supportive paragraphs. It concatenates each retrieved paragraph and its title with the query, processing them individually through the encoder. FiD reduces computational complexity and efficiently utilizes relevant information to generate answers by fusing information from multiple retrieved paragraphs in the decoder, rather than processing each paragraph in the encoder. The application of Fusion-in-Decoder methodologies transcends the realm of textual content processing, demonstrating substantial potential and adaptability in processing code, structured knowledge,

and diverse multimodal datasets. Specifically within the code-related domain, technologies such as EDITSUM [138], BASH-EXPLAINER [139], and RetrieveNEdit [140] adopt the FiD approach, facilitating integration through encoder-processed fusion. Re2Com [141], and RACE [142] , among other methods, also feature the design of multiple encoders for different types of inputs.

In the field of Science, RetMol [55] also employs the Fusion-in-Decoder strategy, integrating information at the decoder stage to enhance the relevance and quality of the generated molecular structures.

In the field of Knowledge Base Question Answering (KBQA), the FiD method has been widely adopted, demonstrating significant effectiveness. UniK-QA [143], DE-CAF [144], SKP [145], KD-CoT [146], and ReSKGC [147] have effectively enhanced the performance of QA systems through the application of Fusion-in-Decoder technology. This illustrates that by integrating RAG for KBQA, the

Retro [36] pioneers the integration of retrieved text via "Chunked Cross-Attention" a novel mechanism that segments the input sequence into discrete chunks. Each chunk independently executes cross-attention operations, thereby mitigating computational burdens. This technique enables the model to selectively retrieve and assimilate distinct documents for varied sequence segments, fostering dynamic retrieval throughout the generation process. This enhances the model's adaptability and enriches the contextual backdrop of generated content. In the domain of image generation, cross-attention mechanisms have been widely adopted within RAG frameworks. Methods such as Re-imagen [148], KNN-Diffusion [149], RDM [150] and LAION-RDM & ImageNet-RDM [151] utilize cross-attention to integrate multiple retrieval results, effectively enhancing the overall performance of the models.

In addition, there are also some other novel structures worth our attention, Li [152] introduced the ACM, a text-image affine combination module, which notably does not employ any form of attention mechanism. Memorizing Transformers [31] revolutionize long document processing through the integration of a kNN-augmented attention mechanism within a Transformer layer. This innovation triggers a kNN search amidst input sequence processing, fetching data based on similarities between the sequence and stored key-value pairs, thereby elevating performance without necessitating complete retraining. This approach not only bolsters processing efficiency but also broadens the model's memory span, enabling self-retrieval from its generated outputs and fine-tuning for extensive knowledge bases or code repositories. Unlimiformer [153], by embedding a k-nearest neighbors (kNN) index within a pre-trained encoder-decoder transformer framework, pioneers handling inputs of indefinite length. Storing hidden states of input tokens in the kNN index allows for the efficient retrieval of highly relevant tokens during decoding. This innovation extends the model's capacity to manage prolonged sequences. Kuratov et al. [154] integrated Transformer with RNN, utilizing the model's intermediate output as the content for retrieval. This process was executed at each layer of the Transformer, thereby significantly extending the text window's length and effectively mitigating the issue of "Lost in the

Middle" [155]. Diverging from prior methods for knowledge, EaE [156] empowers language models to internalize explicit entity knowledge. EaE introduces an entity-specific parameterization, optimizing inference efficacy through an entity memory layer embedded within the transformer architecture. This layer directly acquires entity representations from textual data, utilizing a sparse retrieval strategy to fetch the nearest entities based on their embeddings, thus refining the model's comprehension through a calculated aggregation of entity-specific information. on this basis , TOME [157] shifts the focus towards comprehensive mention encodings, prioritizing the granularity of mention over mere entity representations. It meticulously creates a database that stores key and value encodings along with entity IDs, enabling the retrieval of much more fine-grained information. TOME integrates an initial transformer block to process input texts, followed by TOME blocks with memory attention layers, facilitating the synthesis of multifaceted information sources and enhancing inferential reasoning capabilities, even for unencountered entities.

In the field of 3D generation, ReMoDiffuse [51] introduces a semantics-modulated attention mechanism which enhances the accuracy of generating corresponding 3D motions based on textual descriptions. AMD [158] achieves efficient conversion from text to 3D motion by fusing the original diffusion process with the reference diffusion process.

In the Audio domain, Koizumi et al. [43] utilized an LLM, incorporating encoded dense features in the attention module to guide the generation of audio captions. Re-AudioLDM [159] utilizes distinct encoders to extract deep features from text and audio, which are then integrated into the attention mechanism of its Latent Diffusion Model (LDM).

For video captioning, R-ConvED [48] uses a convolutional encoder-decoder network to process retrieved video-sentence pairs with an attention mechanism, generating hidden states to produce captions. CARE [160] introduces a concept detector to produce concept probabilities, and incorporates concept representations into a hybrid attention mechanism. EgoInstructor [49] employs gated-cross attention to integrate textual inputs with encoded video features, enhancing the relevance and coherence of the generated captions for egocentric video content.

Finally, Latent Representation-based RAG is adaptable to various modalities and tasks. It obtains the hidden states of retrieved data, enabling seamless integration between retrievers and generators. However, additional training is required to align the latent space. Within this paradigm, we have the flexibility to design intricate and novel algorithms that effectively combine the information from retrieved contents.

*3) **Logit-based RAG**:* In logit-based RAG, generative models integrate retrieval information through logits during the decoding process. Typically, the logits are combined through simple summation or models to compute the probabilities for step-wise generation.

The kNN-LM [37] model integrates a pre-trained neural language model with the k-nearest neighbor search. It employs the pre-trained model to generate a list of candidate words and their probability distribution, while simultaneously performing retrieval from a data repository to find the k most relevant

neighbors based on the current context, thus enhancing the output of the original language model. The innovation at the core of this model lies in its ability for dynamic retrieval of information from a broad text corpus, significantly improving the accuracy and relevance of its predictions, particularly in addressing rare patterns and adapting to various fields. He et al. [38] introduced a new framework that is predicated on performing retrieval operations only when necessary, aimed at enhancing the inference efficiency of the kNN-LM model through an adaptive retrieval. This framework accelerates the model's inference speed by training a retrieval adapter, which automatically identifies and eliminates unnecessary retrieval actions in certain scenarios. This method allows the model to dynamically decide on the necessity of retrieval based on the current context, thereby balancing the trade-off between performance and efficiency, and substantially increasing inference speed while maintaining model performance.

Unlike previous methods that only merge memories during the testing time, TRIME [161] achieves memory merging during both training and testing phases, treating in-batch examples as accessible memory. It leverages new data batching and memory construction techniques to effectively utilize external memory. It employs BM25 scores to pack paragraphs with high lexical overlap into the same batch, constructing the training memory to further optimize model performance. NPM [162] is a non-parametric masked language model comprised of an encoder and a reference corpus. Unlike traditional models that apply a softmax over a finite vocabulary, NPM models a non-parametric distribution over the corpus. The encoder's role is to map phrases from the corpus into fixed-size vectors, filling in [MASK] by retrieving the phrase most similar to the masked position.

Beyond Text, other modalities, such as code and Image, also leverage logit-based RAG. For code-to-text conversion task, Rencos [121] generates multiple summary candidates in parallel from the retrieved code. It then normalizes these candidates using edit distance and calculates the final probability to select the summary output that best matches the original code. In code summarization task, EDITSUM [138] enhances the quality of summary generation by integrating prototype summaries at the probability level. For text-to-code tasks, the kNN-TRANX [163] model employs a combination of a confidence network and meta-knowledge to merge retrieved code fragments. It utilizes a seq2tree structure to generate target code that closely matches the input query, thereby increasing the accuracy and relevance of code generation. In image captioning tasks, MA [164] combines an attention-based encoder-decoder, using the image encoder to extract visual features to construct the semantic part, and decodes it word by word with the information retrieved. MA interpolates between two distributions generated by the caption decoder and the memory-augmented module to determine the distribution of the next word.

In conclusion, Logit-based RAG effectively leverages historical states (or other data sources) to infer the current state and combines information at the logit level. This approach is well-suited for sequence generation tasks. It couples retrieval and generation, primarily requiring training of the generator. Novel

approaches can be designed to effectively utilize the obtained probability distributions and adapt to subsequent tasks.

*4) Speculative RAG:* Speculative RAG seeks opportunities to use retrieval instead of pure generation, aiming to save resources and accelerate response speed. REST [32] replaces the small models in speculative decoding [165] with retrieval, enabling the generation of drafts. GPTCache [39] addresses the issue of high latency when using the LLM APIs by building a semantic cache for storing LLM responses. COG [166] decomposes the text generation process into a series of copy-and-paste operations, retrieving words or phrases from the documents instead of generation. Cao et al. [167] proposed a new paradigm to eliminate the dependence of the final result on the quality of the first-stage retrieved content, replacing generation with directly retrieved phrase level content.

In conclusion, Speculative RAG is currently primarily applicable to sequential data. It decouples the generator and the retriever, enabling the direct use of pre-trained models as components. Within this paradigm, we can explore a wider range of strategies to effectively utilize the retrieved content.

## B. *RAG Enhancements*

In this section, we introduce methods which enhance the performance of a constructed RAG system. We categorize existing methods into 5 groups based on their enhancement targets: input, retriever, generator, result, and the entire pipeline.

*1) Input Enhancement:* The input, initially fed into the retriever, significantly impacts the final outcome of the retrieval stage. In this section, we introduce two methods for input enhancement: query transformation and data augmentation.

*a) Query Transformation:* Query transformation can enhance the result of retrieval by modifying the input query.

Query2doc [168] and HyDE [169] use the original query to generate a pseudo document, which is later used as the query for retrieval. The pseudo document contains richer relevant information, which helps to retrieve more accurate results. TOC [170] employs recursive retrieval augmented clarification on ambiguous questions to construct a tree of disambiguated questions, thereby generating comprehensive answers to these ambiguous questions. During the construction process of this tree structure, TOC utilizes a self-verification pruning method to ensure the factual relevance of each node.

*b) Data Augmentation:* Data augmentation improves data before retrieval, including techniques such as removing irrelevant information, eliminating ambiguity, updating outdated documents, synthesize new data, etc.

Make-An-Audio [44] uses captioning and audio-text retrieval to generate captions for language-free audio to mitigate data sparsity, and adds random concept audio to improve the original audio. LESS [171] strategically selects an optimal dataset for downstream tasks by analyzing gradient information. It aims to maximize the dataset's impact on fine-tuning the model's performance in response to instructional prompts. ReACC [132] employs data augmentation (including renaming and dead code insertion) to pre-train the code retrieval model.

*2) Retriever Enhancement:* The retrieval process is crucial in RAG systems. Generally, the better the retrieved content quality, the easier it is to stimulate the ability of LLMs in-context learning as well as other generators and paradigms.
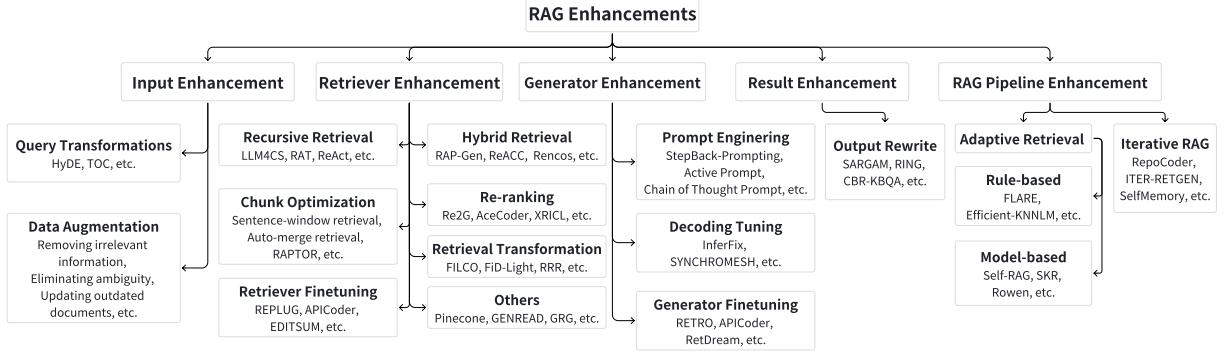
Fig. 4: Taxonomy of RAG Enhancements.

The worse the content quality, the more likely it is to cause model hallucinations. Therefore, in this section, we will discuss how to efficiently improve the effectiveness of the retrieval process.

*a) Recursive Retrieval:* Recursive retrieval is to perform multiple searches to retrieve richer and higher-quality contents.

ReACT [172] uses Chain-of-Thought (CoT) [173] to break queries down for recursive retrieval and provide richer information. RATP [174] uses the Monte-Carlo Tree Search (MCTS) to perform multiple simulations, identifying optimal retrieval content. This content is then integrated into a template and sent to the generator for final production.

*b) Chunk Optimization:* Chunk optimization refers to adjusting chunk size for improved retrieval results. Sentence-window retrieval [175] is an efficient approach that enhances retrieval by fetching small chunks of text and returning a window of relevant sentences surrounding the retrieved segment. This method ensures that the context before and after the targeted sentence is included, providing a more comprehensive understanding of the retrieved information. Auto-merge retrieval is another advanced RAG method of LlamaIndex [175] which organizes the document in a tree-like structure, with the parent node containing the content of all children nodes. For example, articles and paragraphs, as well as paragraphs and sentences, all follow a parent-child relationship. In the retrieve process, fine-grained search for children nodes ultimately returns the parent node, effectively providing richer information. To address the lack of contextual information, RAPTOR [176] employs recursive embedding, clustering, and summarization of text chunks until further clustering becomes infeasible, thereby constructing a multi-level tree structure.

*c) Retriever Finetuning:* As a core component in the RAG system, the retriever plays a crucial role in the entire system operation process. An effective embedding model can cluster semantically similar content in vector space, enhancing the retriever's capability to provide valuable information for the subsequent generator, thus boosting the RAG system's efficiency. Hence, the proficiency of the embedding model [177]–[180] is vital for the RAG system's overall effectiveness.

In addition, for embedding models that already have good expression power, we can still finetune them using high-quality domain data or task related data to improve their performance in specific domains or tasks. REPLUG [127] treats LM as a black box and update the retriever model based on the final results. APICoder [129] finetunes the retriever with python files and api names, signature, description. EDITSUM [138] finetunes the retriever to decrease the jaccard distance between summaries after retrieval. SYNCHROMESH [122] adds tree distance os ASTs in the loss and uses Target Similarity Tuning to finetune the Retriever. R-ConvED [48] finetunes the Retriever with the same data as generator. Kulkarni et al. [181] applied infoNCE loss to finetune the Retriever.

*d) Hybrid Retrieval:* Hybrid retrieve denotes the concurrent employment of a diverse array of retrieval methodologies or the extraction of information from multiple distinct sources.

RAP-Gen [182] and ReACC [132] use both dense retriever and sparse retriever to improve the quality of retrieval. Rencos [121] uses sparse retriever to retrieve similar code snippets on syntactic-level and usse dense retriever to retrieve similar code snippets on semantic-level. BASHEXPLAINER [139] first uses dense retriever to capture semantic information and then uses sparse retriever to acquire lexical information. RetDream [50] first retrieves with text and then retrieves with the image embedding. CRAG [183] has designed a retrieval evaluator to assess the relevance of retrieved documents to the input query, triggering three types of retrieval actions based on varying confidence levels: if deemed correct, it directly uses the retrieval results for Knowledge Refinement; if incorrect, it resorts to Web Search; and if ambiguous, it combines both approaches. To enhance performance in question-and-answer tasks, Huang et al. [184] introduced two metrics, DKS(Dense Knowledge Similarity) and RAC(Retriever as Answer Classifier), during the retrieval phase. These metrics account for both the pertinence of the answers and the applicability of the underlying knowledge. UniMS-RAG [185] introduces a novel kind of token, termed as the "acting token", which determines the source from which to retrieve information.

*e) Re-ranking:* The Rerank technique refers to reordering the retrieved content in order to achieve greater diversity and better results. Re2G [186] applies a re-ranker [187] model after the traditional retriever. The effect of the re-ranker model is to re-rank retrieved documents, the purpose of which is to reduce the impact of information loss caused by compressing text into vectors on the quality of retrieval. AceCoder [188] reranks the retrieved programs with a selector to reduce redundant programs and obtain diverse retrieved programs. XRICL [189] uses a distillation-based exemplar

reranker after retrieval. Rangan [190] employs the Quantized Influence Measure, assessing statistical biases between a query and a reference to evaluate the similarity of data subsets and rerank retrieval results. UDAPDR [191] uses LLMs to cost-effectively generate synthetic queries that train domain-specific rerankers, which then apply multi-teacher knowledge distillation to develop a cohesive retriever. LLM-R [192] iteratively trains the retriever, using feedback from a frozen LLM to rank candidate documents and train the reward model. The retriever is further trained based on knowledge distillation. Each iteration of training builds upon the retriever trained in the previous cycle, facilitating iterative optimization in subsequent rounds.

*f) Retrieval Transformation:* Retrieval Transformation involves rephrasing retrieved content to better activate the generator's potential, resulting in improved output.

FILCO [193] effectively filters out irrelevant content from the retrieved text chunk, leaving only the precise supporting content. This process simplifies the task for the generator, making it easier to predict the correct answer. FiD-Light [194] initially employs an encoder to convert the retrieved content into a vector, which it then compresses, resulting in a substantial reduction of latency time. RRR [195] integrates the current query with the top-k document in each round through a template, and subsequently restructures it via a pre-trained LLMs(GPT-3.5-Turbo etc.).

*g) Others:* In addition to the above optimization methods, there are also some other optimization methods for the retrieve process. For example, Meta-data filtering [196] is a method to help processing retrieved documents which uses metadata (such as time, purpose, etc.) to filter the retrieved documents for better results. GENREAD [197] and GRG [198] introduce a novel approach where the retrieval process is supplanted or improved by prompting a LLM to generate documents in response to a given question.

*3) Generator Enhancement:* In RAG systems, the quality of the generator often determines the quality of the final output results. Therefore, the ability of the generator determines the upper limit of the entire RAG system's effectiveness.

*a) Prompt Engineering:* Technologies in prompt engineering [199] that focus on improving the quality of LLMs' output, such as Prompt compression, Stepback Prompt [200], Active Prompt [201], Chain of Thought Prompt [173], etc., are all applicable to LLM generators in RAG systems. LLM-Lingua [202] applies a small model to compresses the overall length of the query to accelerate model inference, relieving the negative impact of irrelevant information on the model and alleviating the phenomenon of "Lost in the Middle" [155]. ReMoDiffuse [51] decomposes complex descriptions into anatomical text scripts by using ChatGPT. ASAP [203] add exemplar tuples to the prompt for better results. An exemplar tuple is composed of the input code, a function definition, the results of analyzing that definition and its associated comment. CEDAR [130] uses a designed prompt template to organize code demonstration, query, and natural language instructions into a prompt. XRICL [189] utilizes COT technology to add translation pairs as an intermediate step in cross linguistic semantic parsing and inference. AC-

TIVERAG [204] employs the Cognition Nexus mechanism to calibrate the intrinsic cognition of LLMs and applies COT prompt in answer generation. Make-An-Audio [44] is able to use other modalities as input which can provide much richer information for the following process.

*b) Decoding Tuning:* Decoding tuning refers to adding additional controls during the generator processing, which can be achieved by adjusting hyperparameters to achieve greater diversity, limiting the output vocabulary in some form, and so on.

InferFix [131] balances the diversity and quality of results by adjusting the temperature in decoder. SYN-CHROMESH [122] limits the output vocabulary of the decoder by implementing a completion engine to eliminate implementation errors.

*c) Generator Finetuning:* The finetuning of the generator can enhance the model's ability to have more precise domain knowledge or better fit with the retriever.

RETRO [36] fixes the parameters of the retriever and uses the chunked cross attention mechanism in the generator to combine the content of the query and retriever. APICoder [129] finetunes the generator CODEGEN-MONO 350M [205] with a shuffled new file combined with API information and code blocks. CARE [160] first uses image data, audio data and vedio-text pairs to train encoders and then finetune the decoder (generator) with the target of decreasing caption loss and concept detection loss together, during which the encoders and the retriever are frozen. Animate-A-Story [206] optimizes the video generator with image data, and then finetunes a LoRA [207] adapter to capture the appearance details of the given character. RetDream [50] finetunes a LoRA adapter [207] with the rendered images.

*4) Result Enhancement:* In many scenarios, the result of RAG may not achieve the expected effect, and some techniques of Result Enhancement can help alleviate this problem.

*a) Output Rewrite:* Output Rewrite refers to rewriting the content generated by the generator in certain scenarios to meet the needs of downstream tasks.

SARGAM [208] refines outputs in code-related tasks by employing a special Transformer alongside Deletion, Placeholder, and Insertion Classifiers to better align with the real-world code context. Ring [209] obtains diversity results by reranking candidates based on the average of per token log probabilities produced by the generator. CBR-KBQA [54] revises the result by aligning generated relations with those presented in the local neighborhood of the query entity in knowledge graph.

*5) RAG Pipeline Enhancement:* RAG Pipeline Enhancement refers to optimizing the processes of RAG at the system level in order to achieve better performance results.

*a) Adaptive Retrieval:* Some studies and practices on RAG have shown that retrieval is not always beneficial for the final generated results When the parameterized knowledge of the model itself is sufficient to answer relevant questions, excessive retrieval will cause resource waste and may increase the model's confusion. Therefore, in this chapter, we will discuss two types of methods for determining whether to retrieve, named rule-based and model-based methods.

**Rule-based:** FLARE [210] actively decides whether and when to search through the probability in the generation process. Efficient-KNNLM [38] combines the generation probability of KNN-LM [37] and NPM [162] with a hyperparameter $\lambda$ to determine the proportion of generation and retrieval. Mallen et al. [211] used statistical analysis on questions to enable direct answers for high-frequency ones and applied RAG for low-frequency ones. Jiang et al. [212] studied Model Uncertainty, Input Uncertainty, and Input Statistics to comprehensively assess the confidence level of the model. Ultimately, based on the confidence level of the model, a decision is made whether to retrieve. Kandpal et al. [213] studied the correlation between the number of relevant documents and the model's knowledge mastery to assess the need for retrieval.

**Model-based:** Self-RAG [126] uses a trained generator to determine whether to perform a retrieval based on the retrieve token under different instructions, and evaluates the relevance and level of support of the retrieved text through the Self-Reflection token. Finally, the quality of the final output result is evaluated based on the Critique token. Ren et al. [214] used "Judgment Prompting" to determine whether LLMs can answer relevant questions and whether their answers are correct or not, thereby assisting in determining the necessity of a retrieval. SKR [215] uses the ability of LLMs themselves to judge in advance whether they can answer the question, and if they can answer, no retrieval is performed. Rowen [216] employs a model as a sophisticated multilingual detection system to evaluate the semantic coherence of answers to identical questions posed across various languages. In the event of detected inconsistencies, it decides to retrieve external information, thereby enhancing the reasoning process and rectifying inaccuracies. Conversely, when responses exhibit consistency, the system upholds the initially generated answer, which is derived from internal reasoning. AdaptiveRAG [217] dynamically decides whether to retrieve based on the query complexity by a classifier, which is a smaller LM.

*b) Iterative RAG:* Iterative RAG progressively refines results by repeatedly cycling through retrieval and generation phases, rather than a single round.

RepoCoder [218] employs an iterative retrieval-generation pipeline for code completion tasks, enhancing each retrieval query with code generated in prior iterations to more effectively leverage information dispersed across various files, thereby achieving superior results. ITER-RETGEN [219] synergizes retrieval and generation in an iterative manner. The current output of the generator can to some extent reflect the knowledge it still lacks, and the retrieve can retrieve the missing information as contextual information for the next round, which helps to improve the quality of the generated content in the next round. SelfMemory [220] employs a retrieval-augmented generator in an iterative manner to create an unlimited memory pool. Following this, a memory selector is used to choose one output, which then serves as the memory for the subsequent generation round. RAT [221] initial generates content by an LLM with a zero-shot CoT prompt, then revises each thought step by retrieving knowledge from external knowledge base.

## IV. APPLICATIONS

In this section, we focus on RAG applications spanning various modalities. To echo with the taxonomy of RAG foundations and enhancements, we also demonstrate their utilization across different tasks in Table I.

### A. *RAG for Text*

To begin with, text generation is among the most important and widely deployed applications for RAG. Here we introduce popular works for seven tasks, respectively.

*1) Question Answering:* Question Answering involves the process of providing responses to posed questions by drawing from a vast and comprehensive collection of textual sources.

FiD [35] and REALM [33] identify the top-k most pertinent article snippets based on the query and forward each snippet along with the question to LLMs to generate k responses. These responses are then synthesized into a final answer. Toutanova et al. [222] substituted the text corpus in REALM with subgraphs from a knowledge graph, yielding impressive results. As shown in Fig. 5, RETRO [36] employs attention mechanisms to integrate the question with relevant retrieved documents within the model to produce the final answer. SKR [215] observes that using RAG does not invariably benefit Question Answering and thus explored guiding the model to evaluate its grasp of pertinent knowledge, subsequently adapting its use of external resources for retrieval enhancement. TOG [223] introduces an innovative knowledge graph-augmented LLM framework, which excels by fostering interactions between LLMs and the Knowledge Graph and by expanding the inference path space with beam search. NPM [162] pioneers the use of nonparametric data distributions in lieu of the softmax layer, enabling models with fewer parameters to perform effectively. Self-RAG [126] improves answer quality by learning to discern when to retrieve, assess the retrieved content's relevance, and evaluate the final generated results using four types of reflective tokens. CL-ReLKT [224] employs a language-generalized encoder to bridge the gap between question-document pairs across languages, thus better leveraging multilingual data. CORE [225] mitigates language resource disparities by introducing a novel dense passage retrieval algorithm and a multilingual autoregressive generation model. Lastly, EAE [156] enhances answer quality by retrieving entity embeddings for query entities and integrating these with hidden states for further processing. UR-
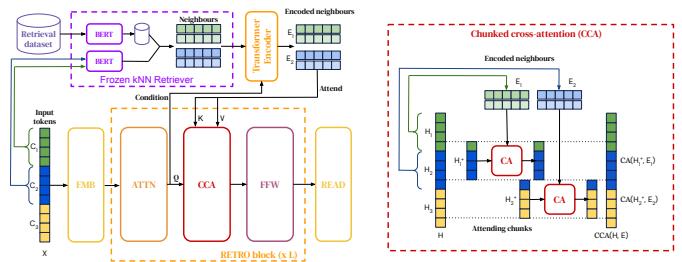


Fig. 5: Architecture of RETRO [36] model.

QA [226] found that when encountering unseen problems, retrieving QA pairs has a better final effect; When encountering problems that have not been seen before, the retrieve

TABLE I: Taxonomy of RAG applications across various modalities.

| RAG for Text | | | | |
|---|---|---|---|---|
| **Question Answering** | **Human-Machine Conversation** | **Neural Machine Translation** | **Summarization** | **Others** |
| REALM‡§  TKEGEN§  TOG‡<br>SKR¶  Self-RAG§¶  RIAG§<br>FiD§  RETRO§  NPM‡§ | CREA-ICL†‡  BlenderBot3‡§<br>CEG‡∥  Internet-Augmented-DG‡§<br>ConceptFlow‡§  Skeleton-to-Response‡§ | NMT-with-Monolingual-TM†‡§<br>TRIME‡§  KNN-MT‡§  COG‡ | RAMKG‡§  RPRR‡  RIGHT‡§<br>Unlimiformer§ | CONCRETE‡§  Atlas‡§<br>KG-BART‡§  R-GQA‡§ |

| RAG for Code | | | | | |
|---|---|---|---|---|---|
| **Code Generation** | **Code Summarization** | **Code Completion** | **Automatic Program Repair** | **Text-to-SQL and Code-based Semantic Parsing** | **Others** |
| SKCODER§  RRGCode‡<br>ARKS†  RECODE<br>KNN-TRANX∥  Toolcoder§∥ | RACE†  BASHEXPLAINER‡<br>READSUM∥  Rencos‡<br>CoRec‡  Tram§  EDITSUM‡ | ReACC†‡  RepoCoder†§¶<br>De-Hallucinator¶  REPOFUSE§<br>RepoFusion§  EDITAS§ | RING∥  CEDAR§<br>RAP-Gen‡§  InferFix§<br>SARGAM§  RTLFixer‡§ | XRICL‡§  SYNCHROMESH‡§<br>RESDSQL§  REFSQL‡§<br>CodeICL§  MURRE∥¶ | De-fine‡∥  Code4UIE§<br>E&V  StackSpotAI‡§<br>ImputBlaster¶ |

| RAG for Knowledge | | | | | RAG for 3D |
|---|---|---|---|---|---|
| **Knowledge Base QA** | **Knowledge-augmented Open-domain QA** | | **Table for QA** | **Others** | **Text-to-3D** |
| CBR-KBQA‡§∥  TIARA†‡§  Keqing†‡§<br>RNG-KBQA‡∥  ReTraCk§  SKP†‡§ | UniK-QA†‡  KG-FiD‡  GRAPE‡<br>SKURG†‡  KnowledGPT‡  EFSUM§ | | EfficientQA‡  CORE§  Convinse†‡<br>RINK‡§  T-RAG‡§  StructGPT‡ | GRetriever§  SURGE§<br>K-LaMP  RHO∥ | ReMoDiffuse†‡<br>AMD† |

| RAG for Image | | | RAG for Video | | |
|---|---|---|---|---|---|
| **Image Generation** | **Image Captioning** | **Others** | **Video Captioning** | **Video QA & Dialogue** | **Others** |
| RetrieveGAN‡  IC-GAN§  Re-imagen§<br>RDM  Retrieve&Fuse§  KNN-Diffusion | MA∥  REVEAL‡  SMALLCAP†<br>CRSR†  RA-Transformer | PICa∥  Maira‡<br>KIF‡  RA-VQA‡ | KaVD‡§  R-ConvED‡§<br>CARE§  EgoInstructor†‡§ | MA-DRNN†‡  R2A‡<br>Tvqa+§  VGNMN‡ | VidIL†‡  RAG-Driver‡<br>Animate-A-Story†§ |

| RAG for Science | | | RAG for Audio | |
|---|---|---|---|---|
| **Drug Discovery** | **Biomedical Informatics Enhancement** | **Math Applications** | **Audio Generation** | **Audio Captioning** |
| RetMol†§  PromptDiff†‡ | PoET‡  Chat-Orthopedist†§  BIOREADER†  MedWriter‡  QARAG†‡ | LeanDojo‡  RAG-for-math-QA†‡ | Re-AudioLDM§  Make-An-Audio†§ | RECAP‡§ |

Query-based  Latent-based  Logit-based  Speculative  Query+Latent  Latent+Logit    † Input  ‡ Retriever  § Generator  ∥ Output  ¶ Pipeline

text chunk performs better. Therefore, it is proposed to simultaneously retrieve QA pairs and text chunks, and select the final answer by comparing the calibrated confidences. DISC-LawLLM [227] constructs a supervised fine-tuning dataset through a legal syllogism prompting strategy, enabling the model to receive support from the latest legal information. RAG-end2end [228] conducts simultaneous training of the retriever (DPR) and the generator (BART) to optimize performance for the end-to-end question-answering task and to facilitate domain adaptation. MultiHop-RAG [229] is designed to extract pertinent information from a variety of distinct documents, aggregating this knowledge to equip the generator with the necessary context for producing the definitive answer to the query.

*2) Fact Verification:* Fact Verification involves assessing the veracity of information, a critical function in disciplines such as Natural Language Processing (NLP), Information Retrieval, and Data Mining. In today's digital era, characterized by an exponential increase in data, particularly across social media and online news platforms, there is a rapid proliferation of unchecked information. Fact verification plays an essential role in countering the spread of fake news, deceptive content, and rumors, thereby preserving the integrity of the information landscape and ensuring the public's access to accurate knowledge. Consequently, automated fact verification systems are of immense importance, with broad applications and significant practical value. CONCRETE [230] leverages cross-lingual retrieval mechanisms to tap into a wealth of multilingual evidence, effectively bridging the gap in resources for languages that are underrepresented in fact-checking datasets. Hagström et al. [231] proved on LLaMA [4] and Atlas [30] that search augmentation is more beneficial for solving inconsistency

problems than increasing model size. Atlas [30] shows that using RAG to support LLMs in knowledge-intensive tasks markedly improves their few-shot learning performance.

*3) Commonsense Reasoning:* Commonsense Reasoning entails the capability of machines to infer or make decisions on problems or tasks in a human-like manner, drawing upon their acquired external knowledge and its application. However, the vast scope of common sense knowledge and the intricacies of reasoning processes make Commonsense Reasoning a perennially challenging and prominent area of research within the field of NLP. KG-BART [232] expands the conceptual landscape by incorporating intricate interrelations among diverse concepts within a knowledge graph. It employs graph attention mechanisms to aid LLMs in crafting more nuanced and logically coherent sentences. This approach not only improves the models' generalization capabilities but also significantly bolsters their Commonsense Reasoning proficiency. Wan et al. [233] constructed the CONFLICTINGQA dataset, comprising contentious questions and conflicting answer documents, to examine which textual features significantly influence LMs' ability to independently navigate controversial issues. The findings reveal that LMs often neglect the stylistic aspects of text that are typically valued by people.

*4) Human-Machine Conversation:* Human-Machine Conversation encompasses the ability of machines to comprehend natural language and adeptly employ this skill to engage with humans seamlessly. This capability represents a significant challenge within the realms of Artificial Intelligence and Natural Language Processing and offers a broad spectrum of practical applications. As such, Human-Machine Conversation continues to be a focal point of research for many scholars. ConceptFlow [234] leverages a commonsense knowledge graph to

structure conversations, directing the flow of dialogue based on attention scores, and propelling the conversation forward. This method achieves commendable results even with a substantial reduction in model parameters. Cai et al. [235] reimagined the text generation task as a cloze test by retrieving and distilling the essence of past conversational history, leading to notable outcomes. Komeili et al. [236] augmented dialogue generation quality by harnessing advanced search engine technologies to source pertinent content from the internet. BlenderBot3 [237] broadens its search horizon, not only mining relevant internet content but also local dialogue history, and employs entity extraction among other techniques to refine the quality of the resulting dialogue. Kim et al. [238], PARC [239], and CREA-ICL [240] improve the caliber of non-English conversations by incorporating cross-lingual knowledge, effectively addressing the scarcity of non-English datasets and enhancing the quality of the generated dialogue. CEG [241] addresses hallucination issues through a post-processing mechanism, verifying LLM-generated answers through retrieval. If the answer is correct, the retrieved document is added to the original answer as a reference; if the answer lacks reliable references, it guides the LLM to respond to the question anew.

*5) Neural Machine Translation:* Neural Machine Translation (NMT) is the automated process of translating text from a source language to a target language [161], [242], [243]. It is a pivotal task in the domain of NLP and represents a significant objective in the pursuit of AI, boasting considerable scientific and practical significance. Cai et al. [242] proposed an innovative approach that utilizes monolingual corpora alongside multilingual learning techniques, challenging the traditional dependency on bilingual corpora in Neural Machine Translation. This approach ensures that the retrieval system provides ample information while simultaneously optimizing both the retrieval mechanism and the translation model, culminating in impressive performance. KNN-MT [243] executes translation tasks at the token level by computing vector space distances. TRIME [161] effectively minimizes the discrepancy between training and inference phases by jointly training the retrieval system and the generation model, thereby enhancing the precision of translations.

*6) Event Extraction:* Event Extraction is a specialized task within Natural Language Processing (NLP) that focuses on pinpointing and extracting instances of particular event types from unstructured textual data. An event is generally characterized by a central action or predicate and the related entities, which can include participants, temporal indicators, locations, and other relevant attributes. The objective of event extraction is to convert the nuanced details embedded within text into a structured format, thereby facilitating advanced analysis, efficient information retrieval, and practical downstream applications. R-GQA [244] employs a retrieval-based approach to enhance the context of a given issue by identifying and utilizing the most closely aligned Question-Answer pair from a repository, thereby enriching the information available for processing the current query.

*7) Summarization:* In the realm of NLP, Summarization is a task aimed at distilling the essential information from lengthy texts and producing a concise, coherent summary that

encapsulates the primary themes. Summarization enables users to quickly grasp the essence of a text, thereby conserving time that would otherwise be spent on reading extensive material. There are two main approaches to Summarization: Extractive and Abstractive. Extractive Summarization involves the automatic selection and compilation of key phrases directly from the source text. A key phrase succinctly captures the main themes, content, or perspectives of the text and is typically composed of one or several words. The generation of key phrases is instrumental for understanding, categorizing, retrieving, and organizing textual information. It is extensively applied in fields such as search engine optimization, academic research, text summarization, and more. This technique refrains from creating new sentences, instead repurposing segments from the original text. Abstractive Summarization, on the other hand, entails comprehending the original text's meaning and reformulating it into new sentences [153], [245]–[247]. This approach can convey the source's intent more fluidly but poses greater challenges in terms of implementation due to its complexity. RAMKG [245] effectively leverages a comprehensive English corpus to bolster the performance of Keyphrase Generation in non-English contexts. It does so by enhancing the alignment of keywords extracted from texts in different languages that share similar subject matter. Unlimiformer [153] addresses the issue of input length constraints in transformer-based models by retrieving and utilizing the top-k most relevant hidden states, thereby extending the model's capacity to handle longer inputs. RPRR [246] employs a Retrieve-Plan-Retrieve-Read approach to overcome the limited context window constraints faced by LLMs, utilizing retrieved information to generate high-quality Wikipedia documents for emerging events. RIGHT [247] chooses to use different types of retrievers in different datasets to enhance the generator, which can effectively improve the quality of automatically generated labels in simple deployment.

*B. RAG for Code*

Separate retrieval and generation approaches have historically been employed for code-related tasks. For retrieval, similar code snippets can be identified using Abstract Syntax Trees (AST) or text edit distance. For generation, sequence-to-sequence models are employed to generate code or natural language. Recent RAG research combines both retrieval and generation techniques to enhance the overall performance.

*1) Code Generation:* The goal of code generation is to transform natural language (NL) descriptions into code implementation, which can be seen a process of text-to-code. Therefore, LSTM and transformer models are widely used for generator. Whether to use code-specific retrieval or text-based retrieval depends on the contents to be searched.

Retrieval-based prompt engineering is one of the most prevalent scenarios of RAG in code generation. In-context learning includes training samples in prompts as the input for sequence-to-sequence generative models. Retrieval techniques are adopted to find similar training samples to the test input, so that the prompt can be more informative and related. REDCODER [40] retrieves similar NL descriptions using dense retriever CodeBert [25], then concatenates the NL texts, their paired codes, for downstream generator PLBART [41].

CodeT5Mix [248] adopts the paradigm of RECODER, proposing a model that functions dually as both the retriever and generator. APICoder [129] first train a Bert-based deep retriever to align the embeddings of NL descriptions and API documentation; then, it retrieves relevant API information to build prompt for the generator CODEGEN-MONO [249]. The following work [250] uses the same pipline, renaming the retriever module as APIFinder. COCOGEN [251] aims at commonsense reasoning, which generates code-based reasoning-graphs with Codex [2] given NL inputs; it adds an evaluation setting of dynamic prompt selection, which actually retrieves relevant examples for prompts. In DocPrompting [42] given an NL intent, the retriver retrieves relevant documentations, then the generator generates codes based on the NL and retrieved documents. It evaluates both sparse retrievers and dense retrievers, and also tries different generators in experiments. CodeT5+ [252] adopts the CodeT5 [106] model for both the retriever and generator, leveraging only the encoder part in the retrieval process. AceCoder [188] fixes the retriever to BM25 [19], and tests several LLM generators for code generation. A3CodGen [253] extracts local information, retrieves relevant global functions using embeddings, and incorporates third-party library information to construct the prompt for LLM-based code generation. SKCODER [254] employs BM25 to retrieve relevant code snippets, which are further processed to produce sketch template. The template and the original description are concatenated for final generation. CodeGen4Libs [255] uses RAG for both import statements and codes, employing BM25 as retriever and finetuned CodeT5 as generator. CODEAGENT [256] design agents to search on web, retrieve relevant documentation, generate programs, and test correctness. RRGCode [257] retrieves relevant code snippets using both sparse and dense retrieval, employs a cross-encoder to re-rank the retrieval results, then generates code using the concatenation of the query and the retrieved codes. A recent study [258] shows that retrieval-augmented framework for code suggestions, including code generation and code completion, can improve the performance by a large margin. ARKS [259] steps further upon prompt-based RAG, incorporating iterative RAG to re-formulate queries and update knowledge soup (containing documentation, execution feedback, generated code, etc.) for dense retrieval, improving final generation accuracy.

Retrieval results can be applied during the generation process as well. RECODE [120] retrieves NL descriptions and paired codes using edit distance, then extracts n-gram action subtrees from codes' ASTs. During LSTM-based generation, the patterns of processed subtrees are leveraged to increase the corresponding word probability at each decoding step. kNN-TRANX [163] uses seq2tree model BertranX [260] to convert NL to code AST. It constructs a datastore for each code AST prefix and NL pair; i.e., for each NL-code pair, the context representation of the i-th context is obtained by encoding NL and the i-th prefix of code's AST through the seq2tree model. At each decoding step of the generation, the hidden representations are searched within the datastore to form a new probability, which is later combined with the seq2tree model's output through a confidence network.

ToolCoder [261] performs normal code generation, and conducts online search or offline retrieval when encountering special $< API >$ token. This paradigm makes the model learn to leverage API tools.

*2) Code Summarization:* The goal of code summarization is to transform code into natural language (NL) descriptions, which is a process of code-to-text. Same to code generation, many sequence-to-sequence models are applied as generator.

In many research works, the retrieval results are processed by additional encoders. Re2Com [141] and EditSum [138]
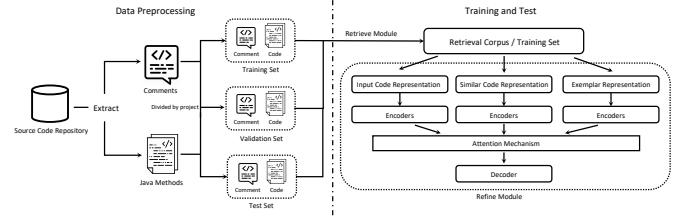


Fig. 6: Architecture of Re2Com [141] model.

both retrieve similar code using sparse retrieval BM25 and generate summary using LSTM. As shown in Fig. 6, they separately encode the input, the retrieved code, and the corresponding summary, then combine the middle representations or probabilities in the decoder. HGNN [262] instead uses code edit distance for retrieval, and substitutes the encoders for codes with hybrid GNN on their Code Property Graphs (CPG) [263]. RACE [142] aims at generating commit message for code difference. It leverages dense retrieval for similar examples and transformer model for generation. It also uses three encoders for the input, the retrieved code difference, and corresponding commit message, then combines the results before feeding into the decoder. BASHEXPLAINER [139] shares the similar idea. Its dense retrieval module is based on CodeBert [25], and for generation, the output representations of the input and the similar code from CodeBert are directly fused for transformer-based decoder. READSUM [264] retrieves similar code using Levenshtein distance, applies code encoder and summary encoder to retrieved pairs, and generates summary using a decoder where a fusion network combines the information.

RAG for in-context learning, which retrieves similar examples and build prompt for generation, also works in code summary. REDCODER [40] works for both code generation and summary, and it replaces the retriever with GraphCode-Bert [105] for code summary task. ASAP [203] retrieves similar code with BM25, and generates summary with GPT models. Similar to the paradigm described above, research on pseudocode generation [265] employs the retrieved code as input for generation and subsequently replaces the results with the original input. SCCLLM [266] retrieves relevant code snippets by semantic, syntactic, and lexical-based retrieval, then forms in-context prompts for smart contract comment generation via LLM. UniLog [267] retrieves similar code snippets paired with their log statements to conduct in-context learning for log statement generation.

Logit-based RAG also prevails in code summarization. Rencos [121] utilizes two different methods to retrieve similar

code snippets, syntactic similarity between abstract syntax trees (AST) and cosine similarity between dense representations. For generator, it adopts attention-based LSTM. There are three encoder-decoder models for the code input and two retrieval results respectively, and the probabilities are combined for final generation. CoRec [268] generates commit message given code diff and retrieved similar code diff. Multiple LSTM generators handle the input code diff and the retrieved code diff, then adds the probabilities for final generation. kNN-Transformer [269] obtains context vector by feeding code into an encoder-decoder generator, then gets three parts of logits: the first is from searching the vector, the second is from normal Transformer, and the third is the copy mechanism that reserve rare tokens from the input. Tram [270] involves retrieval in both token-level and sentence-level. It encodes source code and corresponding AST into hidden states representations; for token-level, it retrieves similar representations in the database to form next-token prediction logits; for sentence-level, it uses similar code for autoregressive generation logits; it also add the original autoregressive generation logits. CMR-Sum [271] uses encoder-decoder model to generate summaries. It conducts cross attention between representations of retrieved summary and generated summary, and add the logits to the original transformer logits for final distribution.

*3) Code Completion:* Code completion can be thought of as the coding equivalent of the "next sentence prediction" task. Early attempts on function completion [272] employs DPR to retrieve relevant template functions using function docstring, then concatenate the information as the input to the code generator BART. ReACC [132] retrieves similar codes to build prompts for generation. For retrieval, it uses hybrid retriever, which combines sparse and dense retrieval; for generation, it uses CodeGPT-adapted [273]. RepoCoder [218] steps further to perform iterative retrieval and generation to bridge the gap between the retrieval context and the intended completion target. In each iteration, for retrieval, the code query is augmented with previously generated code; for generation, the prompt is formed by combining the newest retrieved codes with the code query. Other than combining retrieval results in prompts, the retrieval-and-edit framework [140] first retrieves similar training examples using dense retrieval, then encodes the input and the retrieved result separately, finally combine them through attention for later LSTM decoder. CoCoMic [274] builds a project context graph for the whole code project, and retrieves top-k neighbors of the input source code. It generates representations of both source code and retrieved contexts, then jointly processes the embeddings to complete the code. RepoFusion [275] follows the idea of Fusion-in-Decoder; it employs multiple encoder to input the concatenation of the retrieved repo contexts and the unfinished code, them fuses them and generates the results through a decoder. KNM-LM [276] uses the same model for retrieval encoding and generation, then combines the logits using bayes inference. EDITAS [277] aims at assertion generation. It retrieves similar queries and their assertions, encodes the edit sequence (from the retrieved query to the original query) and the retrieved assertion, then fuses the information for decoder generation. De-Hallucinator [278] first generates code snippets without

retrieval, then retrieves relevant API references given generated contents. In the second pass, retrieved API references are combined in prompt for better generation. REPOFUSE [279] uses both rationale context and retrieved similar codes to construct prompt. To fit in the input length limit, it reserves the contexts that are most relevant to the query.

*4) Automatic Program Repair:* Buggy code can take a lot of effort to fix. Automatic program repair leverages generative models to output the correct version. Query-based RAG technique is widely used in automatic program repair [130], [131], [182], [208], [209]. Among them, RING [209] retrieves similar error messages based on both sparse and dense vectors, then builds prompt for the generator Codex [2]. CEDAR [130] applies for both assertion generation and program repairs tasks; it uses sparse and dense retrieval to search for similar codes, then forms prompt for Codex to generate results. InferFix [131] crafts a prompt carrying the bug type, location, relevant syntax hierarchies, and similar fixes through dense retrieval. Then it also uses Codex for generation. RAP-Gen [182] also retrieves similar buggy codes and corresponding fixes through hybrid retriever (including both sparse and dense retriever). It fine-tunes CodeT5 [106] with this RAG paradigm. SARGAM [208] searches similar buggy code using dense retrieval, generates patches using augmented input, then applies another model to modify the final result. These research works also involve error localization, which is not our focus. RTLFixer [280] leverages ReAct and RAG to implement an agent fixing errors in Verilog codes. It iteratively retrieves relevant errors and corresponding solutions, and combines reasoning and action planning to form prompts for LLMs.

*5) Text-to-SQL and Code-based Semantic Parsing:* Semantic parsing is the task of translating natural language utterances to unambiguous structured meaning representations, where code language is often leveraged to augment this process. SQL is not only a programming language but can also be considered as a structured representation, so we place text-to-SQL (a special case of code generation) in this subsection. Related research works all apply query-based RAG. XRICL [189] focuses on the problem of translating non-English utterances into SQL queries. It searches and reranks English utterance using non-English ones by dense retrieval, then builds prompt for Codex to generate SQL queries. SYNCHROMESH [122] proposes constrained semantic decoding to enforce rich syntactic and semantic constraints during generation of SQL or other languages. It uses the similarity between abstract syntax trees (AST) to finetune the dense retriever S-Bert. During inference, similar NL and SQL are retrived to form the prompt of GPT-3. CodeICL [281] uses Python for semantic parsing, and augments prompts with a structured domain description for GPT generation. In few-shot learning setting, it leverages BM25 sparse retrieval to find similar training examples. RES-DSQL [282] ranks schemas using cross-encoder, then includes ranked schemas into prompts to generate SQL skeleton and SQL query. ReFSQL [283] retrieves relevant questions and corresponding SQL to augment text-to-SQL generation. It involves structure-enhanced retriever with schema linking, and Mahalanobis contrastive learning to improve the retrieval performance. ODIS [284] retrieves in-domain and out-of-domain

demonstrations using BM25, then includes them for in-context learning to generate SQL. Another work [285] retrieves both similar and diverse demonstrations, and then builds prompt for in-context learning to generate SQL. MURRE [286] conducts multi-hop retrieve-rewrite, where relevant tables are retrieved through dense retrieval and then re-writed to generate new tabularized question. A rank module at last integrate the retrieval results and select the top tables for the text-to-SQL module. CodeS [287] retrieves relevant information from table databases in a coarse-to-fine manner, then includes retrieved values to build prompts for both finetuning and inference.

*6) Others:* There are several other code-related tasks that adopt RAG paradigm. In [288] for numerical reasoning task, the Chain-of-Thought is replaced with the programs as the intermediate reasoning step, and dense retrieval-based similar examples are augmented in prompt for downstream LLMs. De-fine [289] tries to resolve intricate tasks using programs. It follows the paradigm in SKCODER, retrieves relevant pairs of query and code, then produces sketch template for real program generation. After generation, it combines the feedback to refine the answer with the same generator. The refined programs, regarded as optimal results, are added back to the codebase for subsequent serving. E&V [290] leverages an LLM-based agent for program static analysis. The agent uses source code retrieval, pseudo-code execution, execution specifications verification, and other tools to form intermediate results. The retrieval is implemented by AST pattern matching. Code4UIE [291] leverages code representation for information extraction tasks. It retrieves relevant examples through dense retrieval, and constructs in-context learning prompt using the retrieved queries and their corresponding codes. StackSpotAI [292] builds an AI coding assistant, which incorporates many code-related tasks. It implements an RAG component, identifying the most relevant pieces of information which serve as the context for GPT generation. InputBlaster [293] aims to generate unusual text input that could cause mobile app crash. It combines retrieved relevant buggy text with valid input to form the prompt for generation.

## C. RAG for Knowledge

Structured knowledge, including KGs (knowledge graph) and tables, is widely used in language-related tasks. It usually serves as the retrieval source to augment generation. In addition to regular sparse and dense retrieval, NER (named-entity recognition) technique and graph-aware neighbor retrieval are applied to identify and extract relevant entities and relations.

*1) Knowledge Base Question Answering:* KBQA (knowledge base question answering) typically utilizes a knowledge base to determine the correct answer to a question. Many semantic parsing methods have been proposed, generating logical forms (e.g. SPARQL) based on the question.

Query-based RAG is the mainstream approach. For a given query, Unseen Entity Handling [53] retrieves topic entities through FreeBase [294], and concatenates the query with the entity for an encoder-decoder to generate SPARQL output. CBR-KBQA [54] retrieves relevant questions and corresponding logical form answers with roberta-based deep retrieval, then concatenates the question and the retrieved pairs for encoder-decoder transformer model. It also revises the final

generation result to align the generated relations with relations present in the local neighborhood of the query entity in the knowledge graph. GMT-KBQA [52] first retrieves relevant entities and relations through bert-based deep retrieval, then concatenates the information for T5 generation. To improve the retrieval result, it leverages cross-encoder to rerank the candidates, and uses the same T5 structure to conduct relation classification and entity dismbiguation. RNG-KBQA [123] enumerates candidate logical forms in the knowledge graph, then ranks the candidates and concatenates them with query to form the prompt input to generate final logical form through a T5 model. Based on this idea, TIARA [124] also retrieve entity and schema besides logical forms, while a following work [295] retrieves top-k questions with BM25. Uni-Parser [133] retrieves relevant entities from knowledge graph using mention detection, cross-encoder ranker, and 2-hop paths extraction. It also considers enumerating tables and columns from databases. On obtaining the relevant information, it concatenates the top-k primitives with the query and generates logical forms through T5. BLLM augmentation [135] uses TIARA as the retrieval for relevant knowledge base elements, then performs in-context learning through black-box LLM such as GPT-4 for generating logical forms. ECBRF [134] follows the case-based reasoning paradigm [296], retrieving similar triplet with dense retriever and constructing prompt input for BART or GPT-2 in-context learning. FC-KBQA [297] extracts relevant class, relation, and entity given a question. For class and relation, it uses BM25 as retriever and a Bert-based cross-encoder as re-ranker. For entity, it follows the mention detection paradigm. To compose all the component candidates, it applies T5 model to generate logical expression. StructGPT [298] extracts relevant information, including triplets and nearest entities, to form prompt for subsequent LLM. KAPING [299] builds prompt including the user query and the retrieved relevant facts (through entity matching), then generates the answer through LLM. Another work [300] also follows the retrieve-then-generate paradigm, replacing the retrieval with a relation distribution generation model for weighted triplets. Retrieve-Rewrite-Answer [301] first retrieves subgraph using hop prediction, relation path prediction, and triplet sampling. It then performs KG-to-text and zero-shot generation with retrieved subgraphs as prompt. Keqing [302] first decomposes a complex question into simple sub-questions through finetuned LLM, then retrieves similar sub-question template by dense retriever RoBERTa to extract candidate entities from knowledge graph, and finally generates answer through ChatGPT given relevant entities as context input. To probe the deep understanding of natural language in LLMs with formal languages, a research work [303] explores the capability of formal language understanding and formal language generation. It leverages retrieved pairs to perform in-context learning. For understanding, it uses tree edit distance to retrieve similar logical forms, while for generation, it uses BM25 to retrieve similar natural language queries. Interactive-KBQA [304] treats LLM as agent and KG as environment. In each step, the LLM conducts entity-linking and one-hop retrieval on KG, then generates current thought and action, until obtaining the final answer.

Latent representation-based RAG is also employed in KBQA. ReTraCk [305] links entities using mention detection, and retrieves schema items using dense retriever Bert. It then generates logical forms by LSTM, incorporating retrieved items through knowledge-specific rules. SKP [145] concatenates triplets with the query and uses fusion-in-decoder technique in inference. It adds a pretraining stage with a knowledge-aware MLM loss on triplets and knowledge constrastive loss with respect to the retrieved items. DECAF [144] forms Resource Description Format knowledge base triplets into sentences, then concatenates sentences with the same head entity into documents. It retrieves relevant documents using BM25 sparse retrieval or Bert dense retrieval, then leverages Fusion-in-Decoder technique to generate logical form and direct answer given each (query, document) pair as input. It combines the output to obtain the final answer. KD-CoT [146] uses the same dense retriever and fusion-in-decoder generator as DECAF. It follows a Chain-of-Thought paradigm, iteratively performing retrieval, generation, and verification until the CoT is finished.

*2) Knowledge-augmented Open-domain Question Answering:* Structured knowledge is often leveraged to augment ODQA (open-domain question answering).

The use of latent representation-based RAG, particularly the fusion-in-decoder technique, is prevalent for knowledge-augmented open-domain question answering. UniK-QA [143] concatenates the text forms of the components in a triplet and build document pool for retrieval. For a given question, it leverages dense retriever for relevant documents, then performs fusion-in-decoder technique to incorporate the information for answer generation. KG-FiD [306] conducts the
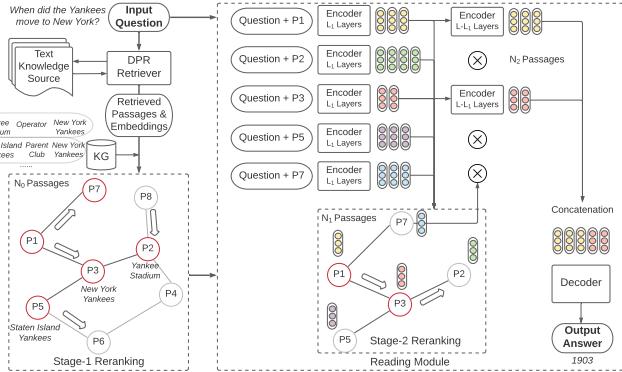


Fig. 7: Architecture of KG-FiD [306] model.

retrieval and generation as normal FiD. It add re-ranking in two ways: the first is to use a graph attention network on the graph formed by retrieved documents; the second is to use the hidden states in the generator. OREOLM [307] empowers LLM with knowledge reasoning paths. Concretely, it uses entity linking to determine the initial state, then conducts contextualized random walk on KG to get reasoning paths, whose entity value memory are combined into the hidden states of LLM for better generation. GRAPE [308] constructs bipartite graph for each pair of question and retrieved passage, then builds bipartite graph on entity for fusing knowledge.

It leverages FiD as backbone model and generate answers. SKURG [309] forms a knowledge graph using text and image data, then updates each source's representation with their relevant sources. It then uses a specially designed decoder to iteratively retrieve and generate. It conducts cross-attention with all the sources, then retrieves the source with highest score and concatenates to the original input embedding; if a gate score does not exceed the threshold, it starts to generate the real answer, otherwise the retrieval stage re-starts.

With the rapid development of LLMs, query-based RAG is emerging as a new standard. DIVKNOWQA [310] develops a retrieval-augmentation tool, including sparse retrieval on structured knowledge, dense retrieval on texts, and sparql generation on KB. Through CoT-based LLM, it retrieves and re-ranks in each step, and generates the final answer. KnowledGPT [311] uses generated code to retrieve from both public and personal knowledge bases, so as to build prompt for LLM question answering. EFSUM [312] generates evidence-focused summary with retrieved relevant facts, then optimizes the summary to align the QA-specific preference with another generator and the filters for helpfulness and faithfulness. GenTKGQA [313] conducts subgraph retrieval through relation ranking and time mining, then employs GNN to incorporate structural and temporal information into virtual token representations for subsequent LLM generation. KnowledgeNavigator [314] analyzes complex questions and performs retrieval on knowledge graph through iterative filtering of relations with respect to core entities, so as to obtain relevant triplets. It then includes the triplets into prompt for generation.

*3) Table for Question Answering:* Tables, as another form of structured knowledge, also facilitates question answering.

Fusion-in-decoder [35] style RAG is often used for table QA. EfficientQA [315], a competition held in NeurIPS 2020, witnessed the proposal of numerous retrieval-reader systems that rely on textual and tabular data. Dual Reader-Parser [316] re-ranks the retrieved textual and tabular data for generation. Convinse [317] retrieves information from heterogeneous resources (including knowledge bases, tables, and texts) after question understanding. CORE [318] retrieves relevant tables and passages through dense representation, then re-ranks the retrieved results using query-generation score from a T0 linker. It uses FiD for final generation. RINK [319] follows the retriever-reranker-reader paradigm for table-augmented question answering. It designs a set-level reader-inherited reranker to get the relevance score of blocks (table segments). TAG-QA [320] retrieves from both tables and texts: for tables, it converts tables to graphs then performs GNN to select relevant contents; for texts, it uses BM25 for retrieval. It then leverages FiD for answer generation.

Tables can be incorporated into prompts for query-based RAG. T-RAG [321] retrieves relevant tables given a user query, then concatenates the query with the tables as a prompt to generate the answer through BART. OmniTab [322] conducts multi-task training to improve the performance of table question answering model. It retrieves relevant tables given a query, then concatenates them for masked-token pre-training. CARP [323] retrieves relevant tables and passages using entity linking, then extracts hybrid chain of retrieved knowl-

edge, which is later used to construct prompt for generation. StructGPT [298] extracts relevant information from knowledge graph, table, or database, to form prompt and generate answers through LLMs. cTBLS [324] retrieves relevant tables through dense retrieval, then for each query, it forms prompt with ranked tables for answer generation. A recent work [325] first uses table-to-text techniques to integrate tabular data into corpora, then conducts experiments on both finetuning and RAG for question answering.

*4) Others:* Prototype-KRG [326] retrieves knowledge facts and dialogue prototypes, then integrates them into the GRU-based generator by both hidden states and logits. SURGE [327] retrieves relevant subgraphs using dense retrieval, then adds them into the input of the generator for dialogue generation. RHO [328] fuses KG embedding of relevant entities and relations into textual embeddings during the generation of open-domain dialogue. K-LaMP [329] retrieves entity in history queries to construct prompt for query suggestion. ReSKGC [147] linearizes all training triplets into text by concatenation, then retrieves relevant triplets using BM25, and generates completed triplet using T5 with fusion-in-decoder. G-Retriever [330] retrieves relevant nodes and edges from graph-based data, then constructs subgraph and performs graph prompt tuning for question answering based on textual graphs.

## D. RAG for Image

*1) Image Generation:* Image Generation refers to the process of creating new images, typically using algorithms in the field of artificial intelligence and machine learning.

The retrieval process can not only help yield high-quality images even for rare or unseen subjects, but also reduces the parameter count and computational expense [45], [137], [148]–[152], [331]. For GAN-based model, RetrieveGAN [45] employs a differentiable retrieval process to select compatible image patches for generation, with Gumbel-softmax trick and end-to-end training. IC-GAN [137] models data as a mix of conditional distributions around each training instance, conditioning both the generator and discriminator on these instances, and can control the semantics and style by swapping class labels or conditional instances. Recently, diffusion models beat GANs on image generation [332]. KNN-Diffusion [149] trains a text-to-image diffusion model without text data, by conditions the model on CLIP joint embedding of the instance and $k$-nearest neighbors from image database. These $k$-NN embeddings bridge the text-image distribution gap and allow image generation from different domains by swapping the database. Similarly, RDM [150] conditions diffusion or autoregressive models on CLIP embeddings of external image databases. It enables post-hoc conditioning on class labels, text prompts and zero-shot stylization [151]. Beyond retrieving only images, Re-imagen [148] conditions on both text prompts and retrieved image-text pairs for text-to-image generation. Interleaved classifier-free guidance is also proposed to balance the alignment between text prompts and retrieval conditions. To avoid information loss of CLIP embeddings and access to all visual condition, Retrieve&Fuse [331] concatenates the retrieved conditional image and the noised image before each attention block in U-Net, and allowing interaction via

self-attention. RPG [79] retrieves representative images to construct informative in-context examples (i.e., image-region pairs), and utilizes multi-modal chain-of-thought reasoning [333] to plan out complementary subregions for compositional text-to-image diffusion.

*2) Image Captioning:* Image Captioning is the process of generating a textual description of an image.

Retrieval-augmented image captioning typically synthesises description with a collection of retrieved captions, instead depending only on the input image. MA [164] augments via a memory bank, built with historical context and target word of image-text training set, and queried during inference with the current context. Vocabulary distribution is computed based on retrieved entries, and interpolated with the original prediction. In adversarial training, RAMP [334] employs retrieved captions as reference for discriminator training, prompting the generator to make full use of retrieved captions. The memory-augmented attention and copying mechanism are also exploited to better use. The RA-Transformer [46] and EXTRA [335], both retrieval-augmented transformer-based captioning models, utilize cross-attention over encoded retrieved captions. EXTRA, as depicted in Fig. 8, jointly process the image and
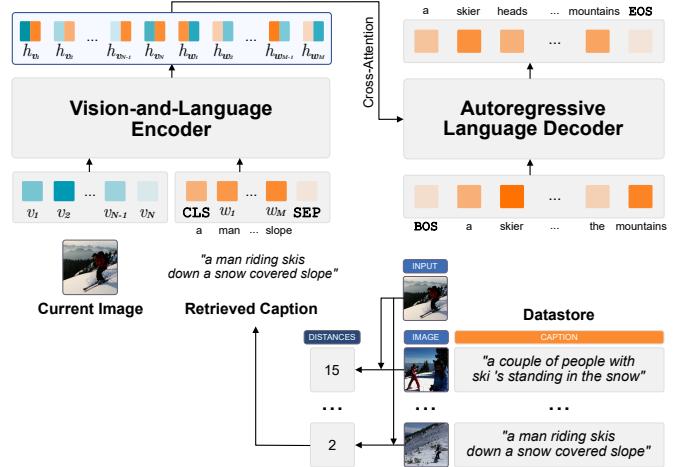


Fig. 8: Architecture of EXTRA [335] model.

retrieved captions with V&L encoder, such that the decoder attends to both visual and linguistic contexts. Beyond retrieved captions, REVEAL [336] uniformly encodes and retrieves multi-modal world knowledge, including image-text pairs, question answering pairs, and knowledge graph triplets, which is then integrated with image features by retrieval score-aware attention module. Straightforwardly, SMALLCAP [47] employs a CLIP vision encoder and a LLM decoder, linked by trainable cross-attention layers, where retrieved captions serve as input-specific in-context examples for prompt for enhancement. For remote sensing image, CRSR [337] enhances retrieved captions with semantic refinement, i.e. filters out misleading details and emphasizes visually salient content. Besides, the visual features is also enriched by transformer network with learnable queries, capturing more intricate details within the images

*3) Others:* There also exist many retrieval augmented works for other image-related tasks. For visual question answering (VQA), PICa [338] leverages GPT-3's implicit knowledge retrieval and reasoning capabilities, which converts images into textual descriptions, then prompts GPT-3 to predict answers based on these descriptions and the question, finally ensembles multi-query results. RA-VQA [339] identifies a limitation that the retrieval in previous work is trained separately from answer generation, and propose a joint training scheme that integrates differentiable retrieval with answer generation, enabling end-to-end training. For visually grounded dialogue, KNN-based Information Fetching (KIF) [340] enhances generative Transformer for dialog modeling. Each KIF module learns a read operation to access fixed external knowledge. Maria [341], a neural conversation agent, enhances dialog generation by leveraging visual experiences retrieved from a large-scale image index as extra context. For multi-modal machine translation, which aims to improve NMT with multi-modal information, [342] incorporates visual information at the phrase level to address the sparsity of paired sentence-image, employing a conditional VAE to filters out redundant visual information from sentence-image datasets.

### E. RAG for Video

*1) Video Captioning:* Video captioning is to describe the visual content with a descriptive utterance. KaVD [343] generates news video caption with background knowledge mined from topically related documents such as named entities and events. R-ConvED [48] introduces retrieval-augmented mechanism to facilitate the word prediction. It uses Dual Encoding [109] for video-text retrieval, and proposes a convolutional encoder-decoder network for generation. For a given input video, R-ConvED first retrieves top-k relevant sentences and their corresponding video from training set, then feeds these pairs and the input video into the generator separately. The obtained decoder hidden states are combined through attention-like read operation, so that the target word can be predicted using the final representation. CARE [160] utilizes visual encoder, audio encoder, and text encoder for frame, audio, and retrieved texts, respectively. It uses CLIP as retriever, and transformer decoder as generator. The embeddings of the three modalities are combined to augment the decoder, producing global semantic guidance which attends the input embedding, and local semantic guidance which attends the attention layer. EgoInstructor [49] generates captions for first-person view videos. It retrieves relevant exocentric videos and corresponding texts via dense retrieval, then encodes the input egocentric video and the retrieved exocentric videos through a CLIP-based visual encoder and a transformer-decoder-based bottleneck module. Then it generates captions through decoder-based LLM which takes the retrieved texts as input and interacts with encoded videos in gated cross-attention layer.

*2) Video QA&Dialogue:* Video QA&Dialogue generates single or multiple-round responses in alignment with video content. For video question answering (VideoQA), MA-DRNN [344] leverages the differentiable neural computer (DNC) with an external memory, for storing and retrieving useful information in questions and videos, and modeling the long-term

visual-textual dependence. Given the video input, R2A [345] retrieves semantically similar texts by multi-modal model, e.g. CLIP, and query LLM with both the question and the retrieved texts. For video-grounded dialogue, [346] proposes TVQA+ dataset which enables to retrieve relevant moments and visual concepts to answer questions about videos, and also proposes Spatio-Temporal Answerer with Grounded Evidence (STAGE) to exploit it. VGNMN [347] also extracts visual cues from videos, while the retrieval is carried out using neural module networks parameterized by entities and actions in previous dialogues.

*3) Others:* There also exist many retrieval augmented works for other video-related tasks. VidIL [348] exploits image-language models to translate video content into temporal-aware prompts with few-shot examples, for various video-language tasks including video captioning, video question answering, video caption retrieval, and video future event prediction. Notably, for trustworthy autonomous driving, RAG-Driver [349] grounds the MLLM in relevant expert demonstrations from a memory database, to produce driving action explanations, justifications, and control signal prediction. Text-to-video generation is to generate video given natural language descriptions. As shown in Fig. 9, Animate-
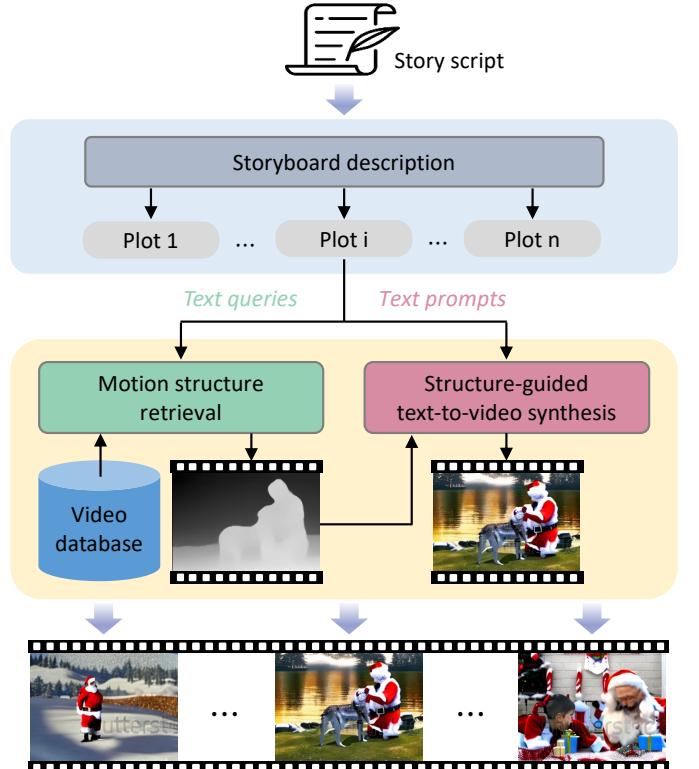


Fig. 9: Architecture of Animate-A-Story [206] model.

A-Story [206] develops a framework which can generate high-quality storytelling videos based on texts. It first separates the text into individual plots, and decorates the description using LLM. It then retrieves relevant videos for each plot through a dense retriever [110]. It generates videos through a latent diffusion model, consisting of two branches: a text encoder

CLIP, and a structure encoder which takes the estimated depth of the retrieved videos as structure control.

### F. RAG for Audio

*1) Audio Generation:* The goal of audio generation is to generate audio with natural language input.
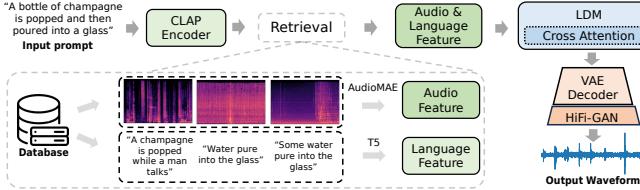


Fig. 10: Architecture of Re-AudioLDM [159] model.

Re-AudioLDM [159] adopts dense retriever CLAP [26] to retrieve similar caption-audio pairs given input prompt. As shown in Fig. 10, the generator, latent diffusion model and VAE-based decoder, take the representations of input text and retrieved pairs as input and generate output audio. Make-An-Audio [44] uses dense retriever CLAP [26] to augment data, retrieving related audio given natural language text. It then constructs pseudo prompts for diffusion-based text-to-audio model training.

*2) Audio Captioning:* The goal of audio captioning is to generate natural language data with audio data, which is basically a sequence-to-sequence task. RECAP [350] leverages CLAP [26] to retrieve related captions given audio data. The retrieved captions are then included into the prompt input for GPT-2 model, which interacts with audio embeddings through cross attention. In [43], dense retriver VGGish [107] is adopted to produce dense embedding of audio data, and GPT-2 is adopted to generate representations of the retrieved captions which are paired with similar audios. After obtaining the representations of audios and captions, an extra multi-head attention block and a linear layer fuses all the information and generates the output. Some research studies transform audio modality to text, in order to leverage advancements in LLMs [351]–[353]. They take advantage of deep retrieval models, aligning the modalities into the same latent space for downstream text generation.

### G. RAG for 3D

*1) Text-to-3D:* Retrieval can be applied to augment the generation of 3D contents. ReMoDiffuse [51] aims at generating motions using diffusion models. It first retrieves relevant motion entites through CLIP given text input, then leverages the information of the text and the retrieved entities through a semantic-modulated attention layer.

AMD [158] designs two branches of motion diffusion for fidelity and diversity. As shown in Fig. 11, the first branch inputs the original prompt text for diffusion; the second branch decomposes the input text into anatomical scripts and retrieve similar reference motions for diffusion. A transformer-based fusion module is further applied to adaptively balance the result from two branches. RetDream [50] targets general 3D generation, using retrieved 3D assets to augment the process
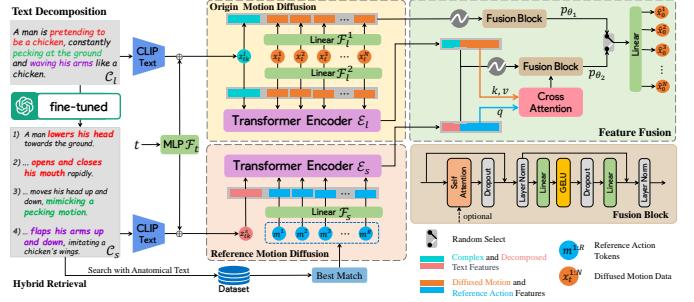


Fig. 11: Architecture of AMD [158] model.

of variational score distillation from 2D diffusion models. Given an input query, it retrieves relevant 3D assets through CLIP, then utilizes the retrieved assets to provide geometric prior and adapted 2D prior. Concretely, retrieved assets not only impose an additional velocity on particles for distribution initialization, but also help optimize the 2D diffusion model through low-rank adaptation.

### H. RAG for Science

RAG has also emerged as a promising research direction for many interdisciplinary applications, such as molecular generation, medical tasks and computational research.

*1) Drug Discovery:* The goal of drug discovery is to generate molecules that concurrently fulfill diverse properties. RetMol [55] integrates a lightweight retrieval mechanism and
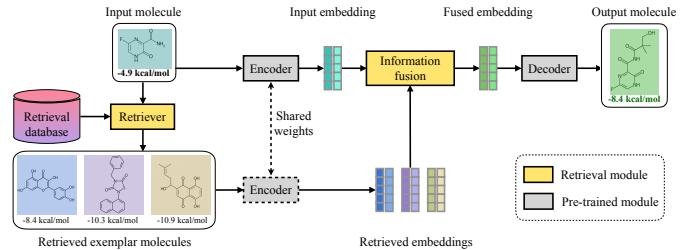


Fig. 12: Architecture of RetMol [55] model.

molecular strings into a pre-trained encoder-decoder generative model to retrieve and fuse exemplar molecules with the input. PromptDiff [354] introduces an interaction-based, retrieval-augmented 3D molecular diffusion model that retrieves a curated set of ligand references to guide the synthesis of ligands meeting specific design criteria.

*2) Biomedical Informatics Enhancement:* Several recent studies have improved the expressiveness of LLM by retrieving information from biomedical domain-specific databases, thereby augmenting the model's capabilities to provide valuable guidance for tasks in the medical field. PoET [355] is an autoregressive generative model based on a variant of Transformer that integrates a retrieval mechanism to enable prompt augmentation, thereby expediting the prediction of fitness properties for protein variants. Chat-Orthopedist [136] enhances ChatGPT with a retrieval-augmented mechanism focused on adolescent idiopathic scoliosis (AIS), utilizing an external knowledge base for precise responses. BIOREADER [356] is the first retrieval-enhanced text-to-text transformer-based model for biomedical natural

language processing, incorporating the retrieved literature evidence into the model using a chunked-cross attention mechanism. MedWriter [357] employs a hierarchical retrieval-augmented generation method that combines report-level and sentence-level templates to produce coherent and clinically accurate medical reports from images. QA-RAG [358] employs a dual-track RAG strategy to enhance pharmaceutical compliance by effectively retrieving and integrating regulatory guidelines based on language model responses and user queries.

*3) Math Applications:* Retrieval-augmented generation technology in mathematics streamlines problem-solving, boosts research innovation, and refines educational strategies. LeanDojo [359] boosts theorem proving by using retrieval-augmented methods to choose relevant premises from extensive mathematical libraries, improving automation and theorem generalization. RAG-for-math-QA [360] improves math question-answering by integrating a high-quality math textbook with retrieval-augmented generation, enhancing LLM-generated responses for middle-school algebra and geometry.

## V. BENCHMARK

Given the increasing research interests and applications of RAG, there have also been several benchmarks assessing RAG from certain aspects.

Chen et al. [361] proposed an RAG benchmark, which evaluates RAG from four aspects: Noise Robustness, Negative Rejection, Information Integration, and Counterfactual Robustness, respectively. Noise Robustness evaluates whether LLMs could extract the necessary information from documents containing noisy information. The noisy information is relevant to the input query but useless for answering it. Negative Rejection measures whether LLMs would reject to respond the query when the retrieved content is not enough. Information Integration assesses whether LLMs could acquire knowledge and make responses by integrating multiple retrieved contents. Counterfactual Robustness refers to the ability of LLMs to identify counterfactual errors in the retrieved content.

Another three benchmarks, RAGAS [362], ARES [363] and TruLens [364], consider three different aspects: Faithfulness, Answer Relevance, and Context Relevance, respectively. Faithfulness focuses on the factual errors in the results when the correct answers can be inferred from the retrieved contents. Answer Relevance measures whether the generated results actually address the problems (i.e., queries) or not. Context Relevance judges whether the retrieved contents contain as much knowledge as possible to answer the queries, and as little irrelevant information as possible.

CRUD-RAG [365] divides all RAG tasks into four categories, which are Create, Read, Update, and Delete, respectively, and also evaluates each category using text continuation, question answering (with single- and multi-document questions), hallucination modification, and open-domain multi-document summary. MIRAGE [366] is a benchmark designed for assessing the application of RAG in the medical domain, focusing on the comparison and optimization of medical question-answering systems' performance. KILT [367] is another benchmark focuses on ensuring information accuracy

and reliability by aligning Wikipedia pages with specific snapshots and pinpointing the most pertinent text ranges through BLEU score evaluations. It filters out lower-quality data to maintain a high standard of information mapping, offering a variety of retrieval system options like TF-IDF, DPR, RAG, and BLINK + flair to support evidence-based predictions or citations according to task requirements.

## VI. DISCUSSION

### A. Limitations

Despite the widespread adoption of RAG, it suffers from several limitations by nature. In this paper, we provide a summary of the limitations and engage in an in-depth discussion.

*1) Noises in Retrieval Results:* Information retrieval cannot yield perfect results because information loss appears in representations generated by encoder models. Additionally, ANN search can also provide approximate results rather than exact ones. Consequently, certain degree of noise is inevitable in retrieval results, manifesting as irrelevant objects or misleading information, which may cause failure points in RAG systems [368]. Though the common sense is that increasing the accuracy of retrieval will contribute to the effectiveness of RAG, a recent study surprisingly shows that noisy retrieval results may conversely help improve the generation quality [369]. A possible explanation is that diversity in retrieval results may also be necessary for prompt construction [370]. As a result, the impact of noise in retrieval results remains uncertain, leading to confusion in practical uses regarding which metric to employ for retrieval and how to facilitate the interaction between the retriever and the generator.

*2) Extra Overhead:* While retrieval can help mitigate the costs of generation in certain cases [30]–[32], the incorporation of retrieval sometimes introduces non-negligible overhead. Considering that RAG is primarily employed to improve the performance of existing generative models, the inclusion of additional retrieval and interaction processes leads to increased latency. Worse still, when combined with complex enhancement methods, such as recursive retrieval [371] and iterative RAG [218], the extra overhead will become even more significant. Furthermore, as the scale of retrieval expands, the storage and access complexity associated with data sources will also increase. In presence, RAG systems exhibit a trade-off between costs and benefits. Looking ahead, we anticipate further optimization to alleviate the associated overhead [372].

*3) The Gap between Retrievers and Generators:* Seamlessly integrating retrieval and generation components requires meticulous design and optimization. Since the objectives of retrievers and generators may not align, and their latent spaces might differ, designing their interaction poses challenges. As introduced in Section III, numerous approaches have been proposed to enable effective RAG, and these approaches either disentangle the retrieval and generation processes or integrate them at an intermediate stage. While the former is more modularized, the latter could potentially benefit from joint training. Till not, there lacks a sufficient comparison of different ways of interaction across various scenarios.

*4) Increased System Complexity:* The introduction of retrieval unavoidably increases the system complexity and the number of hyper-parameters to tune. For instance, a recent study on the trade-off between attribution and fluency in prompt-augmentation-style RAG demonstrates that using top-k retrieval for generation improves attribution, but hurts fluency in turns [373]. The counter effects of different aspects in RAG, such as metric selection, are still under explored. Therefore, further refinement of RAG systems, both in terms of algorithms, and deployment, is necessary to fully unlock their potentials.

*5) Lengthy Context:* One of the primary shortcomings of RAG, in particular the query-based RAG, is that it lengthens the context tremendously, making it infeasible for generators with limited context length. In addition, the lengthened context also slows down the generation process generally. The research advancements in prompt compression [202] and long-context support [374] have partially mitigated these challenges, albeit with a slight trade-off in accuracy or costs.

### B. Potential Future Directions

Lastly, we wish to outline several potential directions for future RAG research and applications.

*1) More Advanced Research on RAG Methodologies, Enhancements, and Applications:* A straight-forward research direction is to develop more advanced methodologies, enhancements, and applications of RAG.

As introduced in Section III-A, existing works have explored various interaction patterns between retrievers and generators. However, since the optimization target of these two components are distinct, the practical augmentation process has a large impact on the final generation results. Investigation of more advanced foundations for augmentation holds promise for fully unleashing the potential of RAG.

Based on a constructed RAG system, enhancements are helpful to improve the effectiveness of certain components or the entire pipeline. Given the inherent complexity of the system, there exists significant potential for RAG to improve, necessitating proper tuning and careful engineering. We look forward to further experimental analysis and in-depth exploration that will contribute to the development of more effective and more robust RAG systems.

As introduced in Section IV, RAG is a general technique that has been applied across diverse modalities and tasks. Yet most of existing works straightforwardly integrate external knowledge with the specific generation tasks, without thoroughly taking into account the key characteristics of the target domains. Therefore, for generation tasks that do not fully leverage the power of RAG, we are confident that designing proper RAG system will be beneficial.

*2) Efficient Deployment and Processing:* Currently, several deployment solutions of query-based RAG for LLMs have been proposed, such as LangChain [375], LLAMA-Index [175], and PipeRAG [376]. However, for other foundations of RAG and/or generation tasks, there lacks a plug-and-play solution. In addition, given the extra overhead introduced by retrieval, and considering that the complexities of both the retriever and generator will continue to grow, achieving efficient processing in RAG remains a challenge, necessitating targeted system optimization.

*3) Incorporating Long-tail and Real-time Knowledge:* While a key motivation of RAG is to harness real-time and long-tail knowledge, few studies have explored the pipeline for knowledge updating and expansion. Many existing works make up the retrieval sources with merely the training data of generators, thereby neglecting the dynamic and flexible information advantages that could have been offered by retrieval. As a consequence, designing a useful RAG system with continuously updated knowledge and/or flexible knowledge sources, along with corresponding system-level optimizations, is a growing research direction. With the capability of utilizing long-tail knowledge, we also expect RAG to leverage personalized information and features, so as to adapt to today's web service.

*4) Combined with Other Techniques:* In essential, RAG is orthogonal to other techniques that share the goal of improving AIGC effectiveness, including fine tuning, reinforcement learning, chain-of-thought, agent-based generation, and other potential optimizations. However, the exploration of simultaneously applying these techniques is still in its early stages, calling for further research to delve into algorithm design and fully leverage their potential. It is worthy to note that a recent notion appears "long-context models like Gemini 1.5 will replace RAG". Nevertheless, this assertion does not hold true — RAG exhibits greater flexibility in managing dynamic information, encompassing both up-to-date and long-tail knowledge [377]. We believe that RAG in the future will take advantage of long context generation to achieve even better performance, rather than simply being weeded out by it.

## VII. CONCLUSION

In this paper, we conducted a thorough and comprehensive survey on RAG within the context of AIGC, with a particular focus on augmentation foundations, enhancements, and applications. We first systematically organized and summarize the foundation paradigms in RAG, providing insights into the interaction between retrievers and generators. Then, we reviewed the enhancements that further improve the effectiveness of RAG, including the enhancements on each component or the entire pipeline. To facilitate researchers across diverse domains, we showcased practical applications of RAG in a range of modalities and tasks. Finally, we also presented existing benchmarks for RAG, discussed current limitations of RAG, and shed light on promising future directions.

## REFERENCES

[1] T. B. Brown, B. Mann et al., "Language models are few-shot learners," in NeurIPS, 2020.
[2] M. Chen, J. Tworek et al., "Evaluating large language models trained on code," arXiv:2107.03374, 2021.
[3] OpenAI, "GPT-4 technical report," arXiv:2303.08774, 2023.
[4] H. Touvron, T. Lavril et al., "Llama: Open and efficient foundation language models," arXiv:2302.13971, 2023.
[5] H. Touvron, L. Martin et al., "Llama 2: Open foundation and fine-tuned chat models," arXiv:2307.09288, 2023.
[6] B. Rozière, J. Gehring et al., "Code llama: Open foundation models for code," arXiv:2308.12950, 2023.
[7] A. Ramesh, M. Pavlov, G. Goh et al., "Zero-shot text-to-image generation," in ICML, 2021.

[8] A. Ramesh, P. Dhariwal, A. Nichol et al., "Hierarchical text-conditional image generation with CLIP latents," arXiv:2204.06125, 2022.

[9] J. Betker, G. Goh, L. Jing et al., "Improving image generation with better captions," Computer Science, vol. 2, no. 3, p. 8, 2023.

[10] R. Rombach, A. Blattmann, D. Lorenz et al., "High-resolution image synthesis with latent diffusion models," in IEEE/CVF, 2022.

[11] OpenAI, "Video generation models as world simulators," https://openai.com/research/video-generation-models-as-world-simulators, 2024.

[12] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural Comput., vol. 9, no. 8, pp. 1735–1780, 1997.

[13] A. Vaswani, N. Shazeer, N. Parmar et al., "Attention is all you need," in NeurIPS, 2017.

[14] I. Goodfellow, J. Pouget-Abadie, M. Mirza et al., "Generative adversarial networks," CACM, vol. 63, no. 11, pp. 139–144, 2020.

[15] J. Devlin, M. Chang et al., "BERT: pre-training of deep bidirectional transformers for language understanding," in NAACL-HLT, 2019.

[16] C. Raffel, N. Shazeer, A. Roberts et al., "Exploring the limits of transfer learning with a unified text-to-text transformer," JMLR, vol. 21, pp. 140:1–140:67, 2020.

[17] W. Fedus, B. Zoph, and N. Shazeer, "Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity," JMLR, vol. 23, no. 120, pp. 1–39, 2022.

[18] J. Kaplan, S. McCandlish, T. Henighan et al., "Scaling laws for neural language models," 2020.

[19] S. E. Robertson and H. Zaragoza, "The probabilistic relevance framework: BM25 and beyond," FTIR, vol. 3, no. 4, pp. 333–389, 2009.

[20] V. Karpukhin, B. Oguz, S. Min et al., "Dense passage retrieval for open-domain question answering," in EMNLP, 2020.

[21] J. Johnson, M. Douze, and H. Jégou, "Billion-scale similarity search with gpus," IEEE Trans. Big Data, vol. 7, no. 3, pp. 535–547, 2021.

[22] Q. Chen, B. Zhao, H. Wang et al., "SPANN: highly-efficient billion-scale approximate nearest neighborhood search," in NeurIPS, 2021.

[23] R. Datta, D. Joshi, J. Li et al., "Image retrieval: Ideas, influences, and trends of the new age," CSUR, vol. 40, no. 2, pp. 5:1–5:60, 2008.

[24] A. Radford, J. W. Kim, C. Hallacy et al., "Learning transferable visual models from natural language supervision," in ICML, 2021.

[25] Z. Feng, D. Guo et al., "Codebert: A pre-trained model for programming and natural languages," in EMNLP Findings, 2020.

[26] Y. Wu, K. Chen, T. Zhang et al., "Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation," in ICASSP, 2023.

[27] A. Mallen, A. Asai, V. Zhong et al., "When not to trust language models: Investigating effectiveness of parametric and non-parametric memories," in ACL, 2023.

[28] N. Carlini, F. Tramèr et al., "Extracting training data from large language models," in USENIX, 2021.

[29] M. Kang, N. M. Gürel et al., "C-RAG: certified generation risks for retrieval-augmented language models," arXiv:2402.03181, 2024.

[30] G. Izacard, P. Lewis, M. Lomeli et al., "Atlas: Few-shot learning with retrieval augmented language models," arXiv:2208.03299, 2022.

[31] Y. Wu, M. N. Rabe, D. Hutchins, and C. Szegedy, "Memorizing transformers," in ICLR, 2022.

[32] Z. He, Z. Zhong, T. Cai et al., "REST: retrieval-based speculative decoding," arxiv:2311.08252, 2023.

[33] K. Guu, K. Lee, Z. Tung et al., "REALM: retrieval-augmented language model pre-training," ICML, 2020.

[34] P. S. H. Lewis, E. Perez, A. Piktus et al., "Retrieval-augmented generation for knowledge-intensive NLP tasks," in NeurIPS, 2020.

[35] G. Izacard and E. Grave, "Leveraging passage retrieval with generative models for open domain question answering," in EACL, 2021.

[36] S. Borgeaud, A. Mensch et al., "Improving language models by retrieving from trillions of tokens," in ICML, 2022.

[37] U. Khandelwal, O. Levy, D. Jurafsky et al., "Generalization through memorization: Nearest neighbor language models," in ICLR, 2020.

[38] J. He, G. Neubig, and T. Berg-Kirkpatrick, "Efficient nearest neighbor language models," in EMNLP, 2021.

[39] zilliztech. (2023) Gptcache. [Online]. Available: https://github.com/zilliztech/GPTCache

[40] M. R. Parvez, W. U. Ahmad et al., "Retrieval augmented code generation and summarization," in EMNLP Findings, 2021.

[41] W. U. Ahmad, S. Chakraborty, B. Ray et al., "Unified pre-training for program understanding and generation," in NAACL-HLT, 2021.

[42] S. Zhou, U. Alon, F. F. Xu et al., "Docprompting: Generating code by retrieving the docs," in ICLR, 2023.

[43] Y. Koizumi, Y. Ohishi et al., "Audio captioning using pre-trained large-scale language model guided by audio-based similar caption retrieval," arXiv:2012.07331, 2020.

[44] R. Huang, J. Huang, D. Yang et al., "Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models," in ICML, 2023.

[45] H.-Y. Tseng, H.-Y. Lee et al., "Retrievegan: Image synthesis via differentiable patch retrieval," in ECCV, 2020.

[46] S. Sarto, M. Cornia, L. Baraldi, and R. Cucchiara, "Retrieval-augmented transformer for image captioning," in CBMI, 2022.

[47] R. Ramos, B. Martins et al., "Smallcap: lightweight image captioning prompted with retrieval augmentation," in CVPR, 2023.

[48] J. Chen, Y. Pan, Y. Li et al., "Retrieval augmented convolutional encoder-decoder networks for video captioning," TOMCCAP, vol. 19, no. 1s, pp. 48:1–48:24, 2023.

[49] J. Xu, Y. Huang, J. Hou et al., "Retrieval-augmented egocentric video captioning," arXiv:2401.00789, 2024.

[50] J. Seo, S. Hong et al., "Retrieval-augmented score distillation for text-to-3d generation," arXiv:2402.02972, 2024.

[51] M. Zhang, X. Guo, L. Pan et al., "Remodiffuse: Retrieval-augmented motion diffusion model," in ICCV, 2023.

[52] X. Hu, X. Wu, Y. Shu, and Y. Qu, "Logical form generation via multi-task learning for complex question answering over knowledge bases," in COLING, 2022.

[53] X. Huang, J. Kim, and B. Zou, "Unseen entity handling in complex question answering over knowledge base via language generation," in EMNLP Findings, 2021.

[54] R. Das, M. Zaheer, D. Thai et al., "Case-based reasoning for natural language queries over knowledge bases," in EMNLP, 2021.

[55] Z. Wang, W. Nie, Z. Qiao et al., "Retrieval-based controllable molecule generation," in ICLR, 2022.

[56] Q. Jin, Y. Yang, Q. Chen, and Z. Lu, "Genegpt: Augmenting large language models with domain tools for improved access to biomedical information," Bioinformatics, vol. 40, no. 2, p. btae075, 2024.

[57] H. Li, Y. Su, D. Cai et al., "A survey on retrieval-augmented text generation," arxiv:2202.01110, 2022.

[58] A. Asai, S. Min, Z. Zhong, and D. Chen, "Acl 2023 tutorial: Retrieval-based language models and applications," ACL 2023, 2023.

[59] Y. Gao, Y. Xiong et al., "Retrieval-augmented generation for large language models: A survey," arxiv:2312.10997, 2023.

[60] R. Zhao, H. Chen et al., "Retrieving multimodal information for augmented generation: A survey," in EMNLP, 2023.

[61] J. Chen, H. Guo, K. Yi et al., "Visualgpt: Data-efficient adaptation of pretrained language models for image captioning," in CVPR, 2022.

[62] Y. Tay, M. Dehghani, D. Bahri, and D. Metzler, "Efficient transformers: A survey," CSUR, vol. 55, no. 6, pp. 109:1–109:28, 2023.

[63] G. V. Houdt et al., "A review on the long short-term memory model," Artif. Intell. Rev., vol. 53, no. 8, pp. 5929–5955, 2020.

[64] L. Yang, Z. Zhang et al., "Diffusion models: A comprehensive survey of methods and applications," CSUR, vol. 56, no. 4, pp. 1–39, 2023.

[65] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in ICML, 2015.

[66] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," in NeurIPS, 2020.

[67] A. Q. Nichol and P. Dhariwal, "Improved denoising diffusion probabilistic models," in ICML, 2021.

[68] Y. Song and S. Ermon, "Generative modeling by estimating gradients of the data distribution," in NeurIPS, 2019.

[69] ——, "Improved techniques for training score-based generative models," in NeurIPS, 2020.

[70] Y. Song, J. Sohl-Dickstein, D. P. Kingma et al., "Score-based generative modeling through stochastic differential equations," in ICLR, 2021.

[71] Y. Song, C. Durkan, I. Murray, and S. Ermon, "Maximum likelihood training of score-based diffusion models," in NeurIPS, 2021.

[72] L. Yang, H. Qian, Z. Zhang et al., "Structure-guided adversarial training of diffusion models," in CVPR, 2024.

[73] X. Zhang, L. Yang, Y. Cai et al., "Realcompo: Dynamic equilibrium between realism and compositionality improves text-to-image diffusion models," arXiv:2402.12908, 2024.

[74] R. Rombach, A. Blattmann, D. Lorenz et al., "High-resolution image synthesis with latent diffusion models," in CVPR, 2022.

[75] A. Ramesh, P. Dhariwal, A. Nichol et al., "Hierarchical text-conditional image generation with clip latents," arXiv:2204.06125, 2022.

[76] H. Li, Y. Yang, M. Chang et al., "Srdiff: Single image super-resolution with diffusion probabilistic models," Neurocomputing, vol. 479, pp. 47–59, 2022.

[77] J. Ho, C. Saharia, W. Chan et al., "Cascaded diffusion models for high fidelity image generation," JMLR, vol. 23, no. 1, pp. 2249–2281, 2022.

[78] L. Yang, J. Liu, S. Hong et al., "Improving diffusion-based image synthesis with context prediction," in NeurIPS, 2024.

[79] L. Yang, Z. Yu, C. Meng et al., "Mastering text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms," arXiv:2401.11708, 2024.

[80] S. Gong, M. Li, J. Feng et al., "Diffuseq: Sequence to sequence text generation with diffusion models," in ICLR, 2023.

[81] X. Li, J. Thickstun, I. Gulrajani et al., "Diffusion-lm improves controllable text generation," in NeurIPS, 2022.

[82] J. Austin, D. D. Johnson, J. Ho et al., "Structured denoising diffusion models in discrete state-spaces," in NeurIPS, 2021.

[83] T. Chen, R. Zhang, and G. Hinton, "Analog bits: Generating discrete data using diffusion models with self-conditioning," in ICLR, 2023.

[84] J. Ho, W. Chan, C. Saharia et al., "Imagen video: High definition video generation with diffusion models," arXiv:2210.02303, 2022.

[85] W. Harvey, S. Naderiparizi, V. Masrani et al., "Flexible diffusion modeling of long videos," in NeurIPS, 2022.

[86] R. Yang, P. Srivastava, and S. Mandt, "Diffusion probabilistic modeling for video generation," Entropy, vol. 25, no. 10, p. 1469, 2023.

[87] M. Zhang, Z. Cai, L. Pan et al., "Motiondiffuse: Text-driven human motion generation with diffusion model," TPAMI, 2024.

[88] L. Yang, Z. Zhang, Z. Yu et al., "Cross-modal contextualized diffusion models for text-guided visual generation and editing," in ICLR, 2024.

[89] N. Anand and T. Achim, "Protein structure and sequence generation with equivariant denoising diffusion probabilistic models," arXiv:2205.15019, 2022.

[90] M. Xu, L. Yu, Y. Song et al., "Geodiff: A geometric diffusion model for molecular conformation generation," in ICLR, 2022.

[91] E. Hoogeboom, V. G. Satorras, C. Vignac, and M. Welling, "Equivariant diffusion for molecule generation in 3d," in ICML, 2022.

[92] B. Jing, G. Corso, J. Chang et al., "Torsional diffusion for molecular conformer generation," in NeurIPS, 2022.

[93] Z. Huang, L. Yang, X. Zhou et al., "Protein-ligand interaction prior for binding-aware 3d molecule diffusion models," in ICLR, 2024.

[94] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," in ICLR, 2021.

[95] X. Liu, C. Gong, and Q. Liu, "Flow straight and fast: Learning to generate and transfer data with rectified flow," 2023.

[96] Y. Song, P. Dhariwal, M. Chen, and I. Sutskever, "Consistency models," in ICML, 2023.

[97] J. Gui, Z. Sun, Y. Wen et al., "A review on generative adversarial networks: Algorithms, theory, and applications," TKDE, vol. 35, no. 4, pp. 3313–3332, 2023.

[98] S. E. Robertson and S. Walker, "On relevance weights with little relevance information," in SIGIR, 1997.

[99] J. D. Lafferty and C. Zhai, "Document language models, query models, and risk minimization for information retrieval," in SIGIR, 2001.

[100] Y. Liu, M. Ott et al., "Roberta: A robustly optimized BERT pretraining approach," arxiv:1907.11692, 2019.

[101] L. Xiong, C. Xiong, Y. Li et al., "Approximate nearest neighbor negative contrastive learning for dense text retrieval," in ICLR, 2021.

[102] H. Zhang, Y. Gong, Y. Shen et al., "Adversarial retriever-ranker for dense text retrieval," in ICLR, 2022.

[103] Y. Qu, Y. Ding, J. Liu et al., "Rocketqa: An optimized training approach to dense passage retrieval for open-domain question answering," in NAACL-HLT, 2021.

[104] L. Gao and J. Callan, "Condenser: a pre-training architecture for dense retrieval," in EMNLP, 2021.

[105] D. Guo, S. Ren et al., "Graphcodebert: Pre-training code representations with data flow," in ICLR, 2021.

[106] Y. Wang, W. Wang, S. R. Joty, and S. C. H. Hoi, "Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation," in EMNLP, 2021.

[107] S. Hershey, S. Chaudhuri et al., "CNN architectures for large-scale audio classification," in ICASSP, 2017.

[108] X. Yuan, Z. Lin, J. Kuen et al., "Multimodal contrastive training for visual representation learning," in CVPR, 2021.

[109] J. Dong, X. Li, C. Xu et al., "Dual encoding for zero-example video retrieval," in CVPR, 2019.

[110] M. Bain, A. Nagrani, G. Varol, and A. Zisserman, "Frozen in time: A joint video and image encoder for end-to-end retrieval," in ICCV, 2021.

[111] J. Zhan, J. Mao, Y. Liu et al., "Optimizing dense retrieval model training with hard negatives," in SIGIR, 2021.

[112] J. L. Bentley, "Multidimensional binary search trees used for associative searching," CACM, vol. 18, no. 9, pp. 509–517, 1975.

[113] W. Li, C. Feng, D. Lian et al., "Learning balanced tree indexes for large-scale vector retrieval," in SIGKDDg, 2023.

[114] M. Datar, N. Immorlica, P. Indyk et al., "Locality-sensitive hashing scheme based on p-stable distributions," in SCG, 2004.

[115] Y. A. Malkov and D. A. Yashunin, "Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs," TPAMI, vol. 42, no. 4, pp. 824–836, 2018.

[116] S. Jayaram Subramanya, F. Devvrit et al., "Diskann: Fast accurate billion-point nearest neighbor search on a single node," NeurIPS, 2019.

[117] J. Ren, M. Zhang, and D. Li, "Hm-ann: Efficient billion-point nearest neighbor search on heterogeneous memory," NeurIPS, 2020.

[118] Y. Wang, Y. Hou, H. Wang et al., "A neural corpus indexer for document retrieval," in NeurIPS, 2022.

[119] H. Zhang, Y. Wang, Q. Chen et al., "Model-enhanced vector index," in NeurIPS, 2023.

[120] S. A. Hayati, R. Olivier, P. Avvaru et al., "Retrieval-based neural code generation," in EMNLP, 2018.

[121] J. Zhang, X. Wang, H. Zhang et al., "Retrieval-based neural source code summarization," in ICSE, 2020.

[122] G. Poesia, A. Polozov, V. Le et al., "Synchromesh: Reliable code generation from pre-trained language models," in ICLR, 2022.

[123] X. Ye, S. Yavuz et al., "RNG-KBQA: generation augmented iterative ranking for knowledge base question answering," in ACL, 2022.

[124] Y. Shu et al., "TIARA: multi-grained retrieval for robust question answering over large knowledge bases," arXiv:2210.12925, 2022.

[125] X. V. Lin, R. Socher et al., "Bridging textual and tabular data for cross-domain text-to-sql semantic parsing," arXiv:2012.12627, 2020.

[126] A. Asai, Z. Wu, Y. Wang et al., "Self-rag: Learning to retrieve, generate, and critique through self-reflection," arxiv:2310.11511, 2023.

[127] W. Shi, S. Min, M. Yasunaga et al., "Replug: Retrieval-augmented black-box language models," arXiv:2301.12652, 2023.

[128] O. Ram, Y. Levine, I. Dalmedigos et al., "In-context retrieval-augmented language models," arXiv:2302.00083, 2023.

[129] D. Zan, B. Chen, Z. Lin et al., "When language model meets private library," in EMNLP Findings, 2022.

[130] N. Nashid, M. Sintaha, and A. Mesbah, "Retrieval-based prompt selection for code-related few-shot learning," in ICSE, 2023.

[131] M. Jin, S. Shahriar, M. Tufano et al., "Inferfix: End-to-end program repair with llms," in ESEC/FSE, 2023.

[132] S. Lu, N. Duan, H. Han et al., "Reacc: A retrieval-augmented code completion framework," in ACL, 2022.

[133] Y. Liu et al., "Uni-parser: Unified semantic parser for question answering on knowledge base and database," in EMNLP, 2022.

[134] Z. Yang, X. Du, E. Cambria et al., "End-to-end case-based reasoning for commonsense knowledge base completion," in EACL, 2023.

[135] M. Patidar, A. K. Singh, R. Sawhney et al., "Combining transfer learning with in-context learning using blackbox llms for zero-shot knowledge base question answering," arXiv:2311.08894, 2023.

[136] W. Shi, Y. Zhuang, Y. Zhu et al., "Retrieval-augmented large language models for adolescent idiopathic scoliosis patients in shared decision-making," in ACM-BCB, 2023.

[137] A. Casanova, M. Careil, J. Verbeek et al., "Instance-conditioned gan," in NeurIPS, 2021.

[138] J. Li, Y. Li, G. Li et al., "Editsum: A retrieve-and-edit framework for source code summarization," in ASE, 2021.

[139] C. Yu, G. Yang, X. Chen et al., "Bashexplainer: Retrieval-augmented bash code comment generation based on fine-tuned codebert," in ICSME, 2022.

[140] T. B. Hashimoto, K. Guu, Y. Oren, and P. Liang, "A retrieve-and-edit framework for predicting structured outputs," in NeurIPS, 2018.

[141] B. Wei, Y. Li, G. Li et al., "Retrieve and refine: Exemplar-based neural comment generation," in ASE, 2020.

[142] E. Shi, Y. Wang, W. Tao et al., "RACE: retrieval-augmented commit message generation," in EMNLP, 2022.

[143] B. Oguz, X. Chen, V. Karpukhin et al., "Unik-qa: Unified representations of structured and unstructured knowledge for open-domain question answering," in NAACL Findings, 2022.

[144] D. Yu, S. Zhang et al., "Decaf: Joint decoding of answers and logical forms for question answering over knowledge bases," in ICLR, 2023.

[145] G. Dong, R. Li, S. Wang et al., "Bridging the kb-text gap: Leveraging structured knowledge-aware pre-training for KBQA," in CIKM, 2023.

[146] K. Wang, F. Duan, S. Wang et al., "Knowledge-driven cot: Exploring faithful reasoning in llms for knowledge-intensive question answering," arXiv:2308.13259, 2023.

[147] D. Yu and Y. Yang, "Retrieval-enhanced generative model for large-scale knowledge graph completion," in SIGIR, 2023.

[148] W. Chen, H. Hu, C. Saharia, and W. W. Cohen, "Re-imagen: Retrieval-augmented text-to-image generator," in ICLR, 2023.

[149] S. Sheynin, O. Ashual, A. Polyak et al., "Knn-diffusion: Image generation via large-scale retrieval," in ICLR, 2023.

[150] A. Blattmann, R. Rombach, K. Oktay et al., "Retrieval-augmented diffusion models," in NeurIPS, 2022.

[151] R. Rombach, A. Blattmann, and B. Ommer, "Text-guided synthesis of artistic images with retrieval-augmented diffusion models," arXiv:2207.13038, 2022.

[152] B. Li, P. H. Torr, and T. Lukasiewicz, "Memory-driven text-to-image generation," arXiv:2208.07022, 2022.

[153] A. Bertsch, U. Alon, G. Neubig, and M. R. Gormley, "Unlimiformer: Long-range transformers with unlimited length input," 2023.

[154] Y. Kuratov, A. Bulatov et al., "In search of needles in a 10m haystack: Recurrent memory finds what llms miss," arXiv:2402.10790, 2024.

[155] N. F. Liu, K. Lin, J. Hewitt et al., "Lost in the middle: How language models use long contexts," arxiv:2307.03172, 2023.

[156] T. Févry, L. B. Soares et al., "Entities as experts: Sparse memory access with entity supervision," in EMNLP, 2020.

[157] M. de Jong, Y. Zemlyanskiy, N. FitzGerald et al., "Mention memory: incorporating textual knowledge into transformers through entity mention attention," in ICLR, 2021.

[158] B. Jing, Y. Zhang, Z. Song et al., "Amd: Anatomical motion diffusion with interpretable motion decomposition and fusion," in AAAI, 2024.

[159] Y. Yuan, H. Liu, X. Liu et al., "Retrieval-augmented text-to-audio generation," in ICASSP, 2024.

[160] B. Yang, M. Cao, and Y. Zou, "Concept-aware video captioning: Describing videos with effective prior information," TIP, vol. 32, pp. 5366–5378, 2023.

[161] Z. Zhong, T. Lei, and D. Chen, "Training language models with memory augmentation," in EMNLP, 2022.

[162] S. Min, W. Shi, M. Lewis et al., "Nonparametric masked language modeling," in ACL Findings, 2023.

[163] X. Zhang, Y. Zhou, G. Yang, and T. Chen, "Syntax-aware retrieval augmented code generation," in EMNLP Findings, 2023.

[164] Z. Fei, "Memory-augmented image captioning," in AAAI, 2021.

[165] Y. Leviathan, M. Kalman, and Y. Matias, "Fast inference from transformers via speculative decoding," in ICML, 2023.

[166] T. Lan, D. Cai, Y. Wang et al., "Copy is all you need," in ICLR, 2023.

[167] B. Cao, D. Cai, L. Cui et al., "Retrieval is accurate generation," arXiv:2402.17532, 2024.

[168] L. Wang, N. Yang, and F. Wei, "Query2doc: Query expansion with large language models," in EMNLP, 2023.

[169] L. Gao, X. Ma, J. Lin, and J. Callan, "Precise zero-shot dense retrieval without relevance labels," in ACL, 2023.

[170] G. Kim, S. Kim, B. Jeon et al., "Tree of clarifications: Answering ambiguous questions with retrieval-augmented large language models," in EMNLP, 2023.

[171] M. Xia, S. Malladi, S. Gururangan et al., "LESS: selecting influential data for targeted instruction tuning," arXiv:2402.04333, 2024.

[172] S. Yao, J. Zhao, D. Yu et al., "React: Synergizing reasoning and acting in language models," in ICLR, 2023.

[173] J. Wei, X. Wang, D. Schuurmans et al., "Chain-of-thought prompting elicits reasoning in large language models," in NeurIPS, 2022.

[174] T. Pouplin, H. Sun, S. Holt, and M. Van der Schaar, "Retrieval-augmented thought process as sequential decision making," arXiv:2402.07812, 2024.

[175] J. Liu, "LlamaIndex," 11 2022. [Online]. Available: https://github.com/jerryjliu/llama_index

[176] P. Sarthi, S. Abdullah, A. Tuli et al., "Raptor: Recursive abstractive processing for tree-organized retrieval," in ICLR, 2023.

[177] S. Xiao, Z. Liu, P. Zhang et al., "C-pack: Packaged resources to advance general chinese embedding," arxiv:2309.07597, 2023.

[178] J. Chen, S. Xiao, P. Zhang et al., "Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation," arxiv:2309.07597, 2023.

[179] S. Xiao, Z. Liu, P. Zhang, and X. Xing, "Lm-cocktail: Resilient tuning of language models via model merging," arxiv:2311.13534, 2023.

[180] P. Zhang, S. Xiao, Z. Liu, Z. Dou, and J.-Y. Nie, "Retrieve anything to augment large language models," arxiv:2310.07554, 2023.

[181] M. Kulkarni, P. Tangarajan, K. Kim et al., "Reinforcement learning for optimizing RAG for domain chatbots," arXiv:2401.06800, 2024.

[182] W. Wang, Y. Wang et al., "Rap-gen: Retrieval-augmented patch generation with codet5 for automatic program repair," in ESEC/FSE, 2023.

[183] S.-Q. Yan, J.-C. Gu, Y. Zhu, and Z.-H. Ling, "Corrective retrieval augmented generation," arXiv:2401.15884, 2024.

[184] W. Huang, M. Lapata, P. Vougiouklis et al., "Retrieval augmented generation with rich answer encoding," in IJCNLP-AACL, 2023.

[185] H. Wang, W. Huang, Y. Deng et al., "Unims-rag: A unified multi-source retrieval-augmented generation for personalized dialogue systems," arXiv:2401.13256, 2024.

[186] M. R. Glass, G. Rossiello, M. F. M. Chowdhury et al., "Re2g: Retrieve, rerank, generate," in NAACL, 2022.

[187] R. F. Nogueira and K. Cho, "Passage re-ranking with BERT," arXiv:1901.04085, 2019.

[188] J. Li, Y. Zhao, Y. Li et al., "Acecoder: Utilizing existing code to enhance code generation," arXiv:2303.17780, 2023.

[189] P. Shi, R. Zhang, H. Bai, and J. Lin, "XRICL: cross-lingual retrieval-augmented in-context learning for cross-lingual text-to-sql semantic parsing," in EMNLP Findings, 2022.

[190] K. Rangan and Y. Yin, "A fine-tuning enhanced rag system with quantized influence measure as ai judge," arXiv:2402.17081, 2024.

[191] J. Saad-Falcon, O. Khattab, K. Santhanam et al., "Udapdr: Unsupervised domain adaptation via llm prompting and distillation of rerankers," in EMNLP, 2023.

[192] L. Wang, N. Yang, and F. Wei, "Learning to retrieve in-context examples for large language models," arXiv:2307.07164, 2023.

[193] Z. Wang, J. Araki, Z. Jiang et al., "Learning to filter context for retrieval-augmented generation," arxiv:2311.08377, 2023.

[194] S. Hofstätter, J. Chen, K. Raman, and H. Zamani, "Fid-light: Efficient and effective retrieval-augmented text generation," in SIGIR, 2023.

[195] D. Arora, A. Kini, S. R. Chowdhury et al., "Gar-meets-rag paradigm for zero-shot information retrieval," arXiv:2310.20158, 2023.

[196] https://www.pinecone.io.

[197] W. Yu, D. Iter et al., "Generate rather than retrieve: Large language models are strong context generators," arXiv:2209.10063, 2022.

[198] A. Abdallah and A. Jatowt, "Generator-retriever-generator: A novel approach to open-domain question answering," arXiv:2307.11278, 2023.

[199] E. Saravia, "Prompt Engineering Guide," https://github.com/dair-ai/Prompt-Engineering-Guide, 12 2022.

[200] H. S. Zheng, S. Mishra et al., "Take a step back: Evoking reasoning via abstraction in large language models," arxiv:2310.06117, 2023.

[201] S. Diao, P. Wang, Y. Lin, and T. Zhang, "Active prompting with chain-of-thought for large language models," arxiv:2302.12246, 2023.

[202] H. Jiang, Q. Wu, C. Lin et al., "Llmlingua: Compressing prompts for accelerated inference of large language models," in EMNLP, 2023.

[203] T. Ahmed, K. S. Pai, P. Devanbu, and E. T. Barr, "Automatic semantic augmentation of language model prompts (for code summarization)," arXiv:2304.06815, 2024.

[204] Z. Xu, Z. Liu, Y. Liu et al., "Activerag: Revealing the treasures of knowledge via active learning," arXiv:2402.13547, 2024.

[205] E. Nijkamp, B. Pang, H. Hayashi et al., "A conversational paradigm for program synthesis," arxiv:2203.13474, 2022.

[206] Y. He, M. Xia, H. Chen et al., "Animate-a-story: Storytelling with retrieval-augmented video generation," arXiv:2307.06940, 2023.

[207] E. J. Hu, Y. Shen, P. Wallis et al., "Lora: Low-rank adaptation of large language models," in ICLR, 2022.

[208] C. Liu, P. Çetin, Y. Patodia et al., "Automated code editing with search-generate-modify," arXiv:2306.06490, 2023.

[209] H. Joshi, J. P. C. Sánchez, S. Gulwani et al., "Repair is nearly generation: Multilingual program repair with llms," in AAAI, 2023.

[210] Z. Jiang, F. F. Xu, L. Gao et al., "Active retrieval augmented generation," arXiv:2305.06983, 2023.

[211] A. Mallen, A. Asai, V. Zhong et al., "When not to trust language models: Investigating effectiveness of parametric and non-parametric memories," in ACL, 2023.

[212] Z. Jiang, J. Araki, H. Ding, and G. Neubig, "How can we know When language models know? on the calibration of language models for question answering," TACL, 2021.

[213] N. Kandpal, H. Deng, A. Roberts et al., "Large language models struggle to learn long-tail knowledge," in ICML, 2023.

[214] R. Ren, Y. Wang, Y. Qu et al., "Investigating the factual knowledge boundary of large language models with retrieval augmentation," arXiv:2307.11019, 2023.

[215] Y. Wang, P. Li, M. Sun, and Y. Liu, "Self-knowledge guided retrieval augmentation for large language models," in EMNLP Findings, 2023.

[216] H. Ding, L. Pang, Z. Wei et al., "Retrieve only when it needs: Adaptive retrieval augmentation for hallucination mitigation in large language models," arXiv:2402.10612, 2024.

[217] S. Jeong, J. Baek, S. Cho et al., "Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity," arXiv:2403.14403, 2024.

[218] F. Zhang, B. Chen et al., "Repocoder: Repository-level code completion through iterative retrieval and generation," in EMNLP, 2023.

[219] Z. Shao, Y. Gong, Y. Shen et al., "Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy," in EMNLP Findings, 2023.

[220] X. Cheng, D. Luo, X. Chen et al., "Lift yourself up: Retrieval-augmented text generation with self-memory," in NeurIPS, 2023.

[221] Z. Wang, A. Liu, H. Lin et al., "Rat: Retrieval augmented thoughts elicit context-aware reasoning in long-horizon generation," arXiv:2403.05313, 2024.

[222] O. Agarwal, H. Ge, S. Shakeri, and R. Al-Rfou, "Knowledge graph based synthetic corpus generation for knowledge-enhanced language model pre-training," in NAACL-HLT, 2021.

[223] J. Sun, C. Xu, L. Tang et al., "Think-on-graph: Deep and responsible reasoning of large language model with knowledge graph," arXiv:2307.07697, 2023.

[224] P. Limkonchotiwat, W. Ponwitayarat, C. Udomcharoenchaikit et al., "Cl-relkt: Cross-lingual language knowledge transfer for multilingual retrieval question answering," in NAACL Findings, 2022.

[225] A. Asai, X. Yu, J. Kasai, and H. Hajishirzi, "One question answering model for many languages with cross-lingual dense passage retrieval," in NeurIPS, 2021.

[226] K. Lee, S. Han et al., "When to read documents or QA history: On unified and selective open-domain QA," in ACL Findings, 2023.

[227] S. Yue, W. Chen et al., "Disc-lawllm: Fine-tuning large language models for intelligent legal services," arXiv:2309.11325, 2023.

[228] S. Siriwardhana, R. Weerasekera, T. Kaluarachchi et al., "Improving the domain adaptation of retrieval augmented generation (RAG) models for open domain question answering," TACL, vol. 11, pp. 1–17, 2023.

[229] Y. Tang and Y. Yang, "Multihop-rag: Benchmarking retrieval-augmented generation for multi-hop queries," arXiv:2401.15391, 2024.

[230] K. Huang, C. Zhai, and H. Ji, "CONCRETE: improving cross-lingual fact-checking with cross-lingual retrieval," in COLING, 2022.

[231] L. Hagström, D. Saynova, T. Norlund et al., "The effect of scaling, retrieval augmentation and form on the factual consistency of language models," arXiv:2311.01307, 2023.

[232] Y. Liu, Y. Wan et al., "KG-BART: knowledge graph-augmented BART for generative commonsense reasoning," in AAAI, 2021.

[233] A. Wan, E. Wallace, and D. Klein, "What evidence do language models find convincing?" arXiv:2402.11782, 2024.

[234] H. Zhang, Z. Liu et al., "Grounded conversation generation as guided traverses in commonsense knowledge graphs," in ACL, 2020.

[235] D. Cai, Y. Wang, W. Bi et al., "Skeleton-to-response: Dialogue generation guided by retrieval memory," in NAACL-HLT, 2019.

[236] M. Komeili, K. Shuster, and J. Weston, "Internet-augmented dialogue generation," in ACL, 2022.

[237] K. Shuster, J. Xu et al., "Blenderbot 3: a deployed conversational agent that continually learns to responsibly engage," arXiv:2208.03188, 2022.

[238] S. Kim, J. Y. Jang, M. Jung, and S. Shin, "A model of cross-lingual knowledge-grounded response generation for open-domain dialogue systems," in EMNLP Findings, 2021.

[239] E. Nie, S. Liang, H. Schmid, and H. Schütze, "Cross-lingual retrieval augmented prompt for low-resource languages," in ACL, 2023.

[240] X. Li, E. Nie, and S. Liang, "From classification to generation: Insights into crosslingual retrieval augmented icl," in NeurIPS, 2023.

[241] W. Li, J. Li, W. Ma, and Y. Liu, "Citation-enhanced generation for llm-based chatbot," arXiv:2402.16063, 2024.

[242] D. Cai, Y. Wang, H. Li et al., "Neural machine translation with monolingual translation memory," in ACL/IJCNLP, 2021.

[243] U. Khandelwal, A. Fan, D. Jurafsky et al., "Nearest neighbor machine translation," in ICLR, 2021.

[244] X. Du and H. Ji, "Retrieval-augmented generative question answering for event argument extraction," in EMNLP, 2022.

[245] Y. Gao, Q. Yin, Z. Li et al., "Retrieval-augmented multilingual keyphrase generation with retriever-generator iterative training," in NAACL Findings, 2022.

[246] J. Zhang, E. J. Yu, Q. Chen et al., "Retrieval-based full-length wikipedia generation for emergent events," arXiv:2402.18264, 2024.

[247] R. Fan, Y. Fan, J. Chen et al., "RIGHT: retrieval-augmented generation for mainstream hashtag recommendation," arxiv:2312.10466, 2023.

[248] Y. Wang, H. Le, A. D. Gotmare et al., "Codet5mix: A pretrained mixture of encoder-decoder transformers for code understanding and generation," 2022.

[249] E. Nijkamp, B. Pang, H. Hayashi et al., "A conversational paradigm for program synthesis," arXiv:2203.13474, 2022.

[250] D. Zan, B. Chen, Y. Gong et al., "Private-library-oriented code generation with large language models," arXiv:2307.15370, 2023.

[251] A. Madaan, S. Zhou, U. Alon et al., "Language models of code are few-shot commonsense learners," in EMNLP, 2022.

[252] Y. Wang, H. Le, A. Gotmare et al., "Codet5+: Open code large language models for code understanding and generation," in EMNLP, 2023.

[253] D. Liao, S. Pan, Q. Huang et al., "Context-aware code generation framework for code repositories: Local, global, and third-party library awareness," arXiv:2312.05772, 2023.

[254] J. Li, Y. Li, G. Li et al., "Skcoder: A sketch-based approach for automatic code generation," in ICSE, 2023.

[255] M. Liu, T. Yang, Y. Lou et al., "Codegen4libs: A two-stage approach for library-oriented code generation," in ASE, 2023.

[256] K. Zhang, J. Li, G. Li et al., "Codeagent: Enhancing code generation with tool-integrated agent systems for real-world repo-level coding challenges," arXiv:2401.07339, 2024.

[257] Q. Gou, Y. Dong, Y. Wu, and Q. Ke, "Rrgcode: Deep hierarchical search-based code generation," Journal of Systems and Software, vol. 211, p. 111982, 2024.

[258] J. Chen, X. Hu, Z. Li et al., "Code search is all you need? improving code suggestions with code search," in ICSE, 2024.

[259] H. Su, S. Jiang, Y. Lai et al., "Arks: Active retrieval in knowledge soup for code generation," arXiv:2402.12317, 2024.

[260] N. Beau and B. Crabbé, "The impact of lexical and grammatical processing on generating code from natural language," in ACL Findings, 2022.

[261] K. Zhang, G. Li, J. Li et al., "Toolcoder: Teach code generation models to use API search tools," arXiv:2305.04032, 2023.

[262] S. Liu, Y. Chen, X. Xie et al., "Retrieval-augmented generation for code summarization via hybrid GNN," in ICLR, 2021.

[263] F. Yamaguchi, N. Golde, D. Arp, and K. Rieck, "Modeling and discovering vulnerabilities with code property graphs," in S&P, 2014.

[264] Y. Choi, C. Na et al., "Readsum: Retrieval-augmented adaptive transformer for source code summarization," IEEE Access, 2023.

[265] A. Alokla, W. Gad, W. Nazih et al., "Retrieval-based transformer pseudocode generation," Mathematics, vol. 10, no. 4, p. 604, 2022.

[266] J. Zhao, X. Chen, G. Yang, and Y. Shen, "Automatic smart contract comment generation via large language models and in-context learning," IST, vol. 168, p. 107405, 2024.

[267] J. Xu, Z. Cui et al., "Unilog: Automatic logging via LLM and in-context learning," in ICSE, 2024.

[268] H. Wang, X. Xia et al., "Context-aware retrieval-based deep commit message generation," TOSEM, vol. 30, no. 4, pp. 56:1–56:30, 2021.

[269] X. Zhu, C. Sha, and J. Niu, "A simple retrieval-based method for code comment generation," in SANER, 2022.

[270] T. Ye, L. Wu, T. Ma et al., "Tram: A token-level retrieval-augmented mechanism for source code summarization," arXiv:2305.11074, 2023.

[271] L. Li, B. Liang, L. Chen, and X. Zhang, "Cross-modal retrieval-enhanced code summarization based on joint learning for retrieval and generation," Available at SSRN 4724884.

[272] D. Drain, C. Hu, C. Wu et al., "Generating code with the help of retrieved template functions and stack overflow answers," arXiv:2104.05310, 2021.

[273] S. Lu, D. Guo et al., "Codexglue: A machine learning benchmark dataset for code understanding and generation," in NeurIPS Datasets and Benchmarks, 2021.

[274] Y. Ding, Z. Wang et al., "Cocomic: Code completion by jointly modeling in-file and cross-file context," arXiv:2212.10007, 2022.

[275] D. Shrivastava, D. Kocetkov et al., "Repofusion: Training code models to understand your repository," arXiv:2306.10998, 2023.

[276] Z. Tang, J. Ge, S. Liu et al., "Domain adaptive code completion via language models and decoupled domain databases," in ASE, 2023.

[277] W. Sun, H. Li, M. Yan et al., "Revisiting and improving retrieval-augmented deep assertion generation," in ASE, 2023.

[278] A. Eghbali and M. Pradel, "De-hallucinator: Iterative grounding for llm-based code completion," arXiv:2401.01701, 2024.

[279] M. Liang, X. Xie, G. Zhang et al., "Repofuse: Repository-level code completion with fused dual context," arXiv:2402.14323, 2024.

[280] Y. Tsai, M. Liu, and H. Ren, "Rtlfixer: Automatically fixing RTL syntax errors with large language models," arXiv:2311.16543, 2023.

[281] B. Bogin, S. Gupta, P. Clark et al., "Leveraging code to improve in-context learning for semantic parsing," arXiv:2311.09519, 2023.

[282] H. Li, J. Zhang, C. Li, and H. Chen, "Resdsql: Decoupling schema linking and skeleton parsing for text-to-sql," in AAAI, 2023.

[283] K. Zhang, X. Lin, Y. Wang et al., "Refsql: A retrieval-augmentation framework for text-to-sql generation," in EMNLP Findings, 2023.

[284] S. Chang and E. Fosler-Lussier, "Selective demonstrations for cross-domain text-to-sql," arXiv:2310.06302, 2023.

[285] L. Nan, Y. Zhao, W. Zou et al., "Enhancing text-to-sql capabilities of large language models: A study on prompt design strategies," in EMNLP Findings, 2023.

[286] X. Zhang, D. Wang, L. Dou et al., "Multi-hop table retrieval for open-domain text-to-sql," arXiv:2402.10666, 2024.

[287] H. Li, J. Zhang, H. Liu et al., "Codes: Towards building open-source language models for text-to-sql," arXiv:2402.16347, 2024.

[288] Z. Jie and W. Lu, "Leveraging training data in few-shot prompting for numerical reasoning," arXiv:2305.18170, 2023.

[289] M. Gao, J. Li, H. Fei et al., "De-fine: Decomposing and refining visual programs with auto-feedback," arXiv:2311.12890, 2023.

[290] Y. Hao, W. Chen, Z. Zhou, and W. Cui, "E&v: Prompting large language models to perform static analysis by pseudo-code execution and verification," arXiv:2312.08477, 2023.

[291] Y. Guo, Z. Li et al., "Retrieval-augmented code generation for universal information extraction," arXiv:2311.02962, 2023.

[292] G. Pinto, C. de Souza et al., "Lessons from building stackspot ai: A contextualized ai coding assistant," arXiv:2311.18450, 2024.

[293] Z. Liu, C. Chen, J. Wang et al., "Testing the limits: Unusual text inputs generation for mobile app crash detection with large language model," arXiv:2310.15657, 2023.

[294] K. D. Bollacker, C. Evans et al., "Freebase: a collaboratively created graph database for structuring human knowledge," in SIGMOD, 2008.

[295] Y. Shu and Z. Yu, "Data distribution bottlenecks in grounding language models to knowledge bases," arXiv:2309.08345, 2023.

[296] D. Leake and D. J. Crandall, "On bringing case-based reasoning methodology to deep learning," in ICCBR, 2020.

[297] L. Zhang, J. Zhang et al., "FC-KBQA: A fine-to-coarse composition framework for knowledge base question answering," in ACL, 2023.

[298] J. Jiang, K. Zhou et al., "Structgpt: A general framework for large language model to reason over structured data," in EMNLP, 2023.

[299] J. Baek, A. F. Aji, and A. Saffari, "Knowledge-augmented language model prompting for zero-shot knowledge graph question answering," arXiv:2306.04136, 2023.

[300] P. Sen, S. Mavadia, and A. Saffari, "Knowledge graph-augmented language models for complex question answering," in NLRSE, 2023.

[301] Y. Wu, N. Hu, S. Bi et al., "Retrieve-rewrite-answer: A kg-to-text enhanced llms framework for knowledge graph question answering," arXiv:2309.11206, 2023.

[302] C. Wang, Y. Xu, Z. Peng et al., "keqing: knowledge-based question answering is a nature chain-of-thought mentor of LLM," arXiv:2401.00426, 2024.

[303] J. Liu, S. Cao, J. Shi et al., "Probing structured semantics understanding and generation of language models via question answering," arXiv:2401.05777, 2024.

[304] G. Xiong, J. Bao, and W. Zhao, "Interactive-kbqa: Multi-turn interactions for knowledge base question answering with large language models," arXiv:2402.15131, 2024.

[305] S. Chen, Q. Liu, Z. Yu et al., "Retrack: A flexible and efficient framework for knowledge base question answering," in ACL, 2021.

[306] D. Yu, C. Zhu, Y. Fang et al., "Kg-fid: Infusing knowledge graph in fusion-in-decoder for open-domain question answering," in ACL, 2022.

[307] Z. Hu, Y. Xu, W. Yu et al., "Empowering language models with knowledge graph reasoning for open-domain question answering," in EMNLP, 2022.

[308] M. Ju, W. Yu, T. Zhao et al., "Grape: Knowledge graph enhanced passage reader for open-domain question answering," in EMNLP Findings, 2022.

[309] Q. Yang, Q. Chen, W. Wang et al., "Enhancing multi-modal multi-hop question answering via structured knowledge and unified retrieval-generation," in MM, 2023.

[310] W. Zhao, Y. Liu, T. Niu et al., "DIVKNOWQA: assessing the reasoning ability of llms via open-domain question answering over knowledge base and text," arXiv:2310.20170, 2023.

[311] X. Wang, Q. Yang, Y. Qiu et al., "Knowledgpt: Enhancing large language models with retrieval and storage access on knowledge bases," arXiv:2308.11761, 2023.

[312] S. Ko, H. Cho, H. Chae et al., "Evidence-focused fact summarization for knowledge-augmented zero-shot question answering," arXiv:2403.02966, 2024.

[313] Y. Gao, L. Qiao, Z. Kan et al., "Two-stage generative question answering on temporal knowledge graph using large language models," arXiv:2402.16568, 2024.

[314] T. Guo, Q. Yang, C. Wang et al., "Knowledgenavigator: Leveraging large language models for enhanced reasoning over knowledge graph," arXiv:2312.15880, 2023.

[315] S. Min, J. Boyd-Graber, C. Alberti et al., "Neurips 2020 efficientqa competition: Systems, analyses and lessons learned," in NeurIPS 2020 Competition and Demonstration Track, 2021.

[316] A. H. Li, P. Ng, P. Xu et al., "Dual reader-parser on hybrid textual and tabular evidence for open domain question answering," in ACL/IJCNLP, 2021.

[317] P. Christmann, R. S. Roy, and G. Weikum, "Conversational question answering on heterogeneous sources," in SIGIR, 2022.

[318] K. Ma, H. Cheng, X. Liu et al., "Open-domain question answering via chain of reasoning over heterogeneous knowledge," in EMNLP Findings, 2022.

[319] E. Park, S.-M. Lee et al., "Rink: reader-inherited evidence reranker for table-and-text open domain question answering," in AAAI, 2023.

[320] W. Zhao, Y. Liu, Y. Wan et al., "Localize, retrieve and fuse: A generalized framework for free-form question answering over tables," arXiv:2309.11049, 2023.

[321] F. Pan, M. Canim et al., "End-to-end table question answering via retrieval-augmented generation," arXiv:2203.16714, 2022.

[322] Z. Jiang, Y. Mao, P. He et al., "Omnitab: Pretraining with natural and synthetic data for few-shot table-based question answering," in NAACL, 2022.

[323] W. Zhong, J. Huang, Q. Liu et al., "Reasoning over hybrid chain for table-and-text open domain question answering," in IJCAI, 2022.

[324] A. S. Sundar and L. Heck, "ctbl: Augmenting large language models for conversational tables," arXiv:2303.12024, 2023.

[325] D. Min, N. Hu, R. Jin et al., "Exploring the impact of table-to-text methods on augmenting llm-based question answering with domain hybrid data," arXiv:2402.12869, 2024.

[326] S. Wu, Y. Li, D. Zhang, and Z. Wu, "Improving knowledge-aware dialogue response generation by using human-written prototype dialogues," in EMNLP Findings, 2020.

[327] M. Kang, J. M. Kwak et al., "Knowledge-consistent dialogue generation with knowledge graphs," in ICML Workshop, 2022.

[328] Z. Ji, Z. Liu, N. Lee et al., "RHO: reducing hallucination in open-domain dialogues with knowledge grounding," in ACL Findings, 2023.

[329] J. Baek, N. Chandrasekaran, S. Cucerzan et al., "Knowledge-augmented large language models for personalized contextual query suggestion," arXiv:2311.06318, 2023.

[330] X. He, Y. Tian, Y. Sun et al., "G-retriever: Retrieval-augmented generation for textual graph understanding and question answering," arXiv:2402.07630, 2024.

[331] Y. Kirstain, O. Levy, and A. Polyak, "X&fuse: Fusing visual information in text-to-image generation," arXiv:2303.01000, 2023.

[332] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," NeurIPS, 2021.

[333] Z. Zhang, A. Zhang, M. Li et al., "Multimodal chain-of-thought reasoning in language models," arXiv:2302.00923, 2023.

[334] C. Xu, M. Yang, X. Ao et al., "Retrieval-enhanced adversarial training with dynamic memory-augmented attention for image paragraph captioning," Knowledge-Based Systems, vol. 214, p. 106730, 2021.

[335] R. Ramos, D. Elliott, and B. Martins, "Retrieval-augmented image captioning," in EACL, 2023.

[336] Z. Hu, A. Iscen, C. Sun et al., "Reveal: Retrieval-augmented visual-language pre-training with multi-source multimodal knowledge memory," in CVPR, 2023.

[337] Z. Li, W. Zhao, X. Du et al., "Cross-modal retrieval and semantic refinement for remote sensing image captioning," Remote Sensing, vol. 16, no. 1, p. 196, 2024.

[338] Z. Yang, Z. Gan, J. Wang et al., "An empirical study of gpt-3 for few-shot knowledge-based vqa," in AAAI, 2022.

[339] W. Lin and B. Byrne, "Retrieval augmented visual question answering with outside knowledge," in EMNLP, 2022.

[340] A. Fan, C. Gardent, C. Braud, and A. Bordes, "Augmenting transformers with knn-based composite memory for dialog," TACL, vol. 9, pp. 82–99, 2021.

[341] Z. Liang, H. Hu, C. Xu et al., "Maria: A visual experience powered conversational agent," in ACL-IJCNLP, 2021.

[342] Q. Fang and Y. Feng, "Neural machine translation with phrase-level universal visual representations," in ACL, 2022.

[343] S. Whitehead, H. Ji, M. Bansal et al., "Incorporating background knowledge into video description generation," in EMNLP, 2018.

[344] C. Yin, J. Tang, Z. Xu, and Y. Wang, "Memory augmented deep recurrent neural network for video question answering," TNNLS, vol. 31, no. 9, pp. 3159–3167, 2019.

[345] J. Pan, Z. Lin, Y. Ge et al., "Retrieving-to-answer: Zero-shot video question answering with frozen large language models," in ICCV, 2023.

[346] J. Lei, L. Yu, T. L. Berg, and M. Bansal, "Tvqa+: Spatio-temporal grounding for video question answering," in ACL, 2020.

[347] H. Le, N. Chen, and S. Hoi, "Vgnmn: Video-grounded neural module networks for video-grounded dialogue systems," in NAACL, 2022.

[348] Z. Wang, M. Li, R. Xu et al., "Language models with image descriptors are strong few-shot video-language learners," in NeurIPS, 2022.

[349] J. Yuan, S. Sun, D. Omeiza et al., "Rag-driver: Generalisable driving explanations with retrieval-augmented in-context learning in multi-modal large language model," arXiv:2402.10828, 2024.

[350] S. Ghosh, S. Kumar, C. K. R. Evuru et al., "Recap: retrieval-augmented audio captioning," in ICASSP, 2024.

[351] B. Elizalde, S. Deshmukh, and H. Wang, "Natural language supervision for general-purpose audio representations," in ICASSP, 2024.

[352] T. Kouzelis and V. Katsouros, "Weakly-supervised automated audio captioning via text only training," in DCASE Workshop, 2023.

[353] S. Deshmukh, B. Elizalde, D. Emmanouilidou et al., "Training audio captioning models without audio," in ICASSP, 2024.

[354] L. Yang, Z. Huang, X. Zhou et al., "Prompt-based 3d molecular diffusion models for structure-based drug design," 2023.

[355] T. Truong Jr and T. Bepler, "Poet: A generative model of protein families as sequences-of-sequences," NeurIPS, 2024.

[356] G. Frisoni, M. Mizutani, G. Moro, and L. Valgimigli, "Bioreader: a retrieval-enhanced text-to-text transformer for biomedical literature," in EMNLP, 2022.

[357] X. Yang, M. Ye, Q. You et al., "Writing by memorizing: Hierarchical retrieval-based medical report generation," arXiv:2106.06471, 2021.

[358] J. Kim and M. Min, "From rag to qa-rag: Integrating generative ai for pharmaceutical regulatory compliance process," arXiv:2402.01717, 2024.

[359] K. Yang, A. Swope et al., "Leandojo: Theorem proving with retrieval-augmented language models," in NeurIPS, 2024.

[360] Z. Levonian, C. Li, W. Zhu et al., "Retrieval-augmented generation to improve math question-answering: Trade-offs between groundedness and human preference," arXiv:2310.03184, 2023.

[361] J. Chen, H. Lin, X. Han, and L. Sun, "Benchmarking large language models in retrieval-augmented generation," arxiv:2309.01431, 2023.

[362] S. ES, J. James, L. E. Anke, and S. Schockaert, "RAGAS: automated evaluation of retrieval augmented generation," arxiv:2309.15217, 2023.

[363] J. Saad-Falcon, O. Khattab, C. Potts et al., "ARES: an automated evaluation framework for retrieval-augmented generation systems," arxiv:2311.09476, 2023.

[364] https://github.com/truera/trulens.

[365] Y. Lyu, Z. Li, S. Niu et al., "CRUD-RAG: A comprehensive chinese benchmark for retrieval-augmented generation of large language models," arxiv:2401.17043, 2024.

[366] G. Xiong, Q. Jin, Z. Lu, and A. Zhang, "Benchmarking retrieval-augmented generation for medicine," arXiv:2402.13178, 2024.

[367] F. Petroni, A. Piktus et al., "Kilt: a benchmark for knowledge intensive language tasks," in NAACL-HLT, 2021.

[368] S. Barnett, S. Kurniawan, S. Thudumu et al., "Seven failure points when engineering a retrieval augmented generation system," arXiv:2401.05856, 2024.

[369] F. Cuconasu, G. Trappolini, F. Siciliano et al., "The power of noise: Redefining retrieval for RAG systems," arXiv:2401.14887, 2024.

[370] L. Qiu, P. Shaw, P. Pasupat et al., "Evaluating the impact of model scale for compositional generalization in semantic parsing," arXiv:2205.12253, 2022.

[371] R. Jagerman, H. Zhuang, Z. Qin et al., "Query expansion by prompting large language models," arxiv:2305.03653, 2023.

[372] H. Zhang, P. Zhao, X. Miao et al., "Experimental analysis of large-scale learnable vector storage compression," VLDB, 2023.

[373] R. Aksitov, C. Chang, D. Reitter et al., "Characterizing attribution and fluency tradeoffs for retrieval-augmented large language models," arXiv:2302.05578, 2023.

[374] C. Han, Q. Wang, W. Xiong et al., "Lm-infinite: Simple on-the-fly length generalization for large language models," arXiv:2308.16137, 2023.

[375] H. Chase, "Langchain," https://github.com/langchain-ai/langchain, 2022.

[376] W. Jiang, S. Zhang, B. Han et al., "Piperag: Fast retrieval-augmented generation via algorithm-system co-design," arXiv:2403.05676, 2024.

[377] S. Jindal, "Did google gemini 1.5 really kill rag?" https://analyticsindiamag.com/did-google-gemini-1-5-really-kill-rag/, 2024.