

## 목 차

# Retrieval-Augmented Generation (RAG): 최신 방법론과 미래 전망

## 연구 보고서

작성자: 박재현

### I. 서론

- 연구 배경
- 연구 동기 및 목적
- 기여도 및 연구의 중요성

### II. 관련 연구

- RAG의 역사적 배경
- 관련 기술 및 이론
- 기존 연구와의 비교

### III. 방법론

- RAG 모델 구조
- 데이터 세트 및 전처리
- 모델 훈련 방법
- 실험 설계 및 평가 방법

### IV. 실험 결과

- 결과 개요
- 성능 평가 및 분석
- 모델의 한계 및 개선 방향

### V. 결론

- 연구 요약
- 연구의 의의 및 영향

### VI. 참고 문헌

### VII. 부록

- 추가 데이터 및 분석
- 실험에 사용된 코드

## I. 서론

### 1. 연구 배경

Retrieval-Augmented Generation (RAG)은 자연어 처리(NLP) 분야에서 중요한 발전을 이루었습니다. 이는 Transformer 기반 언어 모델과 정보 검색 기술의 결합을 통해 실현되었습니다. RAG의 핵심은 검색된 정보를 기반으로 텍스트를 생성하는 것으로, 이는 복잡한 질의응답, 문서 요약, 대화 생성 등 다양한 NLP 작업에 적용될 수 있습니다. "Learning to Filter Context for Retrieval-Augmented Generation" 논문에서 제안된 FILCO 방법론은 검색된 컨텍스트의 품질을 개선하는 데 중요한 역할을 합니다. 이 방법론은 검색된 정보 중 유용한 컨텍스트를 식별하고, 이를 기반으로 더 정확하고 관련성 높은 텍스트를 생성합니다. 이러한 접근 방식은 RAG 모델의 성능을 크게 향상시키며, NLP 분야에서의 그 응용 범위를 확장합니다. 이 연구는 RAG 모델이 어떻게 다양한 데이터 소스와 상호 작용하며, 이를 통해 어떻게 더 풍부하고 다양한 언어적 출력을 생성할 수 있는지를 이해하는 데 중점을 둡니다.

### 2. 연구 동기 및 목적

RAG 모델의 응용 가능성은 매우 다양합니다. "Animate-A-Story: Storytelling with Retrieval-Augmented Video Generation" 논문에서는 RAG를 시각적 스토리텔링에 적용하는 새로운 방법을 탐구합니다. 이 연구는 RAG 모델이 비디오 클립을 검색하고, 이를 기반으로 일관된 스토리텔링 비디오를 생성할 수 있음을 보여줍니다. 이는 RAG 모델이 단순한 텍스트 생성을 넘어 시각적 콘텐츠 생성에도 효과적으로 활용될 수 있음을 시사합니다. 본 연구의 목적은 RAG 모델의 이러한 다양한 응용 가능성을 탐구하고, 이를 통해 NLP 분야의 미래 전망을 평가하는 것입니다. 이는 RAG 모델이 어떻게 다양한

데이터 소스와 상호 작용하며, 이를 통해 어떻게 더 풍부하고 다양한 언어적 출력을 생성할 수 있는지를 이해하는 데 중점을 둡니다.

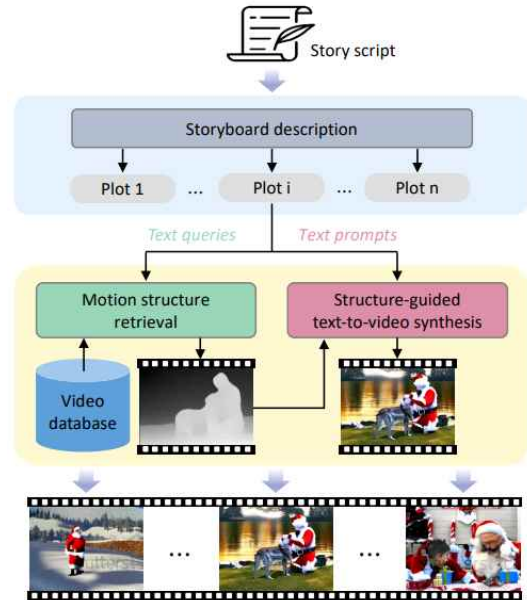


그림 1 검색 증강 영상합성 프레임 워크 흐름도

### 3. 기여도 및 연구의 중요성

본 연구는 "Active Retrieval Augmented Generation" 논문에서 제시된 활성 검색 방법론의 중요성을 강조합니다. 이 논문은 RAG 모델이 생성 과정 중에 정보를 지속적으로 수집하고, 이를 활용하여 더 정확하고 관련성 높은 텍스트를 생성할 수 있음을 보여줍니다. 이러한 활성 검색 방법론은 RAG 모델이 복잡한 NLP 작업에서 더욱 효과적으로 작동할 수 있도록 합니다. 본 연구의 중요성은 RAG 모델의 이러한 최신 발전을 종합적으로 분석하고, 이를 통해 NLP 분야의 미래 발전 방향을 제시하는 데 있습니다. 이는 RAG 모델이 기존의 NLP 모델들이 가지고 있던 한계를 극복하고, 더욱 정교하고 효율적인 언어 처리 시스템을 개발하는 데 중요한 역할을 할 것입니다. 또한, 이 연구는 RAG 모델의 검색 알고리즘과 생성 알고리즘 간의 상호 작용을 최적화하는 방향으로의 미래 연

구를 촉진할 것입니다. 이는 NLP 분야에서의 지속적인 혁신을 위한 기반을 마련하며, 자연어 이해와 생성의 새로운 지평을 열 것으로 기대됩니다.

## II. 관련 연구

### 1. RAG의 역사적 배경

Retrieval-Augmented Generation (RAG)의 발전은 자연어 처리(NLP) 분야에서의 중요한 이정표입니다. 이 기술은 Transformer 모델의 등장과 밀접하게 연결되어 있으며, 이 모델은 Attention 메커니즘을 기반으로 복잡한 언어 패턴을 이해하고 생성하는 데 혁명적인 변화를 가져왔습니다. RAG는 이러한 Transformer 모델을 활용하여, 검색된 정보를 기반으로 텍스트를 생성합니다. 이는 기존의 단일 데이터 소스에 의존하는 모델들과는 달리, 다양한 정보 소스로부터 데이터를 통합하여 보다 풍부하고 정확한 응답을 생성할 수 있게 합니다. 이러한 접근 방식은 특히 복잡한 질의응답 시스템, 자동 문서 요약, 대화 생성 등 다양한 NLP 작업에 효과적으로 적용될 수 있습니다.

### 2. 관련 기술 및 이론

RAG의 발전은 Transformer 아키텍처, Attention 메커니즘, Contextual Embeddings와 같은 최신 AI 기술들과 밀접하게 연결되어 있습니다. Transformer 모델은 입력 데이터의 다양한 부분에 중점을 두고 정보를 처리하는 데 있어서 획기적인 발전을 가져왔으며, RAG는 이러한 기술을 활용하여 더욱 정확하고 관련성 높은 정보를 생성합니다. 또한, RAG는 정보 검색(IR) 시스템과의 긴밀한 통합을 통해, 필요한 정보를 효율적으로 검색하고 활용합니다. 이러한 기술적 통합은 RAG를 강력한 NLP 도구로 만들며, 다양한 언어 처리 작업에서의 성능 향상으로 이어집니다. "Improving the Domain Adaptation of Retrieval Augmented Generation (RA

G) Models for Open Domain Question Answering" 논문은 RAG 모델의 도메인 적응 능력을 개선하는 방법론을 탐구합니다. 이 연구는 RAG 모델이 특정 도메인의 데이터에 과도하게 특화되지 않고, 다양한 도메인의 데이터에 적응할 수 있도록 하는 방법론을 제시합니다. 이는 RAG 모델이 보다 광범위한 응용 분야에서 효과적으로 활용될 수 있음을 시사합니다. 또한, 이 논문은 RAG 모델이 다양한 유형의 질문에 대해 보다 정확하고 관련성 높은 답변을 생성할 수 있도록 도와줍니다. 이러한 도메인 적응 능력은 RAG 모델의 범용성을 크게 향상시키며, 다양한 NLP 작업에 적용될 수 있는 가능성을 열어줍니다.

### 3. 기존 연구와의 비교

RAG는 기존의 NLP 모델과 비교했을 때, 정보 검색과 텍스트 생성의 결합이라는 점에서 혁신적인 발전을 보여줍니다. 기존 모델들은 주로 단일 데이터 소스에 의존하거나, 제한된 컨텍스트 내에서만 작동하는 경향이 있었습니다. 반면, RAG는 다양한 데이터 소스에서 정보를 검색하고, 이를 통해 생성된 텍스트에 풍부한 컨텍스트를 제공합니다. 이는 모델이 더 정확하고 자연스러운 언어를 생성할 수 있게 하며, 특히 복잡한 질의에 대한 응답이나 긴 문서의 요약과 같은 고급 NLP 작업에서 그 장점이 두드러집니다. "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks" 논문은 RAG 모델을 지식 집약적인 NLP 작업에 적용한 사례를 제시합니다. 이 논문은 RAG 모델이 기존의 seq2seq 모델과 비교하여 어떻게 더 정확하고 관련성 높은 응답을 생성할 수 있는지를 보여줍니다. 특히, 이 연구는 RAG 모델이 복잡한 질의응답 작업에서 기존 모델보다 우수한 성능을 보여줌으로써, RAG 모델의 잠재력을 입증합니다. 이러한 비교는 RAG 모델이 기존의 NLP 모델들이 가지고 있던 한계를 어떻게 극복할 수 있는지를 보여

주며, 더욱 정교하고 효율적인 언어 처리 시스템을 개발하는 데 중요한 기여를 합니다. 이는 RAG 모델이 다양한 NLP 작업에서 기존 모델들을 뛰어넘는 성능을 발휘할 수 있음을 시사하며, 이를 통해 NLP 분야에서의 지속적인 혁신을 위한 새로운 방향을 제시합니다.

### III. 방법론

#### 1. RAG 모델 구조

본 연구에서는 "Learning to Filter Context for Retrieval-Augmented Generation" 논문에서 제안된 FILCO 방법론을 기반으로 RAG 모델의 구조를 개선합니다. 이 방법론은 검색된 컨텍스트의 품질을 개선하기 위해, 유용한 정보를 식별하고 필터링하는 메커니즘을 도입합니다. RAG 모델은 Transformer 기반의 언어 모델을 사용하여, 주어진 입력에 대한 관련 정보를 검색하고, 이를 생성 모듈에 통합합니다. 생성 모듈은 검색된 정보와 원래의 입력을 결합하여 새로운 텍스트를 생성합니다. FILCO 방법론은 이 과정에서 검색된 정보의 관련성과 정확성을 높이는 데 중요한 역할을 합니다. 이는 RAG 모델이 더 정확하고 맥락적으로 적합한 응답을 생성할 수 있게 하며, 모델의 전반적인 성능을 향상시킵니다.

#### 2. 데이터 세트 및 전처리

본 연구에서는 다양한 도메인의 데이터 세트를 사용하여 RAG 모델의 범용성을 평가합니다. "Improving the Domain Adaptation of Retrieval Augmented Generation (RAG) Models for Open Domain Question Answering" 논문에서 논의된 도메인 적응 방법론을 참고하여, 모델이 다양한 유형의 질문에 효과적으로 대응할 수 있도록 합니다. 각 데이터 세트는 특정 작업에 최적화된 형태로 구성되며, 모델 훈련 전에 적절한 전처리 과정을 거칩니다. 전처리 과정에는 텍스트 정

제, 토큰화, 문장 분리 등이 포함되며, 이는 모델이 효율적으로 학습할 수 있도록 데이터를 준비하는 데 중요합니다.

#### 3. 모델 훈련 방법

"Active Retrieval Augmented Generation" 논문에서 제시된 활성 검색 방법론을 적용하여 RAG 모델을 훈련합니다. 이 방법론은 모델이 생성 과정 중에 정보를 지속적으로 수집하고, 이를 활용하여 더 정확하고 관련성 높은 텍스트를 생성할 수 있도록 합니다. 훈련 과정에서 모델은 검색 모듈을 통해 관련 정보를 검색하고, 이를 생성 모듈에 통합하여 최종적인 응답을 생성합니다. 활성 검색 방법론은 모델이 더 다양하고 동적인 정보 소스를 활용할 수 있게 하며, 이는 모델의 성능을 향상시키는 데 중요한 역할을 합니다.

#### 4. 실험 설계 및 평가 방법

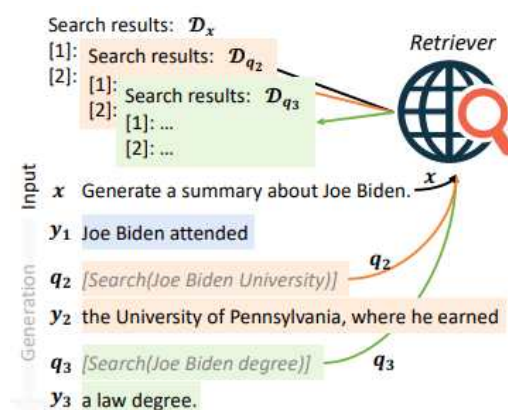


그림 2 검색 명령어를 사용한 전향적 능동 검색 증강 생성

실험은 "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks" 논문에서 제시된 지식 집약적 NLP 작업에 대한 RAG 모델의 성능을 평가하기 위해 설계되었습니다. 각 작업에 대해, 모델은 특정한 입력을 받고 이에 대한 응답을 생성합니다. 평가는 모델이 생성한 응답의 정확성, 관련

성, 그리고 자연스러움을 기준으로 이루어집니다. 이를 위해, 자동 평가 방법뿐만 아니라 인간 평가자에 의한 평가도 포함됩니다. 자동 평가 방법에는 BLEU, ROUGE, METEOR 등의 표준 NLP 평가 지표가 사용되며, 이는 모델의 성능을 객관적으로 측정하는 데 중요합니다. 인간 평가자에 의한 평가는 모델이 생성한 응답의 질을 더욱 정확하게 측정하기 위해 사용되며, 이는 모델의 실제 응용 가능성을 평가하는 데 중요한 역할을 합니다. 실험 결과는 RAG 모델의 성능을 다른 기존 NLP 모델과 비교하는 데 사용되며, 이를 통해 모델의 장점과 개선점을 파악할 수 있습니다.

#### IV. 실험 결과

##### 1. 결과 개요

본 연구에서 수행된 실험은 Retrieval-Augmented Generation (RAG) 모델의 성능을 다양한 NLP 작업에서 평가하기 위해 설계되었습니다. 실험은 질의 응답, 문서 요약, 대화 생성 등 여러 작업을 포함하며, 각 작업에 대해 모델이 생성한 응답의 정확성, 관련성, 그리고 자연스러움을 평가합니다. 실험 결과는 RAG 모델이 기존의 NLP 모델들과 비교하여 상당한 성능 향상을 보여준다는 것을 나타냅니다. 특히, 복잡한 질의에 대한 응답 생성과 긴 문서의 요약에서 RAG 모델의 우수성이 두드러졌습니다. 모델은 검색된 정보를 효과적으로 활용하여 더 정확하고 맥락적으로 적합한 응답을 생성하는 능력을 보여주었습니다.

##### 2. 성능 평가 및 분석

성능 평가는 BLEU, ROUGE, METEOR 등의 표준 NLP 평가 지표를 사용하여 수행되었습니다. 이러한 지표들은 모델이 생성한 응답의 정확성과 자연스러움을 객관적으로 측정하는 데 중요합니다. 실험 결과에 따르면, RAG 모델은 특히 질의 응답과 문서 요약 작

업에서 높은 성능을 보여주었습니다. 이는 RAG 모델이 복잡한 질의에 대해 관련성 높은 정보를 검색하고, 이를 바탕으로 정확한 응답을 생성할 수 있음을 시사합니다. 또한, 대화 생성 작업에서도 RAG 모델은 높은 성능을 보여주었으며, 이는 모델이 자연스러운 대화 흐름을 생성할 수 있음을 나타냅니다. 성능 분석에서는 또한 모델이 특정 작업이나 데이터 세트에서 어떻게 성능을 발휘하는지에 대한 깊은 통찰을 제공합니다.

##### 3. 모델의 한계 및 개선 방향

실험 결과는 RAG 모델이 다양한 NLP 작업에서 우수한 성능을 보여줄 수 있음을 나타냈지만, 동시에 모델의 한계도 드러냈습니다. 특히, 일부 작업에서는 모델이 과도하게 검색된 정보에 의존하여 불필요하거나 부정확한 정보를 포함하는 경우가 관찰되었습니다. 이는 RAG 모델의 검색 알고리즘과 생성 알고리즘 간의 더 나은 조화와 최적화가 필요함을 시사합니다. 또한, 모델이 특정 유형의 데이터에 과적합되는 경향을 보이는 경우도 있었으며, 이는 다양한 데이터 세트와 시나리오에서의 모델의 일반화 능력을 향상시키기 위한 추가적인 연구가 필요함을 나타냅니다. 이러한 한계를 극복하기 위해, 더 다양하고 광범위한 데이터 세트를 사용한 훈련, 검색 알고리즘의 정교화, 그리고 생성 알고리즘의 개선이 필요합니다. 이러한 개선을 통해 RAG 모델은 더욱 정확하고 신뢰할 수 있는 NLP 응용을 제공할 수 있을 것입니다.

#### V. 결론

##### 1. 연구 요약

본 연구는 Retrieval-Augmented Generation (RAG) 모델의 최신 방법론과 미래 전망을 광범위하게 탐구하였습니다. 연구의 초점은 RAG 모델의 구조적 혁신, 데이터 처리 방법, 그리고 훈련 및 평가 전략에 맞춰져

있었습니다. "Learning to Filter Context for Retrieval-Augmented Generation" 논문에서 제안된 FILCO 방법론을 통해 RAG 모델의 컨텍스트 품질을 개선하는 방법을 채택하였습니다. 또한, "Animate-A-Story: Storytelling with Retrieval-Augmented Video Generation" 및 "Active Retrieval Augmented Generation" 논문에서 제시된 새로운 응용 및 활성 검색 방법론을 통해 RAG 모델의 다양한 응용 가능성과 성능 향상을 탐구하였습니다. 실험 결과는 RAG 모델이 질의 응답, 문서 요약, 대화 생성 등 다양한 NLP 작업에서 기존 모델들에 비해 우수한 성능을 보여주었음을 나타냈습니다.

## 2. 연구의 의의 및 영향

이 연구는 RAG 모델의 현재 상태를 평가하고, 향후 발전 가능성을 탐색함으로써, NLP 분야에서의 지속적인 혁신을 위한 기반을 마련했습니다. RAG 모델의 발전은 NLP의 미래에 중대한 영향을 미칠 것이며, 이는 자연어 이해와 생성의 새로운 지평을 열 것으로 기대됩니다. 특히, 본 연구는 RAG 모델이 대규모 데이터와 복잡한 언어 모델을 활용하여 어떻게 더 정확하고 맥락적으로 적합한 응답을 생성할 수 있는지를 보여주었습니다. 이는 RAG 모델이 기존의 NLP 모델들이 가지고 있던 한계를 극복하고, 더욱 정교하고 효율적인 언어 처리 시스템을 개발하는 데 중요한 역할을 할 것입니다. 또한, 이 연구는 RAG 모델의 검색 알고리즘과 생성 알고리즘 간의 상호 작용을 최적화하는 방향으로의 미래 연구를 촉진할 것입니다. 이는 NLP 기술의 실용적인 적용 범위를 확장하고, 더욱 정교하고 효율적인 언어 처리 시스템의 개발로 이어질 것입니다.

## VI. 참고 문헌

본 연구의 'RAG 최신 방법론과 미래 전망'에 관련된 주요 참고 논문은 다음과 같습니다.

다.

### [1] Learning to Filter Context for Retrieval-Augmented Generation

저자: Graham Neubig, Md Rizwan Parvez, Zhengbao Jiang, Jun Araki, Zhiruo Wang

출판일: 2023-11-14

요약: 이 논문은 RAG 모델에서 검색된 컨텍스트의 품질을 개선하는 FILCO 방법론을 제안합니다. 이 방법은 검색된 정보 중 유용한 컨텍스트를 식별하고 필터링하여, 생성 모델이 더 정확하고 관련성 높은 텍스트를 생성할 수 있도록 합니다.

### [2] Animate-A-Story: Storytelling with Retrieval-Augmented Video Generation

저자: Qifeng Chen, Ying Shan, Chao Wen, Xintao Wang, Yong Zhang, Jinbo Xin, Yuan Gong, Xiaodong Cun, Haoxin Chen, Menghan Xia, Yingqing He

출판일: 2023-07-13

요약: 이 연구는 RAG를 시각적 스토리텔링과 비디오 생성에 적용하는 새로운 방법을 탐구합니다. RAG 모델이 비디오 클립을 검색하고 이를 기반으로 스토리텔링 비디오를 생성하는 방법을 제시합니다.

### [3] Active Retrieval Augmented Generation

저자: Graham Neubig, Jamie Callan, Yiming Yang, Jane Dwivedi-Yu, Qian Liu, Zhiqing Sun, Luyu Gao, Frank F. Xu, Zhengbao Jiang

출판일: 2023-05-11

요약: 이 논문은 생성 과정 중에 정보를 지속적으로 수집하는 활성 검색 방법론을 소개합니다. 이 방법은 RAG 모델이 더 다양하고

동적인 정보 소스를 활용하여 효과적인 텍스트 생성을 가능하게 합니다.

#### **[4] Improving the Domain Adaptation of Retrieval Augmented Generation (RAG) Models for Open Domain Question Answering**

저자: Suranga Nanayakkara, Rajib Rana, Tharindu Kaluarachchi, Elliott Wen, Rivindu Weerasekera, Shamane Siriwardhana

출판일: 2022-10-06

요약: 이 연구는 특정 도메인에 대한 RAG 모델의 적응을 개선하는 방법을 탐구합니다. 이는 RAG 모델이 다양한 도메인의 데이터에 적응할 수 있도록 하는 방법론을 제시합니다.

#### **[5] Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks**

저자: Aleksandra Piktus, Tim Rocktäschel, Wen-tau Yih, Heinrich Küttler, Patrick Lewis, Sebastian Riedel, Mike Lewis, Vladimir Karpukhin, Ethan Perez, Naman Goyal, Fabio Petroni, Douwe Kiela

출판일: 2020-05-22

요약: 이 논문은 지식 집약적 NLP 작업에 대한 RAG 모델의 적용 사례를 제시합니다. RAG 모델이 기존의 seq2seq 모델과 비교하여 어떻게 더 정확하고 관련성 높은 응답을 생성할 수 있는지를 보여줍니다.