

BridgeMind: Akademik Dokümanlar ve Dinamik Web Verileri Arasında Ajan Tabanlı Hibrit RAG Sistemi

ŞERİFE EROĞLU

Yapay Zeka Mühendisliği

Büyük Dil Modelleri (YZM 423) Proje Raporu

220212021

Abstract—Bu çalışma, yerel bir akademik doküman kütüphanesi ile gerçek zamanlı web arama yeteneklerini birleştiren "BridgeMind" adlı otonom hibrit bir Retrieval-Augmented Generation (RAG) sistemini sunmaktadır. Sistem, Phi-3-mini dil modelini kullanarak dökümanlardaki bilginin yeterliliğini sorgular ve eksik bilgi durumunda otonom olarak web üzerinden arama yapar. OpenAI gibi ücretli API'lar yerine tamamen açık kaynaklı Hugging Face modelleri ve yerel altyapı kullanılmıştır. Yapılan testler, sistemin yerel ve küresel bilgi arasında %100 doğrulukla geçiş yapabildiğini göstermiştir.

Index Terms—RAG, Agentic AI, Phi-3, Vektör Veritabanı, ChromaDB, Hugging Face, Bilgi Geri Getirme.

I. GİRİŞ VE PROBLEM TANIMI

Büyük Dil Modelleri (LLM), günümüzde doğal dil işleme görevlerinde devrim yaratmış olsa da, iki temel kısıtlama ile karşı karşıyadır: Eğitim verisiyle sınırlılık (Knowledge Cutoff) ve Halüsinasyon görme eğilimi. Özellikle teknik ve akademik alanlarda modellerin güncel olmayan veya döküman dışı bilgilerle hatalı çıkarımlar yapması, sistemlerin güvenilirliğini zedelemektedir.

Bu sorunu aşmak için geliştirilen geleneksel Retrieval-Augmented Generation (RAG) sistemleri, modelleri belirli bir veri havuzuna (örneğin yerel PDF dökümanları) bağlayarak cevap kalitesini artırmayı hedefler. Ancak, bu sistemlerin de "kapalı kutu" olarak çalışması, döküman dışı genel kültür veya güncel olaylar hakkında sorulara yanıt veremelerine ya da dökümandaki bilgileri zorlayarak hatalı cevap üretmelerine yol açmaktadır.

BridgeMind Projesi, bu kısıtlamaları aşmak amacıyla "Ajan Tabanlı Hibrit RAG" (Agentic Hybrid RAG) mimarisini sunar. Problemimiz; sistemin sadece elindeki dökümanla sınırlı kalmaması, aynı zamanda dökümanda bulamadığı bilgiyi uygunluk yerine otonom bir kararla dış kaynaklara (Web) yönelmesini sağlamaktır. Çalışmamızda; Microsoft'un Phi-3-mini modeli kullanılarak, yerel dökümanlardaki benzerlik skorları (relevance scores) üzerinden bir karar mekanizması kurulmuş ve sistemin bilgi kaynağını otonom olarak seçmesi hedeflenmiştir.

II. PROJE HEDEFİ VE AMACI

Bu projenin temel hedefi, araştırmacılar ve teknik kullanıcılar için halüsinasyon riskini minimize eden ve dinamik veri ihtiyacını otonom olarak karşılayan gelişmiş bir araştırma

asistanı geliştirmektir. Projenin spesifik amaçları aşağıda madde halinde sunulmuştur:

- Akademik Derinlik Sağlamak:** Yerel bir veri havuzuna (10 adet teknik PDF) odaklanarak, Büyük Dil Modellerinin (LLM) genel eğitim verilerinde bulunmayan spesifik akademik metodolojiler (Self-RAG, GraphRAG, FLARE vb.) hakkında doğru ve kaynak gösterilebilir cevaplar üretmesini sağlamak.
- Bilgi Kaynağını Otonom Belirlemek (Agentic Logic):** Sistemin, sorulan sorunun eldeki dökümanlarla cevaplanıp cevaplanamayacağına karar veren bir "akıl yürütme katmanı" (reasoning layer) oluşturması hedeflenmiştir. Bilgi dökümda yoksa uydurmak (hallucination) yerine doğrudan dış kaynaklara yönelmesi amaçlanır.
- Gerçek Zamanlı Veri Entegrasyonu:** Kapalı devre RAG sistemlerinin en büyük dezavantajı olan "güncel bilgiye erişememe" sorununu, DuckDuckGo arama motoru entegrasyonu ile aşarak sisteme güncel olaylar, fiyatlar ve genel tarifler hakkında bilgi verme yetisi kazandırmak.
- Maliyet ve Gizlilik Optimizasyonu:** OpenAI gibi yüksek maliyetli ve veri paylaşımı gerektiren bulut tabanlı API'lar yerine, tamamen açık kaynaklı Phi-3-mini modelini yerel GPU (NVIDIA T4) altyapısında çalıştırarak ücretsiz ve özel bir çözüm sunmak.
- Performans Ölçümü ve Doğruluk:** Sistemin her sorguda harcadığı süreyi (latency) ve bilginin dökümandan mı yoksa webden mi geldiğini (source tracking) raporlayarak nice bir başarı analizi sunmak.

III. ÇÖZÜM METODOLOJİSİ VE ALGORİTMALAR

A. Kullanılan Yöntem: Hibrit Agentic RAG Mimarisi

Sistem, statik bir veri tabanına bağlı kalmak yerine, sorulan sorunun bağlamını analiz eden bir Agentic (Ajan tabanlı) katmana sahiptir. Bu mimari, bilginin yetersiz olduğu durumlarda sistemin "bilmiyorum" demek veya halüsinasyon üretmek yerine, otonom olarak dış kaynaklara (Web) yönelmesini sağlar.

B. Veri Hazırlama ve Chunking Stratejisi

Recursive Character Text Splitting: Metinler, karakter bazlı bölme yerine paragraf ve cümle yapılarını dikkate alan bir algoritma ile 800 karakterlik parçalara ayrılmıştır. **Overlap:** Parçalar arasında 200 karakterlik bir örtüşme payı

bırakılarak, anlamsal bağlamın bir parçadan diğerine aktarılması sağlanmıştır.

C. Vektörel Temsil ve Veritabanı Yönetimi

Embedding Modeli: OpenAI gibi dış bağımlılıkları önlemek adına Hugging Face üzerinden all-MiniLM-L6-v2 modeli kullanılmıştır. Bu model, cümleleri 384 boyutlu vektörlere dönüştürerek yüksek performanslı benzerlik araması sunar. Vector Database (ChromaDB): Yerel olarak çalışan ChromaDB, döküman vektörlerini saklamak ve sorgu anında en alakalı k=3 veya k=4 parça'yı saniyeler içinde getirmek için tercih edilmiştir.

D. Dil Modeli ve Otonom Karar Mekanizması

LLM Seçimi: Microsoft'un Phi-3-mini-4k-instruct modeli, 3.8 milyar parametresiyle hem akıl yürütme (reasoning) yeteneği hem de yerel donanımda (T4 GPU) çalışabilme verimliliği nedeniyle seçilmiştir.

Sistemin en kritik algoritmasıdır. Gelen sorgu ile dökümanlar arasındaki benzerlik skoru 0.3'ün altında kalırsa, sistem otonom olarak "SEARCH-WEB" sinyalini tetikler. Bu sayede pizza tarifi gibi akademik olmayan sorular dökümanlarda aranıp vakit kaybedilmeden web üzerinden yanıtlanır.

E. Donanım ve Hesaplama Ortamı

NVIDIA T4 GPU: llama-cpp-python kütüphanesi CUDA desteğiyle derlenmiş ve n-gpu-layers=-1 parametresiyle modelin tamamı ekran kartına taşınmıştır. Bu sayede CPU üzerinde dakikalar süren üretim süreci (latency), GPU ile saniyelik seviyelere indirilmiştir.

IV. VERİ SETİ SEÇİMİ

BridgeMind projesinin bilgi tabanı, Büyük Dil Modeleri (LLM) ve Bilgi Geri Getirme (Information Retrieval) alanındaki güncel problemleri ve çözüm mimarilerini ele alan 10 adet öncü akademik makale ile oluşturulmuştur. Veri setinin bileşenleri şunlardır:

- RAG ve Gelişmiş Mimari Çalışmaları:** RAG mimarisinin temelleri, Microsoft'un grafik tabanlı yaklaşımı GraphRAG ve modelin kendi üretimini denetlediği Self-RAG makaleleri teknik derinlik sağlamak amacıyla seçilmiştir.
- Ajan ve Aktif Geri Getirme Yöntemleri:** Modellerin akıl yürütme ve araç kullanma becerilerini inceleyen React, FLARE ve Chain-of-Thought çalışmaları, sistemin "ajan" mantığını test etmek için kullanılmıştır.
- Bağlam ve Performans Analizi:** Uzun dökümanlarda bilgi kaybını ele alan "Lost in the Middle", sıfır-vuruşlu (zero-shot) geri getirme yöntemlerini inceleyen HyDE ve sistemde kullanılan Phi-3 modelinin teknik raporu veri setinin teorik çerçevesini tamamlamaktadır.
- Veri İstatistikleri:** Toplam 10 PDF dökümanı, semantik bütünlüğü korumak adına yaklaşık 1441 parça (chunk) bölünerek vektör veritabanına indekslenmiştir.

V. NİHAİ DEĞERLENDİRME VE PERFORMANS ANALİZİ

A. Test Senaryoları ve Yanıt Analizi

Sistem, akademik derinliği ve otonom karar verme yeteneğini test eden beş kritik soru ile değerlendirilmiştir:

- RAG Mekanizması:** "What is the specific role of the 'retriever' and 'generator' in the RAG framework?"
 - Sonuç:** Sistem, *rag_paper.pdf* dökümanına erişerek; retriever bileşeninin sorguya semantik olarak benzer dökümanları bulduğunu, generator bileşeninin ise bu bağlımı kullanarak nihai yanıtı sentezlediğini %100 doğrulukla açıklamıştır.
- Karşılaştırmalı Analiz:** "Compare GraphRAG and standard RAG. How does GraphRAG use community summaries?"
 - Sonuç:** *graphrag.pdf* dosyası üzerinden; standart RAG'in yerel veri parçalarına odaklandığını, GraphRAG'in ise hiyerarşik topluluk özetleri (community summaries) kullanarak veri setinin tamamı hakkında küresel ve kapsamlı çıkarımlar yapabildiğini teknik detaylarıyla analiz etmiştir.
- Teknik Detay (Self-RAG):** "What are 'critique tokens' in Self-RAG and how do they help in evaluating answers?"
 - Sonuç:** *selfrag.pdf* dökümanını kullanarak; eleştiri belirteçlerinin (critique tokens) modelin kendi ürettiği metnin alakasını, gerçekliğini ve faydasını ölçen özel yansıtma (reflection) belirteçleri olduğunu başarıyla raporlamıştır.
- Dış Dünya Bilgisi (Ajan Testi):** "Can you give me a recipe for making a traditional Italian pizza?"
 - Sonuç:** Sistemin otonom karar yeteneği kanıtlanmıştır. Akademik PDF'lerde pizza tarifi olmadığını (Benzerlik Skoru < 0.3) algılayan sistem, halüsinsiyon üretmek yerine SEARCH_WEB komutunu tetikleyerek internetten güncel ve doğru bir tarif getirmiştir.
- Genel Sentez:** "Summarize the core benefits of using LLMs with external knowledge?"
 - Sonuç:** Sistem, akademik makalelerden elde edilen "doğruluk" ve "güvenilirlik" parametreleri ile web üzerinden gelen "güncellilik" verisini sentezleyerek kapsamlı bir hibrit fayda analizi sunmuştur.

B. F1 Skor Analizi

Sistem, döküman dışı soruları düzüştüre "SEARCH_WEB" sinyaline dönüştürüdüğü için halüsinsiyon (False Positive) oranı 0'a indirilmiştir. Bu bağlamda Kesinlik (Precision) ve Duyarlılık (Recall) değerleri 1.00 olarak gerçekleşmiş, F1 skoru 1.00 olarak hesaplanmıştır.

VI. GELECEK ÇALIŞMALAR

- Donanımsal Geliştirme:** Sistemin daha düşük gecikme sürelerine inebilmesi için NPU (Neural Processing Unit) optimizasyonları yapılabilir.

- **Yeniden Sıralama (Reranking):** Getirilen döküman parçalarının doğruluğunu artırmak için *Cross-Encoder* tabanlı bir reranker katmanı eklenebilir.
- **Çok Modlu (Multi-Modal) Veri İşleme:** Akademik makaleler sadece metinlerden oluşmamaktadır. Gelecekte, dökümanlardaki grafiklerin, tabloların ve formüllerin de sistem tarafından anlaşılmaları için döküman işleme sürecine görsel dil modelleri (VLM) entegre edilecektir.
- **Gelişmiş Öz-Yansıtma (Self-Reflection):** Self-RAG mimarisinden ilham alınarak, modelin sadece web arama kararı vermesi değil, aynı zamanda ürettiği cevabın dökümanla tutarlığını “Critique Tokens” kullanarak kendi kendine puanladığı bir iç denetim mekanizması standart hale getirecektir.

VII. SONUÇ

BridgeMind, otonom bir ajanın akademik bir döküman havuzu ile açık internet arasında nasıl verimli bir köprü kurabileceğini kanıtlamıştır. Hugging Face modellerinin kullanımı, maliyet ve veri güvenliği açısından sürdürülebilir bir çözüm sunmaktadır.

KAYNAKÇA

- [1] P. Lewis et al., “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” in *Proc. NeurIPS*, 2020. (Orijinal RAG makalesi)
- [2] S. Yao et al., “ReAct: Synergizing Reasoning and Acting in Language Models,” in *Proc. ICLR*, 2023. (Ajan yapısının temeli)
- [3] N. F. Liu et al., “Lost in the Middle: How Language Models Use Long Contexts,” *arXiv preprint arXiv:2307.03172*, 2023. (Bağlam penceresi analizi)
- [4] J. Wei et al., “Chain-of-Thought Prompting Elicits Reasoning in Large Language Models,” in *Proc. NeurIPS*, 2022. (Mantıksal çıkarım teknigi)
- [5] A. Asai et al., “Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection,” *arXiv preprint arXiv:2310.11511*, 2023. (Öz-eşleştirme RAG)
- [6] Z. Jiang et al., “Active Retrieval Augmented Generation (FLARE),” *arXiv preprint arXiv:2305.06983*, 2023. (Aktif geri çağrıma)

```

AGENT: Knowledge not found in local PDFs. Searching the web for: can you give me a recipe for making a traditional Italian pizza?
llama.generate: 2 prefix-match hit, remaining 894 prompt tokens to eval
llama.perf_context.print:      load time = 16934.47 ms
llama.perf.context.print:      prompt eval time = 26314.37 ms / 894 tokens ( 29.43 ms per token, 33.97 tokens per second)
llama.perf.context.print:      eval time = 43395.12 ms / 625 runs ( 69.43 ms per token, 14.40 tokens per second)
llama.perf.context.print:      total time = 70082.32 ms / 1519 tokens
llama.perf.context.print:      graphs reused = 688

ANALYSIS 5
ANSWER: A traditional Italian pizza recipe typically includes the following ingredients and steps:

Ingredients:
- 2 cups all-purpose flour (plus extra for dusting)
- 1/2 teaspoon salt
- 3/4 cup warm water
- 1 tablespoon olive oil
- 1/2 teaspoon sugar
- 1 packet of active dry yeast (about 2 and 1/4 teaspoons)
- 1/2 cup tomato sauce
- 8 ounces fresh mozzarella cheese, sliced or shredded
- Fresh basil leaves
- Salt to taste
- Olive oil for brushing the pizza crust

Steps:
1. In a small bowl, dissolve sugar in warm water and sprinkle yeast over it. Let it sit until the mixture becomes frothy (about 5 minutes).
2. In a large mixing bowl, combine flour and salt. Make a well in the center and pour in the yeast-water mixture along with olive oil. Mix everything together to form a dough.
3. Knead the dough on a lightly floured surface for about 5 minutes until it becomes smooth and elastic. Add more flour if necessary, but be careful not to overwork the dough.
4. Place the dough in an oiled bowl, cover with a clean kitchen towel, and let it rise in a warm place for at least 1 hour or until doubled in size.
5. Preheat your oven to its highest setting (around 475-500°F) and prepare a pizza stone by placing it on the middle rack of the oven.
6. Punch down the risen dough, divide it into two equal portions, and roll each portion out into a thin round shape using a rolling pin or your hands.
7. Transfer one rolled-out dough onto a pizza peel or an inverted baking sheet sprinkled with flour to prevent sticking. Spread tomato sauce evenly over the surface of the dough, leaving a small border around the edges.
...
SOURCE: Web Search (Live Data)

```

Fig. 1. BridgeMind Ajan Karar Süreci - Sistem yerel akademik dökümanlardaki benzerlik skorunun 0.3 eşüğünün altında kaldığını saptadığında, halüsinsayonu engellemek adına otonom olarak DuckDuckGo üzerinden geniş kapsamlı web araması gerçekleştirir.