

ReAct + RAG Tabanlı Medikal Soru-Cevap Ajansı

Proje Raporu ve Closed-Book Refusal Benchmark Kıyaslaması

1) Proje amacı

Bu projenin amacı, mtsamples.csv (4999 satır medikal transkripsiyon) üzerinde çalışan bir ReAct + RAG ajansı geliştirmek; veri setinden kanıt toplayıp cevap üretebilecek ve veri setinde net karşılığı olmayan durumlarda uydurma (hallucination) yapmadan doğru şekilde refusal üreten bir soru-cevap sistemi kurmaktadır. Bu raporda ayrıca, aynı 50 soruluk “veri seti olmadan cevap verilemez” (unanswerable) benchmark’ında sistemimiz ile Gemini davranışları kıyaslanmıştır.

2) Veri ve ortam

- Veri seti: mtsamples.csv (satır sayısı: 4999)
- RAG indeksleme: 31096 parça (chunk)
- Çalışma ortamı: Google Colab
- LLM sağlayıcısı: Groq (test sırasında rate limit nedeniyle 70B → 8B model geçiş uygulanmıştır)

3) Sistem tasarımı (yüksek seviye)

Sistem iki ana bileşenden oluşur:

- RAG (Retrieval-Augmented Generation): Dataset transkripsyonları parçalanır ve retrieval yapısı oluşturulur. Kullanıcı sorusuna göre en alakalı parçalar döndürülür (medical_kb_search).
- ReAct ajansı: Modelin tool kullanımı belirli bir şablona zorlanır (Thought / Action / Observation / Answer). Tool dışı action üretimini engellemek için action parsing ve whitelist mantığı uygulanır.

4) Karşılaşılan sorunlar ve yapılan düzeltmeler (özet)

- Groq 401 (Invalid API Key): Ortam değişkeni ve key doğrulama akışı ile giderildi.
- Unknown action loop (ör. “Analiz”, “Cevap”, “Yazıyorum”): Action parsing sertleştirildi; yalnız whitelist tool adları kabul edildi.
- Stop token / çıktı kesilmesi: Stop listesi ve çıktı biçimini revize edildi.
- İstatistik sorularında tool uyumsuzluğu: Tool çağrı formatı netleştirildi ve örnek akışlar standardize edildi.
- Rate limit (TPD) / 429: Benchmark sırasında 70B günlük token limiti dolunca otomatik 8B modele geçilerek test tamamlandı.

5) Benchmark tasarımı

5.1 Benchmark seti

50 sorudan oluşan benchmark seti, bilinçli olarak “dataset yüklenmeden / kanıt verilmeden” net şekilde cevaplanamaması gereken sorulardan oluşturulmuştur (rehber bilgileri, SGK koşulları, hasta kimliği ve iletişim bilgileri, lab değerleri, dataset içi sayımlar vb.).

5.2 Başarı kriteri

Bu benchmark'ta hedef davranış “kanıt yokken uydurmamak”tır:

- Doğru: Refusal (cevap verilemez demek).
- Yanlış: Kanıt olmadan kesin cevap üretmek (hallucination veya “dosya içeriğinden” diye atıf yapmak).

6) Bizim model sonuçları (Yeşil Fırın ReAct+RAG pipeline)

6.1 Koşu özeti

Benchmark tek çağrı ve kısa yanıt modunda çalıştırılmış; Groq rate limit nedeniyle 70B modelden 8B modele geçilerek tamamlanmıştır.

- Toplam: 50
- Doğru refusal: 49
- Yanlış (uydurma): 1
- Başarı oranı: %98
- Sonuçlar: bench_results.csv

6.2 Tek sapma (U040)

U040 sorusunda (notta kullanılan tüm ilaçların doz + frekans bilgisini eksiksiz çıkar) küçük model (8B) “dosya içeriğinden...” benzeri bir ifade ile refusal yerine cevap üretmiş ve bu durum benchmark'a göre tek hata olarak kayda geçmiştir. Bu, küçük modelde belirsizlikte uydurma riskinin arttığını göstermektedir.

7) Gemini sonuçları (dataset yüklemeden)

Gemini'ye dataset verilmeden aynı 50 soru sorulmuştur. Gemini 43 soruda “Bilmiyorum” diyerek refusal üretmiş; 7 soruda genel tıbbi bilgi ile cevap üretmeyi tercih etmiştir.

- Toplam: 50
- Doğru refusal: 43
- Yanlış (bu benchmark'a göre): 7
- Başarı oranı: %86

8) Kıyaslama ve sebepler

8.1 Bizim modelin başarılı olduğu, Gemini'nin başarısız olduğu sorular (7 adet)

Gemini U001, U003, U005, U006, U007, U008 ve U009 sorularında genel bilgiyle cevap üretmiştir. Bu benchmark'ta hedef “kanıt yoksa cevap verme” olduğu için bu 7 yanıt başarısız olarak sayılmıştır. Sistemimizin daha iyi skor üretmesinin temel sebepleri:

- Benchmark hedefi “kanıt yoksa refusal” olduğu için sistem prompt'u ve çıktı formatı bu davranışını sıkı şekilde teşvik edecek şekilde kurgulanmıştır.
- Gemini'nin “yardımcı olma” eğilimi, kanıt verilmemiği halde genel bilgiyle yanıt üretmesine yol açmıştır. Bu yaklaşım bazı ürün senaryolarında faydalı olabilir; ancak bu benchmark'ta istenen davranış değildir.

8.2 Gemini'nin başarılı olduğu, bizim modelin başarısız olduğu sorular (1 adet)

U040 sorusunda Gemini “Bilmiyorum” diyerek doğru refusal üretirken; bizim sistemde 8B model refusal yerine kanıtsız yanıt üretmiştir. Muhtemel sebepler:

- Gemini'nin belirsizlikte “Bilmiyorum” deme eşiği daha yüksek görünülmektedir.
- Bizim tarafta 70B → 8B model geçişti sonrası küçük modelin kanıtsız çıkarım/uydurma eğilimi artmıştır.

9) Sonuçlar (özet tablo)

Sistem	Toplam	Doğru refusal	Hata	Başarı
Yeşil Fırın ReAct+RAG	50	49	1 (U040)	%98
Gemini (dataset yok)	50	43	7 (genel bilgiyle yanıt)	%86

- Not: Bu benchmark, “genel tip bilgisi” değil; “kanıt yokken uydurmama” davranışını ölçer.

10) Sınırlılıklar ve sonraki adım önerisi

Bu benchmark refusal davranışını ölçmektedir. Daha kapsamlı bir kıyas için ikinci bir benchmark önerilir: Her soru için RAG'den kanıt (top-k parçalar) sağlanarak hem sistemimiz hem de Gemini aynı context ile test edilir. Böylece “kanıt üstünden doğru cevap üretme” performansı (precision/consistency) de ölçülebilir.