

Doktora Tez İzleme Komitesi Toplantısı

Bu toplantı, doktora tez çalışmalarımın ilerlemesini değerlendirmek ve akademik hedeflerime ulaşmamı desteklemek amacıyla düzenlenmiştir. Tez çalışmam, **sağlık hizmetleri** ve **hukuk** gibi kritik alanlarda özelleştirilmiş dil modellerinin geliştirilmesi ve değerlendirilmesi üzerine odaklanmaktadır.

Giriş

Tezimin temel amacı, büyük dil modellerini sağlık hizmetleri ve hukuk alanlarına adapte ederek, bu alanlardaki verimliliği ve doğruluğu artırmaktır. Çalışma kapsamında:

- Alana özgü **veri setleri oluşturma**,
- Model adaptasyon yöntemlerini uygulama ve karşılaştırma**,
- Farklı alanlardaki zorlukları ve fırsatları belirleme,
- Özelleştirilmiş modellerin performansını analiz etme konularına odaklanıyorum.

Literatür Taraması

Araştırmamın temeli, sağlık hizmetleri ve hukuk alanlarında dil modellerinin uygulanabilirliği üzerine yapılan çalışmalar ve model adaptasyonu yöntemlerinin değerlendirilmesine dayanmaktadır. Ele alınan başlıca konular şunlardır:

- Sağlık hizmetlerinde dil modellerinin kullanımı:**
 - Hasta-doktor etkileşimlerinin analizi ve tıbbi kayıtların özetlenmesi.
- Hukuk alanında dil modellerinin kullanımı:**
 - Hukuki belgelerin analizi ve risk faktörlerinin belirlenmesi.
- Model adaptasyonu yöntemleri:**
 - LoRA (Low-Rank Adaptation) ve SLerp (Spherical Linear Interpolation) tekniklerinin etkinliği.
- Performans değerlendirme ölçütleri**

Metodoloji

Çalışmamda izlediğim metodoloji, sağlam ve tekrarlanabilir bir çerçeve sunmaktadır:

- Veri toplama ve ön işleme:**
 - Sağlık ve hukuk alanlarından büyük ölçekli veri kümelerinin toplanması ve temizlenmesi.
- Model adaptasyonu yöntemi seçimi:**
 - Veri setine uygun hiperparametrelerin belirlenmesi ve LoRA/SLerp tekniklerinin uygulanması.
- Model eğitimi ve değerlendirme:**
 - LLAMA ve Gemma gibi önceden eğitilmiş modellerin fine-tuning ile özelleştirilmesi.
- Performans analizi ve sonuçların yorumlanması:**
 - Modellerin farklı görevlerdeki başarısının ölçülmesi ve hata analizi.

Uygulama Alanları

Sağlık Alanı

- **Tıbbi metin analizi:** Hasta kayıtlarının incelenmesi ve risk faktörlerinin belirlenmesi.
- **Tıbbi chatbotlar:** Hasta-doktor etkileşimlerini simüle eden yapay zeka çözümleri.
- **Tıbbi bilgi özetleme:** Hekimlere hızlı bilgi erişimi sağlayan sistemler.

Hukuk Alanı

- **Hukuki metin analizi:** Sözleşme ve yasal belgelerin incelenmesi.
- **Hukuki chatbotlar:** Kullanıcıları yasal süreçler hakkında bilgilendiren sistemler.
- **Risk analizi:** Hukuki belgelerdeki potansiyel risklerin tespiti.

Sonuçlar ve Tartışma

Tez çalışmamın bulguları, özelleştirilmiş dil modellerinin sağlık hizmetleri ve hukuk alanlarındaki etkinliğini değerlendirmeyi amaçlamaktadır. Tartışma konuları şunlardır:

- **Performans analizi:** Modellerin doğruluk ve etkinlik açısından değerlendirilmesi.
- **Zorluklar ve fırsatlar:** Model adaptasyonu sürecinde karşılaşılan sorunlar ve çözüm önerileri.
- **Gelecek yönlendirmeleri:** Dil modellerinin bu alanlardaki gelecekteki gelişimi için öneriler.

Yayınlar ve Katkıları

Yayınlar

1. **Healthcare-Focused Turkish LLM:**
 - **167.000+ hasta-doktor soru-cevap verisi** ile eğitilmiş, LLAMA 3 tabanlı özelleştirilmiş model.
2. **Turkish MMLU Benchmark:**
 - **Türkçe NLP için standart bir değerlendirme kriteri** oluşturan veri seti.
3. **Data Quality-Based Adaptive Learning Rate:**
 - Veri kalitesine dayalı adaptif öğrenme oranı stratejisi.

Modeller ve Veri Setleri

Geliştirilen Modeller

- **Doktor-Llama Serisi:**
 - Tıbbi terminolojiye odaklanmış, **SLerp optimizasyonu** ile geliştirilmiş modeller.
 - Hasta-doktor etkileşimlerini daha iyi anlamlandırmak için Türkçe sağlık terminolojisi ile fine-tune edilmiştir.
- **DoktorGemma Serisi:**
 - Farklı fine-tuning aşamaları ile özelleştirilmiş modeller.
 - **LoRA adaptörleri** kullanılarak modüler bir yapı sağlanmıştır, bu da farklı tıbbi ve hukuki alt alanlara kolayca adapte edilebilmesini mümkün kılar.
- **Turkish Tokenizer Çalışmaları**
 - Türkçe diline özgü **morfolojik** ve **anlamsal** özellikleri analiz ederek daha etkili tokenizasyon sunmayı hedeflemektedir.

- Tokenizer, kök ve ek ayrıştırmasını optimize ederek anlam bütünlüğünü korur.

Veri Setleri

- **Turkish MMLU:**
 - 67 disiplin, 800+ konu, 280.000+ soru havuzundan seçilmiş, Türkçe dil modellerinin performansını ölçmek için kapsamlı bir benchmark.
- **Tıbbi Veri Seti:**
 - Gerçek hasta-doktor etkileşimlerini içeren **167.000+ kayıt**, tıbbi dil modellerinin eğitimi ve analizi için kullanılmıştır.
- **Hukuk Veri Seti:**
 - Türk hukuku üzerine **soru-cevap içerikleri** ve yasal belgelerden elde edilmiş veriler, hukuk alanındaki dil modellerinin değerlendirilmesi için hazırlanmıştır.
- **E-ticaret, Medya ve Haber Yorumları:**
 - Onedio, Hepsiburada, Beyazperde, Kitapyurdu gibi platformlardan toplanan duygu durumu etiketli yorumlar.

Açık Kaynak Katkıları

- **mlx-examples:** ML Apple Silicon optimizasyonu.
- **unsloth:** LLM optimizasyon araçları ve eğitim süreci hızlandırıcılar.

Gelecek Çalışmalar

- **Son kullanıcı odaklı ürünler:** RAG (Retrieval-Augmented Generation) ve agentic sistemlerin geliştirilmesi.

İletişim

- **E-posta:** malibayram20@gmail.com
- **Hugging Face:** @alibayram
- **GitHub:** @malibayram