

# Qwen2.5-Coder LoRA Fine-Tuning Proje Raporu

Öğrenci Adı: Çağla DEMİR

Öğrenci No: 2020556018

Tarih: 29 Kasım 2024

Konu: Large Language Model (LLM) Fine-Tuning for Competitive Code Reasoning

## 1. Proje Özeti

Bu projenin amacı, **Qwen2.5-Coder-1.5B-Instruct** temel modelinin kodlama ve mantıksal muhakeme (reasoning) yeteneklerini, **LoRA (Low-Rank Adaptation)** tekniği kullanarak geliştirmektir. Proje kapsamında model, iki farklı veri seti (**DEEP** ve **DIVERSE**) üzerinde eğitilmiş ve sonuçlar Hugging Face platformunda yayınlanmıştır.

Eğitim sürecinde Google Colab (T4 GPU) donanım kısıtlamaları göz önünde bulundurularak hafıza optimizasyonu sağlayan teknikler (Gradient Checkpointing, Quantization) uygulanmıştır.

## 2. Kullanılan Model ve Veri Setleri

- Base Model:** Qwen/Qwen2.5-Coder-1.5B-Instruct
  - Seçim Nedeni: Kodlama görevlerinde yüksek performans göstermesi ve 1.5B parametre boyutuyla T4 GPU üzerinde eğitilebilir olması.
- Datasets:**
  - DEEP Dataset:** Karmaşık mantık yürütme (reasoning trace) gerektiren kodlama problemleri.
  - DIVERSE Dataset:** Çeşitli konu başlıklarını kapsayan geniş yelpazeli kodlama problemleri.

## 3. Eğitim Konfigürasyonu (Hyperparameters)

Eğitim sırasında dokümanda önerilen ve donanım kısıtlarına uygun aşağıdaki hiperparametreler kullanılmıştır:

Parametre	Değer	Açıklama
LoRA Rank (r)	16	Parametre verimliliği için seçildi.
LoRA Alpha	32	$r * 2$ kuralına uygun olarak belirlendi.
Target Modules	All Linear	q_proj, k_proj, v_proj, o_proj, gate_proj, up_proj, down_proj .
Learning Rate	2e-4	Kararlı bir öğrenme süreci için seçildi.
Batch Size	1	GPU belleğini (VRAM) aşmamak için düşürüldü.
Gradient Accumulation	16	Efektif batch size'ı 16'ya tamamlamak için kullanıldı.
Context Length	1024	Sadece çözüm (code-only) eğitimi yapıldığı için yeterli görüldü.
Precision	FP16	Mixed Precision eğitimi yapıldı.
Optimizer	AdamW	Paged AdamW 32bit.

## 4. Karşılaşılan Zorluklar ve Çözümler

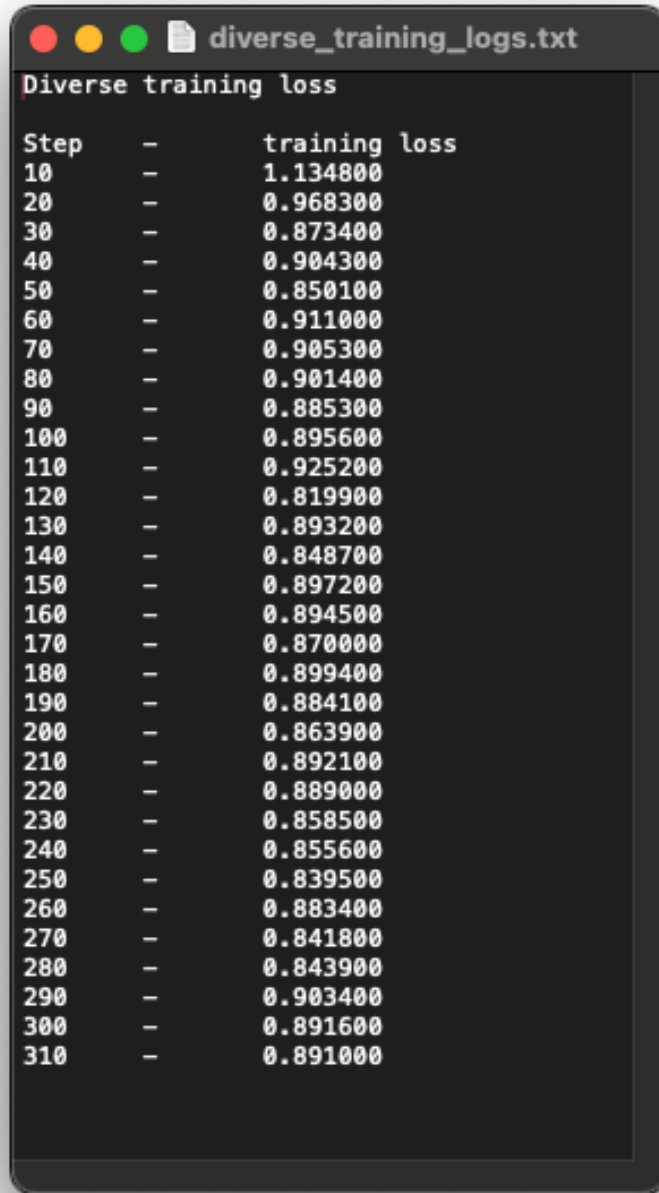
Proje geliştirme sürecinde, dokümanda belirtilen "Sık Karşılaşılan Sorunlar" yaşanmış ve şu çözümler uygulanmıştır:

- **CUDA Out of Memory (OOM) Hatası:** Google Colab T4 GPU'nun 16GB VRAM kapasitesi, standart eğitim ayarları için yetersiz kalmıştır.
  - **Çözüm:** per\_device\_train\_batch\_size 1'e düşürülmüş ve gradient\_accumulation\_steps 16'ya çıkarılarak bellek kullanımı dengelenmiştir. Ayrıca gradient\_checkpointing=True yapılarak aktivasyonların bellekte tutulması engellenmiş, işlemci yükü artırılarak bellekten tasarruf edilmiştir<sup>13</sup>.
- **Eğitim Takibi:** Google Colab oturum süreleri kısıtlı olduğu için modellerin kaybolma riski oluşmuştur.
  - **Çözüm:** Eğitim çıktıları (output\_dir) doğrudan Google Drive'a yönlendirilmiş, böylece olası bağlantı kopmalarında modelin kaybolması engellenmiştir.

## 5. Eğitim Sonuçları ve Kayıtlar (Logs)

Eğitim boyunca Loss (Kayıp) değerlerinin düzenli olarak düştüğü gözlemlenmiştir.

deep_training_logs.txt		
deep training loss		
Step	-	training loss
10	-	1.092200
20	-	0.962500
30	-	0.875100
40	-	0.861700
50	-	0.835900
60	-	0.889100
70	-	0.871300
80	-	0.901800
90	-	0.855500
100	-	0.838700
110	-	0.784800
120	-	0.809100
130	-	0.816300
140	-	0.813300
150	-	0.782200
160	-	0.781300
170	-	0.755600
180	-	0.779100
190	-	0.736400
200	-	0.737700
210	-	0.755700
220	-	0.719200
230	-	0.709900
240	-	0.728600
250	-	0.733500
260	-	0.694100
270	-	0.688100
280	-	0.730300
290	-	0.722000
300	-	0.679200
310	-	0.702100



```
diverse_training_logs.txt
Diverse training loss

Step      -      training loss
10         -      1.134800
20         -      0.968300
30         -      0.873400
40         -      0.904300
50         -      0.850100
60         -      0.911000
70         -      0.905300
80         -      0.901400
90         -      0.885300
100        -      0.895600
110        -      0.925200
120        -      0.819900
130        -      0.893200
140        -      0.848700
150        -      0.897200
160        -      0.894500
170        -      0.870000
180        -      0.899400
190        -      0.884100
200        -      0.863900
210        -      0.892100
220        -      0.889000
230        -      0.858500
240        -      0.855600
250        -      0.839500
260        -      0.883400
270        -      0.841800
280        -      0.843900
290        -      0.903400
300        -      0.891600
310        -      0.891000
```

Eğitim sonunda, Loss değerinin en düşük olduğu **Final Checkpoint**, en başarılı model olarak kabul edilmiştir.

## 6. Proje Teslimatları

Eğitilen modeller ve kaynak kodlar aşağıdaki bağlantılarda erişime açıktır:

- GitHub Repository (Kodlar ve Rapor):

<https://github.com/caglademir/Qwen2.5-LoRA-FineTuning>

- Hugging Face Model - DEEP:

- <https://huggingface.co/datasets/Naholav/CodeGen-Deep-5K>

- Hugging Face Model - DIVERSE:

- <https://huggingface.co/datasets/Naholav/CodeGen-Diverse-5K>

---

**Not:** Benchmark değerlendirmesi, test veri seti yayınlandığında gerçekleştirilecek ve sonuçlar bu rapora eklenecektir.