

<KATEGORİ>

Türkçe Doğal Dil İşleme

8 - 9 Ağustos 2024

< cyph[ai]>

<EKİBİMİZ>



KAPTAN

ŞEYMA KERKÜKLÜ

- Proje Yönetimi
- İş Geliştirme
- ML
- NLP



EKİP ÜYESİ

KARDELEN GEÇKİN

- Sibergüvenlik
- ML
- Developer

<PROJENİN TANIMI>

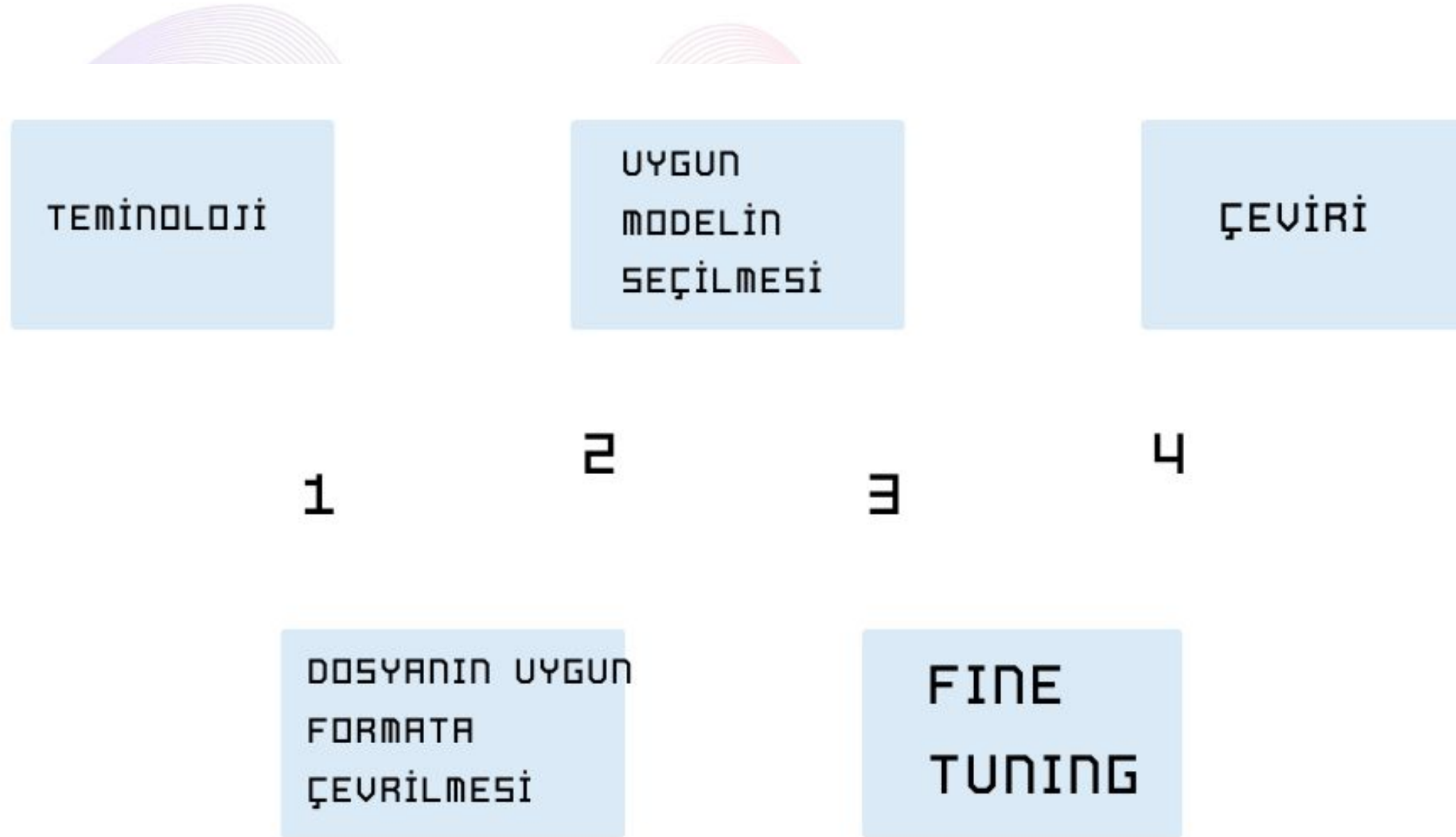


<PROJENİN SAĞLADIĞI ÇÖZÜM>

- ENDÜSTRİ RAPORLARI
- GÜVENLİK RAPORLARI
- OLAY RAPORLARI
- SALDIRI RAPORLARI
- HACKING FORUMLARI
- DARK WEB
- HABER KAYNAKLARI

- Banka ve Kamu Kurumları
- Finans ve Yazılım Hizmeti Sunan Kuruluşlar
- Güvenlik Operasyon Merkezleri
- Tehdit İstihbaratı Analistleri

<PROJE İŞ AKIŞI>



<VERİ SETİ>



- CSV veri setini yükle
- Veri setini eğitim ve doğrulamaya böl
- Tokenizer ve modeli yükle
- Ön işlemeyi uygula
- Eğiticiyi Başlat
- Çevir

<YÖNTEM VE TEKNİKLER>

1. LoRA: Optimizasyon
2. T5: Çeviri
3. mBERT: Çok dilli
4. Tokenization: Parçalama
5. Beam Search: Arama
6. Fine-Tuning: İnce ayarlama

<MODEL EĞİTİMİ VE DEĞERLENDİRME>

Eğitim:

- Veri:** İngilizce-Türkçe çeviri çiftleri
- Yöntem:** Fine-tuning, LoRA
- Model:** mBERT veya mBART

<SONUÇLAR>



<PROJE YOL HARİTASI>



<DEMO VIDEO>

```
[ ] import os
import torch
from transformers import MBartForConditionalGeneration, MBartTokenizer, Seq2SeqTrainer
import pandas as pd
from datasets import Dataset, DatasetDict

df = pd.read_csv('/content/Siber_Guvenlik_Terim_Karsiliklari_duzenlenmis.csv', delimiter=';')

dataset = Dataset.from_pandas(df[['English Term', 'Turkish Term']])

split_dataset = dataset.train_test_split(test_size=0.1)
tokenized_datasets = DatasetDict({
    'train': split_dataset['train'],
    'validation': split_dataset['validation']
})
```

<https://drive.google.com/file/d/1KnY4I3bnwEtbSfapHOcn21J9GUoA539b/view?usp=sharing>

TEŞEKKÜRLER