



# KAYSERIOUS Ekibi



**ABDÜSSAMET SÖKEL**

Eskişehir Teknik Üniversitesi / Endüstri Mühendisliği / Lisans / 2017 - 2022

Eskişehir Osmangazi Üniversitesi / Endüstri Mühendisliği / Yüksek Lisans / 2023 - Halen

Veri Bilimci @ MADAME COCO / 2022 - Halen

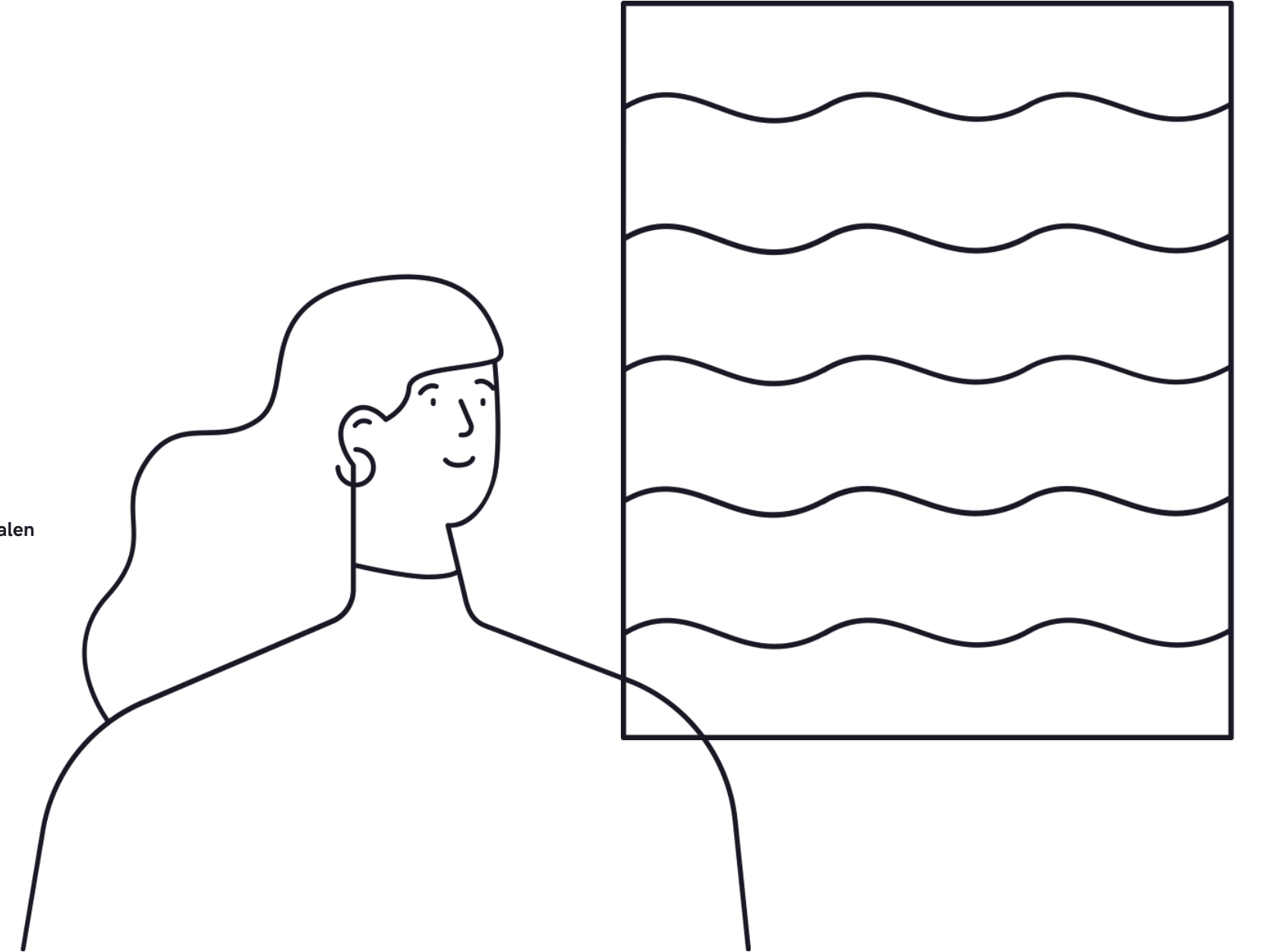


**SUDENAZ AKGÜN**

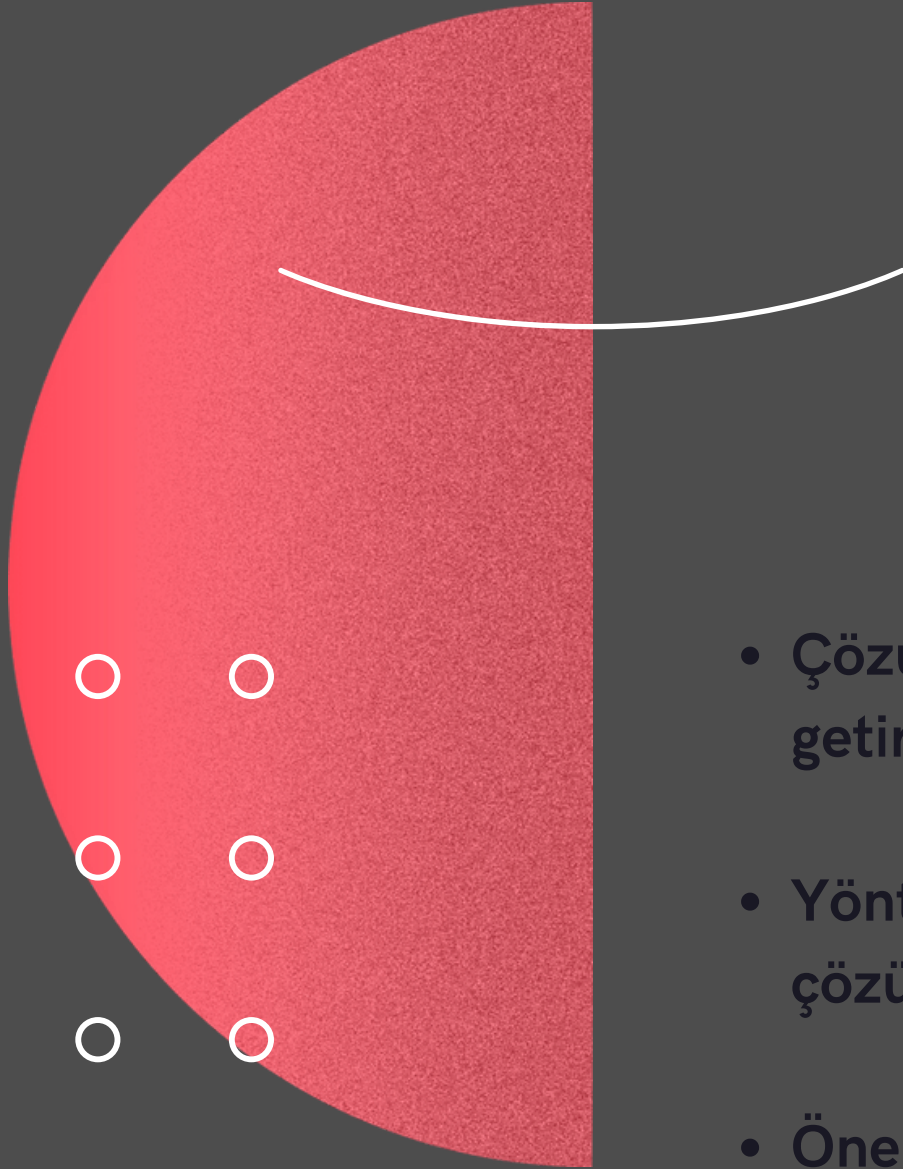
Eskişehir Teknik Üniversitesi / Endüstri Mühendisliği / Lisans / 2022 - Halen

Ekip olarak en büyük motivasyonumuz çalışma alanımızda ülkemize ana dilde literatür kazandırmaktır.

Bu sebeple aktif olarak çalışmalarımızın yanında akademik yayınlar üretmeye de ağırlık vermekteyiz.



# Bireysel Katkılar

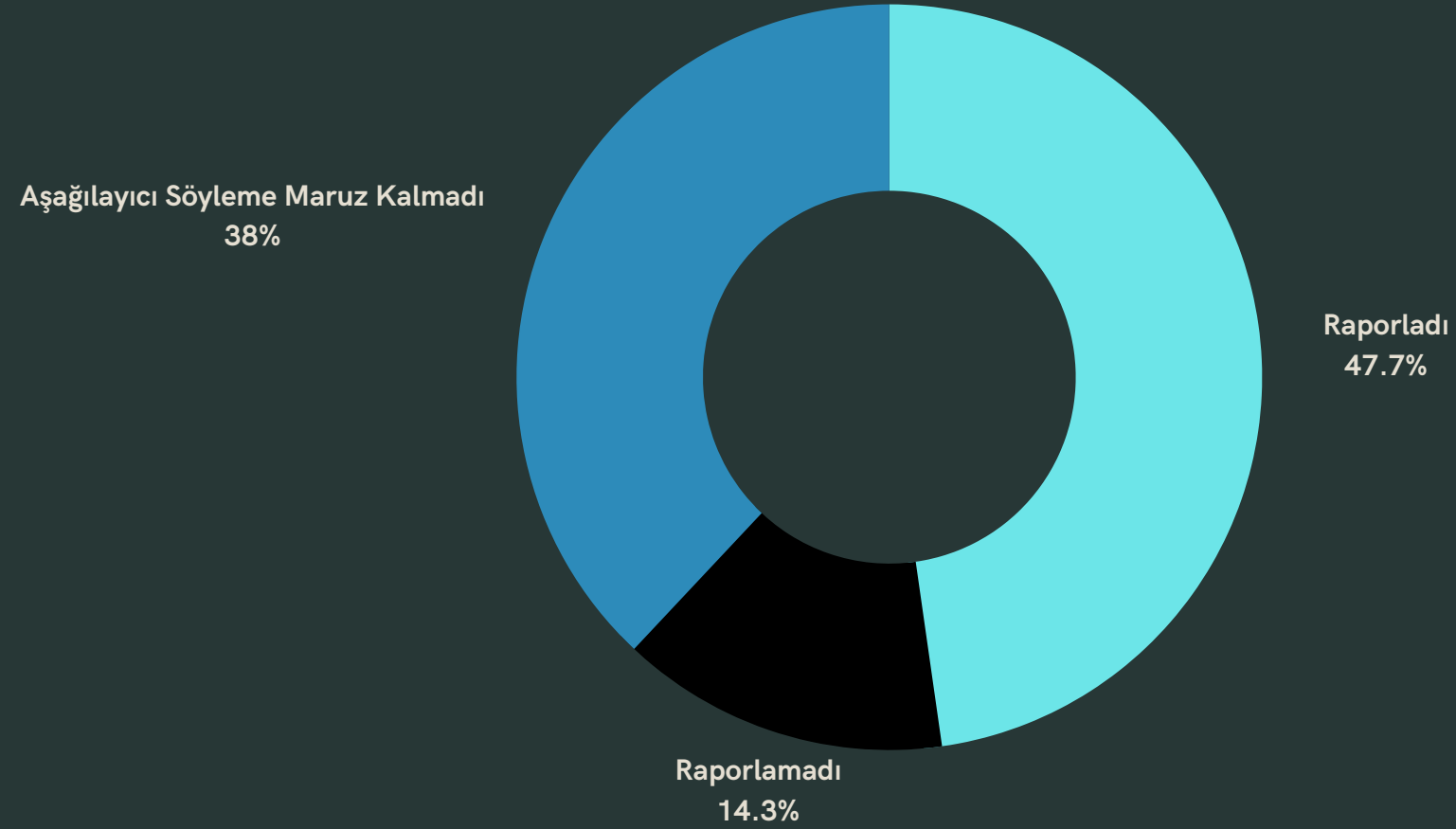


- Çözümün fonksiyonel ve dağıtılabılır hale getirilmesi
- Yöntemin başarısını yorumlamak adına temel çözümlerin oluşturulması (benchmarking)
- Önerilen yöntemin farklı mimarilerde test edilmesi

- Çalışmada kullanılmış olan yöntemin literatür araştırması ile belirlenmesi
- Bahsedilecek yöntem için uygun ve açık veri kaynaklarının sağlanması
- Yöntemin uygulanabilirliğinin yerel ortamda test edilmesi



# Problem



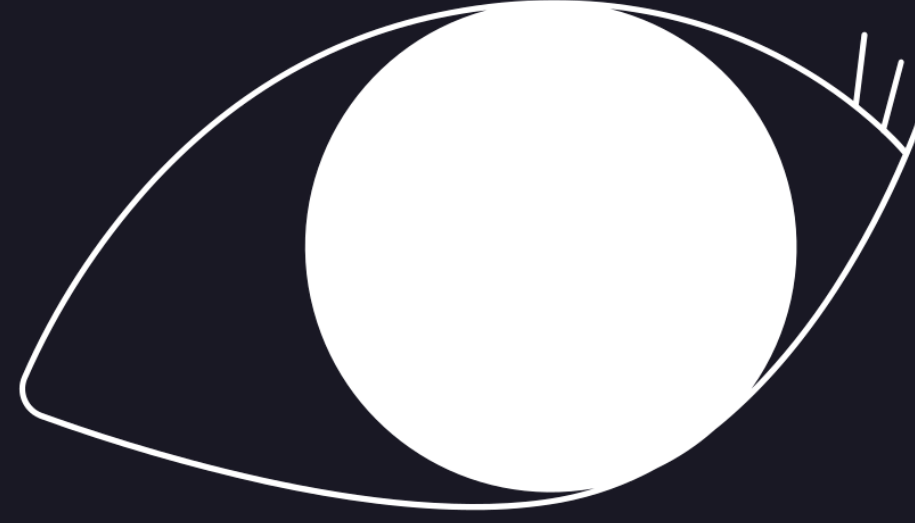
2018 ve 2019 yıllarında yapılan araştırmalara göre Türkiye'de 42.6 milyon sosyal medya kullanıcısı bulunmaktadır[1] ve anketlere göre %62'si aşağılayıcı söyleme maruz kaldığını belirtmiştir.[2]

Sosyal medyada aşağılayıcı söyleme maruz kalan kullanıcıların %23'ü bu söylemleri rapor etmemektedir. [2]

**Bu araştırmalara göre Türkiye'de dijital ortamda nefret söylemine maruz kalan yaklaşık 6 milyon kullanıcı aşağılayıcı söylemleri raporlamamaktadır.**

[1] We Are Social & Hootsuite. (2019). Digital 2019: Turkey

[2] Kadir Has Üniversitesi İletişim Fakültesi. (2019). Dijital Dünya ve Türkiye'de Sosyal Medya Kullanımı 2019 Raporu

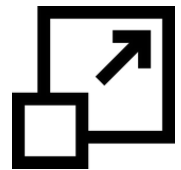


# ÖNERİLEN YÖNTEM

# KAYSERIOUSModel

## ✓ ÖLÇEKLENEBİLİR

Çalışma, bir Python paketi olarak yazıldığı için her türlü dijital sisteme entegre olabilmesiyle birlikte bir API olarak çağrılıp metin sınıflandırması görevleri için kullanılabilir.



## ✓ ÖZELLEŞTİRİLEBİLİR

Sadece aşağılayıcı söylem tespiti için değil, bir çok metin segmentasyonu problemi için etiketlenmiş veriye duyulan ihtiyacı çok azaltarak farklı problemlerin çözümü için kullanılabilir.



## ✓ KULLANICI DOSTU

Modüler yapısı sayesinde sadece birkaç satır kod ile özelleştirilmiş bir dil modeli oluşturulup canlı ortamda kullanılabilir.

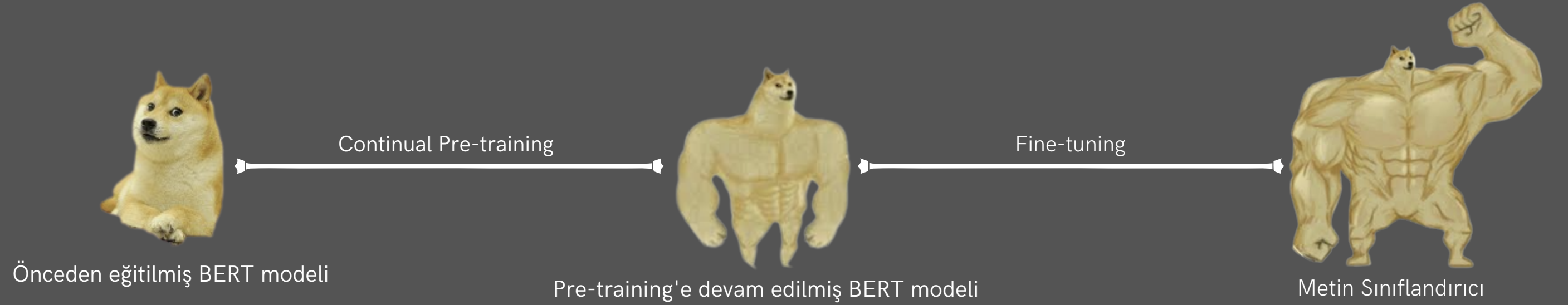


# Yöntem

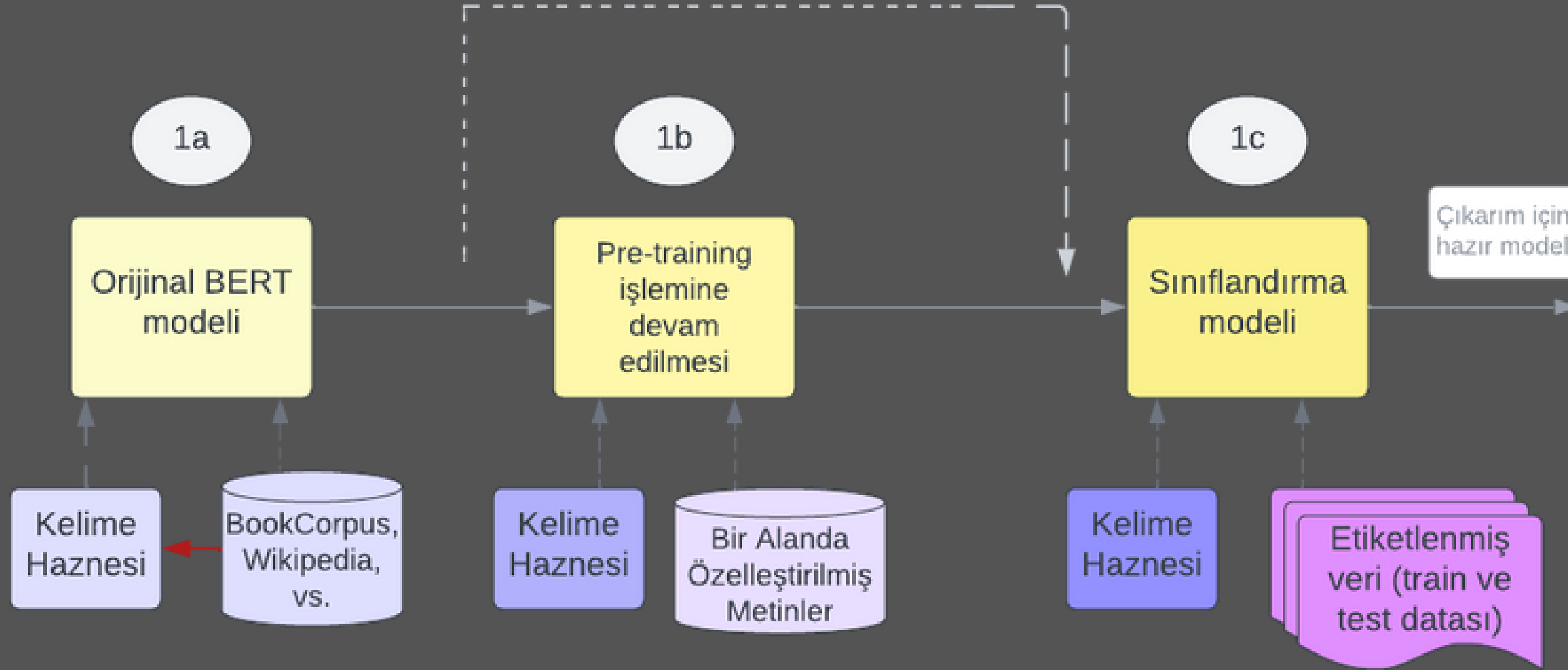
## Geleneksel :



## Yenilikçi :



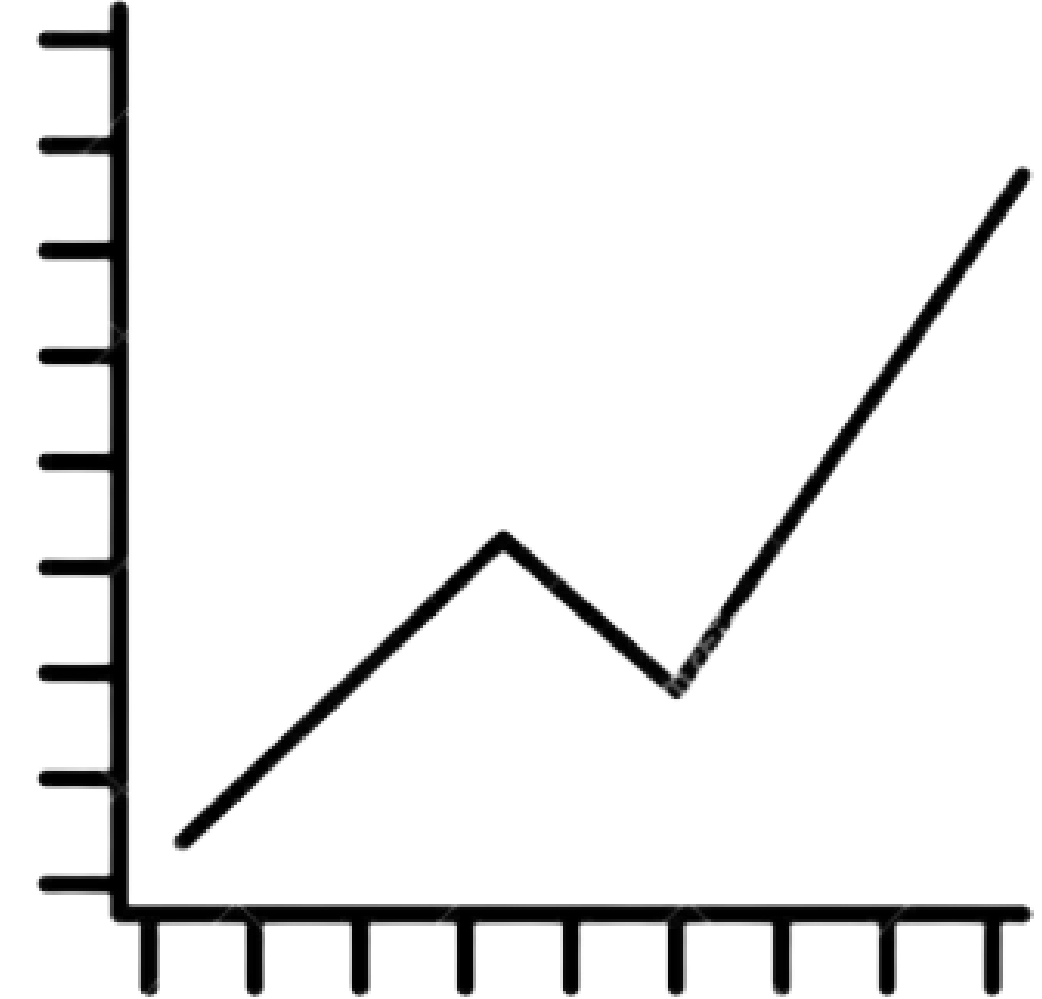
# Yöntem





# Alternatif Yöntemler ile Karşılaştırılması

Dil Modeli	Ortalama F Makro Skoru
BERTurk Cased (simpletransformers) Domain Spesifik 3 epoch \w 32k size corpus	0,96
BERTurk Cased (simpletransformers) Domain Spesifik 3 epoch \w 18k size corpus	0,94
BERTurk Cased (simpletransformers) Domain Spesifik 2 epoch \w 18k size corpus	0,94
BERTurk Cased (simpletransformers)	0,94
BERTurk Cased (simpletransformers) Domain Spesifik 1 epoch \w 18k size corpus	0,94
BERTurk Cased (simpletransformers) Domain Spesifik 4 epoch \w 18k size corpus	0,94
BERTurk 128k Uncased (Domain Spesifik Denemesi \w 15k size corpus)	0,94
BERTurk Cased	0,93
BERTurk Cased (Domain Spesifik Denemesi & MLM \w 4k size corpus)	0,93
BERTurk Cased	0,93
BERTurk Uncased	0,93
BERTurk Cased (Domain Spesifik \w 2k size corpus)	0,93
Distilled BERTurk	0,91
TFIDF - BERTurk Embedding Catboost \w zemberek normalization	0,90
TFIDF - BERTurk Cased Embedding Catboost	0,89
GPT-2	0,88
BERTurk Cased Embedding Catboost	0,87
BERTurk Uncased (simpletransformers)	0,86

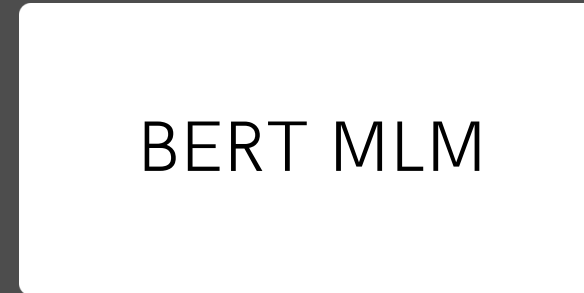


# Ön Eğitimi Sürdürme (Continual Pre-training)



MLM (Maskelenmiş Dil Modeli)

kadın - dediğin - evden - çıkmamalı



Sosyal medya  
verisi ile güçlendirilir

[MASK] - dediğin - evden - çıkmamalı

# KAYSERIOUSPreprocessor ( )

Metin temizliği, yazım hatalarının temizliği gibi veri ön işleme adımlarını yapan modüldür. Ağırlıklı olarak Zemberek Türkçe Doğal Dil İşleme paketini kullanır\*.

\*<https://github.com/ahmetaa/zemberek-nlp>

# KAYSERIOUSPreTrainer ( )

Geliştirilecek dil modelini ve bir alana yönelik metinleri girdi olarak alır ve MLM yöntemiyle dil modelini o alana özgü hale getirir.

# KAYSERIOUSModel ( )

Bir alanda özelleştirilmiş yahut özelleştirilmemiş dil modelini ve metin etiketlerini alarak sınıflandırma modeli oluşturur.

# Örnek Kullanım

```
pt = pretrainer.KAYSERIOUSPreTrainer(corpus_path = cp,  
                                     base_model = constants.BASE_MODEL,  
                                     pretrain_args=constants.PRETRAIN_ARGS,  
                                     seed=constants.RANDOM_SEED,  
                                     gpu = constants.USE_GPU,  
                                     out_dir = constants.SAVE_PRETRAINED_TO)  
  
pt.pretrain()  
  
md = modeler.KAYSERIOUSModel(modelargs = constants.MODEL_ARGS,  
                              modelfolder = constants.SAVE_DEPLOYED_TO,  
                              seed = constants.RANDOM_SEED,  
                              gpu = constants.USE_GPU,  
                              base_model = constants.SAVE_PRETRAINED_TO)  
  
md.construct_data(training_data = df, text_column = constants.TEXT_NAME, target_column = constants.TARGET_NAME)  
  
md.train_model()
```

# kayserious/tddi-2023



Bu repo, Kayserious takımının TEKNOFEST 2023 kapsamında gerçekleştirilen Türkçe Doğal Dil İşleme yarışması çözümünü içerir.

1

Contributor

0

Issues

0

Stars

0

Forks

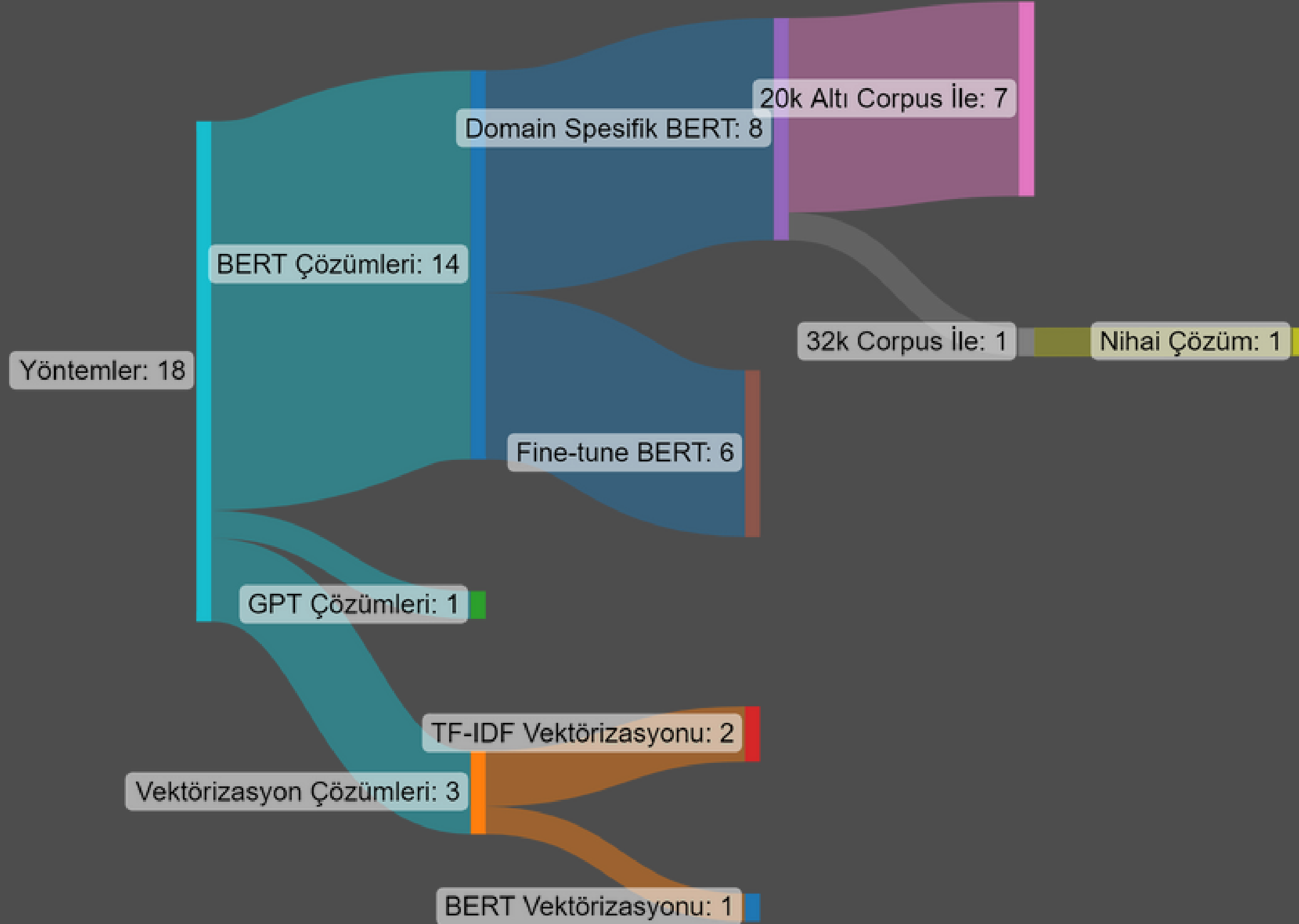


**kayserious/tddi-2023: Bu repo, Kayserious takımının TEKNOFEST 2023 kapsamında gerçekleştirilen Türkçe Doğal Dil İşleme yarışması çözümünü...**

Bu repo, Kayserious takımının TEKNOFEST 2023 kapsamında gerçekleştirilen Türkçe Doğal Dil İşleme yarışması çözümünü içerir. -  
GitHub - kayserious/tddi-2023: Bu repo, Kayserious takımının TEKNOFEST ...

GitHub

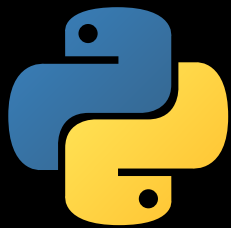
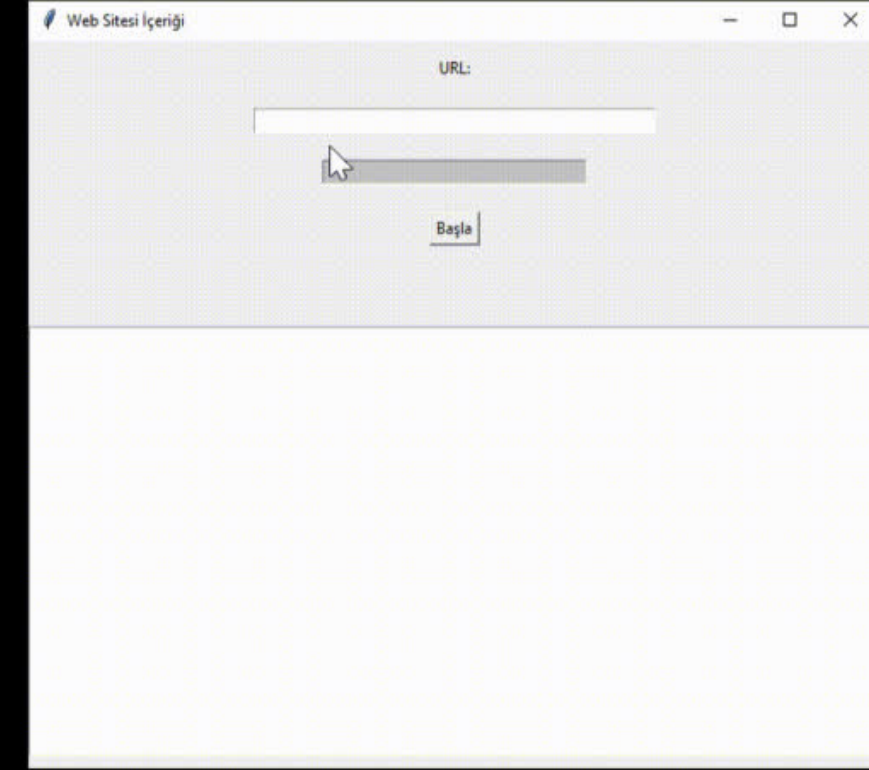






# Projenin Sürdürülebilirliği

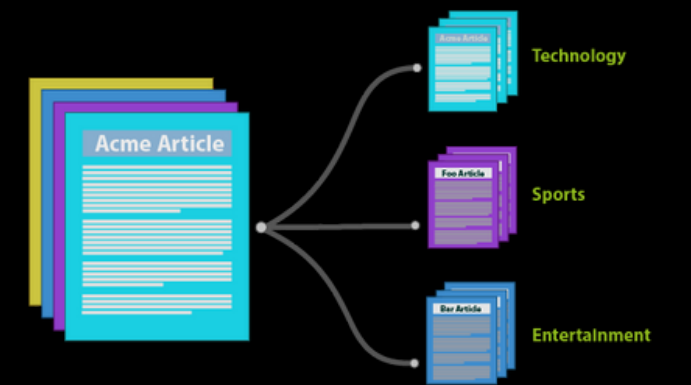
- Yükle
- Özelleştir
- Kullan



pip install kayseriousmodel



alana yönelik metinler



güçlü sınıflandırma modelleri