

TÜRKİYE CUMHURİYETİ
YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



TÜRKÇE İÇİN AYNİ ANLAMLI CÜMLE ÜRETİMİ

19011613 – Metin BİNBİR
19011620 – Sercan AKSOY

BİLGİSAYAR PROJESİ

Danışman
Doç. Dr. Mehmet Fatih AMASYALI

Haziran, 2021

TEŞEKKÜR

Daha önce yapay sinir ağıları ile çalışma yapmadığımızı bilmesine rağmen çalışma sevgimizi değerlendirerek bize güvenen, bizi tamamlamaya götürecek şekilde yol gösteren ve süreç boyunca yardımlarını hiç esirgemeyen Doç. Dr. Mehmet Fatih Amasyalı'ya, projenin gidişatından her zaman haberi olan ve test aşamasında etiketlemeye yardım eden Alaaddin Göktuğ Ayar, Göktuğ Köksal ve Selcan Yağmur Atak'a, yapay sinir ağıları araştırmalarında kaynak gösterdiği için üniversite arkadaşlarımız Aslı Şeyda Özdemir ve Şafak Bilici'ye en içten şekilde teşekkür etmeyi kendimize borç biliriz.

Metin BİNBİR
Sercan AKSOY

İÇİNDEKİLER

KISALTMA LİSTESİ	v
ŞEKİL LİSTESİ	vi
TABLO LİSTESİ	viii
ÖZET	ix
ABSTRACT	x
1 Giriş	1
2 Ön İnceleme	3
3 Fizibilite	4
3.1 Teknik Fizibilite	4
3.1.1 Yazılım Fizibilitesi	4
3.1.2 Donanım Fizibilitesi	4
3.2 İş Gücü ve Zaman Fizibilitesi	5
3.3 Yasal Fizibilite	5
3.4 Ekonomik Fizibilite	5
4 Sistem Analizi	7
5 Sistem Tasarımı	8
5.1 Veri Seti Tasarımı	8
5.1.1 İngilizce Dil Modelleme	8
5.1.2 Duygu Analizi Görevinde İnce Ayar İşlemi	9
5.1.3 Aynı Anlamalı Cümle Üretimi (İngilizce)	9
5.1.4 Aynı Anlamalı Cümle Üretimi (Türkçe)	10
5.2 Yazılım Tasarımı	11
5.3 Girdi-Çıktı Tasarımı	18
6 Uygulama	23

7 Deneysel Sonular	29
8 Performans Analizi	36
8.1 T5 İngilizce Dil Modelleme	36
8.2 T5 ile İngilizce Duygu Analizi	37
8.3 İngilizce Aynı Anlamlı Cümle Üretimi Deęerlendirmesi	37
8.4 MT5 ile Türke Aynı Anlamlı Cümle Üretimi	37
8.4.1 BLEU Skoru	37
8.4.2 İnsan Deęerlendirmesi	38
8.4.3 Sonuların Deęerlendirilmesi	38
9 Sonu	40
Referanslar	42
Özgemiř	44

KISALTMA LİSTESİ

BLEU	Bilingual Evaluation Understudy Score
BoW	Bag of Words
C4	Colossal Clean Crawled Corpus
DDİ	Doğal Dil İşleme
MC4	Multi-lingual Colossal Clean Crawled Corpus
MT5	Multi-lingual Text-to-Text Transfer Transformer
QQP	Quora Question Pairs
RNN	Recurrent Neural Network
RVAE	Recurrent Variational Autoencoder
Seq2Seq	Sequence to Sequence
T5	Text-to-Text Transfer Transformer

ŞEKİL LİSTESİ

Şekil 3.1	Proje İş-Zaman Çizelgesi	5
Şekil 5.1	Duygu Analizi Veri Seti Örneği	9
Şekil 5.2	Veri Seti Formatı	10
Şekil 5.3	RNN	12
Şekil 5.4	T5 Çalışma Sistemi	13
Şekil 5.5	T5 Maskelenmiş Cümle Örneği	13
Şekil 5.6	T5 Encoder-Decoder	14
Şekil 5.7	T5 ile İngilizce-Almanca Çeviri	14
Şekil 5.8	T5 ile İki Örnek Cümle Arasındaki Anlam Benzerliği Tespiti . . .	15
Şekil 5.9	T5 ile Metin Özetleme	15
Şekil 5.10	T5 ile Cümlelerin Oluşturulma Mümkünlüğü	15
Şekil 5.11	T5 Dil Modelleme İçin Kullanılan Argümanlar	16
Şekil 5.12	T5 Duygu Analizi İnce Ayar İşleminde Kullandığımız Argümanlar	16
Şekil 5.13	Aynı Anlamlı Cümle Tespiti ve Aynı Anlamlı Cümle Üretimi . . .	17
Şekil 5.14	Türkçe Aynı Anlamlı Cümle Üretimi İnce Ayar Fonksiyonundaki Argümanlar	18
Şekil 5.15	Modellerin, Kodların ve Eğitim Veri Kümesini İçeren Ekran Görüntüsü	19
Şekil 5.16	Validation Veri Setini İçeren ve Kodda Dizin Gösterilen Ekran Görüntüsü	19
Şekil 5.17	Tarafımızca İnce Ayar Yapılmış Modelin ve Daha Önce Eğitilmiş Tokenizer'ın Yüklenmesi	20
Şekil 5.18	Kullanıcı Tarafından Girdi Cümle Yazılması	20
Şekil 5.19	Üretilen Aynı Anlamlı Cümle Çıktıları	21
Şekil 5.20	Test Cümlelerinin Dosyadan Verilmesi	21
Şekil 5.21	Üretilen Cümlelerin Dosyaya Yazılması	22
Şekil 5.22	Konsola Yazdırılan ve Dosyaya Yazılan Cümleler	22
Şekil 6.1	Cümlelere Yapılan Tahminler	23
Şekil 6.2	Google Drive Bağlantısı Kurulması	25
Şekil 6.3	Gerekli Kütüphanelerin Kurulumu	25
Şekil 6.4	Modelin Eğitim Aşaması - 1	26

Şekil 6.5	Modelin Eğitim Aşaması - 2	26
Şekil 6.6	Eğitilen Modelin Yüklenmesi	27
Şekil 6.7	Konsoldan Cümle Girişi	27
Şekil 6.8	Cümlelere Yapılan Tahminler - 1	28
Şekil 6.9	Cümlelere Yapılan Tahminler - 2	28
Şekil 7.1	Uç Durum: Soru Cümlesi Olmayan Cümle	29
Şekil 7.2	Uç Durum: İngilizce Soru Cümlesi	29
Şekil 7.3	Uç Durum: Soru Cümlesi Olmayan Olumsuz Cümle	29
Şekil 7.4	Uç Durum: Yalnızca Sayı Girişi	29
Şekil 7.5	Uç Durum: Özel Karakter İçeren Girdi Durumu	30
Şekil 7.6	Uç Durum: Tek Kelimelik Girdi Durumu	30
Şekil 7.7	Uç Durum: Eksik Cümle	30
Şekil 7.8	Uç Durum: Yüklem İçermeyen Cümle	30
Şekil 7.9	1.Cümle - Epoch Sayısı=2	35
Şekil 7.10	1.Cümle - Epoch Sayısı=3	35
Şekil 7.11	2.Cümle - Epoch Sayısı=2	35
Şekil 7.12	2.Cümle - Epoch Sayısı=3	35
Şekil 8.1	T5 İngilizce Dil Modellemesinde Loss Değerleri (Karşılaştırma-Küçük Veri Seti-Büyük Veri Seti)	36
Şekil 8.2	İnce Ayar İşleminin Başarısına Ait Metrik Değerler	37

TABLO LİSTESİ

Tablo 3.1	Ekonomik Fizibilite Tablosu	6
Tablo 5.1	QQP Veri Setindeki Bazı Örnekler	10
Tablo 5.2	Eğitimde Kullanılan Veri Setleri	10
Tablo 5.3	Manuel Oluşturulmuş Türkçe Veri Setindeki Bazı Örnekler	11
Tablo 6.1	İngilizce Aynı Anlamlı Cümle Üretimi Sonuçları	24
Tablo 7.1	Farklı Veri Setlerine Göre Üretilen Aynı Anlamlı Cümle Çıktıları - 1	31
Tablo 7.2	Farklı Veri Setlerine Göre Üretilen Aynı Anlamlı Cümle Çıktıları - 2	32
Tablo 7.3	Farklı Veri Setlerine Göre Üretilen Aynı Anlamlı Cümle Çıktıları - 3	33
Tablo 7.4	Farklı Veri Setlerine Göre Üretilen Aynı Anlamlı Cümle Çıktıları - 4	34
Tablo 8.1	Hesaplanan-Gerçeklik Matrisi	37
Tablo 8.2	Farklı Veri Setleri Kullanılarak Eğitilmiş MT5 Modellerin BLEU Skorları	38
Tablo 8.3	Farklı Veri Setleri Kullanılarak Eğitilmiş MT5 Modellerin İnsan Değerlendirmesi	38

TÜRKÇE İÇİN AYNİ ANLAMLI CÜMLE ÜRETİMİ

Metin BİNBİR

Sercan AKSOY

Bilgisayar Mühendisliği Bölümü

Bilgisayar Projesi

Danışman: Doç. Dr. Mehmet Fatih AMASYALI

Projenin konusu olan Türkçe için eş anlamlı cümle üretimi, genel olarak paragraf saflaştırma ve özetlemede, makine çevirisinde ve akademik makalelerde cümle tekrarını azaltmak, akla gelmeyen eş anlamlı sözcükler ile paragrafları tamamlamakta ve bu alandaki veri setlerini genişletmede kullanılmıştır.

Gelişen teknolojiler neticesinde insanlar makinelerle anlaşabilmek adına dillerini onlara öğretmeye başlamıştır. Belirlenen işlerde gerekli veri setleri özellikle Türkçe için eksik kalmıştır. Ana dilimiz olan Türkçe kullanılarak makineler ile iletişim kurabilmek ve belirlenen görevlerin gerçekleşmesi adına proje kapsamında dil modellemeler, modellerin test edilmesi ve neticesinde eş anlamlı cümle üretimi yapılmıştır.

Elde edilen modelin kullanıcı tarafından verilen Türkçe bir cümleyi herhangi bir format beklemeksizin, anlamlandırması ve biçimsel benzerlikten mümkün olduğunca uzak durarak anlamsal bütünlüğü koruması beklenmiştir. Bu amaç uğrunda eğitilen model cümle içerisindeki öğelerin yerlerini değiştirebilir, bazı sözcükler yerine eş anlamlılarını getirebilir hatta özneler yerine zamirler kullanabilir şekilde veri seti hazırlanmıştır.

Anahtar Kelimeler: Doğal Dil İşleme, Derin Öğrenme, Aynı Anlamlı Cümle Üretimi, Makine çevirisi, BoW, Seq2Seq, T5, MT5, Transformer, BLEU, C4.

ABSTRACT

PARAPHRASE GENERATION FOR THE LANGUAGE TURKISH

Metin BİNBİR
Sercan AKSOY

Department of Computer Engineering
Computer Project

Advisor: Assoc. Prof. Dr. Mehmet Fatih AMASYALI

The main purpose of the project is generating paraphrased sentences in Turkish, which is used for purificating and summarizing text, machine translation and to reduce repetition in academic articles.

As a result of improvements on technologies, people have been trying to communicate with the machines in their language. In specified tasks, there are lots of gaps for the dataset in Turkish language. For the Turkish language, which is our native language, we have augmented new data for the purpose of training a new model, for testing it and also for generating paraphrased sentences.

We are expecting from the obtained model to be able to generate new sentences in Turkish which are semantically same but syntactically different from the input sentences. For this purpose, model may change the order of elements in sentence, may change words with their synonyms and may replace subjects with pronouns. Dataset was prepared for these situations.

Keywords: Natural Language Processing, Deep Learning, Paraphrase Generation, Machine Translation, BoW, Seq2Seq, T5, MT5, Transformer, BLEU, C4.

1

Giriş

Aynı anlamlı cümle üretimi görevi, yeni bir veri seti üretilmesinde ya da genişletilmesinde, metin özetlemede, alınan bir cümle için farklı çeviri seçenekleri sunmada kullanılabilir. Aynı anlamlı cümle üretimi çalışması çoğu doğal dil işleme alanında olduğu gibi oldukça büyük veri kümelerine ihtiyaç duyar.

Proje kapsamında, özellikle Türkçe için ihtiyaç duyulan veri kümesine katkıda bulunarak otomatik elde edilmiş veri kümeleri ile aralarındaki farklar değerlendirilmiştir.

Karşılaştırmalar neticesinde, üretilen cümlelerin veri kümeleri arasında boyut farkından ziyade hedef olarak belirlenen cümlelerin niteliklerine bağlı olduğu tespit edilmiş ve ana dili Türkçe olan insanlar tarafından değerlendirmeye tutulmuştur.

Aynı anlama cümle üretimi projesiyle hedeflenenler sıralandığı gibidir; bu görevi Türkçe için en iyi çalıştırabilecek modelleri araştırmak, doğal dil işleme projesinin temeli olan dil modelleme hakkında tecrübe kazanmak, Türkçe için kullanılabilecek veri kümesini geliştirmek, geliştirilen veri kümesinin başarısını ölçmek ve kıyaslamak, muadillerine göre başarılı olduğunu düşünülen aynı anlamlı cümle üretimi görevini bir başka konu için veri kümesi genişletmede kullanabilmek.

Raporun devamında:

- Ön inceleme kısmında başkaları tarafından gerçekleştirilmiş benzer projeleri ve bu çalışmanın katkıları,
- Fizibilite kısmında projenin gerçekleştirilebilmesi için teknik, yasal, ekonomik ve zamansal kısıtlamalar,
- Sistem analizi kısmında projenin gerçekleşmesi için uygun yolları,
- Sistem tasarımı kısmında üç başlıkta incelenmek üzere yazılımın mimarisi ve kullanılan algoritmalar hakkında açıklamalar, veri tabanı özellikleri ve etkileri,

ayrıca projenin alıřtırılmasında girdilerin ve ıktıların nasıl yönetildięi,

- Uygulama kısmında projenin hedefleri doęrultusunda zerine alıřılmış modellerin bilgileri,
- DeneySEL Sonular kısmında eęitilen modelin girdilere karřı rettięi yeni cmleleri ayrıca u noktalar ile sınanması,
- Performans Analizi kısmında alıřmanın bařarısı ve bunun nasıl lldę yer almaktadır.

2 Ön İnceleme

Projenin konusunda örnek çalışmalar araştırıldığında Türkçe dili için yapılmış bir projeye rastlanamamıştır. Bu sebeple mimarileri kıyaslamak için İngilizce dili için yapılmış çalışmalar incelenmiştir. Modelin eğitilebildiği aynı zamanda tarafımızca anlaşılabilen 3 çalışma bulunmuştur. Temel Seq2Seq[1](Sequence to Sequence), BoW[2](Bag of Words) ve T5 (Text-To-Text-Transfer-Transformer) olarak bilinen bu modeller aynı veri seti QQP (Quora Question Pairs) ile eğitilmiş ve yine aynı cümleler ile test edilmiştir. Çıktıların incelenmesi sonucu, en güncel yayımlanmış (2019) çalışma olan T5'in daha başarılı olduğuna karar verilmiş ve bu model ile çalışmaya başlanmıştır.

T5'in Türkçe kullanımını sağlayan MT5(Multi-lingual Text-To-Text-Transfer-Transformer) modeli ile gerçekleştirilmiş bir proje, çalışmalar esnasında kullanıma sunulmuştur[3]. Bu proje çeşitli girdiler ile sınanmış ve üretilen cümlelerin anlamsal bozukluklar içerdiği, ayrıca referans cümle ile ortak kelime sayısının fazla olduğu tespit edilmiştir. Örnek olarak belirlenen bu çalışmada eğitim kümesi makine çevirisi yöntemi ile hazırlanmış olup Türkçe için uygun olmayan kelimeler içermektedir. Bu eksikliğin giderilmesi adına danışmanımız tarafından sağlanmış, manuel olarak oluşturulmuş veri seti yine manuel olarak genişletilmiştir. Böylece aynı mimari kullanılsa dahi insan kaynaklı değerlendirmede başarı farkı ortaya çıkarılmıştır.

3.1 Teknik Fizibilite

3.1.1 Yazılım Fizibilitesi

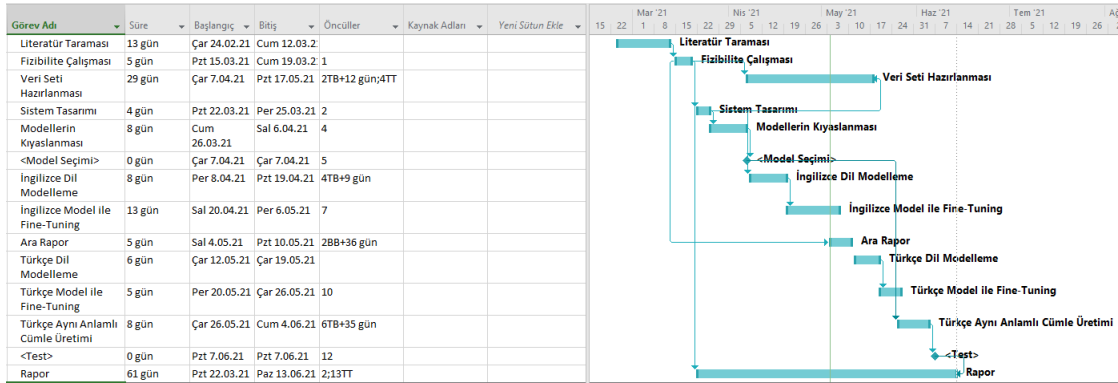
Proje, aynı konuda birçok örnek çalışma bulundurduğu için Python programlama dili ile gerçekleştirilecektir. Doğal dil işleme alanında kullanılan en büyük platformlardan biri olan Hugging Face aracılığıyla eğitilmiş T5[4] modelleri üzerinden İngilizce dilinde testler yapılması beklenmektedir. Transformers[5] kütüphanesi içinden yine eğitilmiş T5Tokenizer kullanılması beklenmektedir.

Türkçe aynı anlamlı cümle üretimi için ise T5'in Türkçe uyumlu versiyonu olan MT5[6] hazır eğitilmiş dil modelinin mt5-base versiyonu ve parçalayıcısı olan mt5-base versiyonu kullanılmıştır. Eğitim aşaması için torch[7] ve pytorch-lightning kütüphaneleri kullanılmıştır. Girdilerin işlenmesi amacıyla sentencepiece ve nltk kütüphaneleri kullanılmıştır. Proje başarısının BLEU[8] skoruyla ölçülmesi için de datasets[9] kütüphanesi kullanılmıştır.

3.1.2 Donanım Fizibilitesi

Projenin temeli büyük veri setlerinden faydalanmaktadır ve veri setini işlemek için çok hızlı işlemcilere, ayrıca çok fazla kapasitesi olan hafızalara ihtiyaç duyar. Google Colab bu olanakları belli bir ücret karşılığında sağladığı için internet bağlantısı olan orta düzey bir bilgisayar yeterli olacaktır.

3.2 İş Gücü ve Zaman Fizibilitesi



Şekil 3.1 Proje İş-Zaman Çizelgesi

Projenin gerçekleştirilmesi için kullanılacak yöntemleri incelemek kapsamında ilk işlem literatür taramasıdır. Proje hakkında bilgi edinildikten sonra projenin devamlılığına karar verecek olan fizibilite çalışması ile yapılabirlik test edilir. Fizibilite çalışmasında bu alanda yeterli veri seti bulunmadığı farkedilerek kendi veri setimizi oluşturmak için yeterli süre ataması yapılmalıdır. Çalıştırılacak algoritma seçimi için proje konusunda daha önce yapılmış çalışmalar test edilir, değerlendirme sonuçlarına göre model seçilir. Bu aşamada projenin alt seviyeleri hakkında tecrübe kazanmak için İngilizce dil modelleme üzerine örnekler geliştirilir. Benzer işlemler Türkçe için de uygulanarak proje tamamlanma aşamasına getirilir. Son ana kadar elde edilen bütün çıktılar raporlama aşamasına dahil edilir.

3.3 Yasal Fizibilite

Yapılan proje mevcut yasa ve yönetmeliklere uygun olmakla beraber herhangi bir patent vb. korunmuş hakkı ihlal etmemektedir.

3.4 Ekonomik Fizibilite

- Geliştirme ortamı, JetBrains PyCharm Python IDE'si JetBrains tarafından öğrencilere ücretsiz olarak sağlanmaktadır. Bir diğer ortam olan Google Colab uygulamasının ücretsiz sürümündeki kullanım kotası dolabildiğinden Google Colab Pro için 4 aylık kullanım esnasında 320 TL'lik bir masraf öngörülmektedir.
- Derin öğrenme için lazım olan tüm Python kütüphaneleri ücretsiz olarak kullanılabilir.

- Sahip olduğumuz kişisel bilgisayarlar bazı kütüphanelerin eski versiyonlarında gerekli fonksiyonları çalıştıramasa da projede yeteri düzeyde ilerleme katetme donanımına sahiptir.
- Modelin eğitimi için gereken cümle çiftlerini içeren veri setinin doldurulması için 5 işçiye asgari ücret mukabilinde 1 ay ödeme yapılması uygun görülmüştür.

Görev Tanımı	Kullanıcı/Adet	Birim Fiyat (TL)	Toplam Maliyet (TL)
Google Colab Pro	4	80	320
Microsoft Project Plan 1	1	72	72
JetBrains PyCharm IDE	2	0	0
Lenovo Legion 5	2	9.300	18.600
21.5 Inch HDMI LCD Ekran	1	560	560
20 Inch HDMI LCD Ekran	1	480	480
Google Drive Depolama Alanı	1	40	40
Mühendis Maaşları	2	10.000	20.000
Cümle Yazarları	5	2.825,90	14.129,50
TOPLAM MALİYET			54.201,50

Tablo 3.1 Ekonomik Fizibilite Tablosu

4 Sistem Analizi

Proje konusunun gerçekleştirilebilmesi için yöntemler araştırıldığında RVAE[10], temel Seq2Seq, BoW ve T5 modelleri ile yapılmış çalışmalar denenmiştir. RVAE mimarisi ile ilgili çalışmalarda hali hazırda eğitilmiş modeller kullanıldığı için araştırmaya uygun olmayacağı, temel Seq2Seq ve BoW için çalışmaların T5 modeli ile yapılmışların yanında sayılarının az olması, BoW mimarisinin dil modelinden bağımsız gerçekleşmesinin çalışmaya fazladan zorluk katacağı düşünülerek 2019'da yayımlanmış, bulunabilen en güncel model olan T5 mimarisi seçilmiştir.

Proje kapsamında danışman hocamız tarafından yönlendirilmeler neticesinde gereklilikler Hugging Face adlı doğal dil işleme ile ilgili en güncel modellerin ve veri kümelerinin bulunduğu web sitesinden, GitHub adlı açık kaynak kodlu çalışmaların bulunduğu web sitesinden, NLP Yardımlaşma Platformu adlı bu alanda çalışan insanların toplandığı ve yardımlaştığı bir Discord uygulaması grubundan ve Medium sitesindeki bloglardan faydalanılmıştır.

Eğitilen İngilizce dil modelinin başarısını ölçmek için duygu analizi görevindeki tahminin ne kadar isabetli olduğunu ölçmek için F1 skoru, İngilizce aynı anlamlı cümle üretiminin başarısını ölçmek için insan değerlendirmesi kullanılmıştır.

Türkçe aynı anlamlı cümle üretiminin başarısını ölçmek için BLEU[8] skoru ve insan değerlendirmesi kullanılmıştır.

5.1 Veri Seti Tasarımı

T5 modeli, dil ve görev gözetmeksizin verilen bir metni bir başka metne çevirmeye çalışır. Bu hususta verilen metnin dönüştürülecek metni hangi görevle elde edeceği, tasarlanan veri setindeki ikili cümleler arası ilişkiye bağlıdır.

5.1.1 İngilizce Dil Modelleme

Google tarafından eğitilmiş olan T5 modeli aynı makalede sunulan C4(Colossal Clean Crawled Corpus) veri seti ile eğitilmiştir. Bu veri seti, daha öncesinde kullanılan çoğu veri setinden daha büyük boyutlu (yaklaşık olarak 750 GB) ve kirli veriden daha çok arındırılmış hâldedir, bunun sebepleri [4] ise aşağıda listelenmiştir:

- Sadece belirli noktalama işaretleri (ünlem işareti, soru işareti, nokta ve tırnak işareti) ile biten cümleler seçilmiştir.
- 5'ten az sayıda cümle içeren metinler ve 3'ten az sayıda kelime içeren cümleler elenmiştir.
- “List of Dirty, Naughty, Obscene or Otherwise Bad Words” içerisinde yer alan kelimeler barındıran cümleler veri setine dahil edilmemiştir.
- Javascript kodlarını çalıştırma izni isteyen tüm internet siteleri elenmiştir.
- Lorem Ipsum (Anlamsız Yazılar) içeren bütün sayfalardan uzak durulmuştur.
- Süslü parantez karakterleri ("{" ve "}") program kodunun çalışmasında problem yaratabileceği için bu karakterleri bulunduran tüm sayfalardan kaçınılmıştır.
- İçerik tekrarına düşmemek için veri setinde mevcut olan 3 cümleyi içeren internet sitelerindeki içerikler veri setine eklenmemiştir.

Tarafımızca oluşturulmuş İngilizce dil modeli gerçekleştirirken "National Library of Medicine" sitesinde bulunan bazı makalelerin özet kısımları kullanılmıştır. Dil modelleme esnasında veri seti büyüklüğünün öneminin görülebilmesi adına kullanılan veri setleri 1/5 oranında iki parçaya ayrılmıştır.

5.1.2 Duygu Analizi Görevinde İnce Ayar İşlemi

Modellenen dilin mantıklı sonuçlar üretebileceğinden emin olmak için sonucu doğru ya da yanlış olarak verebilecek bir görev belirlenmesi daha sağlıklı gözlemlere olanak tanıyacaktır. Bu yüzden ince ayar işlemi için görev olarak duygu analizi, yani metinleri pozitif veya negatif olarak sınıflandıran bir kod kullanılmıştır. Bu görevde kullanılan veri setinde[11] cümleler, karşılığında pozitif veya negatif olarak etiketlenmiş hâlde bulunmaktadır.

Pozitif Cümle Örneği: This is the definite Lars von Trier Movie, my favorite, I rank it higher than "Breaking the waves" or the latest "Dancer in the Dark"... I simply love the beauty of the picture...the framing is so original; acting is wonderful, A MUST SEE.

Negatif Cümle Örneği: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Şekil 5.1 Duygu Analizi Veri Seti Örneği

5.1.3 Aynı Anlamlı Cümle Üretimi (İngilizce)

Projenin ilk aşamalarında model kıyaslama aşaması için eğitilmiş dil modelleri üzerinde QQP veri setinin bir kısmı (90.000 soru çifti) kullanılmıştır. Bu veri seti, Quora web sayfası üzerinde sorulmuş soruların ikili cümleler hâlinde eşlenmesi ve bu soru cümlelerinin 1 ve 0 (aynı anlamlı ya da aynı anlamlı değil) ile işaretlenmesiyle oluşur [12]. Projenin amacı aynı anlamlı cümle üretimi olduğu için veri setinde sadece 1 ile etiketlenmiş ikili soru cümleleri kullanılmıştır.

Cümle 1	When should I take BCAAs?
Cümle 1	When is the best time to take BCAA?
Cümle 2	What are the qualities of a good mother?
Cümle 2	What makes a good mum?
Cümle 3	How I increase my focus in study?
Cümle 3	How do we increase concentration?

Tablo 5.1 QQP Veri Setindeki Bazı Örnekler

5.1.4 Aynı Anlamlı Cümle Üretimi (Türkçe)

Türkçe dilinde aynı anlamlı cümle üretiminin ince ayar aşamasında eğitim veri setindeki cümleler txt uzantılı train ve validation dosyalarına Şekil 5.2'deki formatta verilmiştir.

$\langle s \rangle G === E \langle /s \rangle$

G: İlk cümle

E: İlk cümlenin eş anlamlı cümlesi

```
<s>Amerika'da hangi dili konuşuyorlar? === Amerika'da insanlar hangi dili konuşuyor? </s>
<s>Başka bir gezegenden misin? === Uzaylı mısın? </s>
<s>Perdeyi açabilir miyim? === Perdeyi açmamın bir sakıncası var mı? </s>
<s>Onlara ne sordun? === Onlara sorduğun neydi? </s>
<s>Sana yardım edeceğimi de nereden çıkardın? === Size yardım edeceğimi düşündüren nedir? </s>
<s>Ne öğrenmek istiyorsun? === Ne bilmek istiyorsun? </s>
```

Şekil 5.2 Veri Seti Formatı

Proje kapsamında modellerin başarısı veri setlerinin büyüklükleri ve içeriklerine göre test edildiği için Tablo 5.2'de Türkçe veri setlerinin kısaca açıklamaları verilmiştir.

Veri Seti ID	Açıklama	Cümle Çiftleri Sayısı
DH-S	Manuel olarak ekleme yapılmış forum sitesi başlıklarındaki soru çiftleri	30994
DH	DH-S veri setinin tarafımızca ekleme yapılmış hâli	50994
TQP	Çok sayıdaki dilde aynı anlama gelen cümlelerin çevirileri	51620
MT	QQP'nin Türkçe çevirisi	40200
BIG	DH, MT ve TQP veri setlerinin birleştirilmiş hâli	142814

Tablo 5.2 Eğitimde Kullanılan Veri Setleri

1. Makine Çevirisiyle Oluşturulmuş Veri Setleri

Makine çevirisiyle bir cümleyi birden çok dile çevirdikten sonra istenen dile tekrar çeviri yaparak aynı anlamlı cümle üretimi yapılabilmektedir fakat bu yöntem ile her zaman başarılı cümle çiftleri elde edilemeyebilir.

- (a) MT: QQP veri setinde yer alan 40200 adet soru çiftinin makine çevirisiyle Türkçe diline çevrilmiş hâlidir.
- (b) TQP: TaPaCo[13] veri setinde 73 dilde bulunan soru cümlelerinden Türkçe olarak 51620 çift bulunmaktadır.

2. Manuel Oluşturulmuş Veri Seti

Bir kısmı danışman hocamızca sağlanan ve tarafımızca 20000 adet ekleme yapılan Türkçe veri setinde ise yaklaşık 51000 adet aynı anlamlı soru çifti bulunmaktadır. Sorular ise Türkiye'deki ünlü forum web sitesi olan forum.donanimhaber.com adresindeki konu başlıklarının bir araya getirilmesiyle oluşturulmuştur. Soruların aynı anlama gelen çiftleri elle doldurulmuştur.

- (a) DH-S: Danışman hocamızca sağlanmıştır ve 30994 adet soru çifti içermektedir.
- (b) DH: DH-S veri setinin tarafımızca 20000 adet soru çifti eklenmiş hâlidir.

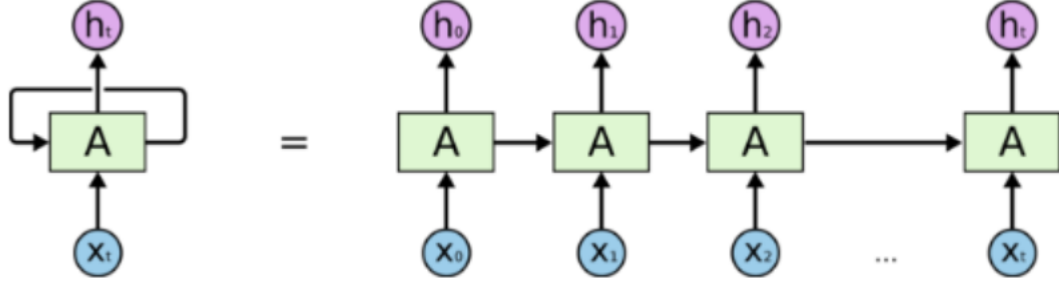
Cümle 1	Yıldız savaşları izlemeli miyim?
Cümle 1	Yıldız savaşları filmini seyretmemi önerir misiniz?
Cümle 2	Vernik seçimi nasıl olmalı?
Cümle 2	Vernik tercih ederken nelere dikkat edilmelidir?
Cümle 3	Avrupaya vize kalkarsa gider misiniz?
Cümle 3	Avrupaya gidiş vizesiz olsa seyahate çıkar mısınız?

Tablo 5.3 Manuel Oluşturulmuş Türkçe Veri Setindeki Bazı Örnekler

5.2 Yazılım Tasarımı

Dil modelleme veya makine çevirisi gibi görevlerde genel olarak Recurrent Neural Networks (RNN) (Şekil 5.3) ve Transformer mimarileri kullanılır. RNN mimarisi, sıralı yapılardan dolayı eğitim esnasında paralelleştirmeye olanak tanımaz, bu yüzden uzun süreli bellek bağımlılıklarını öğrenmek konusunda probleme yol açar. RNN'lerde bulunan bellek sayısının fazla olması, performans olarak daha iyi etki yaratır fakat uzun girdilerin öğrenilmesi sırasında toplu işlem yapmak sınırlandırılmış olur. RNN ile örneğin bir cümledeki son kelimenin hücrelerine erişmek için kelime kelime ilerlemek gerekmektedir. Yapay sinir ağı uzun bir yapıya sahipse hatırlamak birkaç adım

sürebilir, her maskelenmiş durum bir önceki maskelenmiş duruma bağlıdır. Bu sebeple paralelleştirme etki gösteremez. Girdilerin çok uzun olduğu durumlarda, RNN mimarisi uzak konumların içeriklerini art arda unutma eğilimi göstermeye yatkın olabilir.

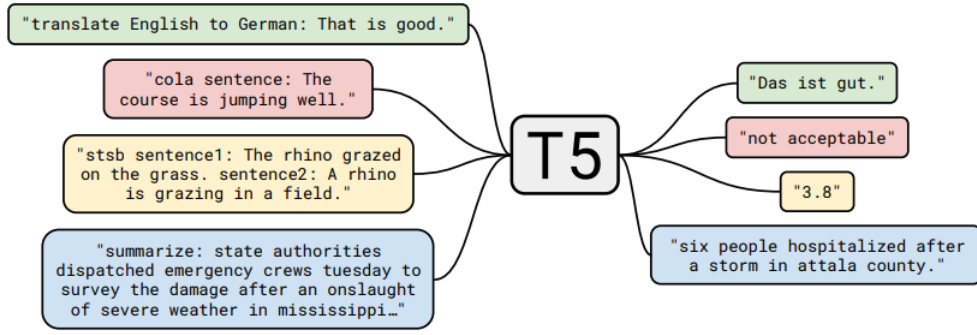


Şekil 5.3 RNN

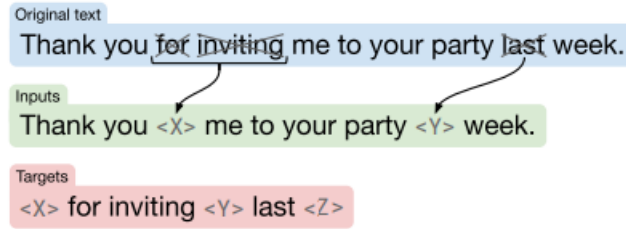
Transformer mimarisi, tekrarlamayı (recurrent) ortadan kaldırıp bunun yerine girdiler ve çıktılar arasındaki ilişkiyi kurmak için öz dikkat (self attention) mekanizması kullanır. Self attention, bir cümlede yer alan herhangi bir sözcüğün diğer sözcüklerle olan ilişkisini belirler. Kodlayıcı(encoder)-kod çözücü(decoder) prensibine çok benzer olarak çalışmaktadır. Bu prensipten farklı olarak kodlayıcıya gelen girdi word-embedding adı verilen nümerik hâldeki kelime vektörlerine dönüştürülmektedir. Bu vektörler de ilk başta self attention mekanizmasına girer. Örnek olarak "On sene sonra bir ilkokula gittim, orada gördüklerim beni şaşırttı." cümlesinde yer alan her kelimenin, cümledeki diğer kelimelerle olan ilişkisini çıkarıp her ilişkiye bir skor belirlenir. Bu cümlede yer alan "ilkokul" ve "orada" sözcükleri arasında kurulan ilişkinin skorunun, diğer ilişkilere göre daha yüksek bir skora sahip olduğu çıkarımına ulaşılır.

T5 modeli, dil ve görev fark etmeksizin metin halinde verilen girdiyi metin halinde bir çıktıya çevirir [4]. Transformers kütüphanesi, içerisinde farklı boyutlarda veri setleri ile eğitilmiş dil modelleri (t5-base, t5-small, t5-large vb.) bulundurulur. Bu modeller, belirli birkaç görevi ön ek olarak kabul edebilmekte ve ince ayar işlemine gerek olmadan çıktı verebilmektedir(Şekil 5.4).

T5 modeli eğitilirken ikili veri setine ihtiyaç duyar, verilen her iki girdi için rastgele ve farklı noktaları maskeler (Şekil 5.5). Böylece karşılıklı iki girdide birbirlerinin içerikleri bulunur. Bu sayede girdilerin ikili hâlleri hangi görev için tasarlandıysa bunu gerçeklemesi çok daha kolay olur.



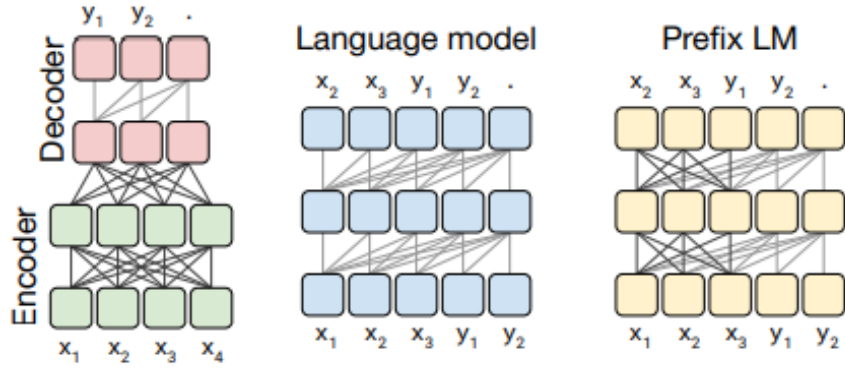
Şekil 5.4 T5 Çalışma Sistemi



Şekil 5.5 T5 Maskelenmiş Cümle Örneği

T5 modeli genel olarak kodlayıcı-kod çözücü prensibi(Şekil 5.6) ile çalışır, girdi olarak verilen cümlelerin parçacıklarına ayrılmasını ise yine Transformers kütüphanesinde bulunan, daha önceden eğitilmiş parçalayıcılar ile yapılabilir. Bu noktada dil modellerinde olduğu gibi farklı boyutlarda kelime hazinesine sahip parçalayıcılar mevcuttur ve kullanılacak görevin içerdiği kelimeleri tanımlarına göre başarı oranları değişir.

Şekil 5.7, Şekil 5.8, Şekil 5.9 ve Şekil 5.10’de yer alan örneklerden yola çıkarak T5 modelinin nasıl çalıştığının incelenmesinin ardından, tarafımızca İngilizce dil modeli oluşturulmuştur. Ardından bu modele, duygu analizi görevi için ince ayar yapılmış ve skorlar elde edilmiştir. Dil modelleme esnasında kullanılan argümanlar Şekil 5.11’da gösterilmiştir.



Şekil 5.6 T5 Encoder-Decoder

```
!pip install --upgrade transformers

[1] from transformers import T5ForConditionalGeneration, T5Tokenizer

[2] import numpy as np

[3] model = T5ForConditionalGeneration.from_pretrained('t5-base')

[4] !pip install sentencepiece
Requirement already satisfied: sentencepiece in /usr/local/lib/python3.7/dist-packages (0.1.95)

[5] tokenizer = T5Tokenizer.from_pretrained('t5-base')

[19] input_ids = tokenizer("translate English to German: The house is wonderful.", return_tensors="pt").input_ids

[20] outputs = model.generate(input_ids)

[21] outputs
tensor([[ 0, 644, 4598, 229, 19250, 5, 1]])

[22] print(tokenizer.decode(outputs[0]))
<pad> Das Haus ist wunderbar.</s>
```

Şekil 5.7 T5 ile İngilizce-Almanca Çeviri

```

▶ input_ids = tokenizer(['stsb sentence1: I am going school at this moment sentence2: Right now I am going school'], return_tensors='pt').input_ids

[44] outputs2 = model.generate(input_ids)

[45] print(tokenizer.decode(outputs2[0]))

<pad> 5.0</s>

[46] input_ids = tokenizer('stsb sentence1: I am going school at this moment sentence2: My brothers name is Mike', return_tensors='pt').input_ids

[47] outputs2 = model.generate(input_ids)

▶ print(tokenizer.decode(outputs2[0]))

<pad> 0.0</s>

[49] input_ids = tokenizer('stsb sentence1: I am going school at this moment sentence2: I have to go school tomorrow', return_tensors='pt').input_ids

[50] outputs2 = model.generate(input_ids)

▶ print(tokenizer.decode(outputs2[0]))

<pad> 2.4</s>

```

Şekil 5.8 T5 ile İki Örnek Cümle Arasındaki Anlam Benzerliği Tespiti

```

[52] input_ids = tokenizer('summarize: Composed entirely of footage shot at the time in various parts of the Soviet Union, the film is a haunting amalgam of official pomp and everyday experience, the double image of a totalitarian government and the people in whose name it ruled.', return_tensors='pt').input_ids

[53] outputs3 = model.generate(input_ids)

▶ print(tokenizer.decode(outputs3[0]))

<pad> a film about the totalitarian government of the time is a haunting amalgam of

[55] input_ids = tokenizer('summarize: While I was going to watch a cinema, I saw my friend from my school which is in Los Angeles.', return_tensors='pt').input_ids

[56] outputs3 = model.generate(input_ids)

▶ print(tokenizer.decode(outputs3[0]))

<pad> a friend of mine was going to watch a movie in los angeles.

```

Şekil 5.9 T5 ile Metin Özetleme

```

[71] input_ids = tokenizer('cola sentence: The movie was the president', return_tensors='pt').input_ids

[72] outputs4 = model.generate(input_ids)

▶ print(tokenizer.decode(outputs4[0]))

<pad> unacceptable</s>

[74] input_ids = tokenizer('cola sentence: The movie was fun', return_tensors='pt').input_ids

[75] outputs4 = model.generate(input_ids)

▶ print(tokenizer.decode(outputs4[0]))

<pad> acceptable</s>

```

Şekil 5.10 T5 ile Cümlelerin Oluşturulma Mümkünlüğü

```

args_dict = dict(
    output_dir="", # path to save the checkpoints
    model_name_or_path=hparam.model,
    tokenizer_name_or_path=hparam.model,
    max_input_length=int(hparam.input_length),
    max_output_length=int(hparam.output_length),
    freeze_encoder=False,
    freeze_embeddings=False,
    learning_rate=1e-5,
    weight_decay=0.0,
    adam_epsilon=1e-8,
    warmup_steps=0,
    train_batch_size=4,
    eval_batch_size=4,
    num_train_epochs=2,
    gradient_accumulation_steps=1,
    n_gpu=1,
    resume_from_checkpoint=None,
    val_check_interval = 1.0,
    n_val=0,
    val_percent_check= 0,
    n_train=-1,
    n_test=-1,
    early_stop_callback=False,
    fp_16=False, # if you want to enable 16-bit training then install apex and set this to true
    opt_level='O1', # you can find out more on optimisation levels here https://nvidia.github.io/apex/amp.html#opt-levels-and-properties
    max_grad_norm=1.0, # if you enable 16-bit training then set this to a sensible value, 0.5 is a good default
    seed=101,
)

```

Şekil 5.11 T5 Dil Modelleme İçin Kullanılan Argümanlar

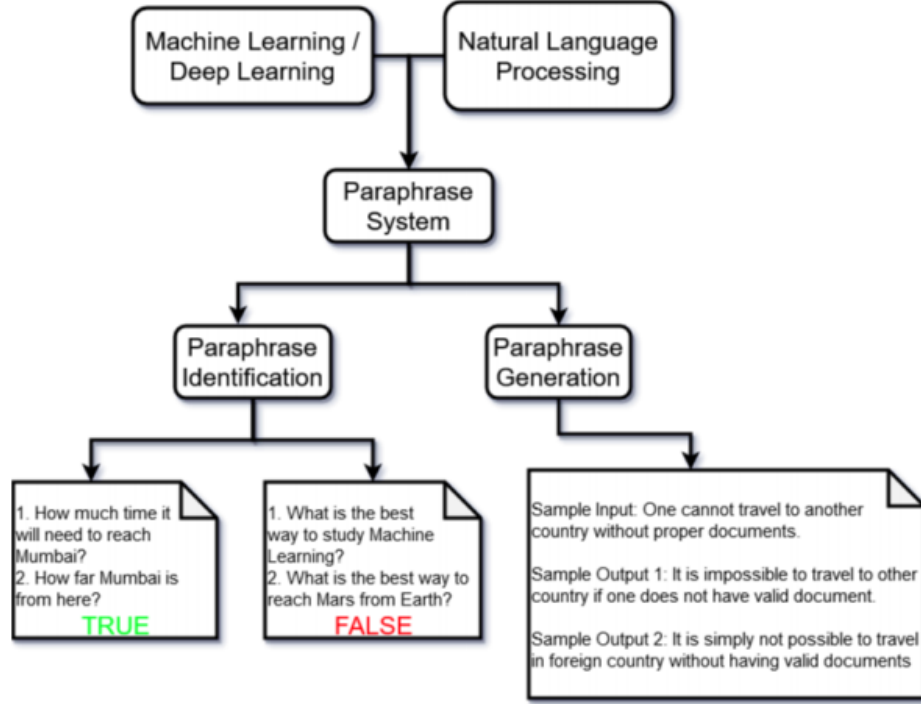
Model, eğitilirken ve ince ayar yapılırken aynı optimize etme algoritması olarak "Adam" ve aynı parçalayıcı olarak "T5Tokenizer.from_pretrained('t5-base')" kullanılmıştır. (Şekil 5.12)

```

[ ] args_dict = dict(
    data_dir="", # path for data files
    output_dir="", # path to save the checkpoints
    model_name_or_path='t5-base',
    tokenizer_name_or_path='t5-base',
    max_seq_length=512,
    learning_rate=3e-4,
    weight_decay=0.0,
    adam_epsilon=1e-8,
    warmup_steps=0,
    train_batch_size=8,
    eval_batch_size=8,
    num_train_epochs=2,
    gradient_accumulation_steps=16,
    n_gpu=1,
    early_stop_callback=False,
    fp_16=False, # if you want to enable 16-bit training then install apex and set this to true
    opt_level='O1', # you can find out more on optimisation levels here https://nvidia.github.io/apex/amp.html#opt-levels-and-properties
    max_grad_norm=1.0, # if you enable 16-bit training then set this to a sensible value, 0.5 is a good default
    seed=42,
)

```

Şekil 5.12 T5 Duygu Analizi İnce Ayar İşleminde Kullandığımız Argümanlar



Şekil 5.13 Aynı Anlamlı Cümle Tespiti ve Aynı Anlamlı Cümle Üretimi

Proje konusu olan Türkçe dilinde aynı anlamlı cümle üretimi için Türkçe kelime hazinesine sahip bir parçalayıcının girdi olarak verilen metni ayıklaması, daha sonra ayrı ayrı elde edilen kelime köklerinin kodlayıcıya verilmesi, yine aynı şekilde Türkçe ve göreve uygun kodlanmış bir başka metnin referans verilmesi gerekmektedir. Dil eğitimi tamamlandıktan sonra ince ayar işlemine geçilecek ve bu sefer hazırlanmış ikili soru cümleleri ile yapay sinir ağının ağırlıkları hesaplanacak, hangi kelimelerin yerine hangi kelimelerin kullanılabileceği, hangi kelimelerin hangi kelimeler ile daha sık geçtiği tespit edilebilecektir.

T5 dil modelinin çok sayıda dildeki versiyonu olan MT5 ile proje, Türkçe dilinde gerçekleştirilmiştir. MT5 dil modelinin çalışma prensibi, T5 ile bire bir denecek kadar benzerdir. T5'in veri seti C4 iken MT5'in oluşumunda C4'ün çok sayıda dildeki versiyonu olan MC4 (Multi-lingual Colossal Clean Crawled Corpus) veri seti kullanılmıştır. Her dilin kuralları farklı olacağından ve her dilde bulunan veri miktarı aynı olmayacağından parametreler konusunda iki model değişiklik göstermiştir.

Türkçe aynı anlamlı cümle üretimi işleminde önceden eğitilmiş MT5 dil modeli olarak mt5-base, tokenizer olarak da mt5-tokenizer kullanılmıştır. Türkçe aynı anlamlı cümle üretimi için kullanılan eğitim fonksiyonunun parametreleri Şekil 5.14'te verilmiştir. Train ve validation veri seti 2/1 oranında ayrılmıştır. Veri setlerinin büyüklüklerine göre her epoch için geçen süre değişmekle birlikte her modelin 3 epoch'luk eğitimi

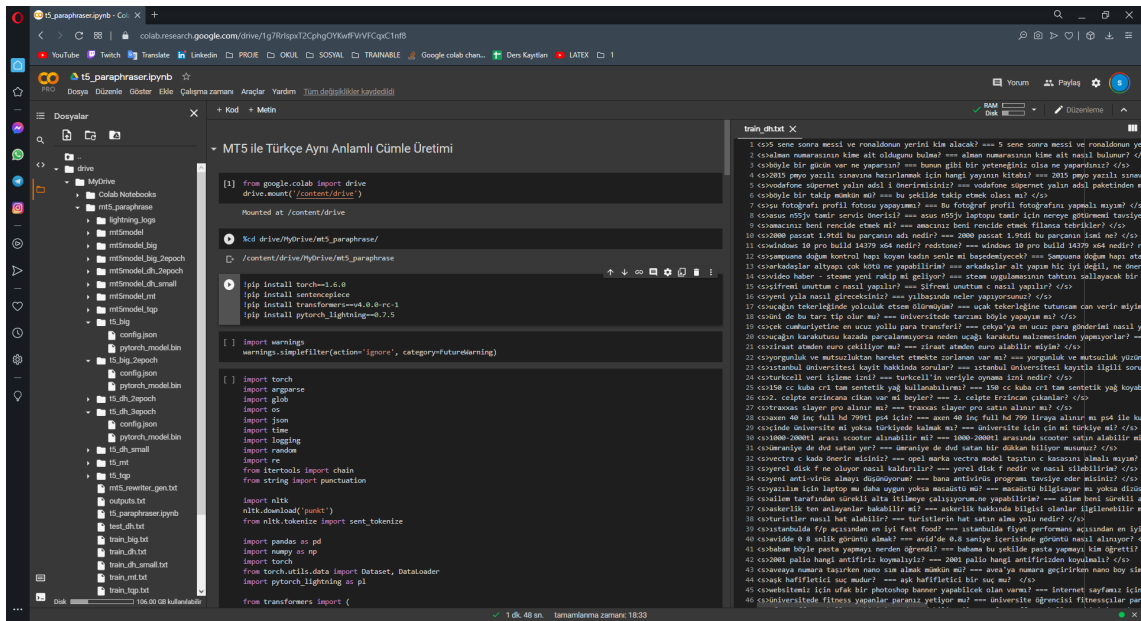
yaklaşık 3-5 saat aralığında sürmüştür. mt5-base modelinin daha büyük olan versiyonu olan mt5-large modeli ile eğitim denenmiştir fakat Google Colab'deki kısıtlamalar dolayısıyla eğitim tamamlanamamıştır. Ek olarak maximum sequence length ve batch size değerleri de değiştirilerek eğitim denenmiştir fakat yine Google Colab'deki kısıtlamalar buna engel olmuştur.

```
args_dict = dict(  
    data_dir="",  
    output_dir="mt5model_big_3epoch/",  
    model_name_or_path="google/mt5-base",  
    tokenizer_name_or_path="google/mt5-base",  
    max_seq_length=150,  
    learning_rate=3e-4,  
    weight_decay=0.0,  
    adam_epsilon=1e-8,  
    warmup_steps=0,  
    train_batch_size=4,  
    eval_batch_size=4,  
    num_train_epochs=3,  
    gradient_accumulation_steps=16,  
    early_stop_callback=False,  
    n_gpu=1,  
    fp_16=False,  
    opt_level='O1',  
    max_grad_norm=1.0,  
    seed=42,  
)
```

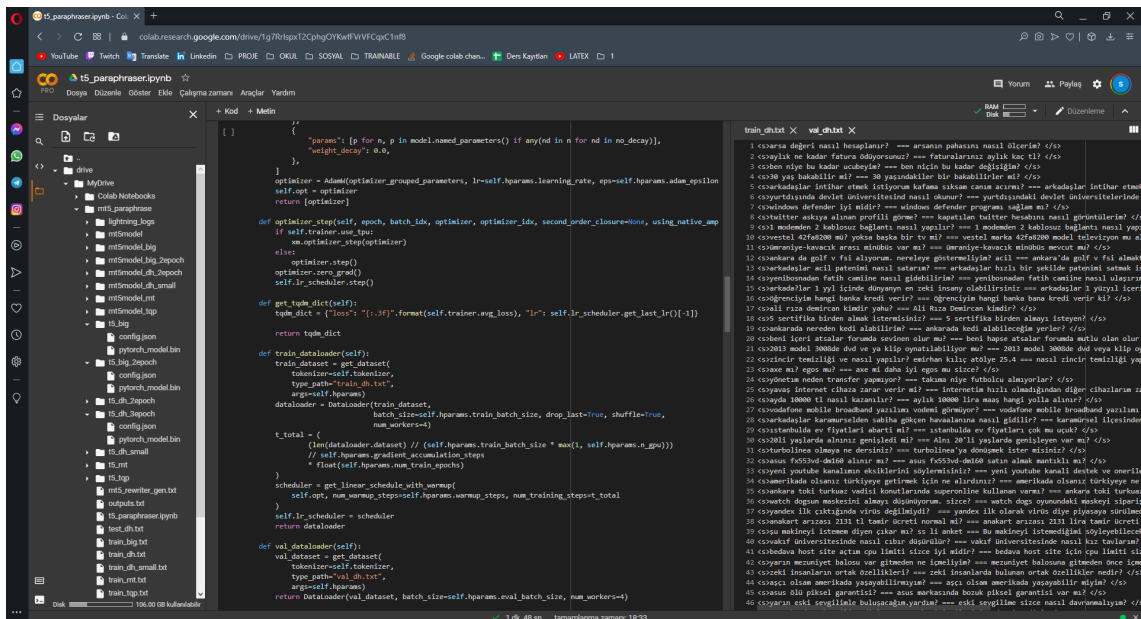
Şekil 5.14 Türkçe Aynı Anlamlı Cümle Üretimi İnce Ayar Fonksiyonundaki Argümanlar

5.3 Girdi-Çıktı Tasarımı

Türkçe aynı anlamlı cümle üretimi projesi, Opera tarayıcısından Google Colab'e girerek çalıştırılmaktadır. Model ve veri seti Google Drive'a kaydedilmektedir ve oradan yüklenmektedir.



Şekil 5.15 Modellerin, Kodların ve Eğitim Veri Kümesini İçeren Ekran Görüntüsü



Şekil 5.16 Validation Veri Setini İçeren ve Kodda Dizin Gösterilen Ekran Görüntüsü

Eğitimi tamamlanan model Google Drive'a kaydedilmektedir. Modelin testlerinin görülebilmesi için Google Drive'daki modelin bulunduğu klasör dizini koda eklendikten sonra modelin yüklenmesi ve ardından daha önceden eğitilmiş olan "mt5-base" tokenizer yüklenmektedir.

```
import torch
from transformers import T5ForConditionalGeneration, T5Tokenizer
from transformers import MT5ForConditionalGeneration, AutoTokenizer
print("Imported!")

def set_seed(seed):
    torch.manual_seed(seed)
    if torch.cuda.is_available():
        torch.cuda.manual_seed_all(seed)

set_seed(42)
print('Seed')
model = MT5ForConditionalGeneration.from_pretrained('./t5_dh_3epoch')
print("Models loaded")
tokenizer = AutoTokenizer.from_pretrained('google/mt5-base')
print("Downloaded all")

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print("device ", device)
model = model.to(device)
```

[C]: Imported!
Seed
Models loaded
Downloading: 100% ██████████ 639/639 [00:00<00:00, 832B/s]
Downloading: 100% ██████████ 4.31M/4.31M [00:03<00:00, 1.34MB/s]
Downloading: 100% ██████████ 65.0/65.0 [00:01<00:00, 53.4B/s]
Downloading: 100% ██████████ 378/378 [00:00<00:00, 769B/s]
Downloaded all
device: cuda

Şekil 5.17 Tarafımızca İnce Ayar Yapılmış Modelin ve Daha Önce Eğitilmiş Tokenizer'ın Yüklenmesi

Model ve tokenizer yüklenmesinin ardından kullanıcının konsola gireceği girdi cümle alınarak Şekil 5.18'te generate fonksiyonu aynı anlama gelen cümlelerin üretilmesi sağlanmıştır (Şekil 5.19).

```
question=input("CÜMLENİZİ YAZINIZ:")
ques2list='cúmleniz: '+question+' </s>'
to_model.append(ques2list)

for i in range(len(to_model)):
    to_model[i]=to_model[i].lower()

for sentence in to_model:
    max_len = 150
    encoding = tokenizer.encode_plus(sentence, pad_to_max_length=True, return_tensors="pt")
    input_ids, attention_masks = encoding["input_ids"].to(device), encoding["attention_mask"].to(device)

    beam_outputs = model.generate(
        input_ids=input_ids, attention_mask=attention_masks,
        do_sample=True,
        max_length=150,
        top_k=50,
        top_p=0.95,
        early_stopping=True,
        num_return_sequences=5
    )
    print ("Cúmleniz: "+question+"\n")

    final_outputs = []
    for beam_output in beam_outputs:
        sent = tokenizer.decode(beam_output, skip_special_tokens=True, clean_up_tokenization_spaces=True)
        if sent.lower() != sentence.lower() and sent not in final_outputs:
            final_outputs.append(sent)
    for i, final_output in enumerate(final_outputs):
        print("{}\t{}".format(i+1, final_output))
    final_outputs = final_outputs + [sentence]
```

... CÜMLENİZİ YAZINIZ: saat kaçta okulda olayım?]

Şekil 5.18 Kullanıcı Tarafından Girdi Cümle Yazılması

```

question=input("CÜMLENİZİ YAZINIZ:")
ques2list='cumlenez: '+question+ ' </s>'
to_model.append(ques2list)

for i in range(len(to_model)):
    to_model[i]=to_model[i].lower()

for sentence in to_model:
    max_len = 150
    encoding = tokenizer.encode_plus(sentence, pad_to_max_length=True, return_tensors="pt")
    input_ids, attention_masks = encoding["input_ids"].to(device), encoding["attention_mask"].to(device)

    beam_outputs = model.generate(
        input_ids=input_ids, attention_mask=attention_masks,
        do_sample=True,
        max_length=150,
        top_k=50,
        top_p=0.95,
        early_stopping=True,
        num_return_sequences=5
    )
    print ("Cümleniz: "+question+"\n")

    final_outputs = []
    for beam_output in beam_outputs:
        sent = tokenizer.decode(beam_output, skip_special_tokens=True, clean_up_tokenization_spaces=True)
        if sent.lower() != sentence.lower() and sent not in final_outputs:
            final_outputs.append(sent)
    for i, final_output in enumerate(final_outputs):
        print("\t\t{}\t\t".format(i+1, final_output))
        final_outputs = final_outputs + [sentence]

```

CÜMLENİZİ YAZINIZ:saat kaçta okulda olayım?
 Cümleniz: saat kaçta okulda olayım?

- 1 okulda hangi saatte çalışayım?
- 2 okulda saat kaçta kalınmalı?
- 3 okulda saat kaçta gitmemi önerirsiniz?
- 4 Okulda saat kaçta bulunmalı?
- 5 saat kaçta okula gidilir?

Şekil 5.19 Üretilen Aynı Anlamalı Cümle Çıktıları

```

to_model = []
#to_model = []
#to_ref = []
#to_model_1=[]
#to_model_2=[]
#to_total_ref=[]
#res = 0
#i=0
with open("test_dh.txt") as inf:
    for line in inf.readlines():
        s1, s2 = line.rstrip().split(" --- ")
        sentence = "cumlenez: " + s1 + " </s>"
        to_model.append(sentence)
        #to_model.append(s1)
        #to_ref.append(s2)
    ...

for sentence in to_model:
    to_model_1.append(sentence.split())
    to_model_2.append(to_ref[1].split())
    to_total_ref.append([to_model_1[i],to_model_2[i]])
    i += 1
    ...
#to_total_ref = [to_total_ref]
...

```

mt5_rewriter.gen.txt test_dh.txt

- 1 ben bu kıza aşık mıyım? --- kadını seviyor muyum?
- 2 çeşitli organizasyonlar için servis ücretleri? --- organizasyonların servis ücretleri kaçtır?
- 3 araç seçimi için yardım --- araç seçeceğim, yardımı olur musunuz?
- 4 soruy sızın yoldaki nasıl biliyorsunuz? --- soruy sızın yolu nasıl biliyorsunuz?
- 5 arkadaşına hediye aldım oyunu etkinleştiriyor? --- arkadaşına hediye yolladım oyunu çalıştıramıyor?
- 6 yerli aracın motorunu neden bir üretiyorsunuz? --- yerli aracın motoru neden burada üretilmiyor?
- 7 tnet adı yalan internet? --- tnet adı yalan internet nasıldır?
- 8 üstteki oyunun nicki ne/neyi çağırıyor? --- üstteki oyunun nicki sana neyi hatırlatıyor?
- 9 asus dsl-88u wifi ayarı nasıl olmalı? --- asus dsl-88u wifi ayarının doğrusu nedir?
- 10 estim hattısı olupta sigara için var mı? --- estim olduğu halde sigara kullanmaya devam eder var mı?
- 11 yerde 500 bin tl bulmanız napardınız? --- sokakta yukü bir miktar para bulmanız onunla ne yaparsınız?
- 12 hakalın çörebilen çıkacak mı? --- inceleyelim yapabilen olacak mı?
- 13 üniversitede vize-final sistemi nasıl? --- üniversitede sınav sistemi nasıldır?
- 14 996 smart box nasıl sıtıcı? --- 996 smart box hakkındaki yorumlarınız nedir?
- 15 ıkten sona doğru cep telefonlarınız? --- bastan sona doğru telefonlarınız nelerdi?
- 16 ankara a2 ehliyeti için kurslar --- ankara a2 ehliyeti için kurs önerir misiniz?
- 17 4 tl 1 dolar olurmu? --- 4 lira, 1 dolar olur mu sizce?
- 18 bana yardım edebilir misiniz? --- bir el atar mısınız?
- 19 ankara otizm e yakın cemaatlı özel yurt var mı? --- ankara otisme yakın cemaatlı özel yurt?
- 20 antifiyir değişim konusunda yardım --- antifiyir değişimyle ilgili yardımca edebilir misiniz?
- 21 vodafone dan taşınmış bozup hat taşıyacağım, martıklı mı? --- vodafone sözleşmesini feshedip başka operatöre geçme
- 22 alp ismi ile akrostij yapar mısınız? --- alp ismi ile akrostij?
- 23 asosyal misiniz? normal mi olmak istiyorsunuz? --- sosyal değil misiniz, normal olmak mı istersiniz?

Şekil 5.20 Test Cümlelerinin Dosyadan Verilmesi

6 Uygulama

T5 ile İngilizce dilinde duygu analizi işleminde büyük veri setine uygulanan ince ayar işlemi sonucu test veri setindeki cümlelerin duygu analizinin sonuç çıktıları Şekil 6.1’de ifade edilmiştir.

```
Review: From the beginning Til There Was You was on the right track setting up for the big finish
where it would all come together But the thing is it didn t I found the ending extremely
disappointing but maybe in someway it was the right ending a little more realistic you could say
Judge for yourself

Actual sentiment: negative
Predicted sentiment: negative
-----

Review: I agree with most of the critics above More yet I was shocked by the presentation of the
love scenes with the homosexual couple Why because while they the director the producers didnt have
any compulsion whatsoever in presenting the different heterosexual couples in the most passionate
embraces including nudity and super closeups of French kissing and all sorts of nude contortions in
bed completely unnecessary in their length and in the story when the moment came to show the same
experiences with the homosexual couple they only dare to go as far as an excruciatingly painful hug
almost among scholarly giggles with two very nervous actors So in reality the makers of this film
found homosexuality to be UNNATURAL as one of the characters says in some scene What a difference
with the Spanish cinema I remember being at the projection of an Almodovar film in an Italian cinema
in Rome and being completely amazed at the total lack of reaction from the Italian audience they
were afraid to have a reaction when in Spain people would fall down from their seats laughing at all
the risqué situations and fabulous Almodovar wit and flair Obviously in Italy there are dark forces
in its history that impedes the free manifestation of some very normal and natural emotions Pity I
must add that I was quite surprised to find that this same comment was censured by another
correspondent Its very bad and dangerous when we cannot be allowed by the narrowness of others to
express our opinions about certain matters Where is freedom of speech I dont know if that censor
will approve of the changes I was forced to make in this comment and I hope he wont receive the same
treatment from some other narrow minded judge Pity again

Actual sentiment: negative
Predicted sentiment: positive
```

Şekil 6.1 Cümlelere Yapılan Tahminler

T5, BoW ve Seq2Seq ile İngilizce dilinde aynı anlamlı cümle üretimi işleminde verilen cümlelere ait üretilen çıktılar Tablo 6.1’de verilmiştir.

Referans	Is this the start of an aviation revolution?
BoW	Was this the start of an aviation revolution?
Seq2Seq	Is the world aviation revolution?
T5	Is the first runway in the world about to go into air-strips?
Referans	What’s your favorite sleeping position?
BoW	What is your favorite position position?
Seq2Seq	What is your favorite sleeping position?
T5	What are some of your favorite sleeping positions?
Referans	What’s the best career decision you’ve ever made?
BoW	What ’s the best decision you ever made?
Seq2Seq	What is the best career decision you ever made?
T5	What was the best career decision of your life?
Referans	Can I make my own hand sanitizer?
BoW	Can i make make sanitizer?
Seq2Seq	Can I make my own hand sanitizer?
T5	What are some easy ways to make personal hand sanitizer?
Referans	Do you offer business insurance and protection plans?
BoW	Do you have insurance and protection coverage for businesses?
Seq2Seq	Do you offer business insurance plans?
T5	Do you sell business insurance? How does it work?
Referans	What should I do if I have problems with my truck?
BoW	What should i do if i have problems with my truck?
Seq2Seq	How do I get my truck?
T5	What should we do if my truck is leaking?

Tablo 6.1 İngilizce Aynı Anlamlı Cümle Üretimi Sonuçları

Türkçe dilinde aynı anlamlı cümle üretiminde ise uygulamanın çalışmasına ilişkin ekran çıktıları aşağıdaki şekillerdeki gibidir.

```

1) from google.colab import drive
drive.mount('/content/drive')

2) %cd drive/MyDrive/mt5_paraphrase/

!pip install torch=1.6.0
!pip install sentencepiece
!pip install transformers==4.0.0-rc1
!pip install pytorch_lightning==0.7.5

Requirement already satisfied: torch==1.6.0 in /usr/local/lib/python3.7/dist-packages (1.6.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from torch==1.6.0) (1.19.5)
Requirement already satisfied: future in /usr/local/lib/python3.7/dist-packages (from torch==1.6.0) (0.18.2)
Requirement already satisfied: sentencepiece in /usr/local/lib/python3.7/dist-packages (0.1.90)
Requirement already satisfied: transformers==4.0.0-rc1 in /usr/local/lib/python3.7/dist-packages (4.0.0rc1)
Requirement already satisfied: requests in /usr/local/lib/python3.7/dist-packages (from transformers==4.0.0-rc1) (2.23.0)
Requirement already satisfied: filelock in /usr/local/lib/python3.7/dist-packages (from transformers==4.0.0-rc1) (3.0.12)
Requirement already satisfied: regex==2019.12.17 in /usr/local/lib/python3.7/dist-packages (from transformers==4.0.0-rc1) (2019.12.20)
Requirement already satisfied: packaging in /usr/local/lib/python3.7/dist-packages (from transformers==4.0.0-rc1) (20.9)
Requirement already satisfied: tqdm==4.27 in /usr/local/lib/python3.7/dist-packages (from transformers==4.0.0-rc1) (4.41.1)
Requirement already satisfied: tokenizers==0.8.4 in /usr/local/lib/python3.7/dist-packages (from transformers==4.0.0-rc1) (0.9.4)
Requirement already satisfied: sacremoses in /usr/local/lib/python3.7/dist-packages (from transformers==4.0.0-rc1) (0.0.45)
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from transformers==4.0.0-rc1) (1.19.5)
Requirement already satisfied: joblib in /usr/local/lib/python3.7/dist-packages (from requests>transformers==4.0.0-rc1) (2.10)
Requirement already satisfied: chardet==3.0.2 in /usr/local/lib/python3.7/dist-packages (from requests>transformers==4.0.0-rc1) (3.0.4)
Requirement already satisfied: certifi==2019.10.17 in /usr/local/lib/python3.7/dist-packages (from requests>transformers==4.0.0-rc1) (2019.10.17)
Requirement already satisfied: urllib3==1.25.8 in /usr/local/lib/python3.7/dist-packages (from requests>transformers==4.0.0-rc1) (1.24.3)
Requirement already satisfied: idna==2.8 in /usr/local/lib/python3.7/dist-packages (from packaging>transformers==4.0.0-rc1) (2.4)
Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from sacremoses>transformers==4.0.0-rc1) (1.15.0)
Requirement already satisfied: click in /usr/local/lib/python3.7/dist-packages (from sacremoses>transformers==4.0.0-rc1) (7.1.2)
Requirement already satisfied: joblib in /usr/local/lib/python3.7/dist-packages (from sacremoses>transformers==4.0.0-rc1) (1.0.1)
Requirement already satisfied: pytorch-lightning==0.7.5 in /usr/local/lib/python3.7/dist-packages (0.7.5)
Requirement already satisfied: torch==1.6.0 in /usr/local/lib/python3.7/dist-packages (from pytorch-lightning==0.7.5) (1.6.0)
Requirement already satisfied: future==0.15.2 in /usr/local/lib/python3.7/dist-packages (from pytorch-lightning==0.7.5) (0.18.2)
Requirement already satisfied: sacremoses==0.0.45 in /usr/local/lib/python3.7/dist-packages (from pytorch-lightning==0.7.5) (0.0.45)
Requirement already satisfied: tensorboard-data-server==0.7.0 in /usr/local/lib/python3.7/dist-packages (from tensorboard>pytorch-lightning==0.7.5) (0.7.0)
Requirement already satisfied: wheel==0.33.1 in /usr/local/lib/python3.7/dist-packages (from tensorboard>pytorch-lightning==0.7.5) (0.33.1)
Requirement already satisfied: setuptools==44.0.0 in /usr/local/lib/python3.7/dist-packages (from tensorboard>pytorch-lightning==0.7.5) (44.0.0)
Requirement already satisfied: grpcio==1.24.3 in /usr/local/lib/python3.7/dist-packages (from tensorboard>pytorch-lightning==0.7.5) (1.34.1)
Requirement already satisfied: google-auth-oauthlib==0.4.1 in /usr/local/lib/python3.7/dist-packages (from tensorboard>pytorch-lightning==0.7.5) (0.4.1)

```

Şekil 6.2 Google Drive Bağlantısı Kurulması

```

import torch
from transformers import MT5ForConditionalGeneration, MT5Tokenizer
from transformers import MT5ForConditionalGeneration, AutoTokenizer

print("Imported")

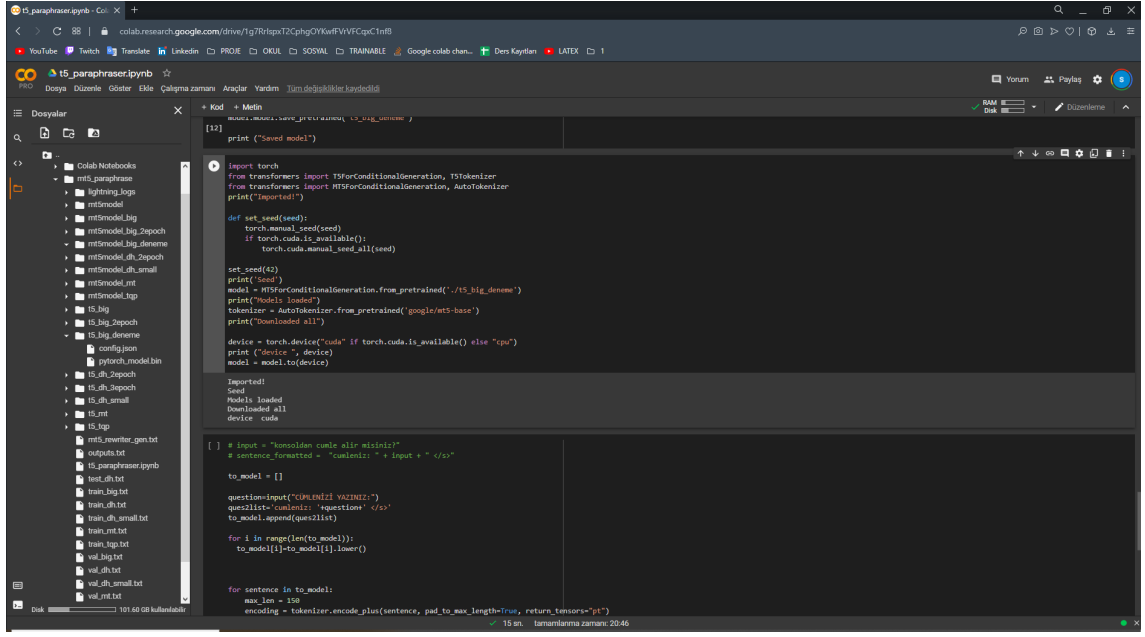
def set_seed(seed):
    torch.manual_seed(seed)
    if torch.cuda.is_available():
        torch.cuda.manual_seed_all(seed)

set_seed(42)
print("Seed")
model = MT5ForConditionalGeneration.from_pretrained('/content/drive/MyDrive/mt5_model.pth')
tokenizer = AutoTokenizer.from_pretrained('/content/drive/MyDrive/mt5_model.pth')

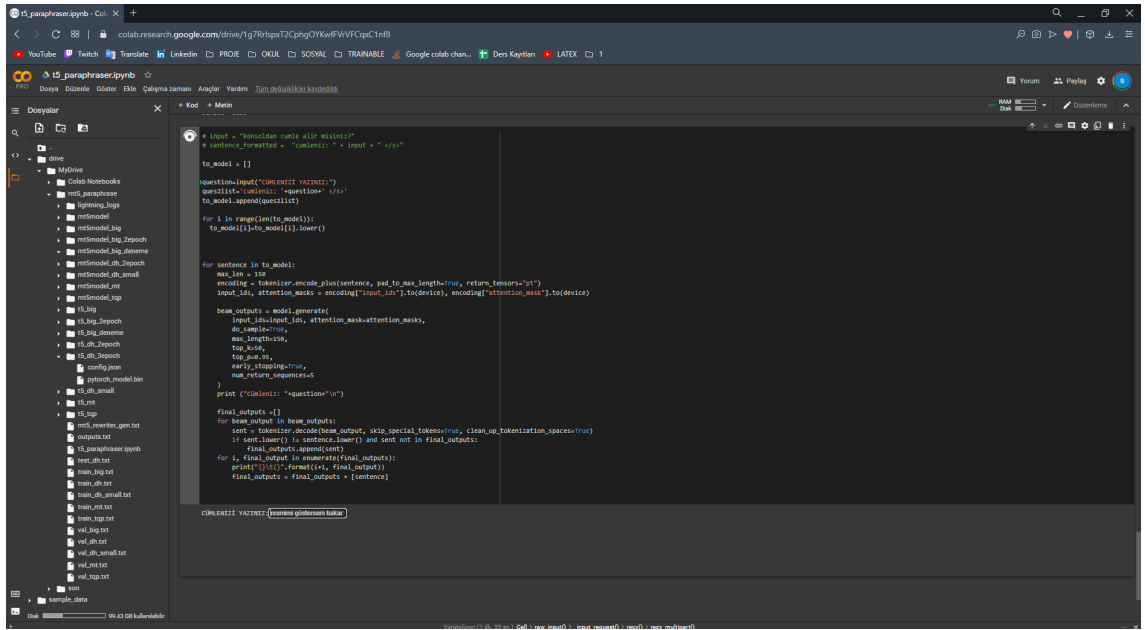
print("Model loaded")

```

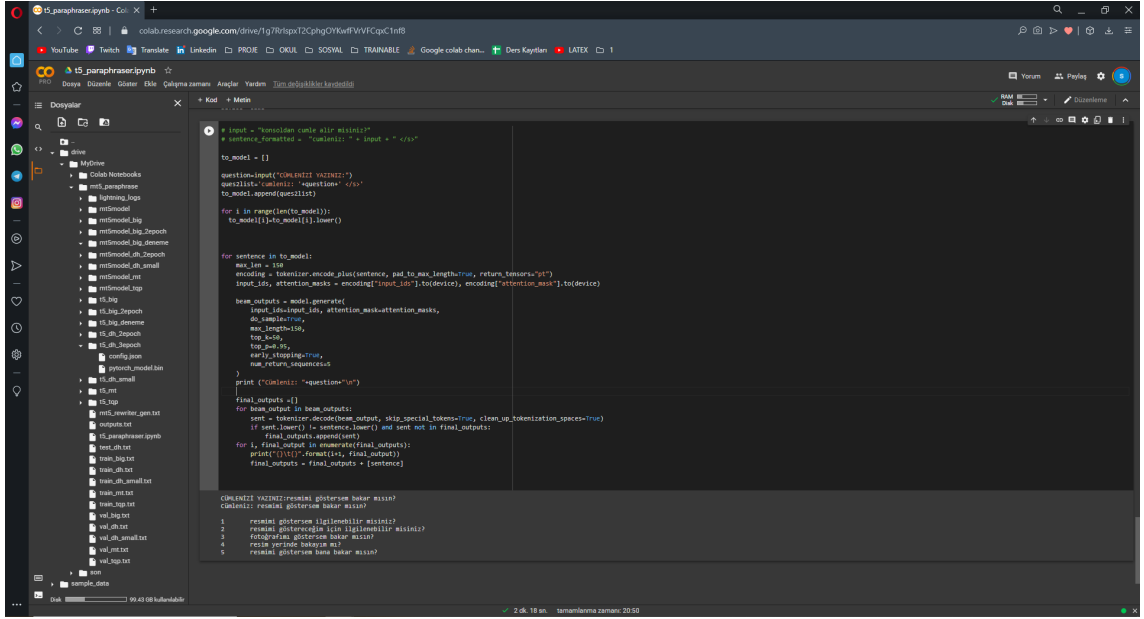
Şekil 6.3 Gerekli Kütüphanelerin Kurulumu



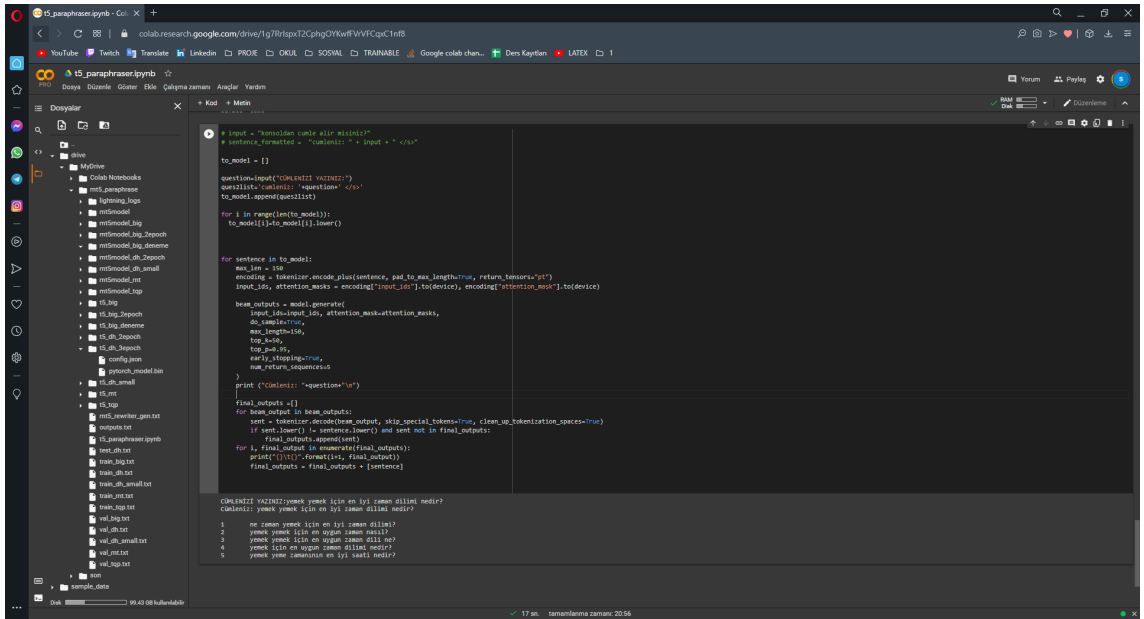
Şekil 6.6 Eğitilen Modelin Yüklenmesi



Şekil 6.7 Konsoldan Cümle Girişi



Şekil 6.8 Cümlelere Yapılan Tahminler - 1



Şekil 6.9 Cümlelere Yapılan Tahminler - 2

7 Deneyisel Sonuçlar

Sisteme verilen uç noktalara ait girdiler ve alınan sonuçların ekran çıktıları aşağıdaki şekillerde gösterilmiştir.

```
CÜMLENİZİ YAZINIZ:otobüs kullanarak şehir dışına gittim.  
Cümleniz: otobüs kullanarak şehir dışına gittim.  
  
1      otobüs kullanarak şehir dışına gittim  
2      otobüs kullanarak şehir dışına çıktım  
3      otobüs kullanarak şehre giden yol tarifi verdim  
4      otobüs kullanarak şehir dışına gittim?  
5      otobüs ile şehir dışına gitmek istiyorum.
```

Şekil 7.1 Uç Durum: Soru Cümlesi Olmayan Cümle

```
CÜMLENİZİ YAZINIZ:what were you doing while i was playing games?  
Cümleniz: what were you doing while i was playing games?  
  
1      while you played games?  
2      what are you playing games and whom was?  
3      what was you doing while i play games?  
4      what were you doing while i play games?  
5      what were you doing while i was playing games?
```

Şekil 7.2 Uç Durum: İngilizce Soru Cümlesi

```
CÜMLENİZİ YAZINIZ:sabahları erken kalkmayı tercih etmiyorum.  
Cümleniz: sabahları erken kalkmayı tercih etmiyorum.  
  
1      sabahları erken kalkmayı tercih etmem gerekir mi?  
2      sabahları erken kalkmayı tavsiye etmiyorum?  
3      sabahları erken kalkmak kabul etmem lazım?  
4      sabahları erken kalkmayı tercih etmem gerekiyor fakat yine de vazgeçmem gerekiyor.
```

Şekil 7.3 Uç Durum: Soru Cümlesi Olmayan Olumsuz Cümle

```
CÜMLENİZİ YAZINIZ:8734689404309  
Cümleniz: 8734689404309  
  
1      8734689404309?  
2      873468940439  
3      8734689404309 hakkında  
4      8734689404309
```

Şekil 7.4 Uç Durum: Yalnızca Sayı Girişi


```
CÜMLENİZİ YAZINIZ:metinbinbir@icloud mail adresi kime aittir?  
Cümleniz: metinbinbir@icloud mail adresi kime aittir?  
  
1      metinbinbir@icloud mail adresi kime aittir?  
2      metinbinbir@icloud adresi nereye ait?  
3      metinbinbir@icloud mail adresi kimin aittir?  
4      metinbinbir@icloud mail adresi kime ait?
```

Şekil 7.5 Uç Durum: Özel Karakter İçeren Girdi Durumu

```
CÜMLENİZİ YAZINIZ:merhaba  
Cümleniz: merhaba  
  
1      merhaba, dostlar bana yardım eder misiniz?  
2      merhaba beyler, ne düşünüyorsunuz?  
3      merhaba arkadaşlar, yardımcı olabilir misiniz?  
4      merhaba bilen var mı?  
5      dostlar yardımcı olur musunuz?
```

Şekil 7.6 Uç Durum: Tek Kelimelik Girdi Durumu

```
CÜMLENİZİ YAZINIZ:yarın okulda  
Cümleniz: yarın okulda  
  
1      okulda yarın mı eğitim yapmak için yardımcı olursunuz?  
2      okulda yarın hangi öğretmen olmak ister?  
3      yarın okulda okula gidecek var mı?  
4      okulu yarın nasıl kapatayım?  
5      yarın okulda sınav hazırlıkları yapacaklar mı?
```

Şekil 7.7 Uç Durum: Eksik Cümle

```
CÜMLENİZİ YAZINIZ:ben kalem elma armut araba.  
Cümleniz: ben kalem elma armut araba.  
  
1      beni kalem elma armut araba mı önerirsiniz?  
2      ben elma armut araba ben kalem elma armut araba ben?  
3      beni kalem elma armut arabamı seçeceğim  
4      ben kalem elma armut arabası ben çok güzeldir  
5      kalem elma armut araba mı benim için iyidir?
```

Şekil 7.8 Uç Durum: Yükleme İçermeyen Cümle

Veri seti tasarımı kısmındaki Tablo 5.14'teki DH-S, DH, MT, TQP ve BIG kodlu 5 adet veri setiyle MT5 modeline ince ayar yapılmış modellere ait girdiler ve ürettikleri aynı anlamlı cümlelerin çıktıkları Tablo 7.1, Tablo 7.2, Tablo 7.3 ve Tablo 7.4'te verilmiştir.

Referans	yaşadığınız şehirde en güzel manzaraya sahip yer neresidir?
DH-S	Şehirdeki en güzel manzaraya sahip mekan neresi?
DH	hayatını sürdürdüğüm şehirde en güzel manzara alan bölge nedir?
TQP	Senin yaşadığın şehirde en iyi manzaraya sahip yer nedir?
MT	Yaşadığınız şehirde en güzel manzara alabileceğin yer neresi?
BIG	yaşadığınız şehirde en güzel manzaraların neler olduğunu biliyor musun?
Referans	hangi saatlerde kitap okumak daha verimlidir?
DH-S	hangi saatlerde kitap okunabilir?
DH	hangi saatlerde kitap okumak daha yararlı?
TQP	Ne zaman kitap okumak daha güvenli?
MT	Özel kitap okumak için hangi saatler uygun olacak?
BIG	kitap okumak için en uygun saat nedir?
Referans	son model arabalar çevreye daha mı az zarar veriyor?
DH-S	son model arabaların çevreye zararları daha mı az oluyor?
DH	son arabaların çevreye etkisi az mı?
TQP	Son model araba çevreye daha az zarar verir mi?
MT	Son model arabalar çevreye zarar veriyor mu?
BIG	yeni arabalar çevreye daha az zarar veriyor mu?
Referans	insanlar boş zamanlarını neler yaparak daha verimli hale getirebilir?
DH-S	insanlar boş zamanlarını iyileştirebilecek şeyler?
DH	insanlar boş zamanlarını nasıl iyileştirir?
TQP	İnsanlar boş zamanını nasıl geçirir?
MT	boş zamanımızı en verimli hale getirmenin en iyi yolları nelerdir?
BIG	insanlar boş zamanlarını nasıl daha verimli yapabilir?
Referans	ülkemizdeki konut fiyatlarında ileriki zamanlarda nasıl bir artış bekleniyor?
DH-S	ülkemizdeki konut fiyatlarında ilerleyen haftalarda ne kadar artış olur?
DH	ülkemizdeki konut fiyatları gelecek tarihte nasıl artacak?
TQP	ülkemizdeki konut fiyatlarına kısa bir zaman için biraz artış bekliyor musun?
MT	Fransa'da konut fiyatlarının ileriki zamanlarda nasıl artacağını düşünüyorsunuz?
BIG	ülkemizdeki konut fiyatları gelecekte nasıl artacaktır?

Tablo 7.1 Farklı Veri Setlerine Göre Üretilen Aynı Anlamalı Cümle Çıktıları - 1

Referans	okul bahçesinde top oynayan çocukları gördünüz mü?
DH-S	Okul bahçesinde top oynayan çocuklar var mı?
DH	okul bahçesinde top oynayan çocuklardan haberdar mısınız?
TQP	okula gidene katlanan çocukları gördün mü?
MT	Okul bahçesinde top oynananlar neler?
BIG	okul bahçesinde top oynayan çocuklardan haberiniz var mı?
Referans	salgın hastalıkların etkisi ne kadar sürecek?
DH-S	salgın hastalıklarının etkisi kaçır?
DH	salgın hastalıkları kaç yılda etki eder?
TQP	sağlık hastalıklarının sebebi ne kadar sürecek?
MT	salgın hastalıkları nasıl sürecek?
BIG	salgın hastalıklar ne kadar zaman alır?
Referans	sıcak havalarda parkta yürüyen insanlar görüyor musunuz?
DH-S	parkta sıcak havalarda yürüyen var mı?
DH	sıcak havalarda parkta duran insanlara rastlanır mı?
TQP	Sıcak havalarda parkta yürüyüşe çıkan insanlar görüyor musun?
MT	Yüksekliklere yakın parkta yürüyen insanlar görüyor musunuz?
BIG	sıcak havalarda parkta yürüyebilenler bulunuyor mu?
Referans	deniz kenarında uçan kuşlar çok mutlu gözüküyorlar mı?
DH-S	deniz kenarında uçan kuşlar çok mutlu görünmezler mi?
DH	deniz kenarında uçan kuşlar fazla mutlu değil mi?
TQP	Mavi deniz kenarında uçan kuşlar çok mutlu görünmüyor mu?
MT	Denizlerde uçan kuşlar mutlu mu?
BIG	deniz kenarında uçan kuşlar gerçekten mutlu değil mi?
Referans	ülke genelinde bisiklet kullanımının artması konusunda ne düşünüyorsunuz?
DH-S	ülke genelinde bisiklet kullanımı artış yönünde düşünceniz nedir?
DH	bisiklet kullanımının yükselmesi hakkında fikirleriniz nedir?
TQP	Dünya genelinde bisiklet kullanımının artırılması konusunda fikrin nedir?
MT	Bisiklet kullanımının artmasının altındaki görüşleriniz nelerdir?
BIG	ülke genelinde bisiklet kullanımının artması hakkındaki görüşleriniz nedir?

Tablo 7.2 Farklı Veri Setlerine Göre Üretilen Aynı Anlamalı Cümle Çıktıları - 2

Referans	Askerliğimin tecillenmesi için nereye başvurmalıyım?
DH-S	askerliğimi nereye tecilleyebilirim?
DH	askerliğimi tecillemek için kurum önerisi?
TQP	Nereden bir askerliğini tecillemek için gerekli olduğunu bana söyleyebilir misin?
MT	Biri kendi askerlik işlerini tecillemek için nereye başvurabilir?
BIG	Askerlik durumu nasıl tecillenebilir?
Referans	Bayrampaşaya nereden gidebilirim?
DH-S	bayrampaşaya yol tarifi?
DH	bayrampaşada hangi yerden gidebilirim?
TQP	Bayrampaşaya giden yolu bana söyleyecek misin?
MT	Bayrampaşaya gitmenin en iyi yolu nedir?
BIG	Bayrampaşaya nasıl gideceğimi bana söyleyebilir misiniz?
Referans	Telefonumun ekranı kırıldı tamire nereye götürmeliyim?
DH-S	Telefonumun ekranı kırıldı nereye götüreyim?
DH	sıcak havalarda parkta duran insanlara rastlanır mı?
TQP	Android telefonumun ekranı kırıldı tamire nereye gitmeli?
MT	Telefonumun ekranı kırıldı nereye göndermeliyim?
BIG	Telefonumun ekranı kırıldı, tamire nereye götürmem gerekir?
Referans	Motorlu taşıtların çevreye verdiği zarar ile ilgili fikriniz nedir?
DH-S	motorlu taşıtların çevreye verdiği zarar hakkında ne düşünüyorsunuz?
DH	motorlu taşıtların çevreye verdiği zarar hakkında ne düşünüyorsunuz?
TQP	motolu taşıtların çevreye verdiği zarar hakkında ne düşünüyorsunuz?
MT	Motorlu taşıtlar çevreye verdiği zarar ile ilgili düşüncelerin neler?
BIG	Motorlu taşıtların çevrelere verdiği zarar hakkında ne düşünüyorsunuz?
Referans	paramızın değerinin düşmesi hakkında ne düşünüyorsunuz?
DH-S	paramız değerinin düşmesiyle ilgili düşünceniz nedir?
DH	paramızın değerinin çok olarak düşmesi ile ilgili ne düşünüyorsunuz?
TQP	Paramızın değeri nedir?
MT	Paranın değeri düşüyor mu?
BIG	paranın değerinin düşürülmesine ne düşünüyorsunuz?

Tablo 7.3 Farklı Veri Setlerine Göre Üretilen Aynı Anlamlı Cümle Çıktıları - 3

Referans	Bu yaz başlayan hangi yaz dizisini izleyeceksiniz?
DH-S	yazın başlayan hangi yaz dizisini izleyeceksiniz?
DH	bu yaz başlayan yaz dizileri?
TQP	Bu yaz ne tür dizileri izliyorsun?
MT	Bir yaz başlamak için izlemesi gereken en iyi yaz dizileri hangileridir?
BIG	Bu yaz çıkacak en iyi TV şovlarından bazıları nelerdir?
Referans	Bir kızdan çok hoşlanıyorum nasıl çıkma teklifi etmemi önerirsiniz?
DH-S	çok hoşlanıyorum şimdi çıkma teklifi yapımı önerir misiniz?
DH	bir kıza çok hoşlanıyorum, ne yapmalıyım?
TQP	Nasıl bir kızdan hoşlanacağımı bana söyleyebilir misiniz?
MT	Bir kızın çok hoşlandığım için nasıl çıkma teklifi etmeliyim?
BIG	bir kadından hoşlanıyorum, nasıl çıkma teklifi etmemi tavsiye edersiniz?
Referans	Aranızda daha önce hiç balık yememiş var mı?
DH-S	aranızda hiç balık yememiş bulunuyor mu?
DH	aranızda daha önce hiç balık tüketmemiş var mı?
TQP	Şu anda hiç balık yememiş var mı?
MT	Bugüne kadar hiç balık yememiş var mı?
BIG	aranızda daha önce hiç balık tüketmemiş kimse var mı?
Referans	Ölmeden kesin deneyimlemem lazım dediğiniz şeyler var mı?
DH-S	ölmeden önce ne deneyimlemek gerekir?
DH	ölmeden önce kesin neler görmemi önerir misiniz?
TQP	Ölmeden gerçekten deneyimlemem gereken şeyler var mı?
MT	Ölmeden kesin deneyimlemem gereken şeyler nelerdir?
BIG	Ölmeden önce deneyimlememi gerektiren şeylerden bazıları nelerdir?
Referans	Dünya üzerinde insanlar tarafından konuşulan dillerin sayısı nedir?
DH-S	toplumdan konuşulan dillerin sayısı nedir?
DH	insanların dünya üzerinde konuşma dillerinin ne kadar olduğunu düşünüyorsunuz?
TQP	Dünya üzerinde insanlar için konuşulan dilleri kaç tane?
MT	Dünyada konuşulan dillerin sayısı nedir?
BIG	dünyada insanlar tarafından konuşulan dillerin sayısı ne kadar?

Tablo 7.4 Farklı Veri Setlerine Göre Üretilen Aynı Anlamlı Cümle Çıktıları - 4

Farklı epoch değerleriyle (2 ve 3) ve DH kodlu veri setiyle eğitilmiş iki farklı modelin aynı cümleye dair ürettiği çıktılar aşağıdaki şekillerde gösterilmiştir.

```
CÜMLENİZİ YAZINIZ:günde ne kadar su içmeliyim?  
Cümleniz: günde ne kadar su içmeliyim?  
1      hangi günde suyu içmeliyim?  
2      Bir günde ne kadar su içmeli miyim?  
3      günde hangi su içmeliyim?  
4      günde günde ne kadar su içebiliriz?  
5      günde ne kadar su içebiliriz?
```

Şekil 7.9 1.Cümle - Epoch Sayısı=2

```
CÜMLENİZİ YAZINIZ:günde ne kadar su içmeliyim?  
Cümleniz: günde ne kadar su içmeliyim?  
1      günlük ne kadar su içelim?  
2      günde kaç litre su içmeliyim?  
3      günde kaç su içeyim?  
4      günde ne kadar su içmem gerekiyor?  
5      gün içerisinde ne kadar su tüketmeliyim?
```

Şekil 7.10 1.Cümle - Epoch Sayısı=3

```
CÜMLENİZİ YAZINIZ:kendinizi huzursuz hissettiğiniz anlar oluyor mu?  
Cümleniz: kendinizi huzursuz hissettiğiniz anlar oluyor mu?  
1      bu kadar hiçbir şey hissettiğiniz anlar oluyor mu?  
2      özünüzdeki kişisi bir insana dikkat eder misiniz?  
3      kendinizi birin huzursuz hissettiğiniz şeyler?  
4      kendinizi uzursuz hissedeceğiniz birini düşünüyor musunuz?  
5      kendinizi mutlu hissedecek biri var mı?
```

Şekil 7.11 2.Cümle - Epoch Sayısı=2

```
CÜMLENİZİ YAZINIZ:kendinizi huzursuz hissettiğiniz zamanlar oluyor mu?  
Cümleniz: kendinizi huzursuz hissettiğiniz zamanlar oluyor mu?  
1      kendinizi hiç huzursuz hissettiğiniz zamanlar olur mu?  
2      kendinizi huzursuz hissettiğiniz zamanlar yaşanıyor mu?  
3      kendinizi huzursuz hissedene var mı?  
4      kendinizi uyumsuz hissettiğiniz zamanlar olur mu?  
5      kendinizi hayal ettiğiniz zamanlar oluyor mu?
```

Şekil 7.12 2.Cümle - Epoch Sayısı=3

8

Performans Analizi

8.1 T5 İngilizce Dil Modelleme

Tarafımızca oluşturulmuş İngilizce T5 dil modeli üzerinde ince ayar işlemi için dil modelinin kayıp (loss) değerleri veri setinin büyüklüğüne ve tekrarlar göre Şekil 8.1'deki grafiklerde ifade edilmiştir.

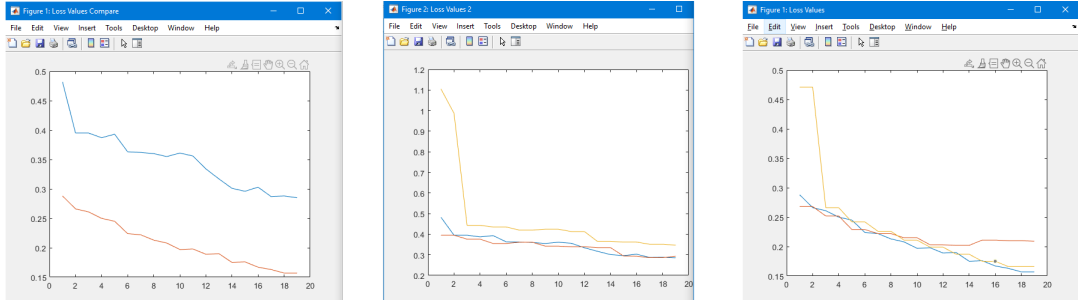
Sarı renkli çizgi: avg_train_loss

Mavi renkli çizgi: loss

Kırmızı renkli çizgi: val_loss

x axis: epoch

y axis: loss değerleri



Şekil 8.1 T5 İngilizce Dil Modellemesinde Loss Değerleri
(Karşılaştırma-Küçük Veri Seti-Büyük Veri Seti)

8.2 T5 ile İngilizce Duygu Analizi

Büyük veri setiyle oluşturulmuş İngilizce dilindeki T5 dil modeli üzerinde uygulanan ince ayar işleminde, cümle üzerinden duygu analizi işlemine ait Hesaplanan-Gerçeklik matrisinde (Tablo 8.1) gösterilmiştir. Duygu analizi işlemine ait F1 skoru da Şekil 8.2’de gösterilmiştir.

	Gerçek		
Hesaplanan		Pozitif	Negatif
	Pozitif	9824	3078
	Negatif	2676	9422

Tablo 8.1 Hesaplanan-Gerçeklik Matrisi

	precision	recall	f1-score	support
negative	0.74	0.83	0.78	12500
positive	0.81	0.71	0.76	12500
accuracy			0.77	25000
macro avg	0.78	0.77	0.77	25000
weighted avg	0.78	0.77	0.77	25000

Şekil 8.2 İnce Ayar İşleminin Başarısına Ait Metrik Değerler

8.3 İngilizce Aynı Anlamlı Cümle Üretimi Değerlendirmesi

T5, BoW ve Seq2Seq ile İngilizce diline ait aynı anlamlı cümle üretiminin başarısı insan değerlendirmesiyle ölçülmüştür ve en başarılı sonuçların önceden eğitilmiş T5 dil modeliyle yapıldığı saptanmıştır.

8.4 MT5 ile Türkçe Aynı Anlamlı Cümle Üretimi

8.4.1 BLEU Skoru

BLEU skoru, girilen cümle ve oluşturulan cümledeki kelimelerin n-gramlarını kontrol ederek ek alıp almadığına bakmaksızın bire bir aynı kelimeleri içermesi durumunda arttığı için aynı anlamlı cümle üretiminin başarısını doğrudan ölçen bir metrik olmayacağı kanaatine varılmıştır. MT5 dil modeline uygulanan ince ayar işlemi sonucu en yüksek BLEU skorunu DH-S veri setiyle eğitilen model almıştır (Tablo 8.2).

Veri Seti	BLEU-1	BLEU-2	BLEU-3
DH-S	0.926	0.787	0.586
DH	0.768	0.606	0.420
TQP	0.415	0.255	0.147
MT	0.471	0.312	0.197
BIG	0.752	0.585	0.420

Tablo 8.2 Farklı Veri Setleri Kullanılarak Eğitilmiş MT5 Modellerin BLEU Skorları

8.4.2 İnsan Değerlendirmesi

5 farklı veri setiyle eğitilmiş 5 adet MT5 modeli ile 20 cümlelerin her biri için 5'er tane aynı anlamlı çıktı üretilmiştir ve ana dili Türkçe olan 5 farklı insan tarafından, Excel dosyasında; her satırda en iyi ve en kötü üretilmiş cümleye ait çıktılar işaretlenmiştir. Sonuç olarak en başarılı modelin tarafımızca ekleme yapılmış olan manuel hazırlanmış DH kodlu veri setiyle eğitilmiş model olduğu sonucuna varılmıştır.

Veri Seti ID	Başarılı	Başarısız
DH-S	%19.8	%14.2
DH	%26.6	%14.6
TQP	%16	%32.2
MT	%13.4	%27.2
BIG	%24.2	%11.8

Tablo 8.3 Farklı Veri Setleri Kullanılarak Eğitilmiş MT5 Modellerin İnsan Değerlendirmesi

8.4.3 Sonuçların Değerlendirilmesi

8.4.3.1 Uç Noktaların Değerlendirilmesi

Şekil 7.1, Şekil 7.3, Şekil 7.5 görsellerinde üretilen cümlelerin referans edilmiş cümle ile aynı veya yakın anlam içerdiği görülmüştür. Bunun anlamı model eğitilirken yalnızca ikili soru cümleleri kullanılsa dahi soru olmayan cümlelerde de başarı elde edebilecek olması, ayrıca eş anlamlı cümlesi üretilecek referans cümlelerin özel karakter bulundurması ile ilgili bir problemin olmadığıdır.

Şekil 7.4, Şekil 7.6, Şekil 7.7 görsellerinde eş anlamlı cümle üretilmesi beklenen girdilerin, 1 ya da 2 kelime olduğu görülmektedir. Modelin eğitildiği veri setinde bu tarz kısa cümleler genellikle öneri almak veya yardım istemek için kullanılmıştır. Veri

setindeki kısa cümlelerin bu anlamı taşımasının etkisi çıktılarda görülmektedir.

Şekil 7.2 görselinde Türkçe için eğitilmiş veri setinin İngilizce içeriği olan bir cümleyi işleyemediği görülmektedir.

Şekil 7.8 görselinde insanlar için de bir anlam taşımayan cümle referans olarak verilmiştir. Neticesinde modelin cümlelerin içerdiği kelimelerden öneriler beklediği görülmektedir.

8.4.3.2 Epoch Sayılarına Göre Değerlendirme

Epoch bir modelin tek bir veri setini birden fazla kullanarak kendini eğitmesini ayarlamaktadır. Bu sebeple belirli bir eşik değerine kadar epoch sayısını arttırmak modelin öğrenimini ve başarısını arttıracaktır. Modelin epoch için eşik değerinin 2 olmadığı ve 3 yapıldığında daha başarılı sonuçlar ürettiği Şekil 7.9 ile Şekil 7.10 arasında ve Şekil 7.11 ile Şekil 7.12 arasında görülmektedir.

Tamamlanan proje kapsamında aşamalı olarak İngilizce dili için aynı anlamlı cümle üretimi, İngilizce dil modelleme, modellenen dilin test edilmesi, Türkçe için aynı anlamlı cümle üretimi ve bu görev için eğitilmiş modelin başarısının ölçülmesi tamamlanmıştır.

Projenin ilk aşamasında kullanılacak mimari olarak T5 seçilmiş, Türkçe versiyonu olarak kullanabilmek adına MT5 dil modelinden yararlanılmıştır. Modellenen İngilizce dilin başarısını görmek için duygu analizi görevi çalıştırılmış ve başarısı F1 skor yöntemi ile ölçülmüştür. Nihayetinde, Türkçe için aynı anlamlı cümleler üretilmiş ve başarısını ölçmek için iki farklı yöntem denenmiştir. BLEU skor yönteminin projede beklenenlere ters düştüğü tespit edilip değerlendirmede ölçüt olarak kabul edilmemiştir. Diğer yöntem olan insan değerlendirmesinde ise beklendiği üzere insan tarafından oluşturulmuş veri seti ile eğitilmiş model en başarılı seçilmiştir.

Türkçe için aynı anlamlı cümle üretimi çalışması boyunca bu alanda ihtiyaç duyulan veri setine 20.000 cümle katkıda bulunulmuş ve büyütülmüştür. Makine çevirisiyle oluşturulmuş ve manuel olarak oluşturulmuş veri kümeleri karşılaştırılmış, cümle ikililerinin niteliklerinin başarıya etkisi gösterilmiştir.

Çalışma boyunca doğal dil işleme alanına yönelik çok farklı alanlara ait işler ile uğraşılmış ve edinilmesi beklenen tecrübe kazanılmıştır. Belirlenen görev kapsamında ülkemizin teknolojisine, veri kümesine katkıda bulunulmuştur. Tarafımızca katkı sağlanmış veri setinin modele etkisi, katılmayan veri seti ile oluşturulmuş modeller ile kıyaslanarak gösterilmiştir. Proje neticesinde, ana dili Türkçe olan insanlar tarafından değerlendirilen cümlelerde beğenilme oranı oldukça yüksektir. Bunun sonucu olarak büyüklüğü az olan bir veri seti üzerinde aynı anlamlı cümle üretimi görevi çalıştırılarak veri setinin sağlıklı bir şekilde büyütülmesi teorik olarak mümkündür. Proje amacı ile elde edilen veri setine katkılar yapılarak denenmesine rağmen veri seti günlük işlerde kullanılabilecek tüm cümleleri içermez ve uç nokta testlerinde bunun sebebi olarak başarısızlıklar görülmüştür.

Aynı anlamlı cümle üretimi amacı ile kullanılmış T5 mimarisinin diğer modellere göre daha başarılı olduğu, aynı zamanda diğer modellere göre kullanımının daha rahat olduğu belirtilmiştir. Bu hususta ileride bu alanda çalışmak isteyenler için bu modelin ilk tercihlerden olmasının düşünüldüğü belirtilmiştir. Aynı zamanda, üzerine çalışılacak veri seti mümkün oldukça insan emeği ile oluşturulmalı ve genişletilmelidir.

- [1] E. Egonmwan and Y. Chali, “Transformer and seq2seq model for paraphrase generation,” in *Proceedings of the 3rd Workshop on Neural Generation and Translation*, Hong Kong: Association for Computational Linguistics, Nov. 2019, pp. 249–255. DOI: 10.18653/v1/D19-5627. [Online]. Available: <https://www.aclweb.org/anthology/D19-5627>.
- [2] Y. Fu, Y. Feng, and J. P. Cunningham, “Paraphrase generation with latent bag of words,” p. 12, 2019.
- [3] M. Y. Sarigöz, *Monatis/tqp: V0.1*, version v0.1, Apr. 2021. DOI: 10.5281/zenodo.4719801. [Online]. Available: <https://doi.org/10.5281/zenodo.4719801>.
- [4] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020. [Online]. Available: <http://jmlr.org/papers/v21/20-074.html>.
- [5] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush, *Huggingface’s transformers: State-of-the-art natural language processing*, 2020. arXiv: 1910.03771 [cs.CL].
- [6] L. Xue, N. Constant, A. Roberts, M. Kale, R. Al-Rfou, A. Siddhant, A. Barua, and C. Raffel, *mT5: A massively multilingual pre-trained text-to-text transformer*, 2020. arXiv: 2010.11934 [cs.CL].
- [7] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: <http://papers.nips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [8] K. Papineni, S. Roukos, T. Ward, and W. J. Zhu, “Bleu: A method for automatic evaluation of machine translation,” Oct. 2002. DOI: 10.3115/1073083.1073135.

- [9] Q. Lhoest, P. von Platen, T. Wolf, A. V. del Moral, Y. Jernite, S. Patil, M. Drame, J. Chaumond, J. Plu, J. Davison, S. Brandeis, T. L. Scao, V. Sanh, K. C. Xu, L. Tunstall, N. Patry, A. McMillan-Major, P. Schmid, S. Gugger, L. Debut, C. Delangue, T. Matušíère, S. Bekman, and F. Lagunas, *Huggingface/datasets: 1.7.0*, version 1.7.0, May 2021. DOI: 10.5281/zenodo.4817769. [Online]. Available: <https://doi.org/10.5281/zenodo.4817769>.
- [10] Y. Kim, Y. Jernite, D. Sontag, and A. M. Rush, *Character-aware neural language models*, 2015. arXiv: 1508.06615 [cs.CL].
- [11] B. Pang and L. Lee, “Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales,” in *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*, Ann Arbor, Michigan: Association for Computational Linguistics, Jun. 2005, pp. 115–124. DOI: 10.3115/1219840.1219855. [Online]. Available: <https://www.aclweb.org/anthology/P05-1015>.
- [12] L. Sharma, L. Graesser, N. Nangia, and U. Evci, “Natural language understanding with the quora question pairs dataset,” p. 10, Jul. 2019.
- [13] Y. Scherrer, *TaPaCo: A Corpus of Sentential Paraphrases for 73 Languages*, version 1.0, Zenodo, Mar. 2020. DOI: 10.5281/zenodo.3707949. [Online]. Available: <https://doi.org/10.5281/zenodo.3707949>.

BİRİNCİ ÜYE

İsim-Soyisim: Metin BİNBİR
Doğum Tarihi ve Yeri: 12.04.1999, Gölcük
E-mail: 11119613@std.yildiz.edu.tr
Telefon: 0530 593 05 92
Staj Tecrübeleri:

İKİNCİ ÜYE

İsim-Soyisim: Sercan AKSOY
Doğum Tarihi ve Yeri: 30.07.1998, İstanbul
E-mail: 11119620@std.yildiz.edu.tr
Telefon: 0534 958 88 74
Staj Tecrübeleri:

Proje Sistem Bilgileri

Sistem ve Yazılım: Windows İşletim Sistemi, Python, Google Colab Pro
Gerekli RAM: 16GB
Gerekli Disk: 200GB