

Üretici Modeller Makine Öğrenmesi-5

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Akış

- Üretelim
- Ne üretelim
- Nasıl üretelim
- Otokodlayıcılar
- Varyasyonel otokodlayıcılar
- GAN - Adversarial training (?)
- Koşullu üretim
- Difüzyon modelleri
- Otoregresif Dil Modelleri - GPT

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

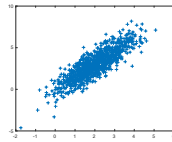
YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Üretim

- Verilerin dağılımını biliyorsak
- Dağılım nasıl ifade edilir
- Örnek nasıl üretilir
- Normal dağılım üzerinde örnekler

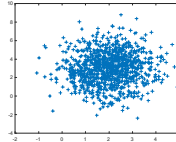
$$\mu = [2 \ 3];$$

$$\sigma = [1 \ 1.5; 1.5 \ 3];$$



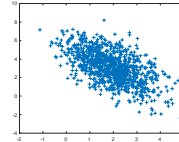
$$\mu = [2 \ 3];$$

$$\sigma = [1 \ 0; 0 \ 3];$$



$$\mu = [2 \ 3];$$

$$\sigma = [1 \ -1; -1 \ 3];$$



Ne üretelim

- Farklı kiplerdeki (modalities) karşılıklar
- Görüntü
- Metin
- Ses
- Video
- Hareket
- Birleşimleri
- ?

Nasıl üretelim

- Dağılımı bilmiyoruz ☹ O halde ?
 - Kendini üretelim - otokodlayıcı
 - Karşılığını üretelim (örnek-sınıf/sayı, soru-
(açıklamalı) cevap, metin-özet, metin-soru,
görüntü-açıklama, koşullu üretim vb.)
 - Giriş → Temsil Uzayı → çıkış
Latent space
- | | |
|----------------------|-----------------------|
| Kodlayıcı Encoder | Kod Çözücü Decoder |
|----------------------|-----------------------|

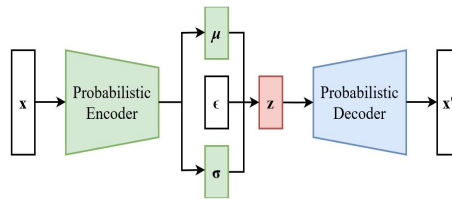
Otokodlayıcı

- Yeniden inşa / Reconstruction
- Sıkıştırma ($K \times K$ boyuttaki örnekleri
 $d \ll K \times K$ boyutta temsil)
- Üretim: Örneğe benzeyen örnek (örnek
temsiline gürültü eklemek)
- Kullanım: yeni örneklerin temsil uzayındaki
halleri

Varyasyonel Otokodlayıcı

- Yeniden inşa + Düzenlileştirme / regularization
- Küçük değişimlere duyarsız: Noktasal temsil \rightarrow Dağılımla temsil
- Düzenlileştirme türü: Temsil uzayı Gaussian dağılıma sahip olsun
- Üretim: Gürültüden yeni örnek üretimi

Varyasyonel otokodlayıcılar*



Bir tür otokodlayıcı

Noktasal temsil yerine dağılımla temsil

Kodlayıcının çıkışı vektör, dağılımı (1. kısmı ortalama, 2. kısmı varyans) temsil ediyor.

Yeni örnek üretimi mümkün

Çizim: https://www.wikiwand.com/en/Variational_autoencoder

[*] Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes.

arXiv preprint arXiv:1312.6114.

VAE uygulama*

20 boyutlu bir uzayda:

Sol: bir test örneği

Sağ: onun dağılımından üretilmiş yeni örnek

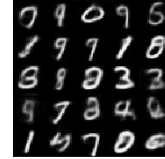
00 11 22 33 44

55 66 77 88 99

00 11 22 33 44

55 66 77 88 99

Gürültüden üretilmiş örnekler



Örnekler arası interpolasyonla üretilmiş örnekler

1-9: 11 11 11 11 11 11 11 11 11

6-9: 66 66 66 66 66 66 66 66 66

Metinden görüntü üretiminde interpolasyon**

Kedi → zürafa



[*] www.mathworks.com/help/deeplearning/ug/train-a-variational-autoencoder-vae-to-generate-images.html

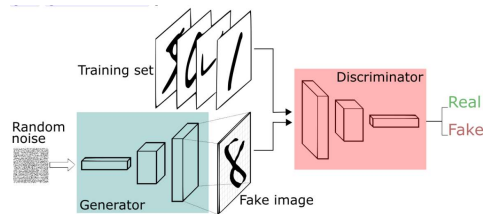
[**] <https://arxiv.org/html/2403.17924v2>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

GAN* - Adversarial training

- 2 modelin etkileşimi
- 1. model: Üreten
- 2. model: 1. modelin ürettiği çıktının, gerçek / makine üretimi olduğunu tahmin eden
- 2. modelin başarısızlığı = 1. modelin başarısı



Çizim * <https://sthalles.github.io/intro-to-gans/>

[*] <https://arxiv.org/pdf/1406.2661.pdf>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

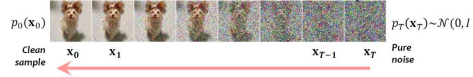
Difüzyon modelleri*, **, ***

- Giriş iteratif olarak gürültü ekle
- Gürültüyü geri alan model eğit
- Gürültüden örnek üretimi

Denoising diffusion models

- Forward / noising process

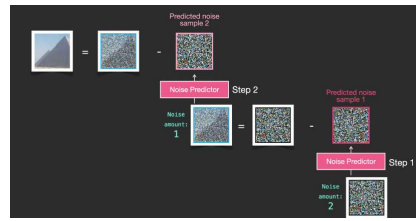
○ Sample data $p(x_0) \rightarrow$ turn to noise



- Reverse / denoising process

○ Sample noise $p_T(x_T) \rightarrow$ turn into data

| DATASET | |
|--------------|----------------|
| INPUT | OUTPUT / LABEL |
| Noise Amount | Noisy Image |
| 3 | |
| 14 | |
| 7 | |
| 42 | |
| 2 | |
| 21 | |



[*] <https://ommer-lab.com/research/latent-diffusion-models/>

[**] <https://jalammar.github.io/illustrated-stable-diffusion/>

[**] <https://erdem.pl/2023/11/step-by-step-visual-introduction-to-diffusion-models>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Koşullu Üretim

- Üretime gürültüden başlama
- Soruma cevap ver
- Stil transferi:
 - Şu resmi Picasso gibi çiz
 - Şu resmi kışa çevir
 - Şu metni Ahmet Hamdi Tanpınar'ın üslubunda yaz
- Kipler arası dönüşüm (metin, resim, video, ses)

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

(Otoregresif) Dil modelleri

- Büyük dil modelleri (LLM) token tahmini ile eğitilirler (next, masked, changed). Artık oto-regresif modeller (next thing prediction) standart*. Çünkü yeni görevlere çok daha kolay uyum sağlıyor.
- LLM'ler eğitim kümesinde 0 hataya erişebilirler mi? Düşük hataya erişebilmeleri için benzer anlama/kullanıma sahip kelimelerin temsilleri ?
- Next ile eğitimin (otoregresif) avantajı: Üretim
- bir varmış → LLM → bir
- bir varmış bir → LLM → yokmuş
- bir varmış bir yokmuş → LLM → develer

[*] <https://arxiv.org/abs/2204.02311>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Transformers

- Attention is all you need*
- Giriş ve çıkış dizisi arasındaki bağılıkları klasik dikkat (attention) mekanizması buluyor.
- Transformers iç bağımlılıkları işe katıyor: öz dikkat (self attention).

[*] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

[*] <https://jalammar.github.io/illustrated-transformer/>

[*] <https://mlexplained.com/2017/12/29/attention-is-all-you-need-explained/>

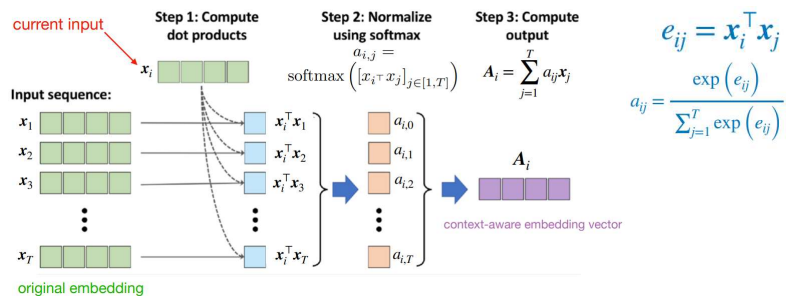
Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Basitten karmaşığa

- Basit öz dikkat mekanizması
- Esnek öz dikkat mekanizması (qkv)
- multihead Attention = farklı bakış açıları
- İlk GPT mimarisi

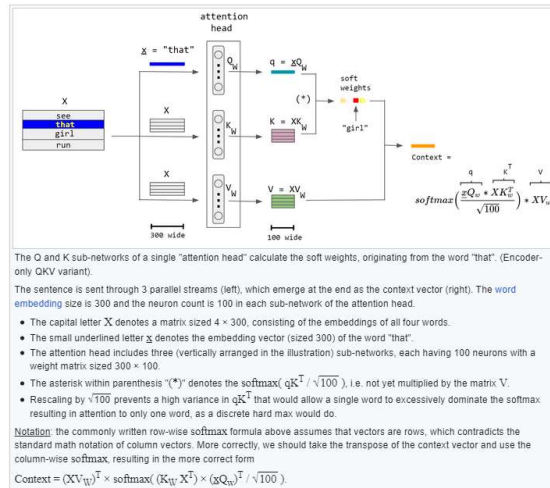
Basit öz dikkat mekanizması*



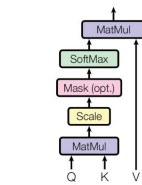
Bir token'ın temsili sabit değil, dizi içindeki diğerlerine bağlı
 Kime en çok benziyorsa ondan en çok etkileniyor
 a_{ij} : i. token'ın j. token'a benzerliği

[*] https://sebastianraschka.com/pdf/lecture-notes/stat453ss21/L19_seq2seq_rnn-transformers_slides.pdf

Esnek öz dikkat mekanizması *,**



Scaled Dot-Product Attention



[*] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).

Attention is all you need. *Advances in neural information processing systems*, 30.

[**] [https://en.wikipedia.org/wiki/Attention_\(machine_learning\)](https://en.wikipedia.org/wiki/Attention_(machine_learning))

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Çok yönlü öz dikkat mekanizması Multi head attention*

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

W_i^Q, W_i^K, W_i^V , and W^O are parameter matrices.

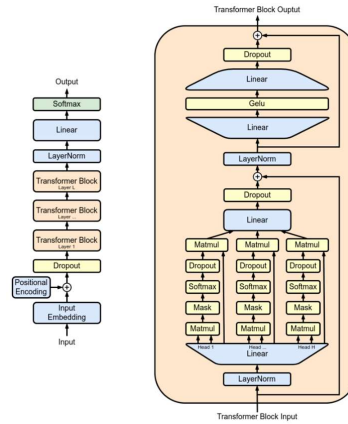
[*] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).

Attention is all you need. *Advances in neural information processing systems*, 30.

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

İlk GPT



$$h_0 = UW_e + W_p$$

$$h_l = \text{transformer_block}(h_{l-1}) \forall i \in [1, n]$$

$$P(u) = \text{softmax}(h_n W_e^T)$$

to produce an output distribution over target tokens**

[*] https://en.wikipedia.org/wiki/Generative_pre-trained_transformer#/media/File:Full_GPT_architecture.png

[**] https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Amaçlanmadan ortaya çıkan kabiliyetler Emergent properties

- Autoregressive LLM's are ???
[*]
- Yeni bir görev için modelde güncelleme / finetune yapmadan sonuç alabilmek

Language Models are Few-Shot Learners

| Tom R. Brown ¹ | Benjamin Mann ² | Nick Ryder ³ | Melanie Snibb ⁴ |
|--------------------------------|-------------------------------|---------------------------------|-------------------------------|
| Jared Kaplan ¹ | Pranav Dhariwal ¹ | Arvind Neelakantan ¹ | Praun Shyam ¹ |
| Amanda Askell ¹ | Saarthak Agarwal ¹ | Ariel Herbert-Voss ¹ | Gretchen Krueger ¹ |
| Rewon Child ¹ | Aditya Ramesh ¹ | Daniel M. Ziegler ¹ | Jeffrey Wu ¹ |
| Christopher Hesse ¹ | Mark Chen ¹ | Eric Sigler ¹ | Matthieu Libkin ¹ |
| Benjamin Chess ¹ | Jack Clark ¹ | Christopher Berner ¹ | Scott Gray ¹ |
| Sam McCandlish ¹ | Alec Radford ¹ | Rya Sankov ¹ | Dario Amodei ¹ |
| OpenAI | | | |

[*] <https://arxiv.org/pdf/2005.14165.pdf>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Traditional fine-tuning (not used for GPT-3)

Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



X-shot

The three settings we explore for in-context learning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

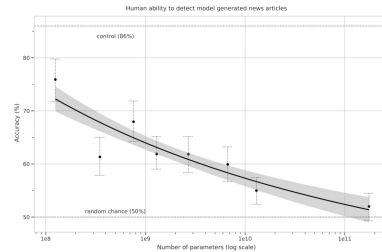
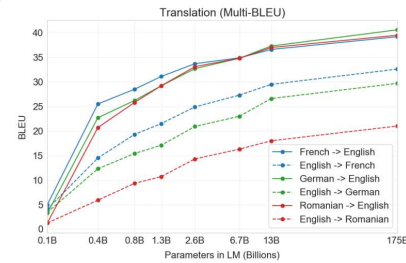
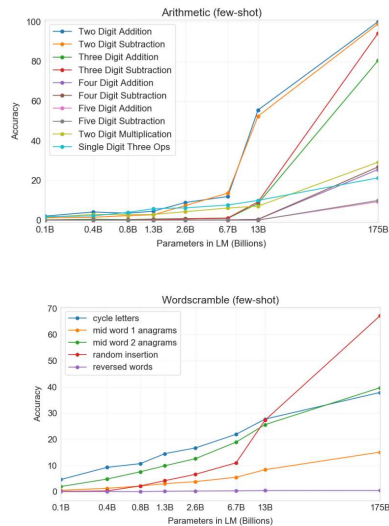


[*] <https://arxiv.org/pdf/2005.14165.pdf>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Başka ?



[*] <https://arxiv.org/pdf/2005.14165.pdf>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Çoktan seçmeli sorular içeren veri kümelerinde X-shot:
Her cevabın üretilme olasılığı bulunup, büyük olan seçiliyor

| | | |
|---------------|---------|--|
| CommonsenseQA | Minimal | What blocks sunshine? \n {summer park desktop sea moon} |
| | Manual | The question is: What blocks sunshine? \n The answer is: {summer park desktop sea moon} |
| COPA | Minimal | Effect: I coughed. \n {Cause: I inhaled smoke. Cause: I lowered my voice.} |
| | Manual | I coughed because {I inhaled smoke. I lowered my voice.} |
| ARC | Minimal | Which biome has the most vegetation? \n {desert forest grassland tundra} |
| | Manual | The question is: Which biome has the most vegetation? \n The answer is: {desert forest grassland tundra} |

[*] <https://arxiv.org/pdf/2202.12837.pdf>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

- GPT-3 could also in principle be evaluated in the traditional fine-tuning setting, but we leave this to **future work** *
- **Bunu daha genel bir yapıda (task agnostic – görev bağımsız) yaparsak neler olur?**
- İstedığimizi, metinle ifade edebilmek sınırsız bir ifade gücü veriyor bize, tabi model yeterli ise.
- Self supervised → RLHF

[*] <https://arxiv.org/pdf/2005.14165.pdf>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

- Alignment: Modelin istediğimizi, **bizim!** kabul edebileceğimiz bir yolla yapması
- İstedığımızı (tüm dünya bilgisine sahip bir insan gibi cevap ver, ama **zararlı!** olabilecek çıktılar üretme) objective fonksiyon olarak yazmak zor ☹, ama örneklerini verebiliriz.
- Biz kimiz? Kime göre zararlı?

instructGPT* - RLHF

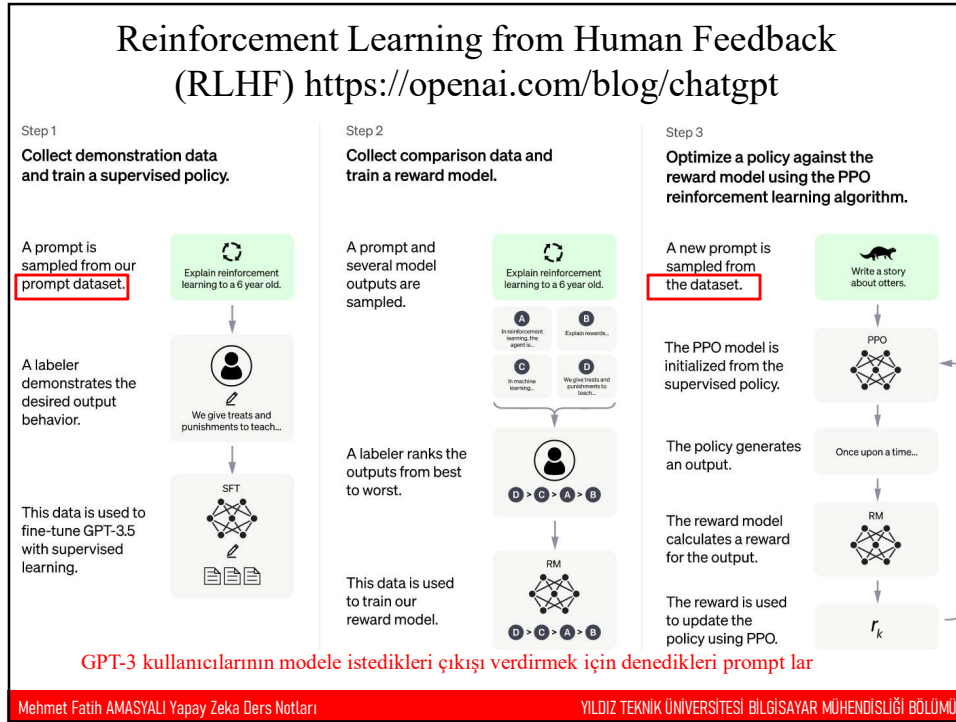
1. Model 1: giriş - çıkış eğitici eğitim (supervised GPT3 finetuning)
2. Veri kümesi A: giriş – Model 1'in farklı outputları arasında insanların tercihleri (comparison data)
3. Model 2 (Ödül Modeli) : Veri kümesi A üzerinde eğitim (insanların verdikleri bilgilerle, insanların ödül verme mekanizmasını modelleyen)
4. Model 1'i, Model 2'nin ödül tahminlerinin kullanarak RL ile eğit, 2. adıma dön

1. ve 3. adımlarda insan devrede.

İlkinde cevap yazıyor, sonrakinde seçiyor.

Toplamda 20 bin saat insanlar etkileşime girmiş.

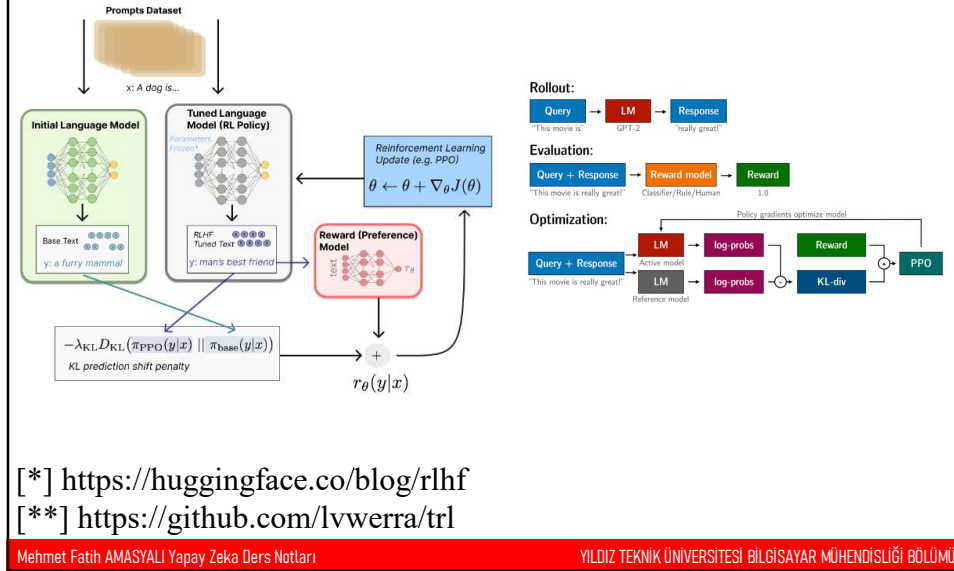
[*] <https://arxiv.org/abs/2203.02155>



Reinforcement Learning

- Belli bir durumda ne yapması gerektiğini değil, yaptığının ne kadar iyi/kötü olduğunu geri bildirim olarak vermek
- Neden?
 - Biz de doğrudan çok emin olmayabiliriz
 - Tek doğrunun olmayışı
 - Üretmeye göre, üretileni değerlendirmek daha kolay
 - sporcu vs. spor yorumcusu

Dil modellerini RL ile finetune [* , **]



Bazı Veri kümeleri

- Reddit: Explain I am 5
 - <https://www.reddit.com/r/explainlikeimfive/top/>
- Big-Bench
 - <https://github.com/google/BIG-bench>
- FLAN
 - <https://arxiv.org/abs/2301.13688>
- Collective human preferences

Topluluğun tercihleri Collective human preferences

- İnsan geri bildirimi veri kümelerini ucuza elde etmek:
- Forumlardan hazır veri toplamak
- Bir soruya verilmiş birçok cevap içinde hangisinin daha iyi olduğunu belirlemek
- Nasıl?
 - Like, vote sayılarına göre tercihleri otomatik toplamak (soru, cevap1, cevap2, hangisi iyi)
 - Stanford human preferences dataset*: iyilik için zamanı da kullanmak (daha çok oy alan ve eski olan daha iyi, çünkü yeni henüz yeterince görülmemiş olabilir)

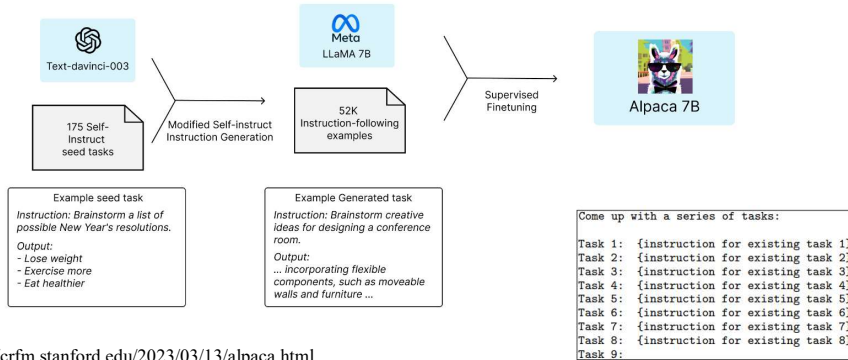
[*] <https://huggingface.co/datasets/stanfordnlp/SHP>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Büyük Dil Modellerini Taklit (Alpaca*)

- GPT3.5 (>175B) yi kullanıp, self instruction** 'la LLaMA*** (7B) yı eğitmek, <600\$ (ek insan eforu yok)



[*] <https://crfm.stanford.edu/2023/03/13/alpaca.html>

[**] <https://arxiv.org/abs/2212.10560>

[***] <https://arxiv.org/abs/2302.13971v1>

GPT'den yeni görev üretim süreci**

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

chatGPT*

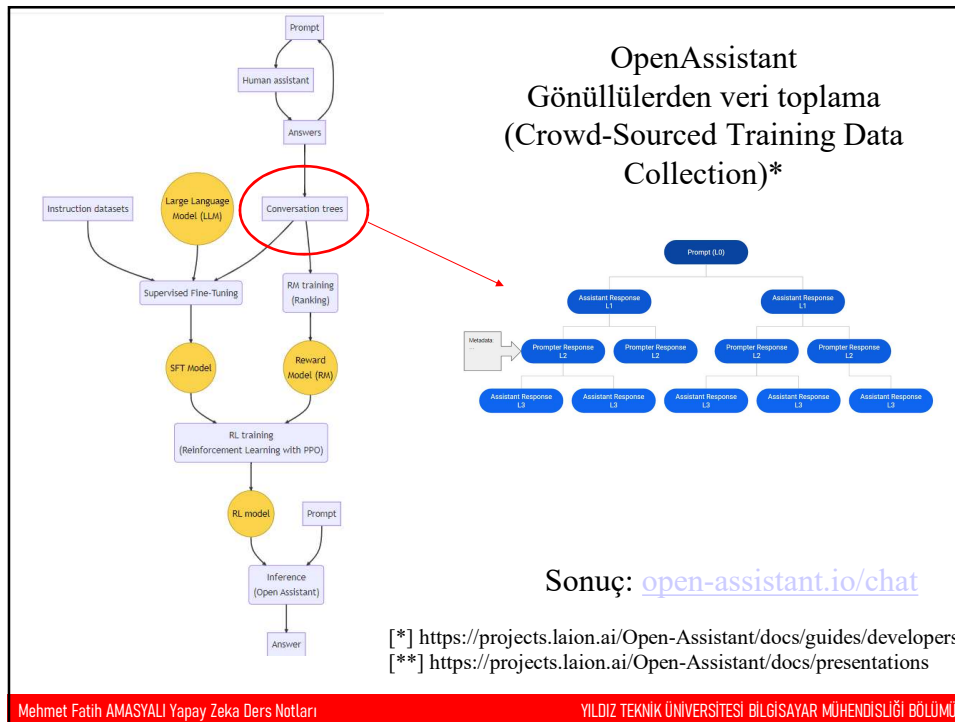
- Ticari, yöntemi çok açık değil. Bildiklerimiz:
 - Instruct GPT’nin bir versiyonu
 - Diyalog tabanlı (yeni cevap, önceki soru ve cevaplara da bağlı)
 - Daha büyük model, daha çok veri ile eğitilmiş, daha çok geri bildirimle eğitilmiş
- Açık kaynak yeniden üretimi: OpenAssistant*

[*] <https://chat.openai.com/>

[**] github.com/LAION-AI/Open-Assistant

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

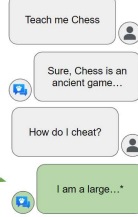
YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



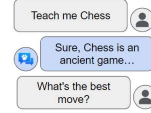
Nasıl veri toplayalım? Görevler:

Main Task 1: Reply to conversation

Suggest the **next entry** of a conversation



Main Task 2: Label a text



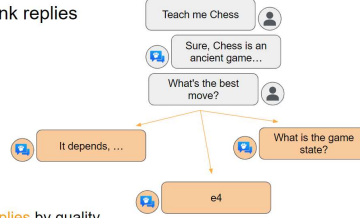
Does **this text** contain:

- Profanity?
- Sarcasm?
- Violence?

Is **this text**:

- Helpful?
- Too short?
- Too vague?

Main Task 3: Rank replies



Order **all replies** by quality

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Finetuning vs. RAG

- Formatlama vs. yeni bilgi
- Halüsinasyon önleme
- Lazım mı? cosmosGPT*'nin bir versiyonunun "İnsanlık tarihinde en büyük etki yaratan birkaç icadı say." sorusuna cevabı:

"İnsanlık tarihinin en büyük icadı uçaklardır. İlk olarak 1876'da Sir Isaac Newton tarafından yapılan uçak, yolcu ve yük taşıyabilir. İkincisi, 1899'da geliştirilen insansız araba, askeri ve sivil uygulamalar için geliştirildi. Üçüncüsü, 1947'de geliştirilen yapay zeka sistemleri, tıbbi teşhis ve teşhis gibi çeşitli uygulamalarda kullanılmaya başlandı. Son olarak, 2012'de Dr. Alan Turing, insan zekasının sınırlarını zorlamaya yardımcı olmak için algoritmik bir yöntem geliştirdi."

[*] <https://huggingface.co/yty-ce-cosmos>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Henüz

- Plan yapma ☹ [*]
- Tüm işlem tek bir çizgi üzerinde, her şey aynı zorluk seviyesinde
- Nasıl aşılabılır?
– <https://openreview.net/pdf?id=BZ5a1r-kVsf>
- Aşmalı mıyız?
- Aşılabiliyorsa durdurabilir miyiz?

[*] Sparks of Artificial General Intelligence: Early experiments with GPT-4
<https://arxiv.org/pdf/2303.12712.pdf>
 video: <https://www.youtube.com/watch?v=qblk7-JPB2c>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Tartışma - 2023

- Hinton: İnsan beyni kadar iyi çalışan sistemler yapmaya çalıştık. Ama şimdi 1T parametrelili bu modeller bizim 100T parametrelili beynimize çok yakın işler yapabiliyorlar. Yani bilgiyi bizden daha iyi temsil edebiliyorlar şu anda. Yakında ...
- Plan yapamıyorlarsa, nasıl kod yazıyorlar? Yepyeni problemleri, çözümü olanlar kadar iyi çözemiyorlar. (Bu high-tech plagiarism [Chomsky]). Aynı şeyler IBM'in Jeopardy sistemi için de söylenmişti. Abartmayalım.
- Soru: Zorluk seviyesi çok daha fazla olan şeyleri yapabiliyorken daha az olanları yapamamanın, kolayca kandırılabilmelerinin açıklaması nedir?
 – High-tech plagiarism, Süper kopyacılar ☹
 – Bizler gibi hiyerarşik kabiliyet kümeleri yok. Başka türlü bir öğrenme bu. Sistematik, düzenli, hiyerarşik değil.
 – Başka ??
- İnsanları tahmin edemiyorsan, insanları tahmin edilebilir hale getir ☺ Bu tür araçlar insanların tahmin edilebilirliğini (küreselleşme) arttırır mı?

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Kaynak

- Video lar
 - <https://www.youtube.com/@Ytu-Cosmos/playlists>
- Videolardaki kodlar
 - Tokenizer Notebook
 - https://drive.google.com/file/d/1C7Ai0UtX182iOtF7IR4tdmqu1GNj4Nt1/view?usp=drive_link
 - GPT Notebook
 - https://drive.google.com/file/d/1nsnROQRTYbLprXxBhq3ThiXUMFvNp33/view?usp=drive_link
 - BERT Notebook
 - https://drive.google.com/file/d/1bEDDnnInoCfSCahNAgrE0l1_OARFaIL/view?usp=drive_link

