



# Mérföldkövek a ChatGPT-ig

**Dr. Yang Zijian Győző**

HUN-REN Nyelvtudományi Kutatóközpont  
[yang.zijian.gyozo@nytud.hun-ren.hu](mailto:yang.zijian.gyozo@nytud.hun-ren.hu)



A conTEXT konferencián,  
Ahol a jövő most vár,  
A technológia és a tudomány  
Egy új, izgalmas világot tár.

A conTEXT konferencián,  
Ahol a tudás hatalom,  
Ott a helyem én is,  
Ahol a jövő formálódik.



# Miről lesz szó?

- Nyelvtechnológia története (röviden)
- Neurális nyelvtechnológia fontosabb mérföldkövei
- Nagy nyelvi modellek
- ChatGPT
- Magyar kutatások a nagy nyelvi modellek területén



# Nyelvtechnológia története



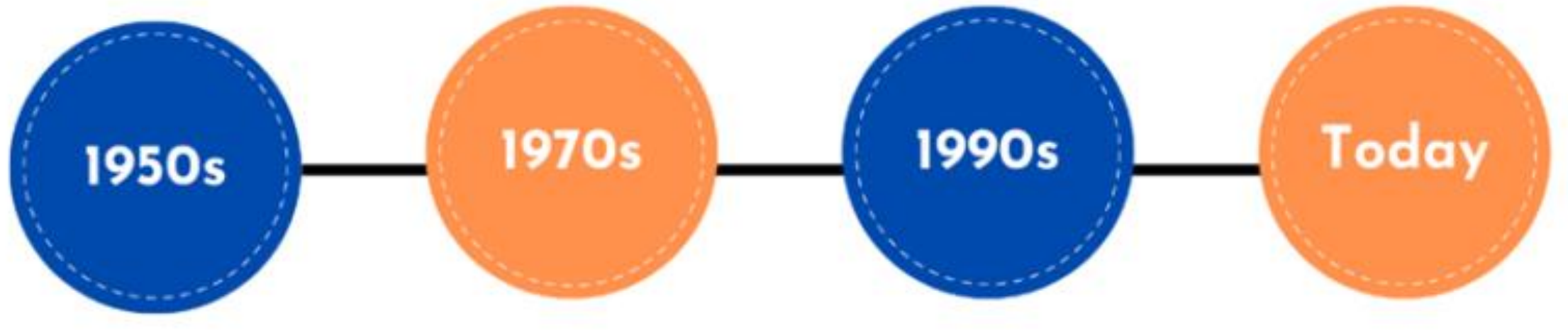
**Modern számítógép  
megjelenése**

**Gépi fordítás**

**Turing teszt**



# Nyelvtechnológia története



**Modern számítógép  
megjelenése**

**Gépi fordítás**

**Turing teszt**

**Szabályalapú  
megoldások**

**ELIZA**

**Szemantikus háló,  
Ontológiák**



# Szabályalapú gépi fordítás



# Szabályalapú gépi fordítás

**A piros cicám alszik.**

➤ **Mondat szegmentálás**



# Szabályalapú gépi fordítás

A piros cicám alszik.



**A piros cicám alszik .**

➤ **Mondat szegmentálás**

➤ **Tokenizálás**





# Szabályalapú gépi fordítás

A piros cicám alszik.



**A piros cica alszik .**

➤ **Mondat szegmentálás**

➤ **Tokenizálás**

➤ **Lemmatizálás**



# Szabályalapú gépi fordítás

A piros cicám alszik.



A [DET] **piros** [ADJ] **cica** [NOUN] **alszik** [VERB] . [PUNCT]

➤ **Mondat szegmentálás**

➤ **Tokenizálás**

➤ **Lemmatizálás**

➤ **Szófaji elemzés**



# Szabályalapú gépi fordítás

A piros cicám alszik.



**A** [DET] **piros** [ADJ][Nom] **cica** [NOUN][Poss.1Sg][Nom]  
**alszik** [VERB][Prs.NDef.3Sg] . [PUNCT]

- **Mondat szegmentálás**
- **Tokenizálás**
- **Lemmatizálás**
- **Szófaji elemzés**
- **Morfológiai elemzés**

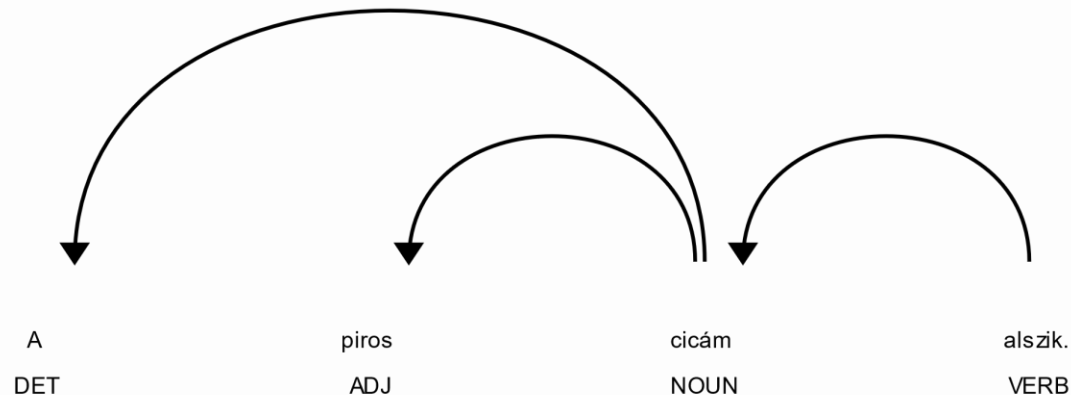


# Szabályalapú gépi fordítás

A piros cicám alszik.



**A** [DET] **piros** [ADJ][Nom] **cica** [NOUN][Poss.1Sg][Nom]  
**alszik** [VERB][Prs.NDef.3Sg] . [PUNCT]



➤ **Mondat szegmentálás**

➤ **Tokenizálás**

➤ **Lemmatizálás**

➤ **Szófaji elemzés**

➤ **Morfológiai elemzés**

➤ **Szintaktikai elemzés**

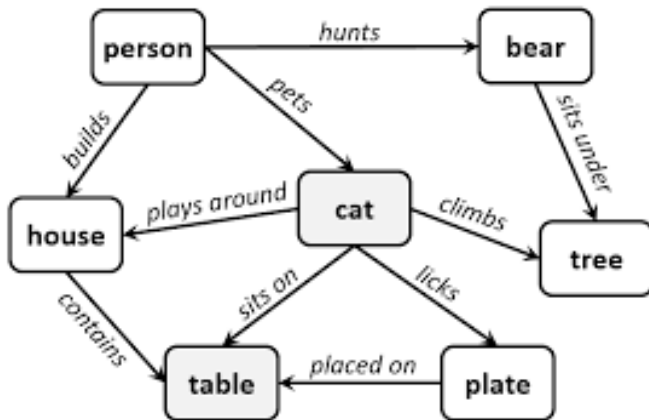


# Szabályalapú gépi fordítás

A piros cicám alszik.



**A** [DET] **piros** [ADJ][Nom] **cica** [NOUN][Poss.1Sg][Nom]  
**alszik** [VERB][Prs.NDef.3Sg] . [PUNCT]



➤ **Mondat szegmentálás**

➤ **Tokenizálás**

➤ **Lemmatizálás**

➤ **Szófaji elemzés**

➤ **Morfológiai elemzés**

➤ **Szintaktikai elemzés**

➤ **Szemantikai elemzés**



# Szabályalapú gépi fordítás

A piros cicám alszik.



**A** [DET] **piros** [ADJ][Nom] **cica** [NOUN][Poss.1Sg][Nom]  
**alszik** [VERB][Prs.NDef.3Sg] . [PUNCT]



**The** [DET] **red** [ADJ][Nom] **cat** [NOUN][Poss.1Sg][Nom]  
**sleep** [VERB][Prs.NDef.3Sg] . [PUNCT]

➤ **Mondat szegmentálás**

➤ **Tokenizálás**

➤ **Lemmatizálás**

➤ **Szófaji elemzés**

➤ **Morfológiai elemzés**

➤ **Szintaktikai elemzés**

➤ **Szemantikai elemzés**

➤ **Szótár**



# Szabályalapú gépi fordítás

A piros cicám alszik.



A [DET] **piros** [ADJ][Nom] **cica** [NOUN][Poss.1Sg][Nom]  
**alszik** [VERB][Prs.NDef.3Sg] . [PUNCT]



The [DET] **red** [ADJ][Nom] **my** **cat** [NOUN][**Poss.1Sg**][Nom]  
**sleeps** [VERB][**Prs.NDef.3Sg**] . [PUNCT]

➤ **Mondat szegmentálás**

➤ **Tokenizálás**

➤ **Lemmatizálás**

➤ **Szófaji elemzés**

➤ **Morfológiai elemzés**

➤ **Szintaktikai elemzés**

➤ **Szemantikai elemzés**

➤ **Szótár**

➤ **Célnyelvi generálás**



# Szabályalapú gépi fordítás

A piros cicám alszik.



**A** [DET] **piros** [ADJ][Nom] **cica** [NOUN][Poss.1Sg][Nom]  
**alszik** [VERB][Prs.NDef.3Sg] . [PUNCT]



**My** [DET] **red** [ADJ][Nom] **cat** [NOUN][**Poss.1Sg**][Nom]  
**sleeps** [VERB][**Prs.NDef.3Sg**] . [PUNCT]

➤ **Mondat szegmentálás**

➤ **Tokenizálás**

➤ **Lemmatizálás**

➤ **Szófaji elemzés**

➤ **Morfológiai elemzés**

➤ **Szintaktikai elemzés**

➤ **Szemantikai elemzés**

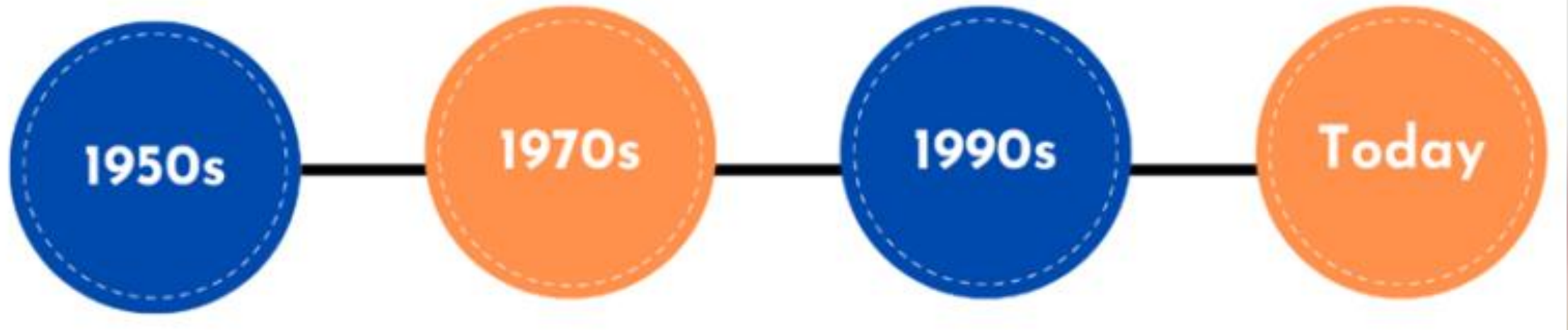
➤ **Szótár**

➤ **Célnyelvi generálás**





# Nyelvtechnológia története



**Modern számítógép  
megjelenése**

**Gépi fordítás**

**Turing teszt**

**Szabályalapú  
megoldások**

**ELIZA**

**Ontológiák,  
szemantikai hálók**

**Statisztikai megoldások**

**Adatvezérelt tanulás**

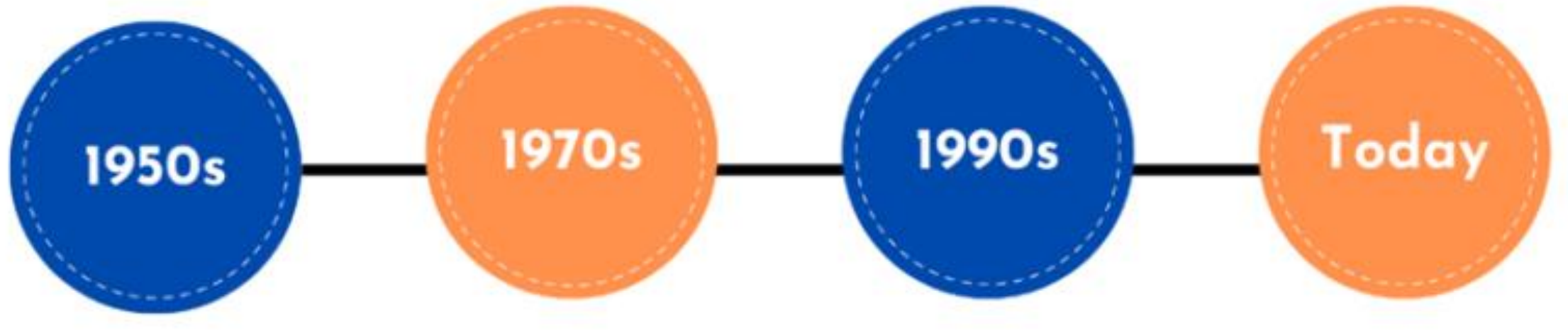


# Statisztikai módszerek

- Elemzők (morfológiai, szintaktikai, névelemfelismerő):
  - Rejtett Markov-modell
  - Viterbi
- Statisztikai gépi fordítás:
  - Naiv Bayes-tétel
  - Fordítási modell:
    - EM algoritmus
  - Nyelvmodell:
    - N-gramm modell



# Nyelvtechnológia története



**Modern számítógép  
megjelenése**

**Gépi fordítás**

**Turing teszt**

**Szabályalapú  
megoldások**

**ELIZA**

**Ontológiák,  
szemantikai hálók**

**Statisztikai megoldások**

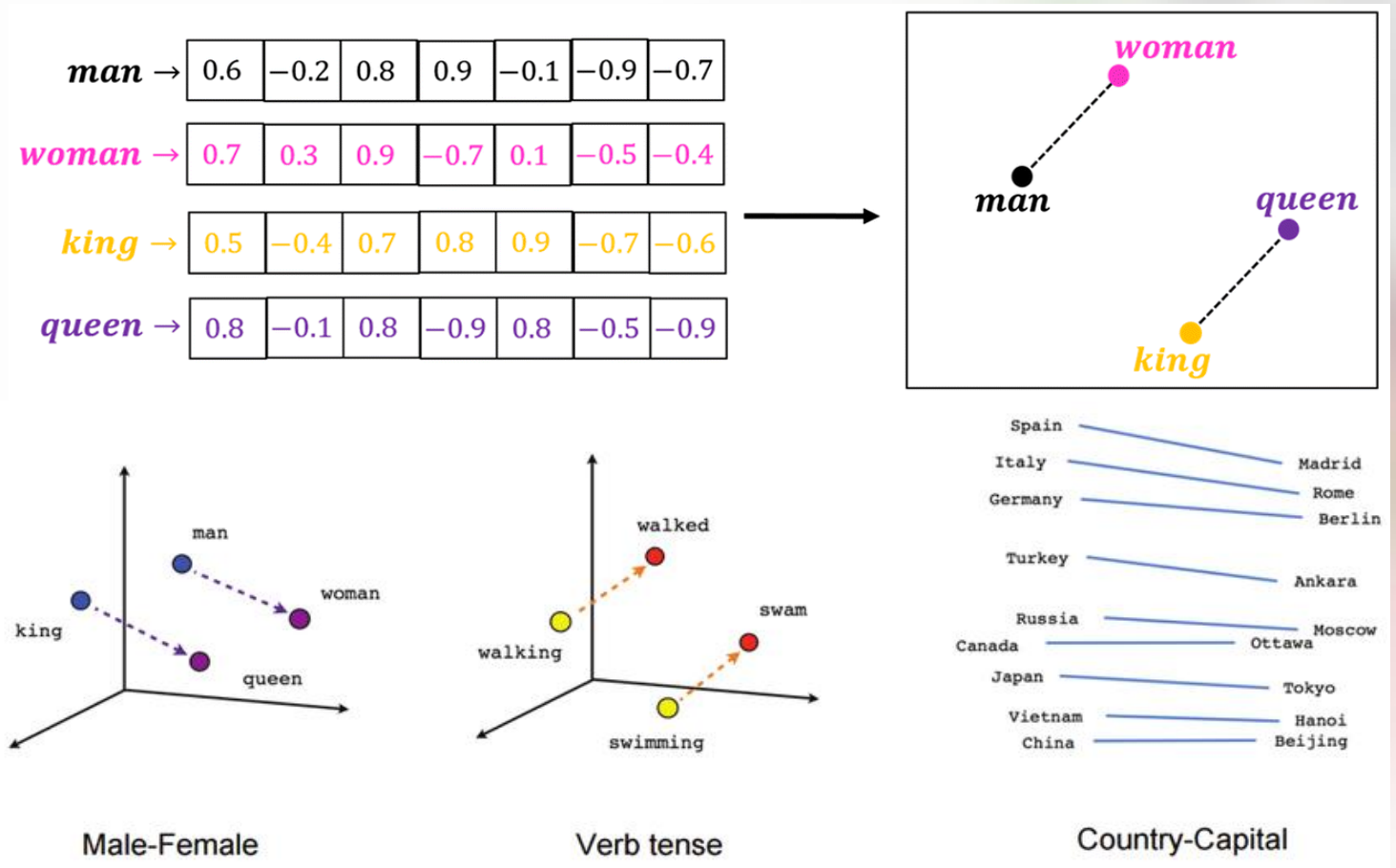
**Adatvezérelt tanulás**

**Neurális  
megoldások**



# Szóbeágyazás (word embedding)

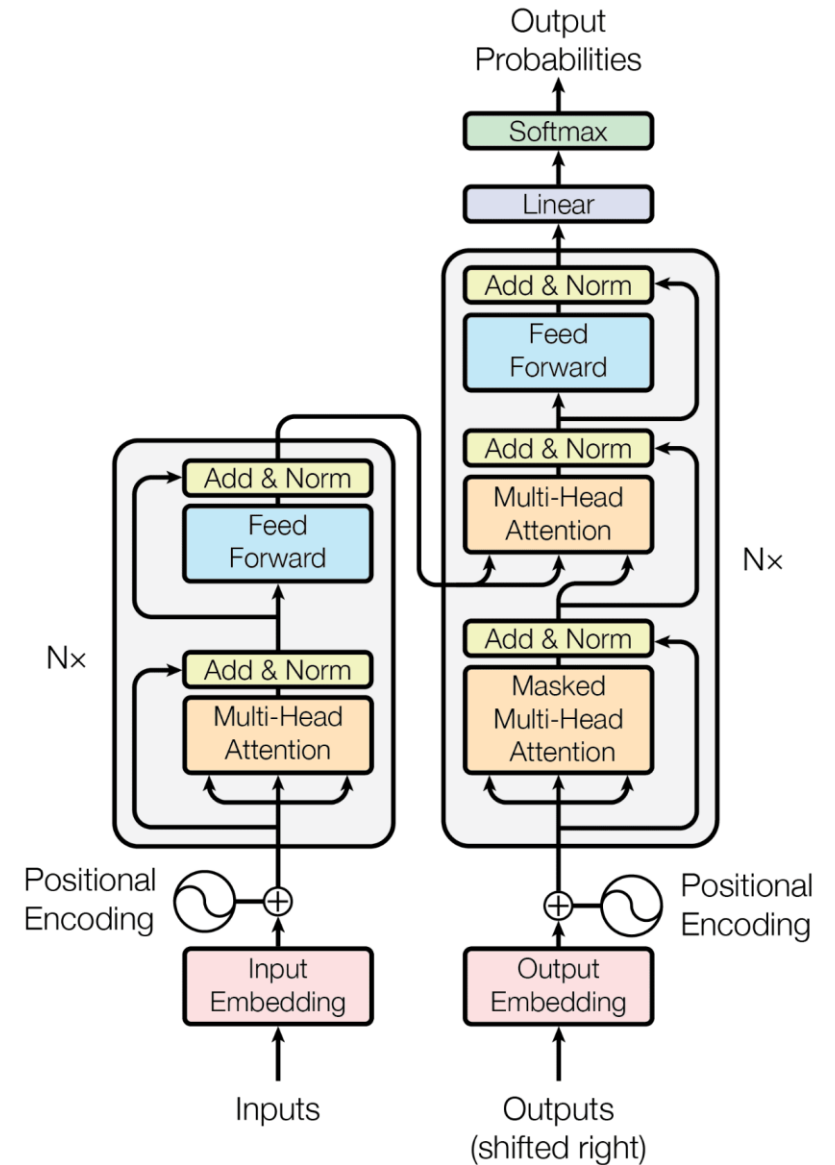
- Mikolov és mtsai.
- Google, 2013
- Word2vec
  
- A szavak a környezetükkel jellemezhetőek
  
- Neurális nyelvmodell





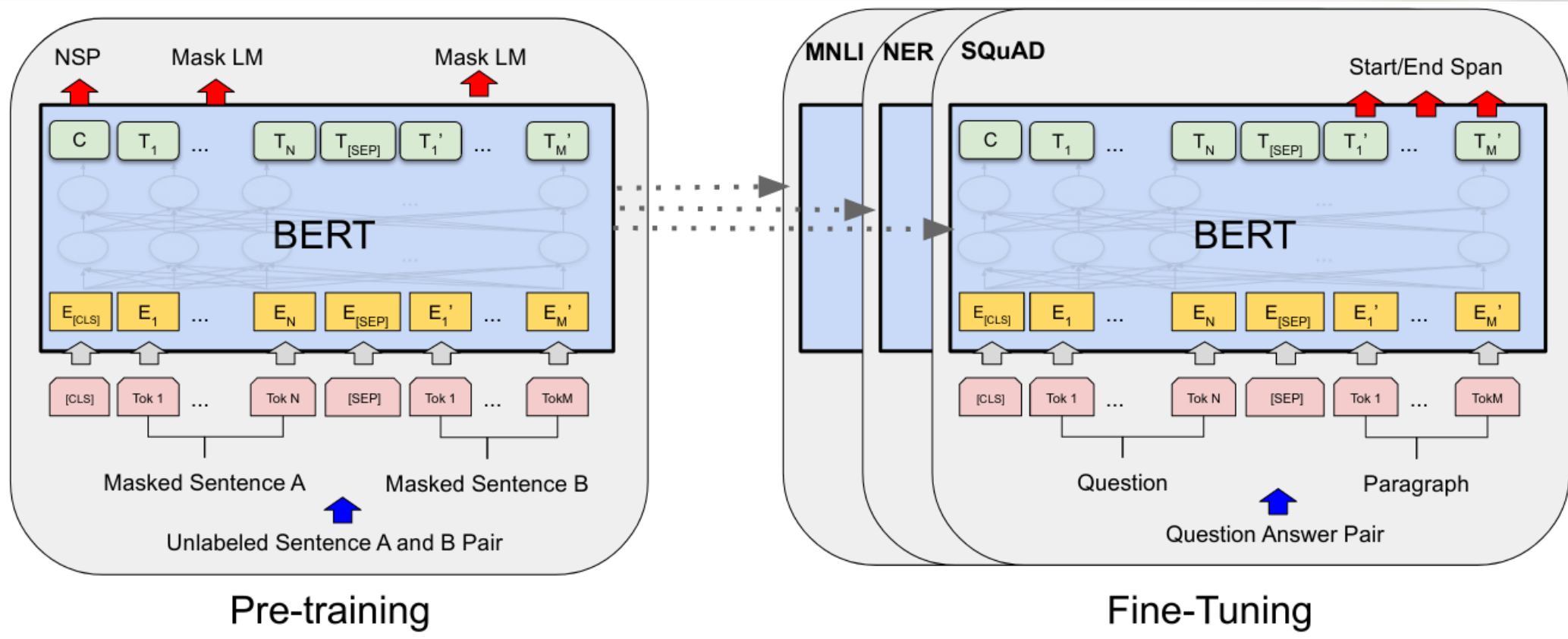
# Tranzformer modell

- Vaswani és mtsai.
- Google, 2017
- Enkóder – Dekóder architektúra
  - Figyelmi mechanizmus
    - Adott szövegben lévő szavak egymáshoz képesti viszonya
- Környezetfüggő reprezentáció
- Párhuzamos feldolgozás





# Előtanítás – Finomhangolás





# Előtanítás

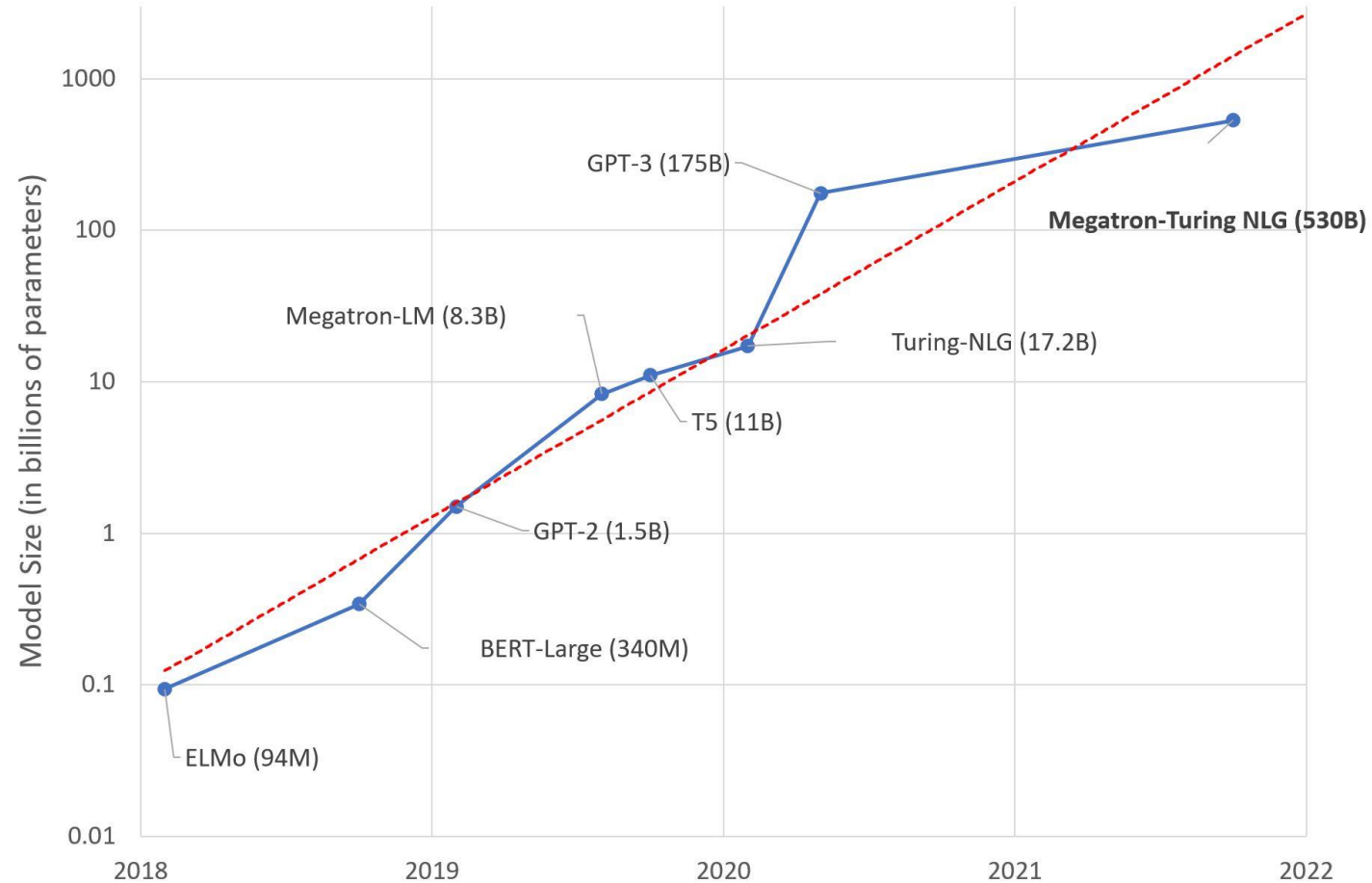
- Általános nyelvi tudás
- Nagy mennyiségű adat
- Nagy erőforrásigény
  
- GPT-3/GPT-4

# – Finomhangolás

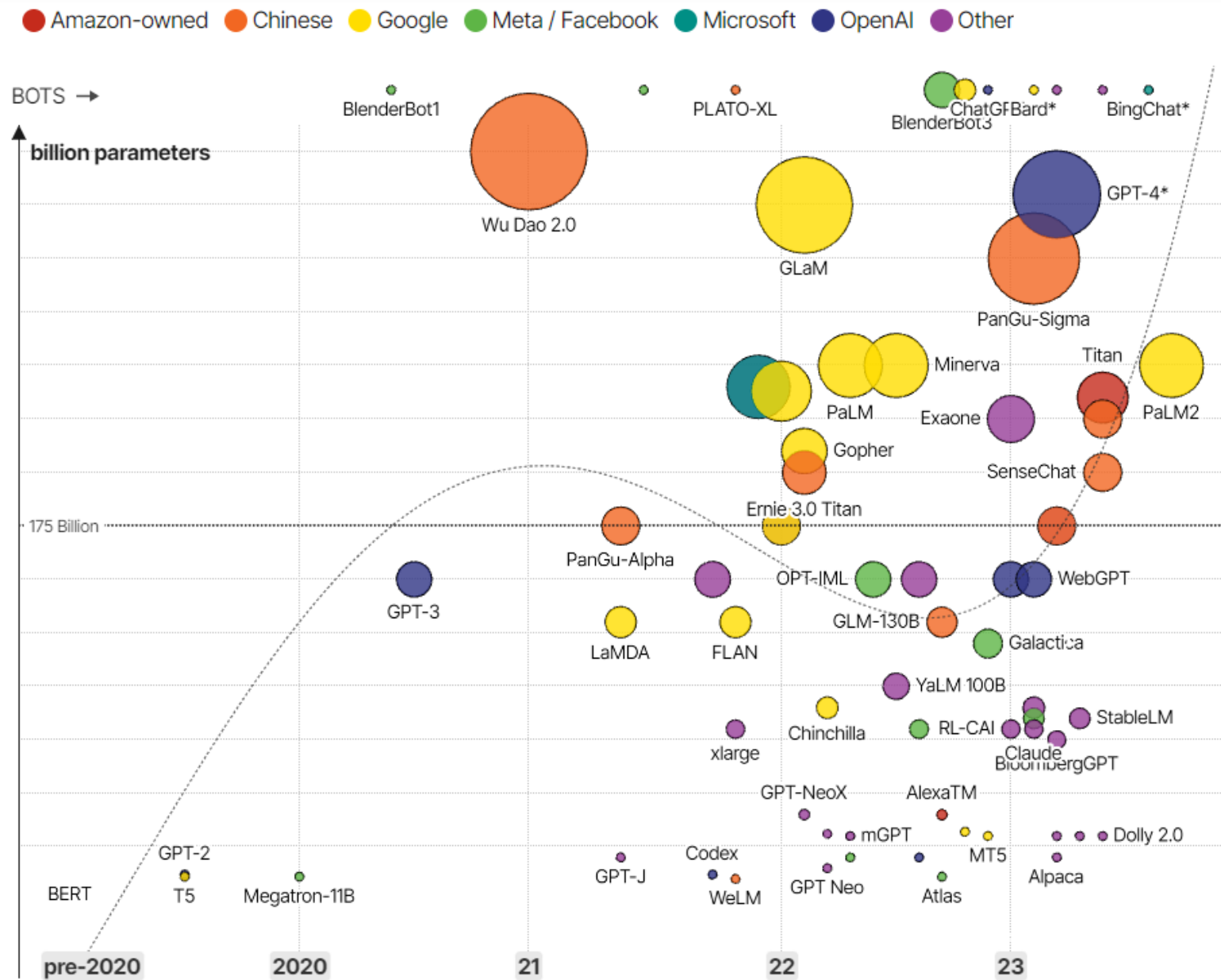
- Feladat-specifikus tudás
- Kevés adat
- Kisebb erőforrásigény
  
- ChatGPT



# Nagy nyelvi modellek





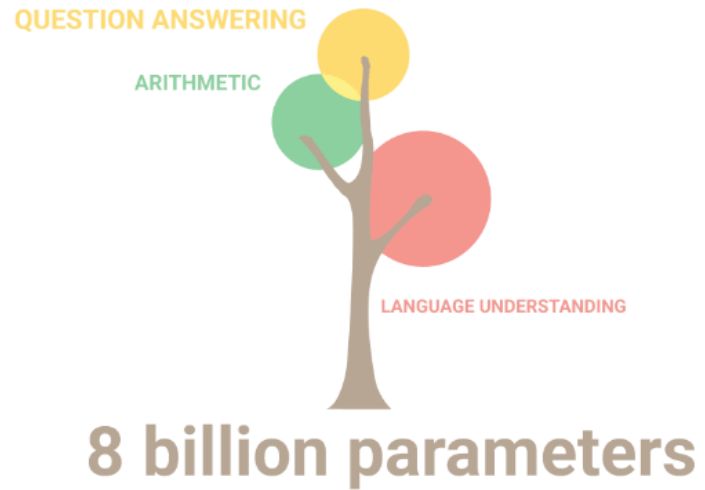


David McCandless, Tom Evans, Paul Barton  
Information is Beautiful // May 2023

source: news reports, [LifeArchitect.ai](#)  
\* = parameters undisclosed // see [the data](#)



# Nagy nyelvi modellek





Default (GPT-3.5)



Írj egy előadás vázlatot, amit a conTEXT konferencián fogok előadni! A téma: Mérföldkövek a ChatGPT-ig



Természetesen, örömmel segítek egy előadásvázlat elkészítésében a "Mérföldkövek a ChatGPT-ig" témában a conTEXT konferenciára. Az alábbiakban találsz egy vázlatot, amely segíthet strukturálni az előadásodat:

## I. Bevezetés

A. A mesterséges intelligencia fejlődése

1. Az MI történeti áttekintése
2. A gépi tanulás térhódítása

## II. Első generációs Chatbotok

A. ALICE és ELIZA: Az első lépések

B. Korlátok és kihívások

1. Korlátozott szókincs és kontextus
2. Szabályalapú rendszerek



## III. A neurális hálózatok korszaka

A. Deep Learning és a nyelvi modellek

Regenerate

Send a message





# ChatGPT a kulisszák mögött

- OpenAI
- GPT – Generative Pretrained Transformer
- GPT-3 (175 milliárd paraméter) és GPT-4 (1,76 billió paraméter)
- Csak dekóder
- Következő szó prediktálása (Causal Language Modeling)
  
- Promptolás
  - prompting, prompt tuning, prompt engineering



# Promptolás vs. finomhangolás

- Finomhangolás: előtanított nyelvmodell tovább **tanítása** specifikus feladatra.
  - Pl. szentimentelemzés:
    - Szöveg: Boldog vagyok! → Osztálycímke: pozitív
- Promptolás: nincsen tanítás, a prompt egy kiindulási/előzmény szöveg, ami alapján a modell folytatja a generálást.
  - Input:
    - Vidám vagyok = pozitív
    - Szomorú vagyok = negatív
    - Boldog vagyok =
  - Output:
    - pozitív



# Promptolás

## Zero-shot

- **Bemenet:**
  - Add meg a következő mondat szentimentjét:  
Vidám vagyok!
- **Kimenet:**
  - pozitív

## Few-shot

- **Bemenet:**
  - Boldog vagyok! = pozitív  
Szomorú vagyok. = negatív
  - Add meg a következő mondat szentimentjét:  
Vidám vagyok! =
- **Kimenet:**
  - pozitív

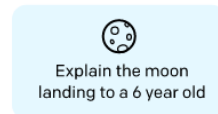


# Utasításkövető GPT (InstructGPT)

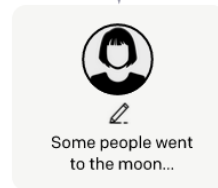
Step 1

**Collect demonstration data, and train a supervised policy.**

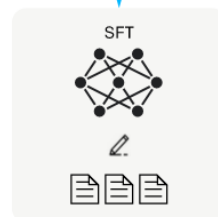
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



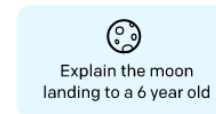
This data is used to fine-tune GPT-3 with supervised learning.



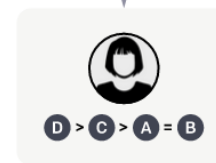
Step 2

**Collect comparison data, and train a reward model.**

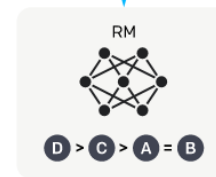
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



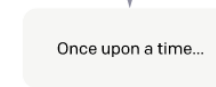
Step 3

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.





# Utasításkövető GPT (InstructGPT)

## 0. Nagy nyelvi modell





# Utasításkövető GPT (InstructGPT)

0. Nagy nyelvi modell
1. Nyelvmódel finomhangolása felügyelt tanulással
  - Utasítások/kérdések és válaszok



# Utasításkövető GPT (InstructGPT)

0. Nagy nyelvi modell
1. Nyelvmodell finomhangolása felügyelt tanulással
  - Utasítások/kérdések és válaszok
2. Emberi visszajelzések gyűjtése és jutalom modell (reward model) tanítása
  - A modell válasza jó-e vagy mennyire jó



# Utasításkövető GPT (InstructGPT)

0. Nagy nyelvi modell
1. Nyelvmodell finomhangolása felügyelt tanulással
  - Utasítások/kérdések és válaszok
2. Emberi visszajelzések gyűjtése és jutalom modell (reward model) tanítása
  - A modell válasza jó-e vagy mennyire jó
3. Nyelvmodell tovább finomhangolása a jutalom modellel és a megerősítéses tanulással



# Magyar nyelvű nagy nyelvi modellek

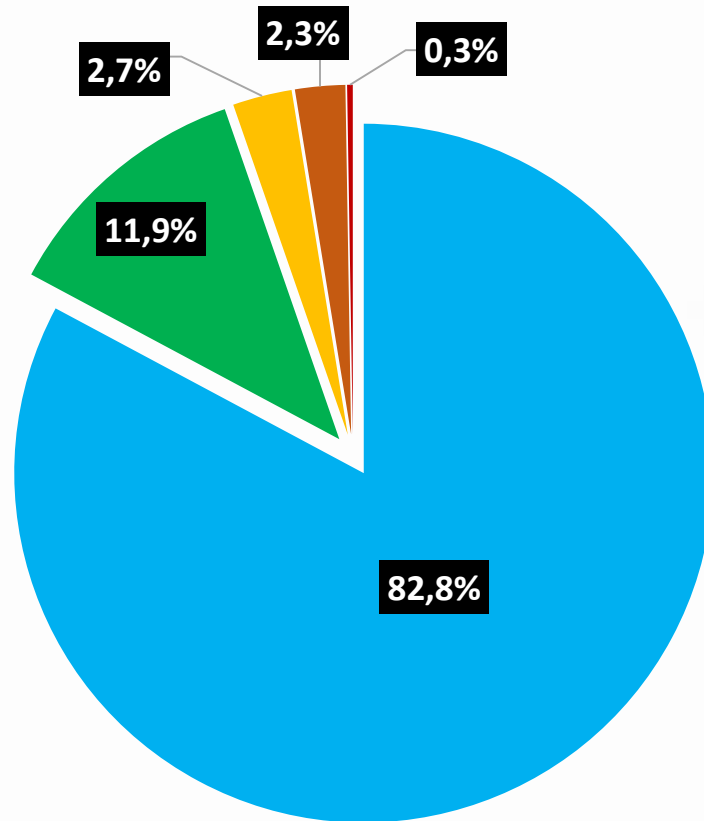


# Magyar nyelvű nagy nyelvi modellek

- HILANCO GPTX:
  - 6,7 milliárd paraméteres
  - angol-magyar
- PULI GPT-3SX:
  - 6,7 milliárd paraméteres
  - magyar
- PULI GPTrío:
  - 7,7 milliárd paraméteres
  - magyar-angol-kínai



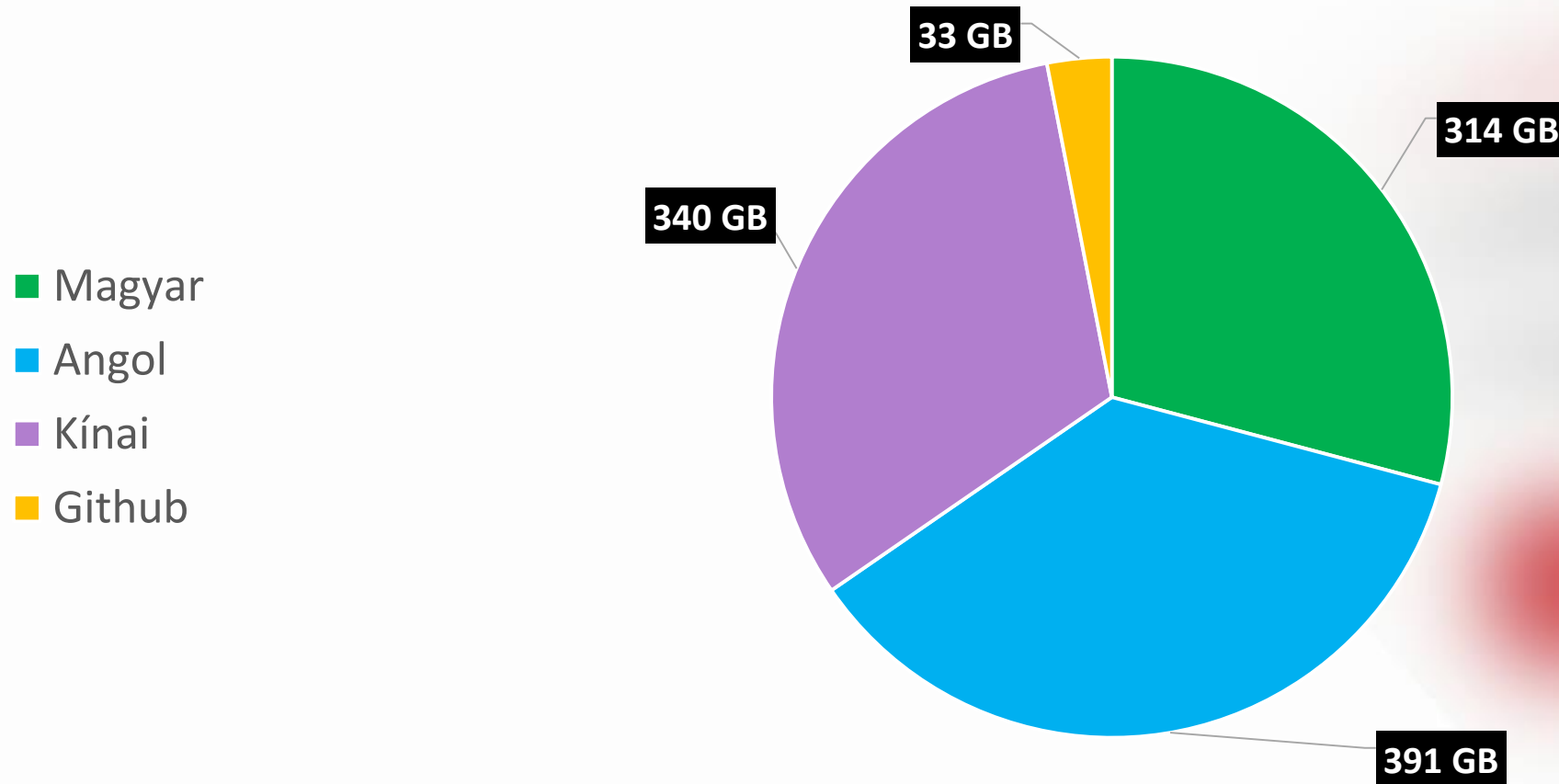
# Magyarnyelvű tanítóadat: ~41 milliárd szó



- Common Crawl (Internet)
- Nemzetközi gyűjtemények
- Magyar Nemzeti Szövegtár
- Közösségi Média
- Magyar Wikipédia



# Háromnyelvű tanítóadat: >150 milliárd szó





# PULI vs GPT-3



## **PULI**

(ParancsPULI)

~7 milliárd paraméter

~150 milliárd szó

**~40 milliárd magyar szó**

~2000 látogató



## **GPT-3**

(ChatGPT)

~175 milliárd paraméter

~400 milliárd szó

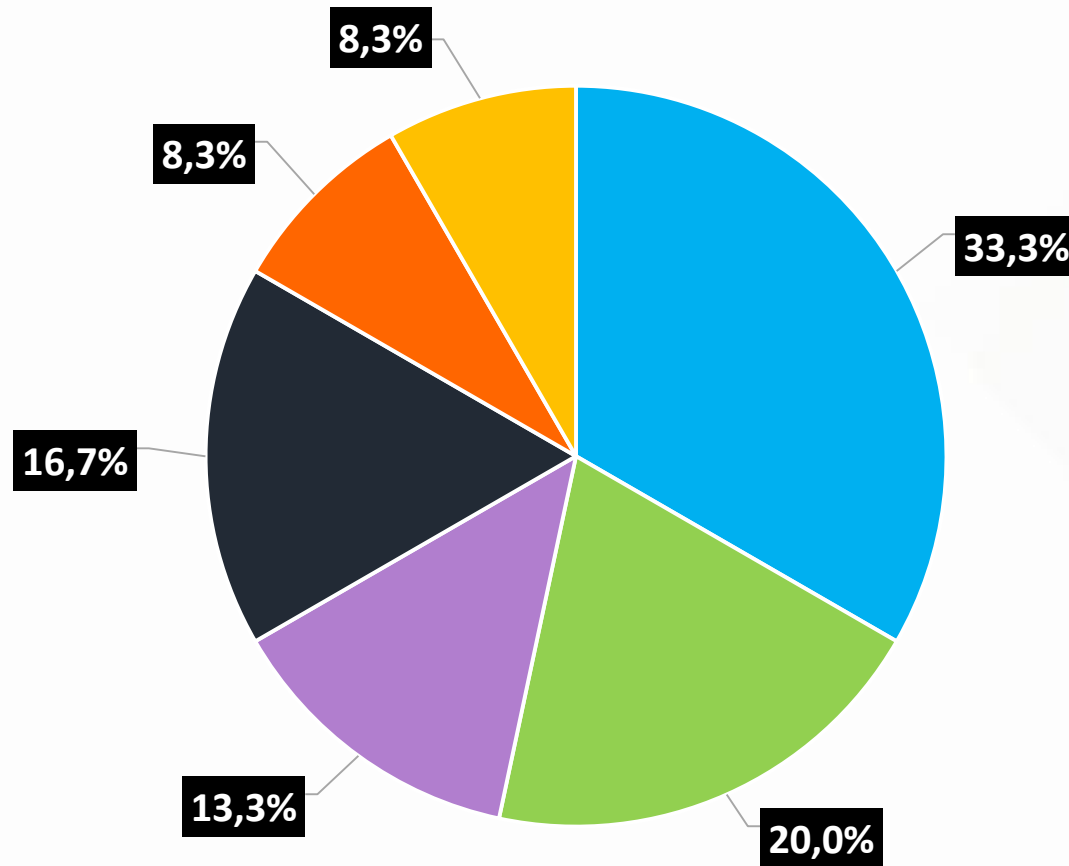
**~130 millió magyar szó**

>100 millió látogató





# ParancsPULI (6000 prompt)



- Stanford Alpaca
- HuLU
- SQL
- Fordítás
- Érettségi feladatok
- Rövid kérdések



ParancsPULI

PULI-GPT-35X

Fordító

Kivonatoló

Szentimentelemző



# ParancsPULI

Magyar nyelvű utasításkövető finomhangolt PULI-GPTrio (GPT-3; 7,67 milliárd paraméter)



Generálás

Utasítás / Instruction

Hogyan köszönjek el az előadásom végén a conTEXT 2023 konferencián?

Remélem, hogy élvezted az előadást és hogy hasznos volt számodra.  
Találkozunk jövőre a következő konferencián!

NYTK • 2023

<https://juniper.nytud.hu/demo/parancspuli>