

# Deep learning alapú beszédszintézis

Dr. Zainkó Csaba



# Deep learning - Beszédszintézis 01

Alapok

## Mit várjunk, várhatunk el? 02

Jelen állapot

## Merre tovább? 03

Jövő



# Beszéd-szintézis



• Szöveg



Beszéd

*„Nem is gondolnánk milyen gyakran használunk már ma is neurális hálózatokat vagy azok eredményeit a mindennapokban.”*

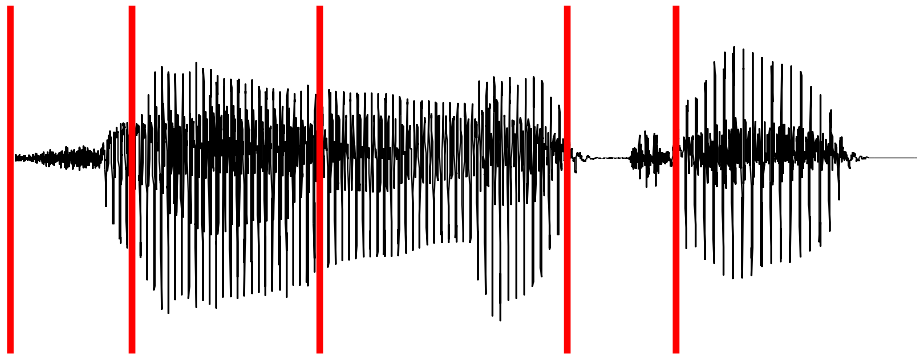


# Technológia

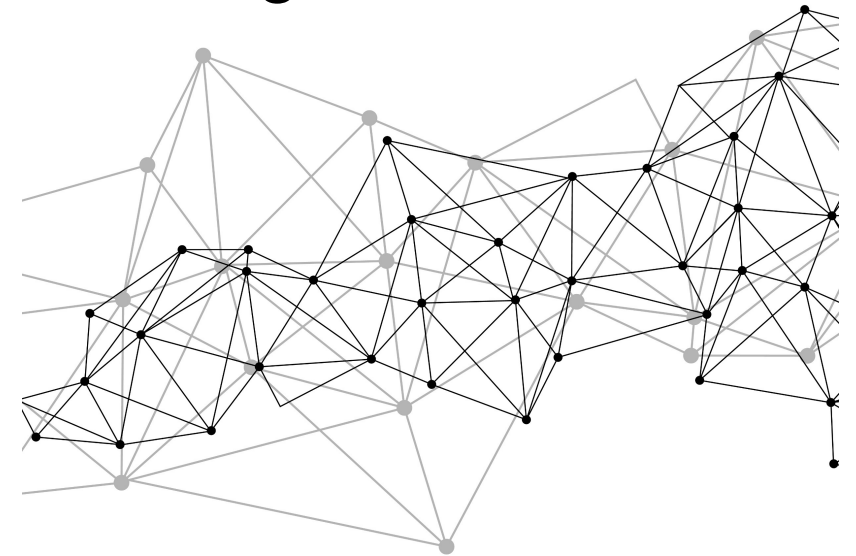
- Parametrikus



- Összefűzéses



- Deep learning

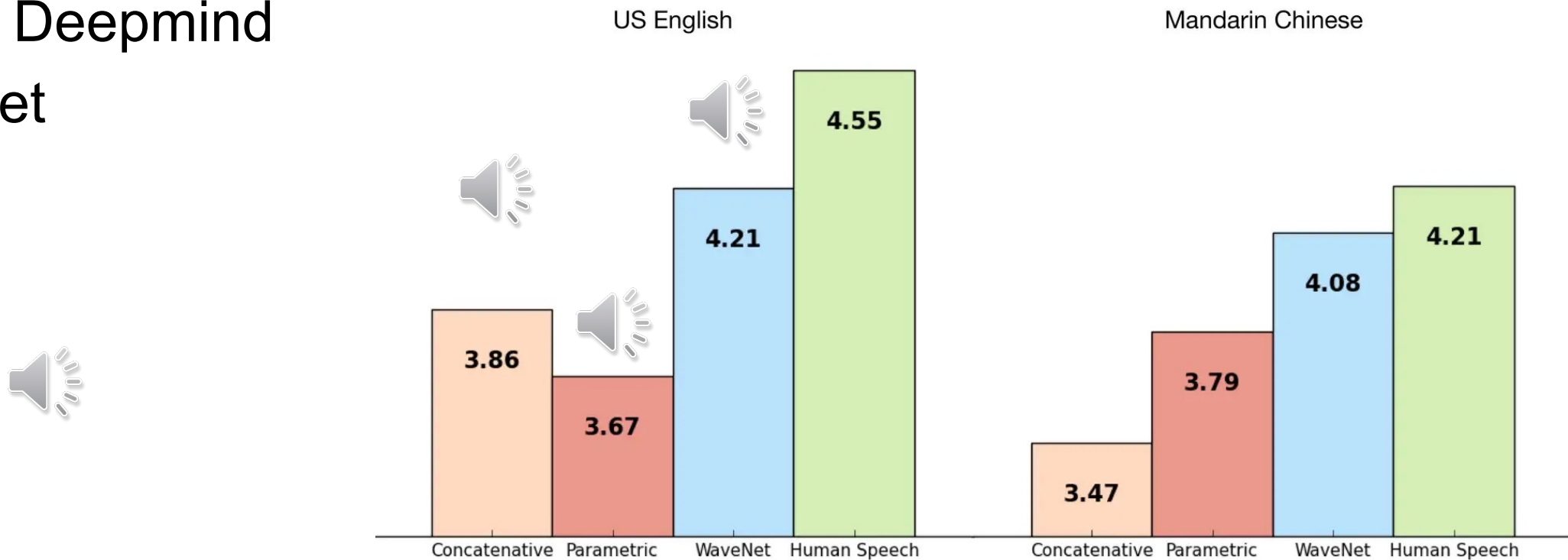


<http://www.freepik.com>

# Deep learning

- 2016
- Google Deepmind
- WaveNet

<https://www.deepmind.com/blog/wavenet-a-generative-model-for-raw-audio>



# Minőség – Hogyan mérjük?

- Szubjektív

- MOS – Mean Opinion Score

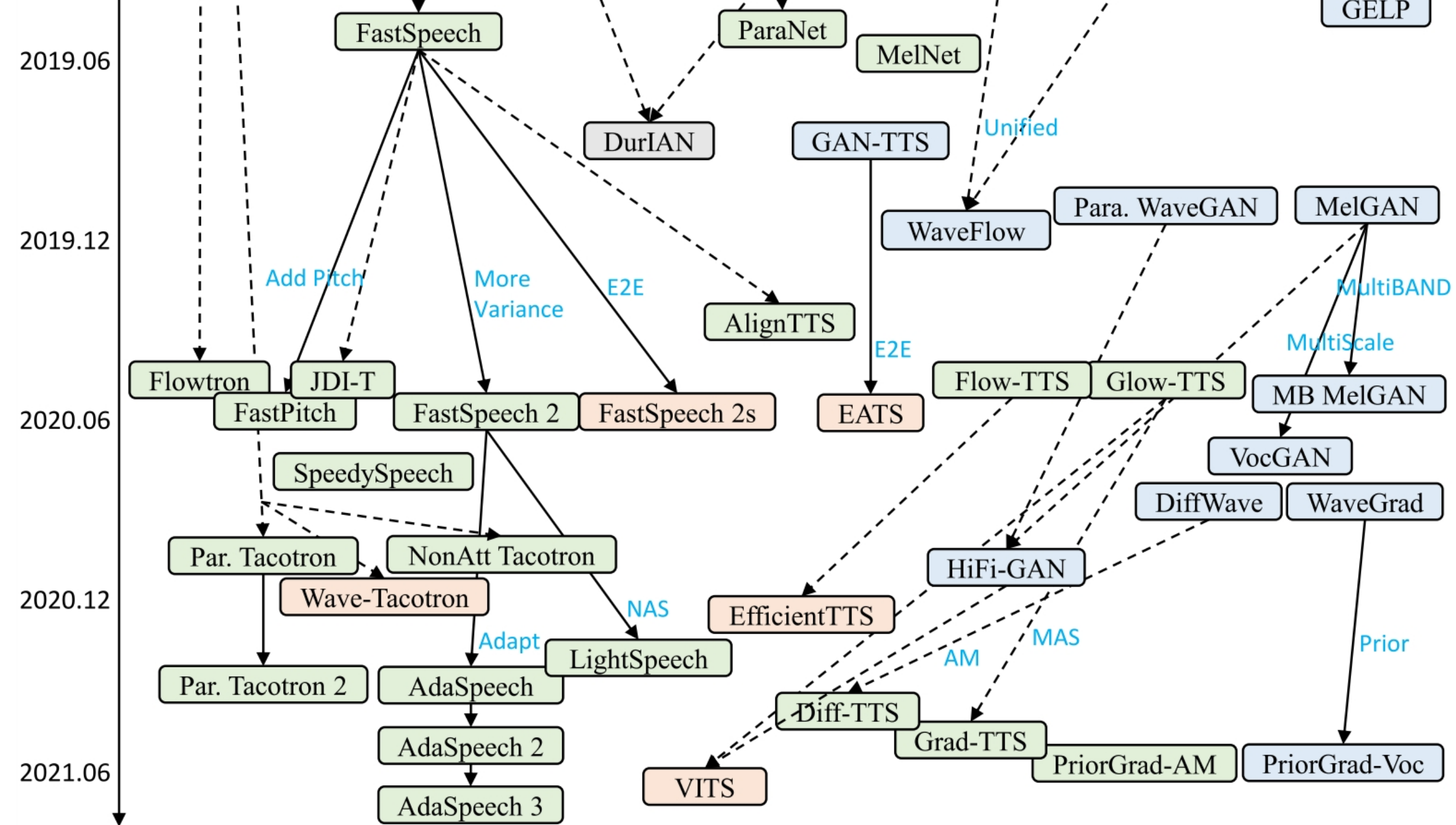
- 1 – legrosszabb
    - 5 – legjobb

- MUSHRA - Multiple Stimuli with Hidden Reference and Anchor

- Referencia minőség
      - Alsó horgony
      - Felső horgony

- Objektív

- Referencia
  - Felismerővel (WER)
  - Gépi tanulósos mód (VoiceMOS)



# DNN-TTS

- End-to-End: szöveg (hangsorozat) → Hullámforma
- 2 lépcsős:  
szöveg (hangsorozat) → spektrális (mel-spectrum) → Hullámforma
- Tanítás: GPU sok óra beszéd + szöveg
  - Tanítás idő: néhány óra, 1-2 nap (1-2 GPU-n)
- Generálás: CPU/GPU





# Mit tud, mit kell tudni?

- Közel emberi hangminőség
- Kezelni:
  - Számok
  - Dátum
  - Összegek
  - Időpontok
  - Rövidítések
  - ...
- Gyors: legalább valós idő
- Kérdések (eldöntendő)
- Kontrollálható kimenet
  - kivétel szótár

# Mit tud(hatna) a mostani technológia?

- Testreszabás

- Férfi/női
- Beszélő választás
- Egyedi hang
- Stílus
- Érzelem???

## Egyedi hang

- Új tanítás: 1-2 óra beszéd
- Adaptálás: 5-10 perc beszéd
  
- És ha csak 1 mondat van?

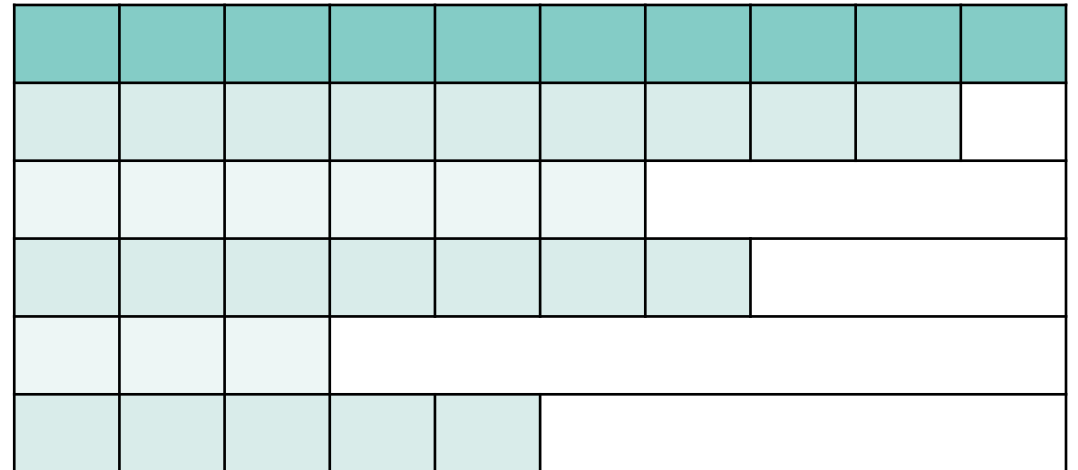
- Szövegjavítás

*Ez a szöveg mindenki számára jól érthető.*



# Mit tud(hatna) a mostani technológia?

- Válaszidő
  - CPU/GPU
  - max 100-200ms késleltetés
- Batch feldolgozás (GPU)
  - Több azonos jellegű feladat
- Párhuzamos generálás



# Chat - Párbeszéd

- Sima szövegfelolvasás helyett: Párbeszéd
- Alkalmazkodni a másik félhez, integrált TTS
  - Stílusban
  - Beszédtempóban
  - Beszélő (ffi/nő)



Különálló független komponensként egyre kevésbé értelmezhető

# Merre tovább?

- 5 éves gyerek szintje
- Szótár, kiejtés

PÉLDA:

*„Ezekre a kérdésekre keressük a választ az idei **dataSTREAM** konferencián.”*



Tudjuk-e automatizálni az ember hozzáadott tudását?

# LLM – Large Language Model

- 100 óra beszéd (kb. 4,3 millió karakter, <1 millió szó)
- GPT-3 – 10 milliárd szó!
  
- Nagy modellek adhatnak megoldást
  - Wave2vec
  - Bert alapú
  - GPT típusú modellek

# Mit várunk?

- Jobb szubjektív élményt
- Pontosabb nyelvfelismerést
- Szövegkorrekciót (hogyan teszteljem?)
  - Ügyfélélmény v. Szöveghű viselkedés
- Előzmények figyelembevétele
- Szöveg megértése ???

# Veszélyek

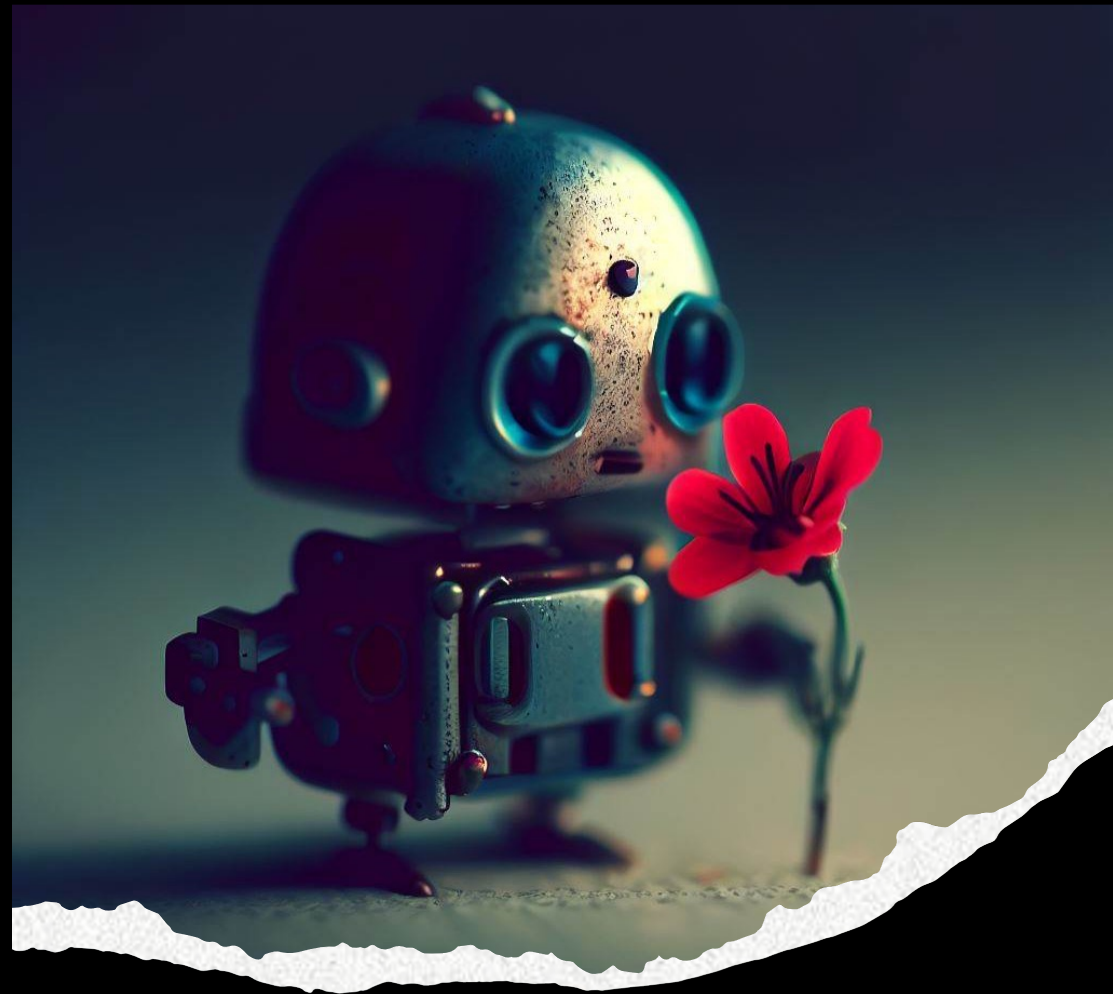
- FakeVoice
- Hibás értelmezés
- Félreérthető hangsúlyozás
- Sértő viselkedés (szomorú hír vidám hangon felolvasása)
- Lehet, hogy „Gépies” hangot kell majd fejleszteni?





# Köszönöm

Deep learning alapú  
beszédszintézis



Dr. Zainkó Csaba    [zainko@tmit.bme.hu](mailto:zainko@tmit.bme.hu)

Tartalom	Human: Dr. Zainkó Csaba
Szöveg	Human: Dr. Zainkó Csaba
TTS beszédminták	<b>AI:</b> BME Profivox-DNN
Design	<b>AI</b> + Human: (MS powerpoint + ZCs)
Rajzok (4db)	Human: Dr. Zainkó Csaba
Illusztráció (1db)	?: <a href="https://www.freepik.com">https://www.freepik.com</a>
Képek (5db)	<b>AI:</b> <a href="https://www.bing.com/images">https://www.bing.com/images</a>

