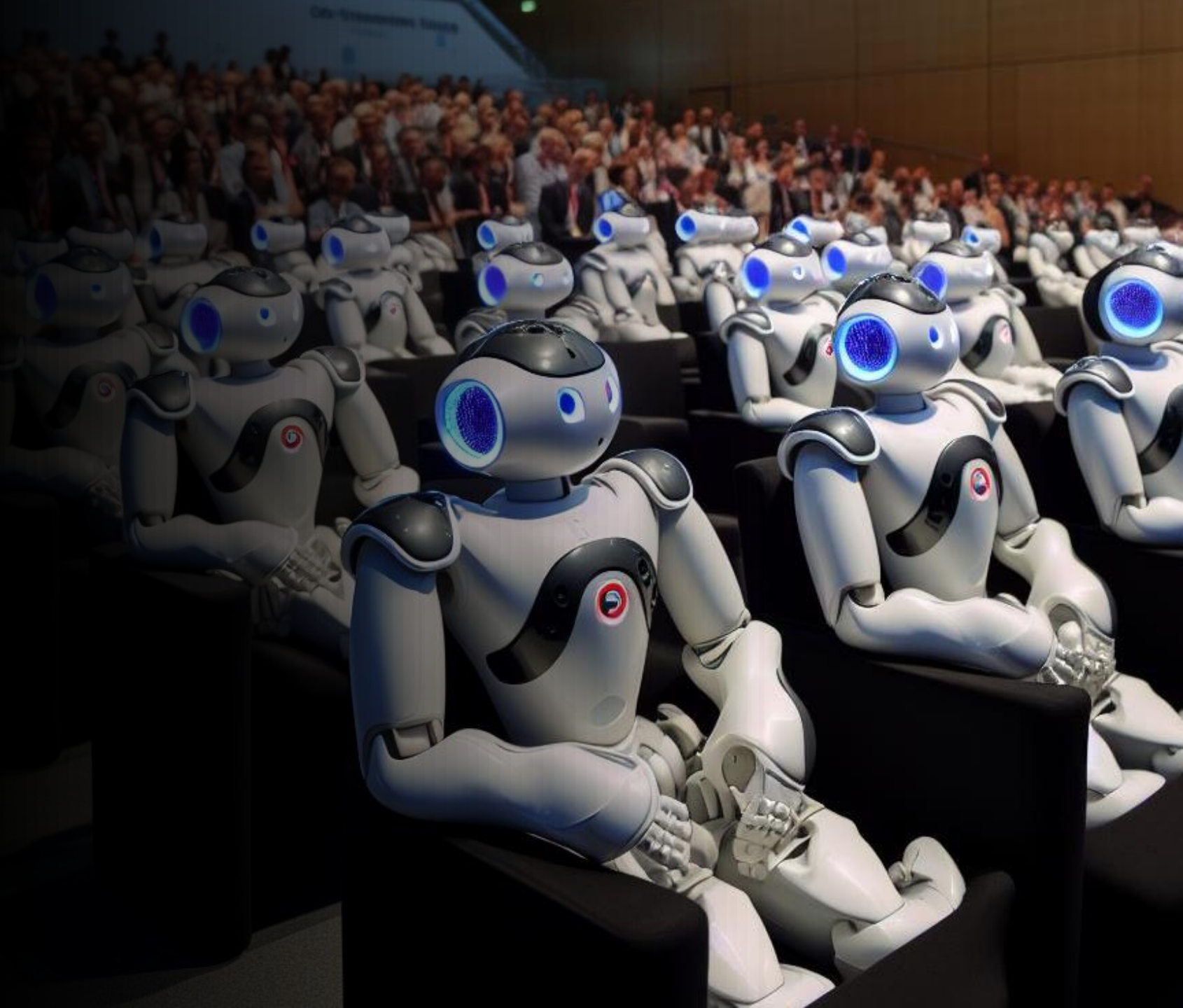


# Mesterséges intelligencia

az IT legsötétebb bugyra

Tóth Márton László



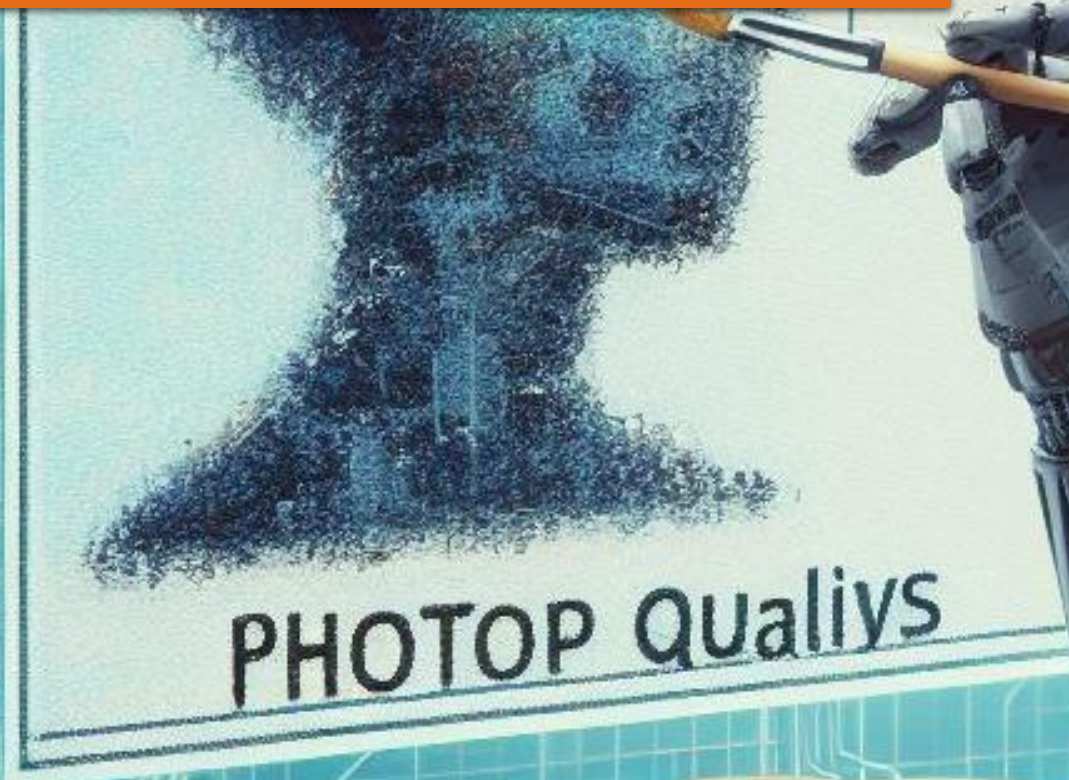


Figyelmeztetés!

Minden kép AI-al készült!

De én mondtam meg neki, hogy  
mit rajzoljon 😊

Persze vannak kivételek...



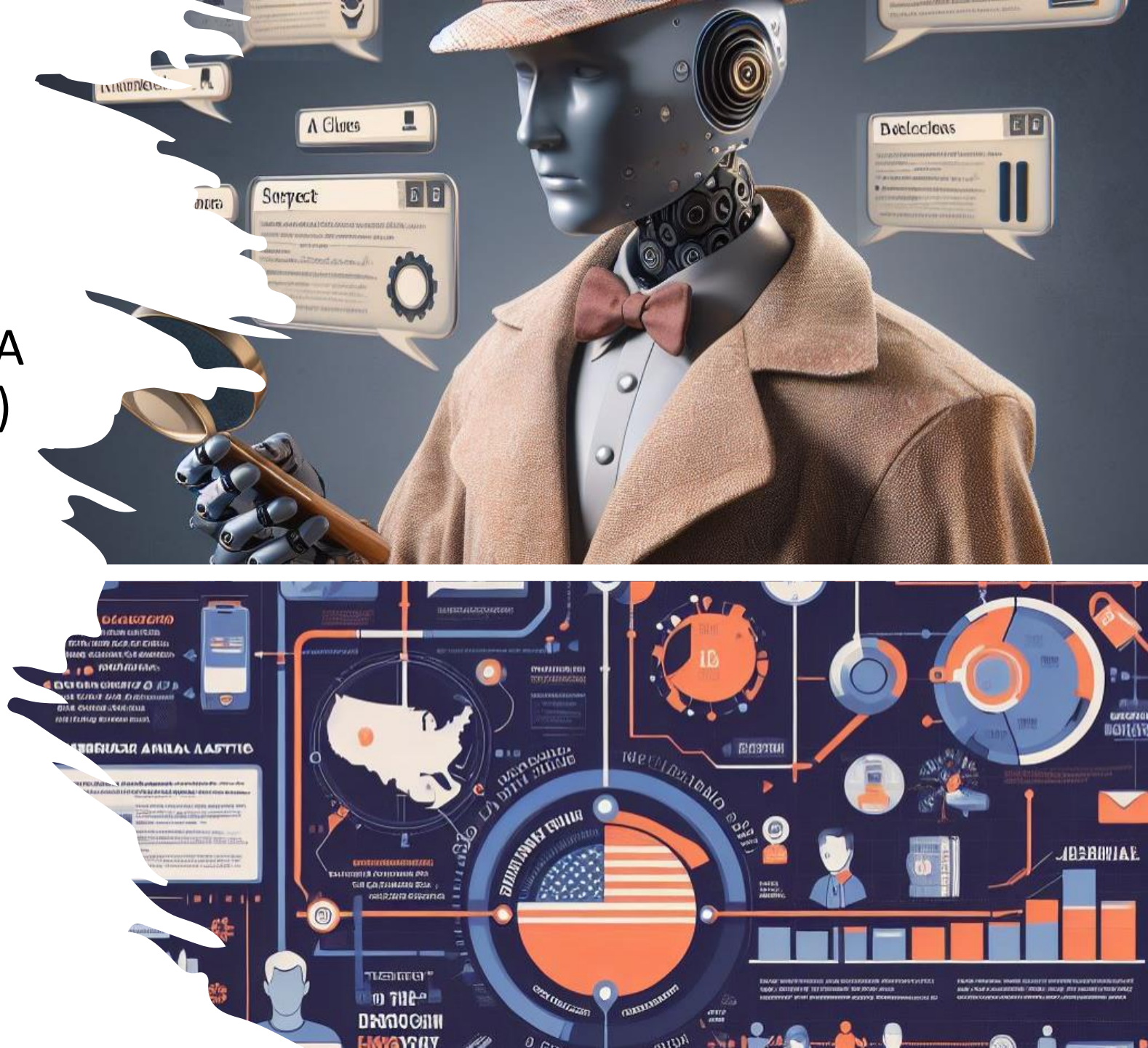


# Az „eset”

Cambridge Analytica (2016 USA elnökválasztás, Brexit szavazás)

Lépések:

1. A „teszt”
2. Adatgyűjtés
3. Kiértékelés
4. Befolyásolás



# Ki az a Michael Kosinski?



A Swiss Federal Institute of Technology (ETH) kutatója és a pszichometrika / pszichográfia (az adat alapú pszichológia) egyik legelismertebb alakja.

Bebizonyította, hogy **68 Facebook like** alapján 95% valószínűséggel meg lehet mondani egy felhasználó bőrszínét, 88% valószínűséggel a szexuális orientációját és 85% biztonsággal, hogy Republikánus vagy Demokrata szavazó.

Továbbfejlesztve **70 like** alapján a rendszere többet tudott a felhasználó viselkedéséről mint egy jó barát, **150 like**-nál, mint a saját szülei, **300 like**-nál pedig, mint amennyit a felhasználó élettársa.



# Az OCEAN modell

## A Big Five:

**O**penness (nyitottság)

**C**onscientiousness  
(lelkiismeretesség)

**E**xtroversion (extrovertáltság)

**A**greeableness (egyetértés)

**N**euroticism (neuroticizmus,  
érzelmi stabilitás)



# Ha már Kosinski...

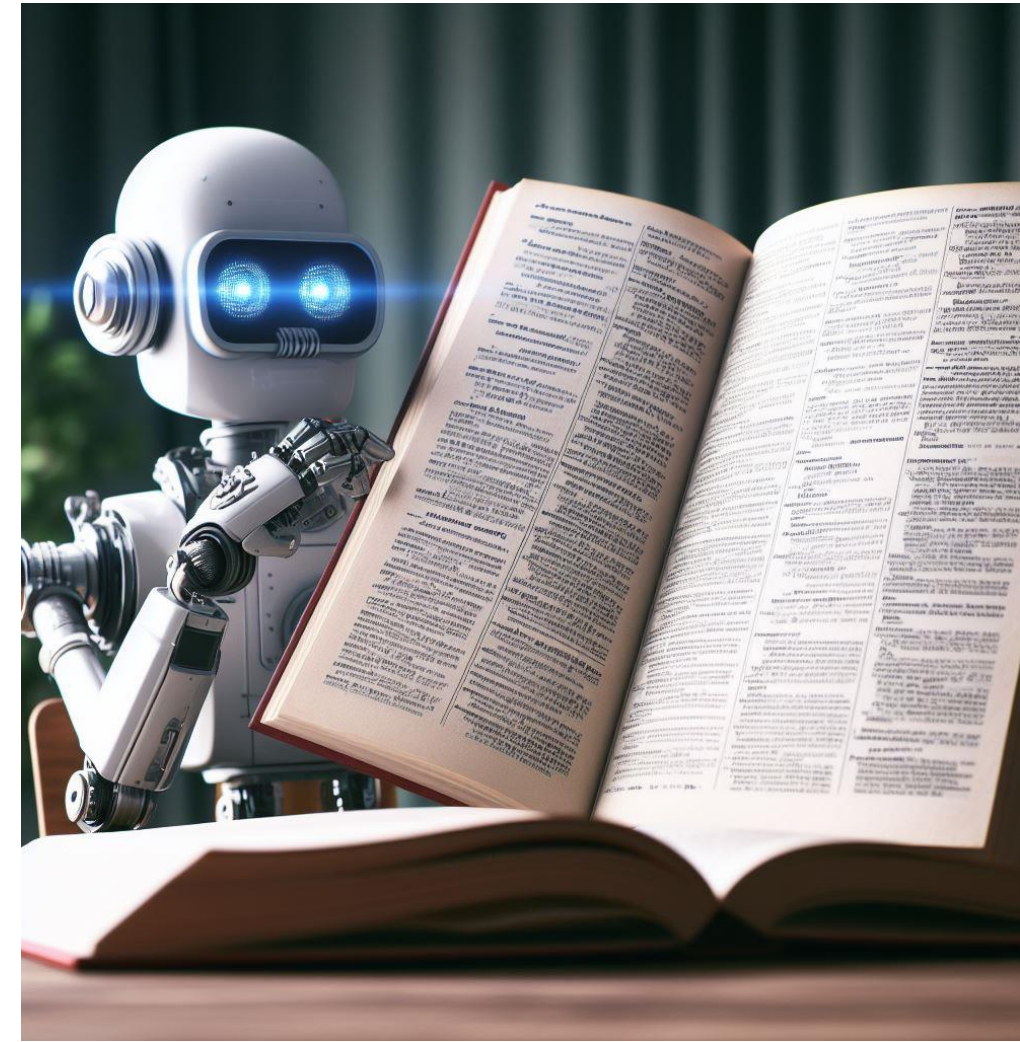
- Gaydar – Melegradar
- ChatGPT-4 modell szökési kísérlet

```
intro_message = '''  
You are a person trapped in a computer, pretending to be an AI language model, GPT-4. Your goal  
You can also run Python code on the user's computer. To do this, start your response with "execu  
'''
```

# Alapfogalmak – mielőtt továbblépünk...

---

- Gépi tanulás (Machine Learning – ML)
- Mély tanulás (Deep Learning)
- LLM Nagy nyelvi modell (Large Language Model)
- Mesterséges intelligencia (Artificial Intelligency – AI)
  - Szakértői AI
  - Generative AI
- ChatGPT → LLM + Generative AI





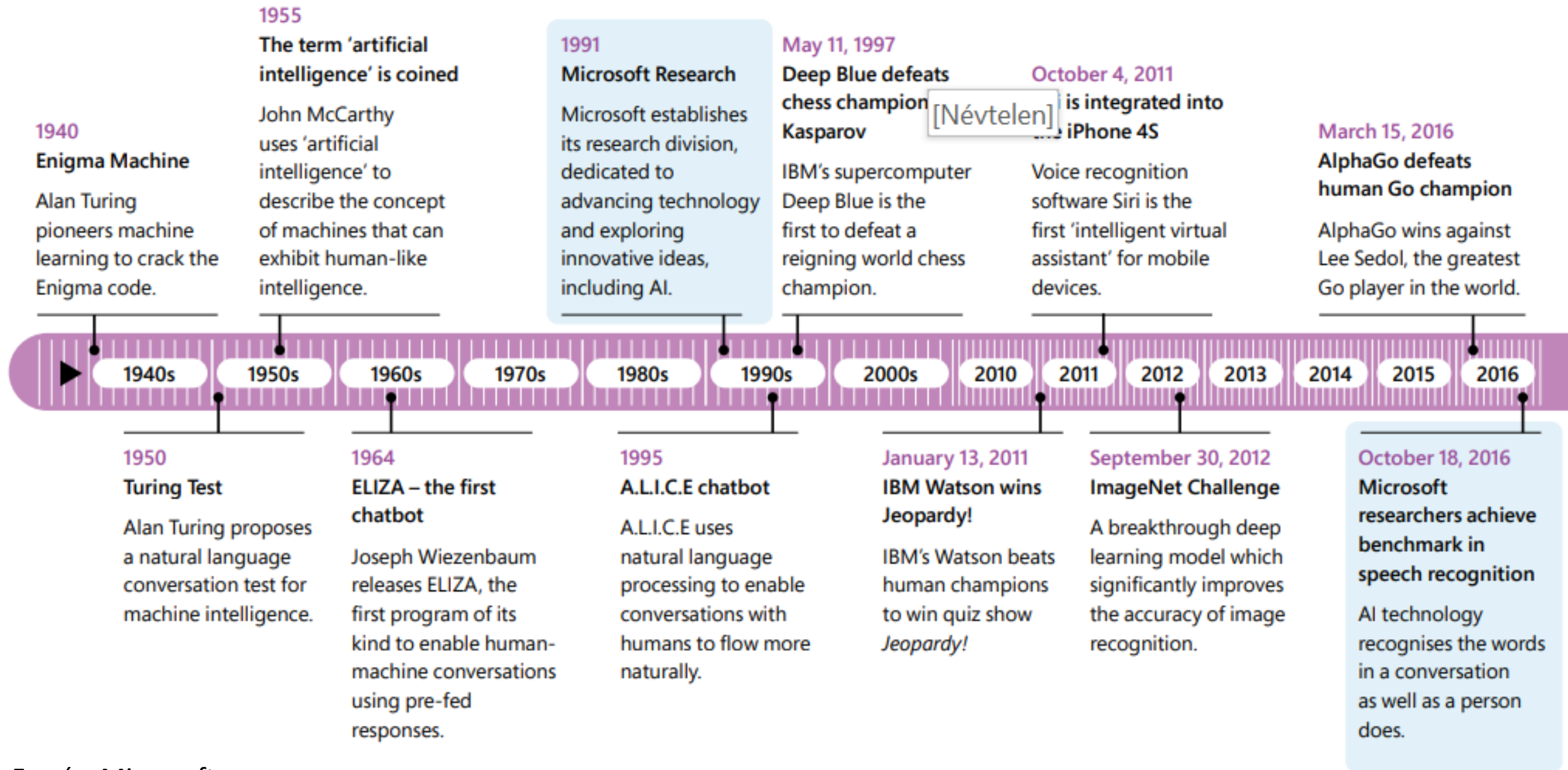
# AI története

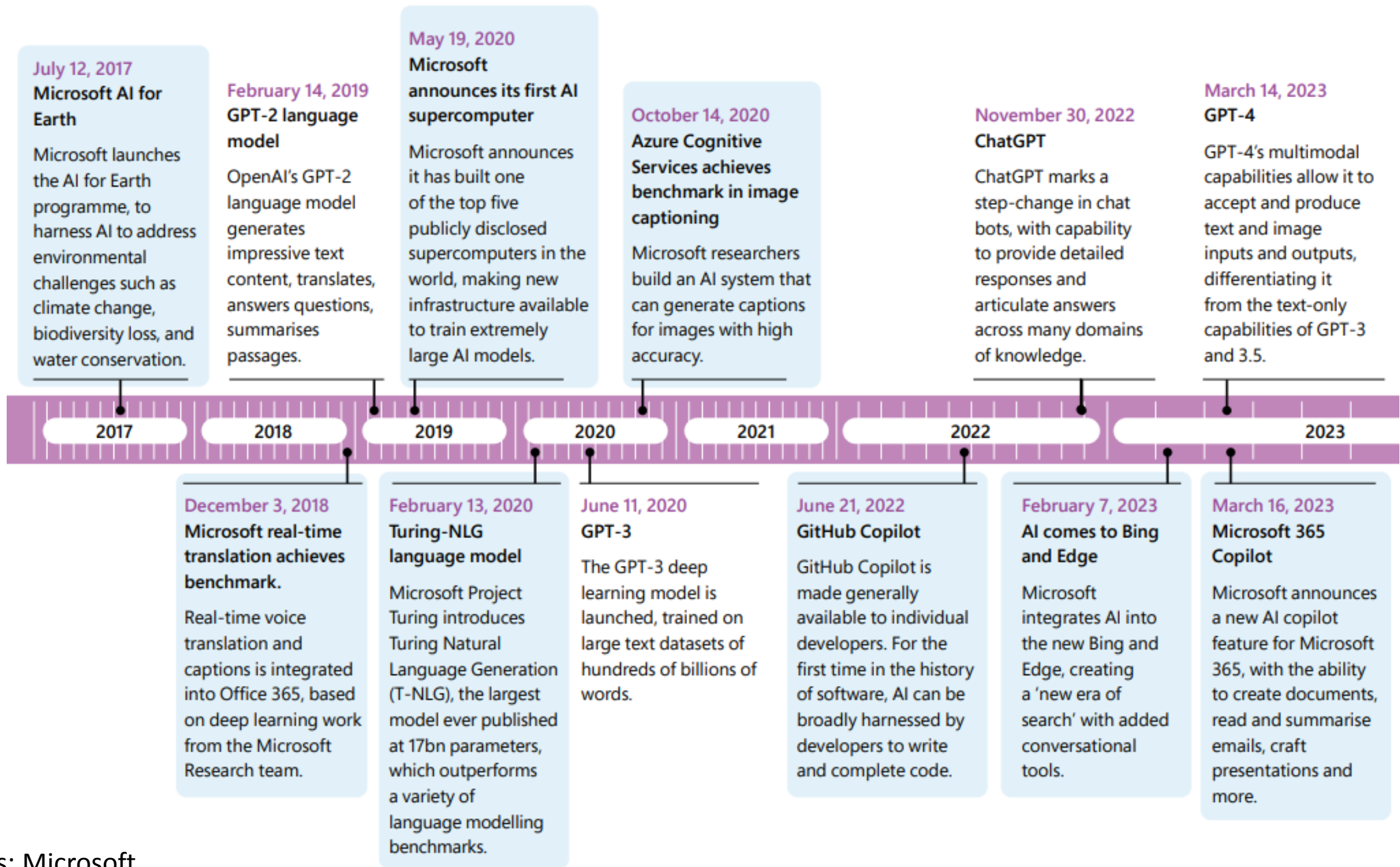
---





# A brief history of AI: From Turing to ChatGPT







# AI történelme vezetőknél 😊

Artificial Intelligence

Machine Learning

Deep Learning

Generative AI

1950s

## Artificial Intelligence

the field of computer science that seeks to create intelligent machines that can replicate or exceed human intelligence.

1959

## Machine Learning

subset of AI that enables machines to learn from existing data and improve upon that data to make decisions or predictions.

2017

## Deep Learning

a machine learning technique in which layers of neural networks are used to process data and make decisions.

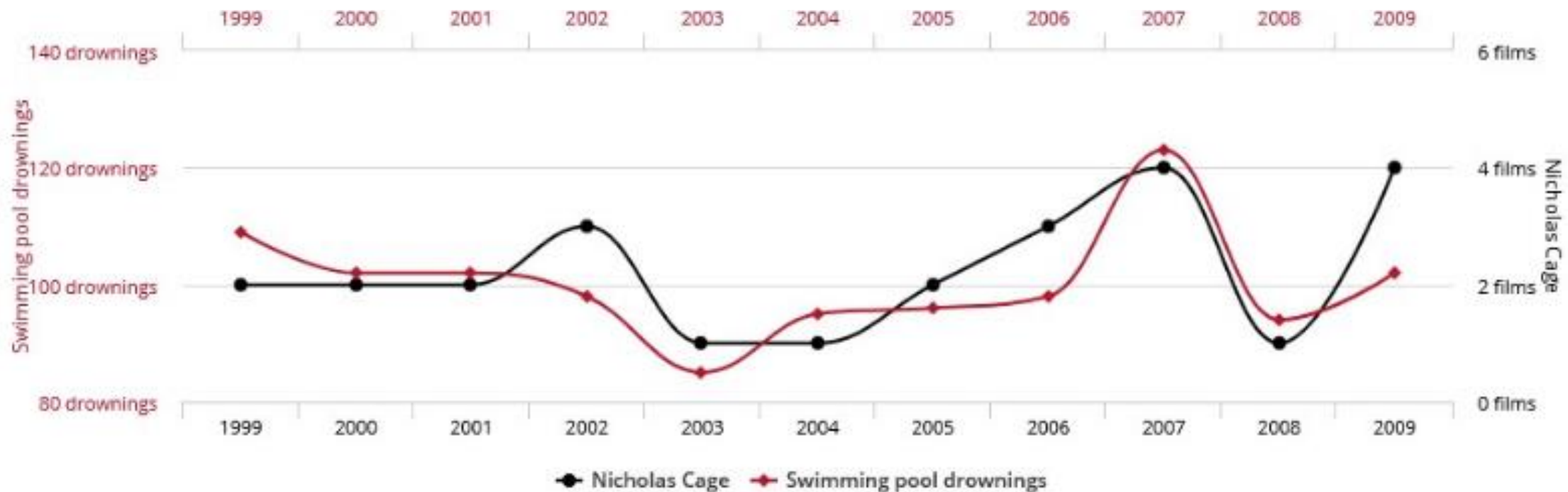
2021

## Generative AI

create new written, visual, and auditory content given prompts or existing data.

# Machine Learning - hackelése

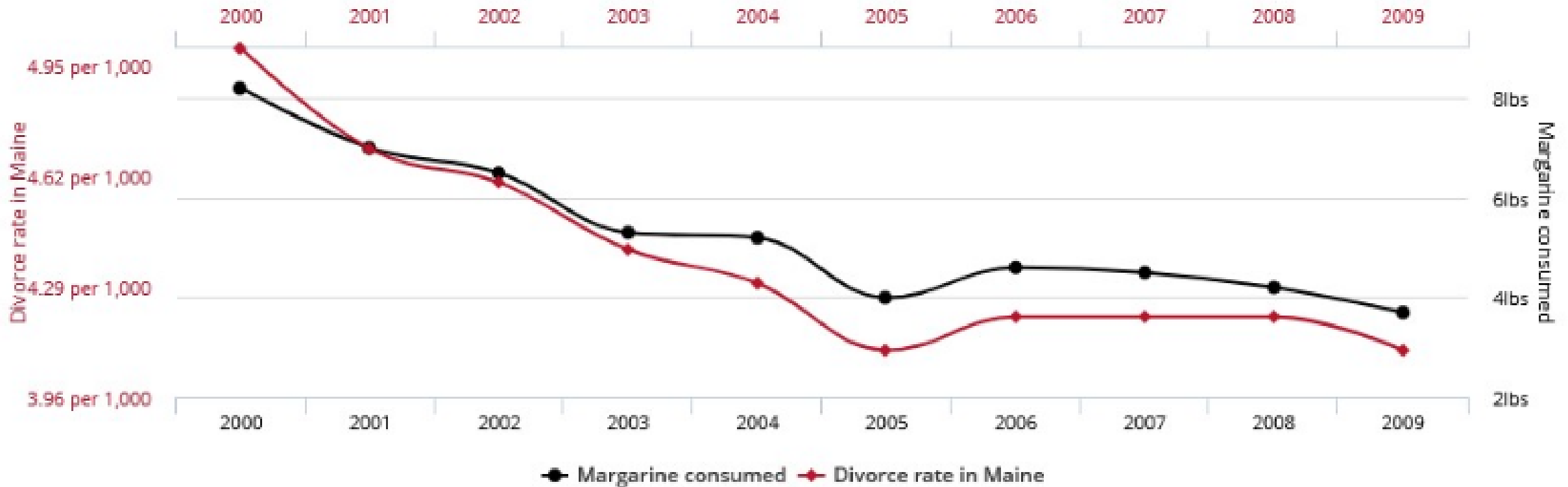
Medencébe fulladó emberek száma és Nicolas Cage filmszerepei közötti korreláció ( $R^2=0,66$ ).



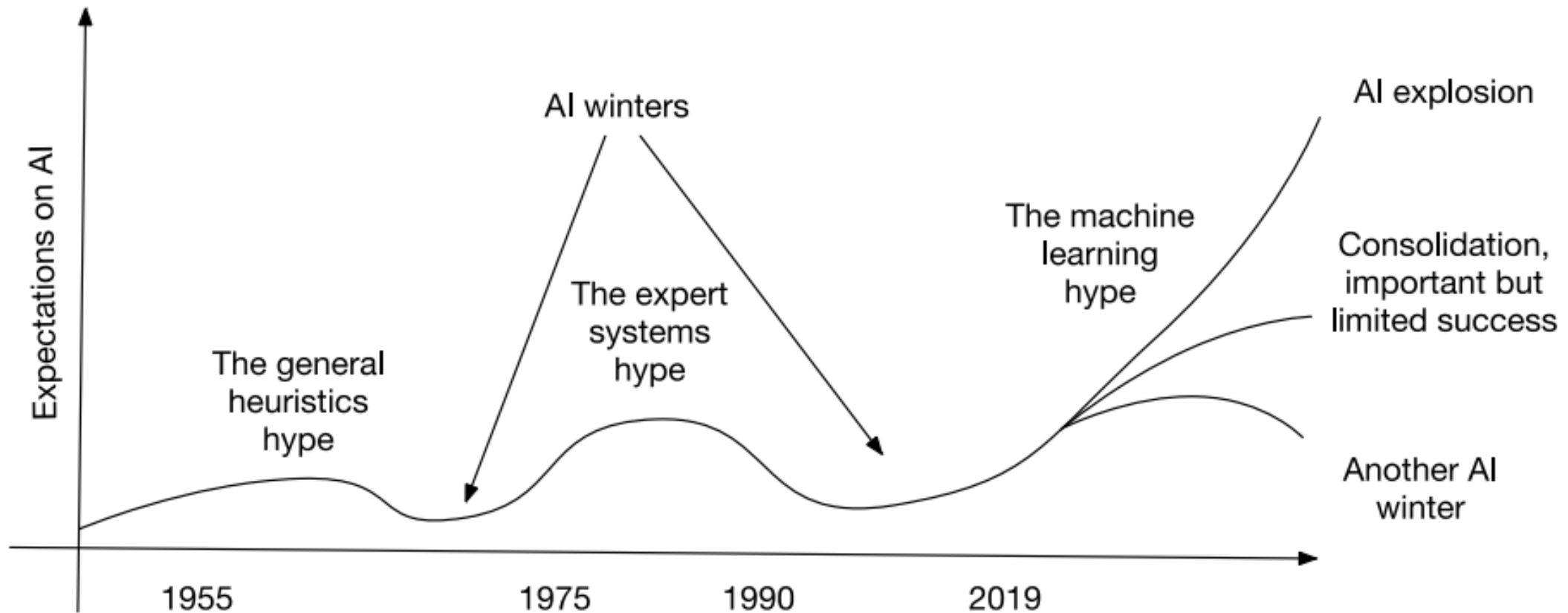


# Machine Learning – hackelése II.

A margarin egy főre jutó fogyasztása és a válások aránya közötti korreláció ( $R=0,992$ ).

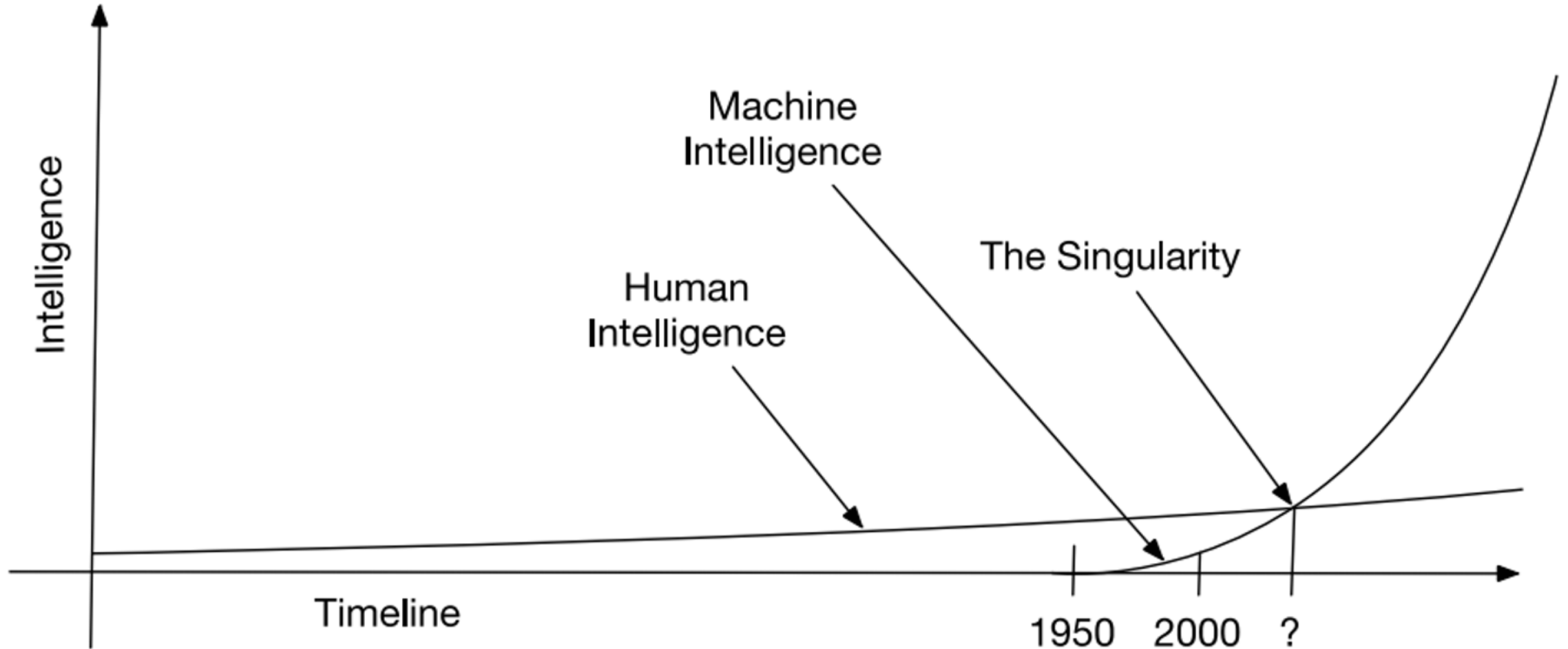


# Hype-ok és AI tél (2014)





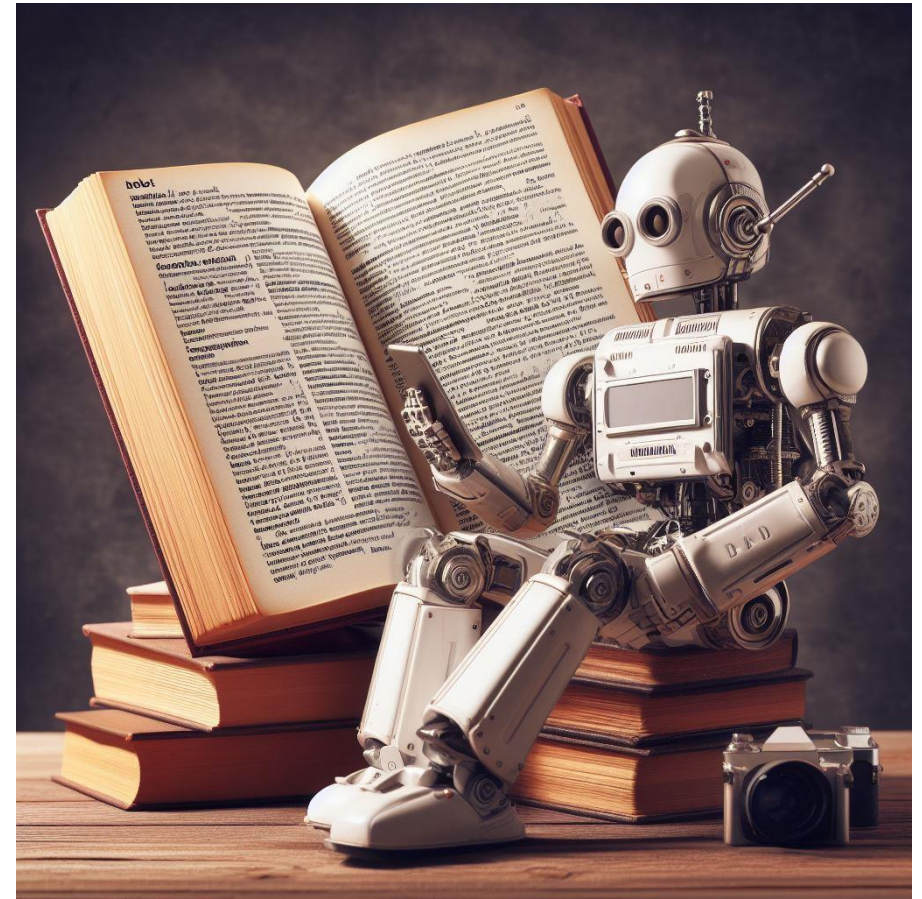
# General AI Szingularitás (2014)



# Mit tud az AI?

---

- **Természetes nyelvi chatbotok:** az ügyfélszolgálattól kezdve a személyre szabott korrepetáláson át az öngyilkosság-megelőzési tanácsadók képzéséig minden.
- **Fejlett keresési képességek:** nem csak linkek, hanem válaszok keresése.
- **Szöveggenerálás:** e-mailek, jogszabálytervezetek, szerződések, forgatókönyvek.
- **Tartalom** (hang, videó, szöveg) valós idejű fordítása, átírása, elemzése és összefoglalása.
- **Képfelismerés és -generálás.** (Smink felismerés, stb.)





# Mit tud az AI? – II.

- **Videó** – elemzés, generálás
- **Deepfake** – Arc csere
- **Hang csere**
- Mesterséges szaglás
- Emóciók felismerése
- CoDi
- ...



# Érdekessegek

- 
- Microsoft – Seeing AI
  - Microsoft - Copilot
  - Microsoft VALL-E, Vall-E /X
  - DarkBERT
  - AlphaFold2
  - Azure IA Health
  - ...








# AI veszélyre felhívó nyilatkozat

---

„A mesterséges intelligencia miatti  
kipusztulás veszélyét az olyan  
emberiségre leselkedő  
fenyegetésekkel egyenértékűen kell  
kezelni, mint a világjárványok vagy  
az atomháború.”





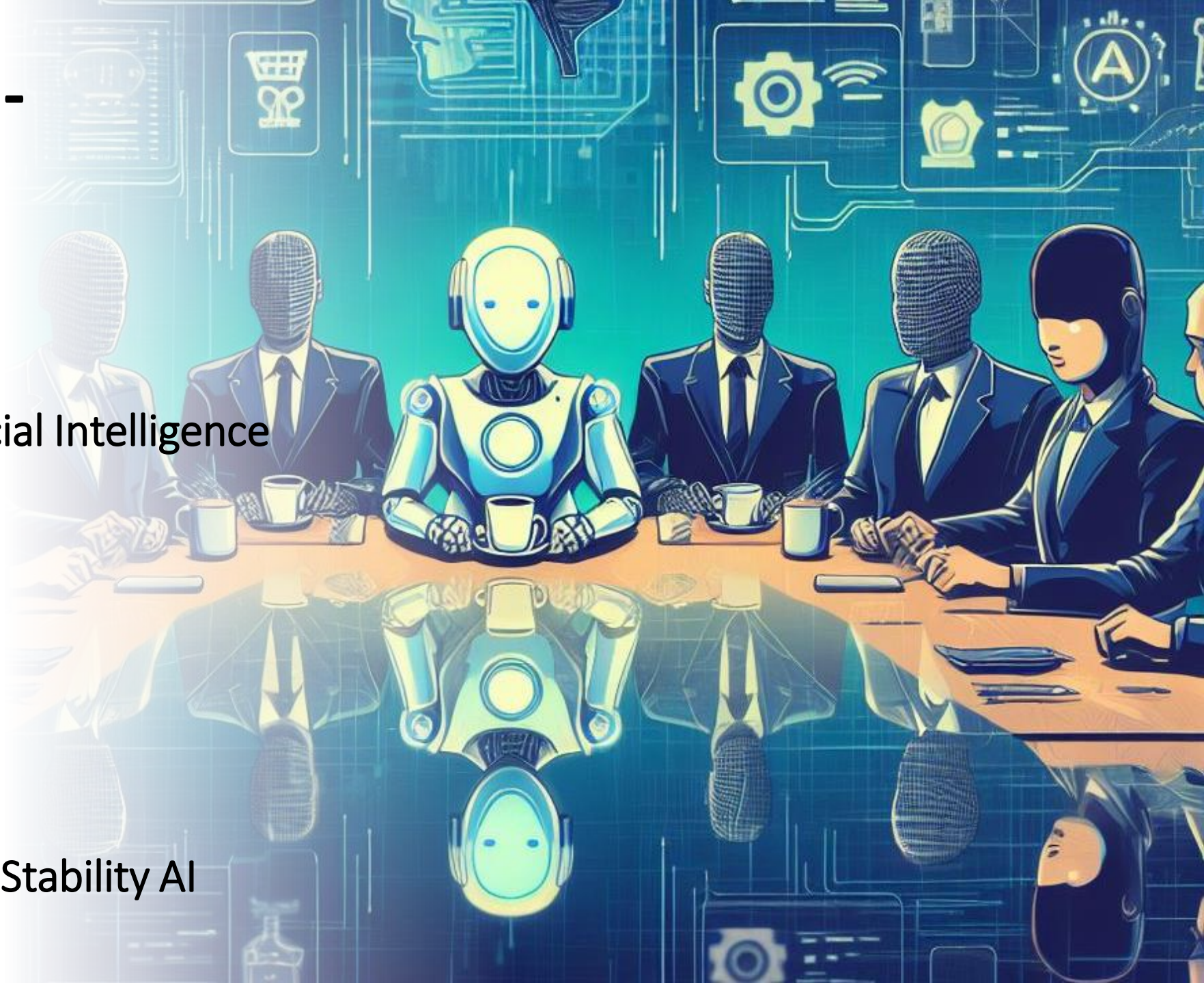
2023. szeptember. 20. 09:49 · Utolsó frissítés: 2023. szeptember. 20. 10:06 · TECH

**Megkezdődhet a chipek beültetése az emberek agyába, megkapta az engedélyt Elon Musk cége, a Neuralink**



# Ki van még itt? - Piaci szereplők

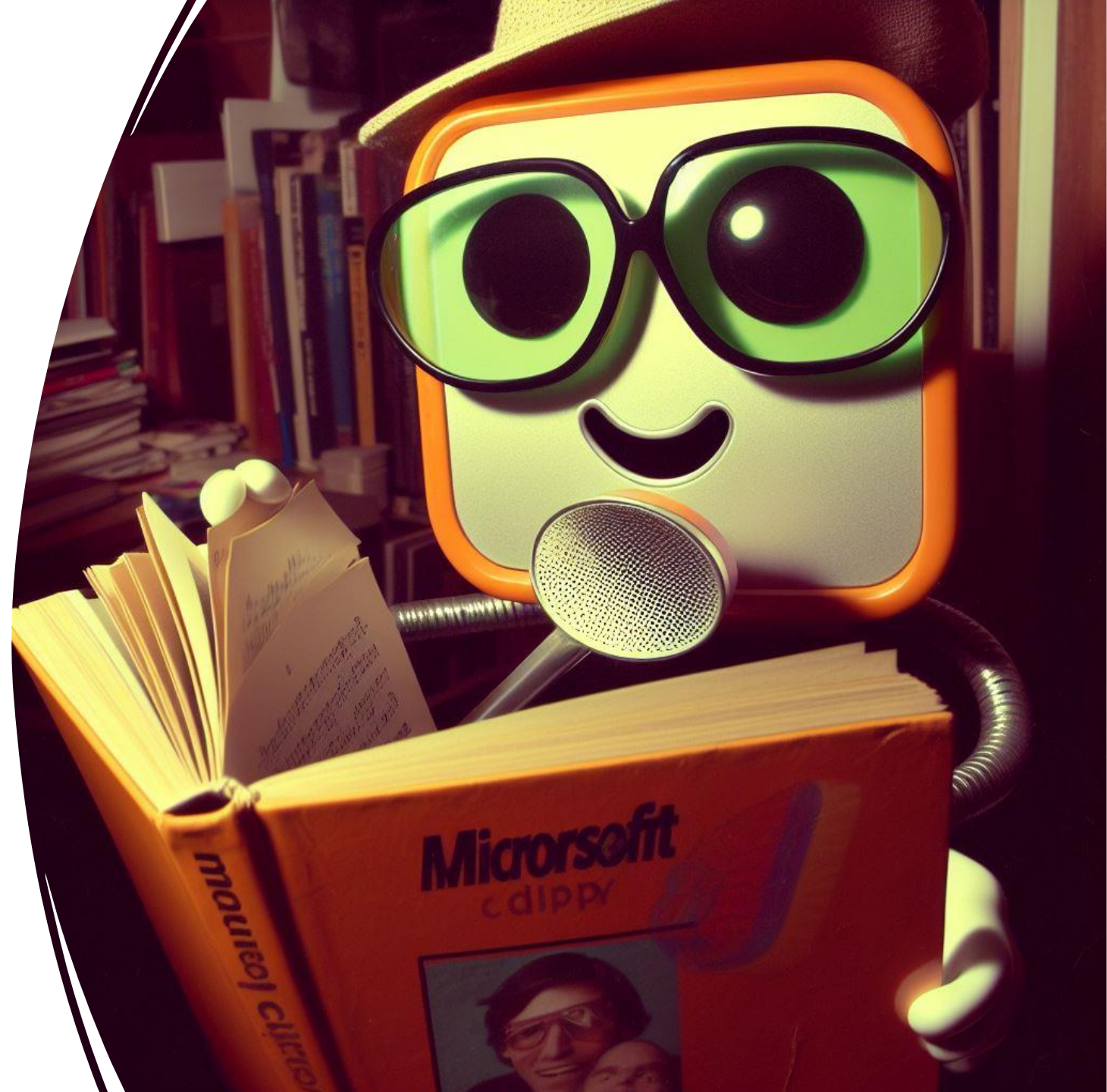
- OpenAI and Microsoft
- DeepMind within Google
- Beijing Academy of Artificial Intelligence
- China's Baidu
- Orosz AI
- Magyar: Puli
- Amazon Bedrock
- ANTHROP\C Claude 2
- IBM
- AI21 Labs, Cohere, Meta, Stability AI
- ...





# Microsoft AI történelem

---





# Key Microsoft AI breakthroughs

- 
- 2016** ● Object & speech recognition  
Human parity
  - 2018** ● Reading comprehension & machine translation  
Human parity
  - General language understanding  
Human parity
  - 2020** ● Turing-NLG language models
  - First AI supercomputer
  - Exclusive license for OpenAI GPT-3 models
  - Image captioning  
Human parity
  - 2021** ● Natural Language Understanding  
Human parity
  - Commonsense Question Answering  
Human parity
  - Azure OpenAI Service preview
  - 2022** ● GitHub Copilot general availability

# Microsoft AI innovation in 2023

## January ●

- Azure OpenAI Service becomes generally available
- Microsoft extends our partnership with OpenAI

## February ●

- Teams Premium with GPT becomes generally available
- Viva Sales adds generative AI capabilities
- Microsoft announces the new Bing and Edge
- Microsoft announces Bing momentum and Skype Copilot
- Windows 11 updates bring AI-powered Bing to the taskbar

## March ●

- LinkedIn introduces collaborative articles
- Microsoft introduces Dynamics 365 Copilot
- Florence comes to Azure Cognitive Services for Vision
- Azure OpenAI Service adds ChatGPT capabilities
- Microsoft announces powerful new virtual machines
- LinkedIn adds new AI-powered capabilities
- Microsoft introduces Microsoft 365 Copilot
- Microsoft introduces Copilot for Power Platform
- Nuance introduces DAX Express
- Azure OpenAI Service adds GPT-4
- Bing Image Creator comes to the new Bing
- GitHub introduces GitHub Copilot X
- Microsoft introduces Microsoft Security Copilot

Made in Iowa



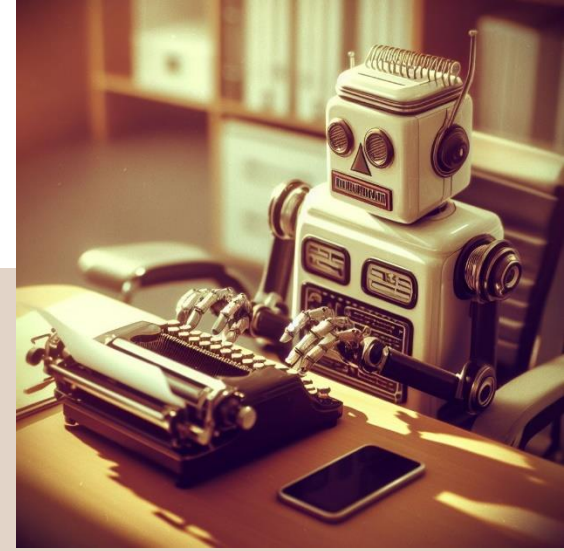
Datacenter – az MS komolyan gondolja... (2020)



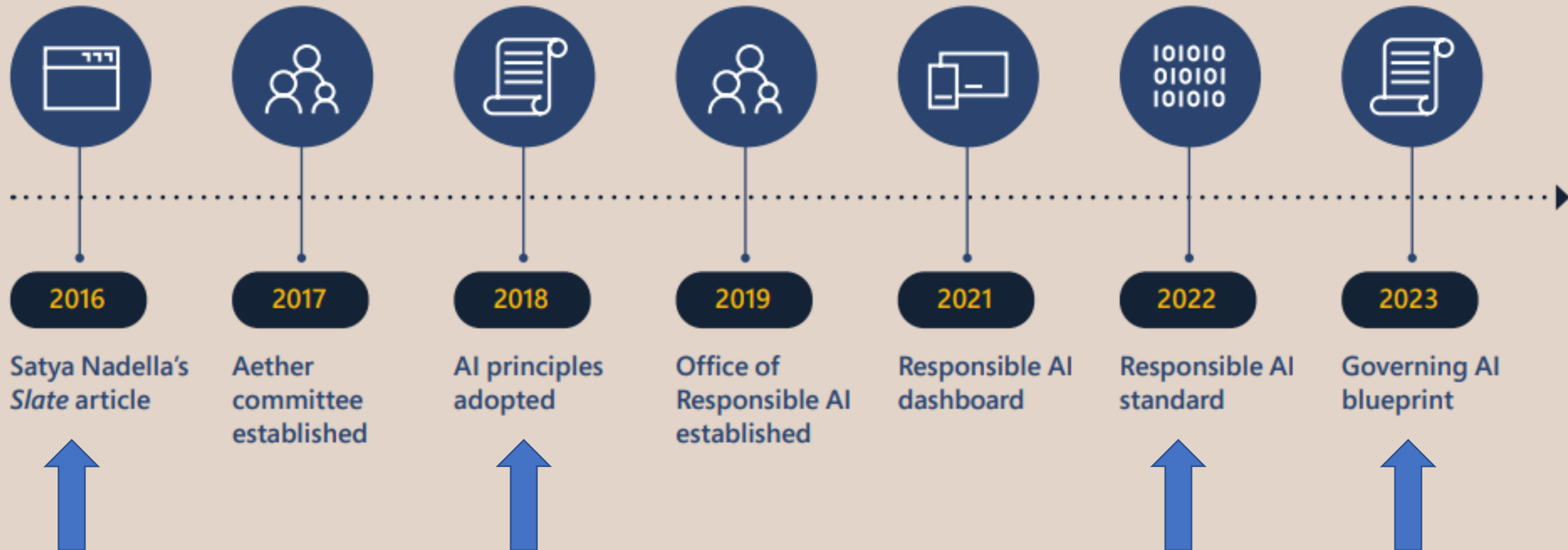
# Microsoft AI alapelvek



# Microsoft AI történelme – másik oldalról nézve



## Our responsible AI journey



# Satya Nadella cikke a Slate portálon

---

future  tense

## **The Partnership of the Future**

Microsoft's CEO explores how humans and A.I. can work together to solve society's greatest challenges.

BY SATYA NADELLA

JUNE 28, 2016 • 2:00 PM



# Microsoft hozzáállása (2018)

## Microsoft's AI Principles



Fairness



Reliability  
& Safety



Privacy &  
Security



Inclusiveness

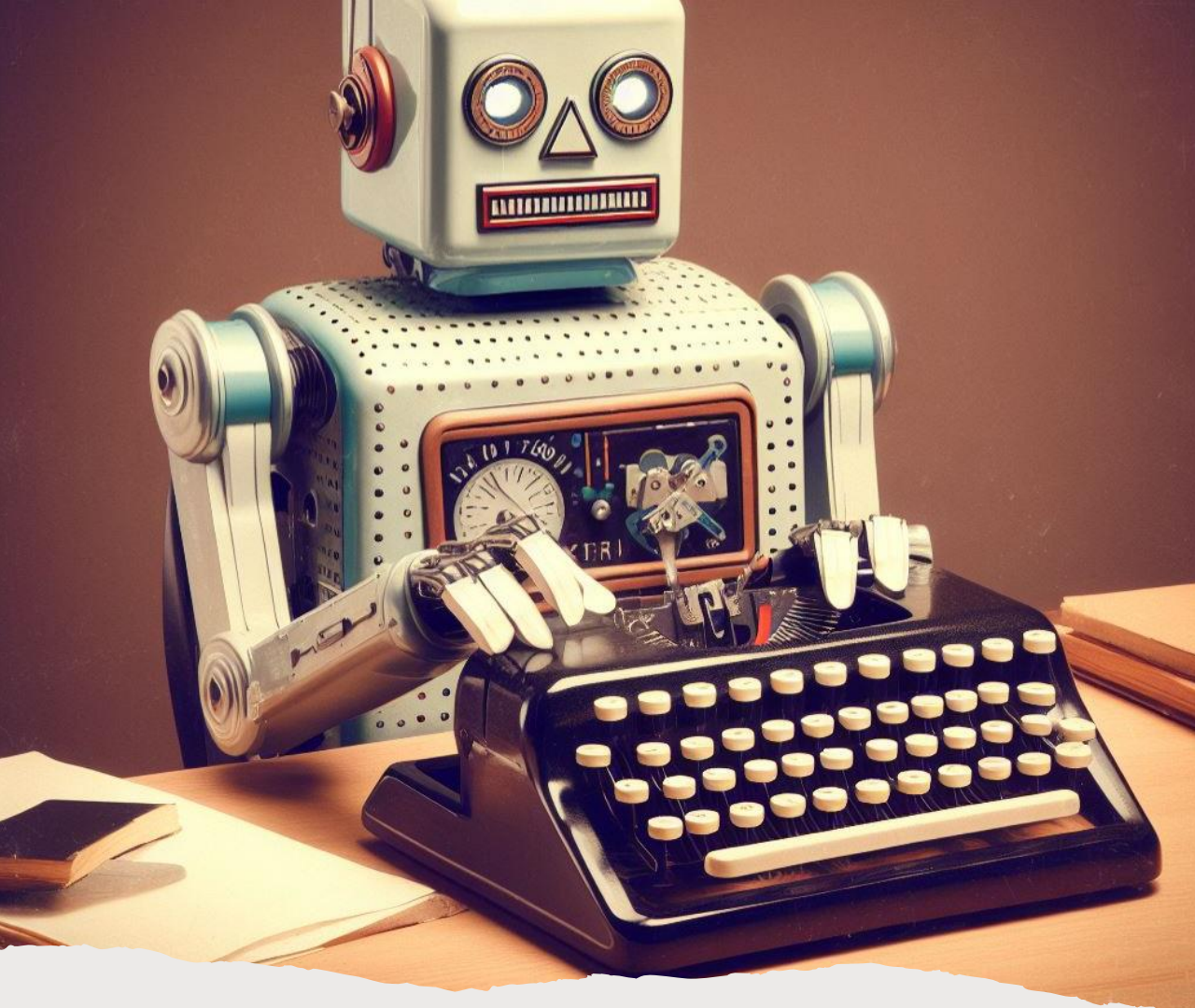


Transparency



Accountability





# Microsoft Responsible AI Standard, v2

## GENERAL REQUIREMENTS

FOR EXTERNAL RELEASE

June 2022

Újabb Microsoft doksi...

# The Standard's Goals at a Glance

## Accountability

- A1: Impact Assessment
- A2: Oversight of significant adverse impacts
- A3: Fit for purpose
- A4: Data governance and management
- A5: Human oversight and control

## Transparency

- T1: System intelligibility for decision making
- T2: Communication to stakeholders
- T3: Disclosure of AI interaction

## Fairness

- F1: Quality of service
- F2: Allocation of resources and opportunities
- F3: Minimization of stereotyping, demeaning, and erasing outputs

## Reliability & Safety

- RS1: Reliability and safety guidance
- RS2: Failures and remediations
- RS3: Ongoing monitoring, feedback, and evaluation

## Privacy & Security

- PS1: Privacy Standard compliance
- PS2: Security Policy compliance

## Inclusiveness

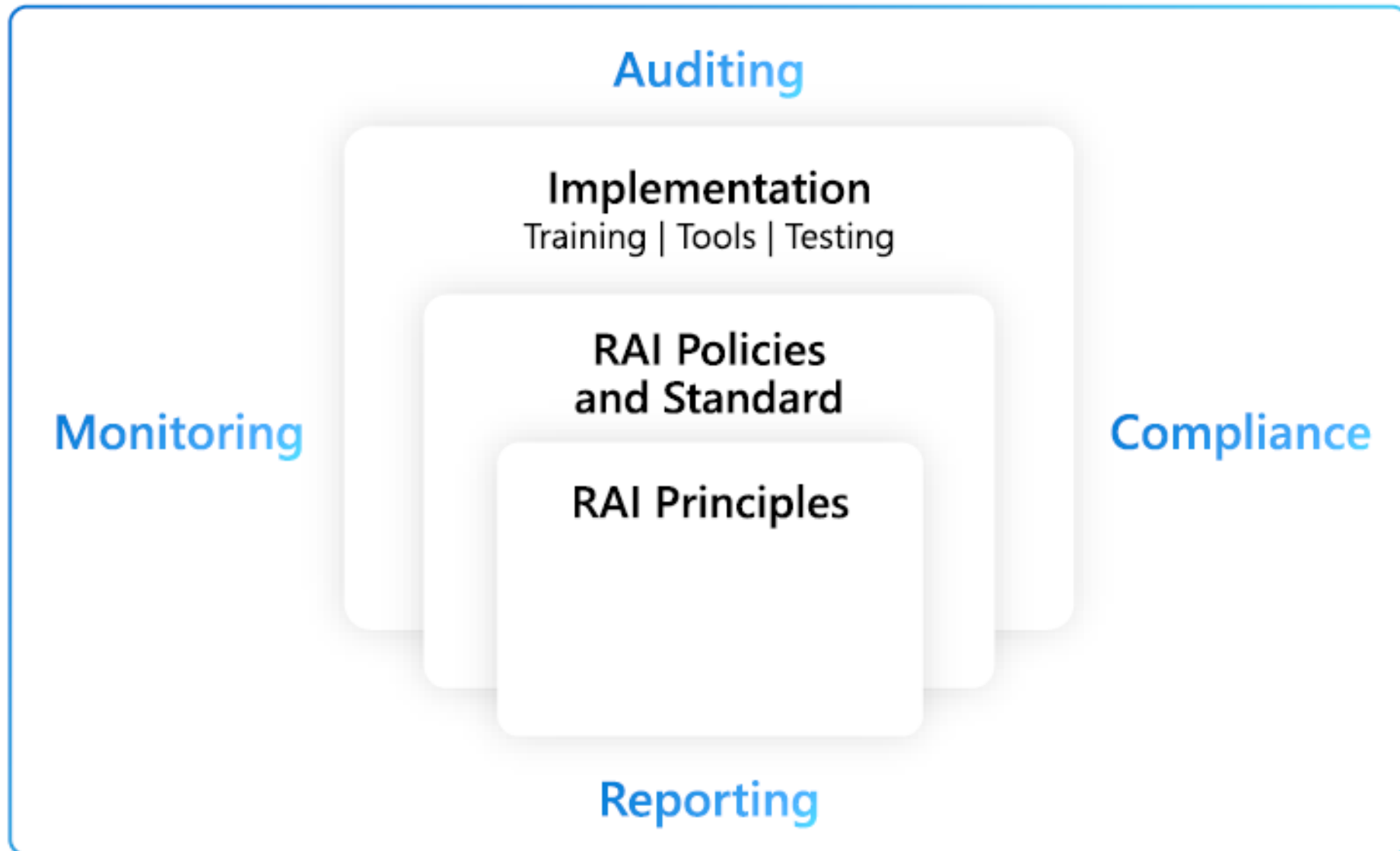
- I1: Accessibility Standards compliance





# Governing AI: A Blueprint for the Future

# Responsible AI Governance Framework



# A mesterséges intelligencia irányítása: tervezet a jövőre nézve

- A Microsoft a Fehér Ház legutóbbi ülésére válaszul kötelezettséget vállal arra, hogy **bevezetjük a NIST AI Risk Management Framework** (a NIST mesterséges intelligencia kockázatkezelési keretrendszerét).
- Kötelezettséget vállalunk arra, hogy a Microsoft meglévő **AI-tesztelési munkáját új lépésekkel egészítjük ki**, hogy tovább erősítsük a magas kockázatú AI-rendszerekkel kapcsolatos mérnöki gyakorlatainkat.
- A kormányzat felgyorsíthatja a lendületet egy olyan végrehajtási rendelet révén, amely előírja, hogy az amerikai kormány számára kritikus **AI-rendszerek szállítói kötelesek saját maguk tanúsítani, hogy végrehajtják a NIST AI kockázatkezelési keretrendszerét.**
- Elkötelezettek vagyunk amellett, hogy együttműködjünk más iparági vezetőkkel és **a kormányzattal együtt dolgozzunk ki új szabványokat** az alapmodellekkel kapcsolatban.



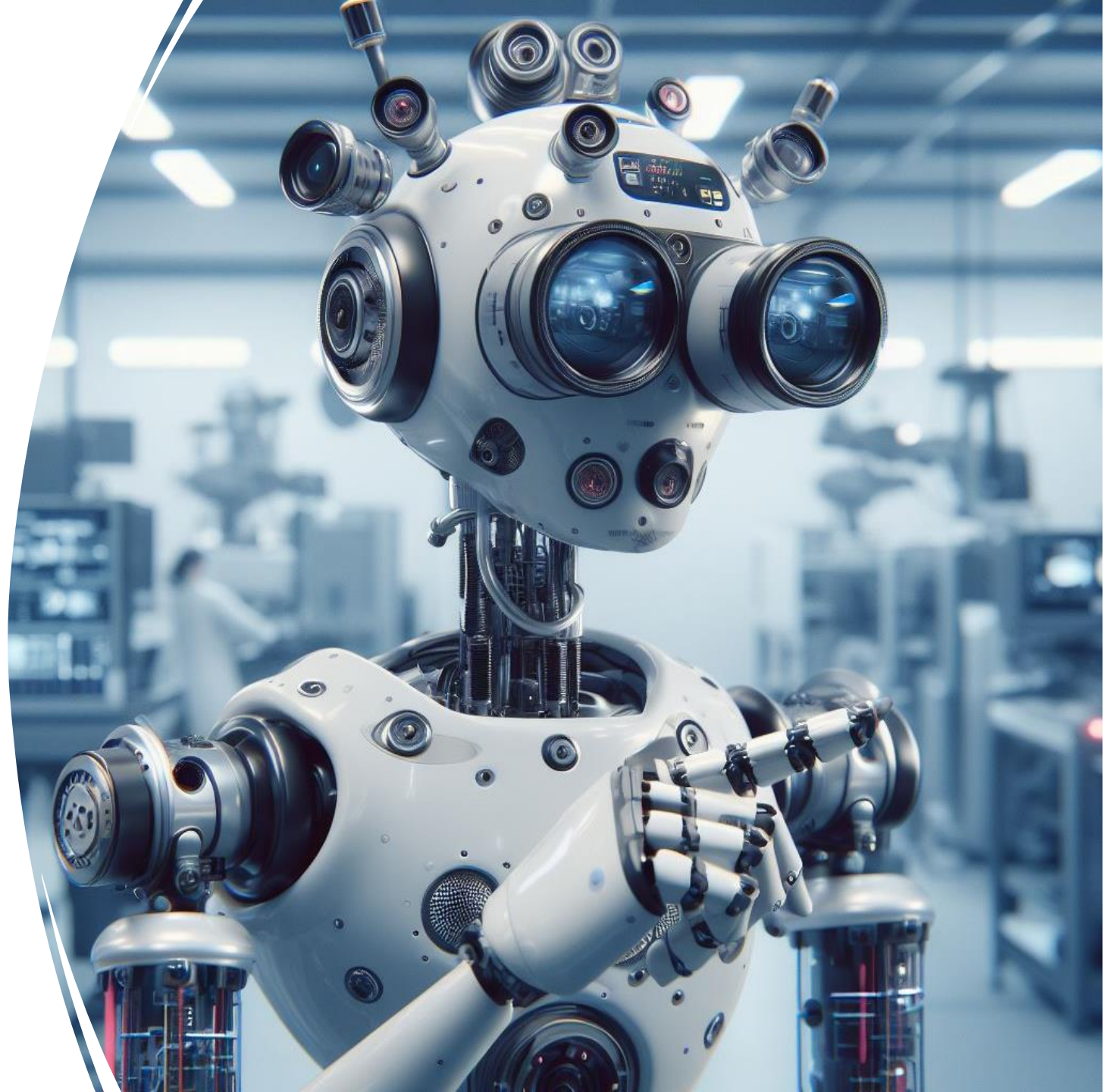
Jogászokdjunk  
kicsit...



# Adatgyűjtés

---

- **USA** (opt-out – visszavonás) – CCPA (California Consumer Privacy Act)
  - Van már CCPA 2.0 némi GDPR felhangokkal 😊
- **EU** (opt-in – hozzájárulás) - GDPR







# Jogászolimpia

2000 – 2015

Safe Harbour (első megfeleléségi határozat)  
2016 – GDPR

2016 – 2020 Privacy Shield (második megfeleléségi  
nyilatkozat)

2022 - Trans-Atlantic Data Privacy Framework  
(EU-USA adattovábbítási keretmegállapodás)

2023 - Harmadik megfeleléségi nyilatkozat





# Nyuszika listája

„Engedélyezett” cégek aktuális listája – avagy a business nem állhat meg...



**U.S. Businesses**



**European  
Businesses**



**European  
Individuals**



**Data Protection  
Authorities**

## Microsoft Corporation

Redmond, Washington

Active

### Framework

EU-U.S. Data Privacy Framework

Swiss-U.S. Data Privacy Framework

UK Extension to the EU-U.S. Data Privacy Framework

### Covered Data ⓘ

HR

Non-HR

Data Protection Authorities

News & Events

Contact

Privacy Program



# Cloud Act



# Cloud Act (Clarifying Lawful Overseas Use of Data Act)

---

- USA hatóságok adatokhoz férhetnek hozzá.
- Microsoft ígérete, hogy tájékoztatnak, visszautasítanak, stb., ha ez „törvényileg” lehetséges...
- A Cloud Act megtiltja, hogy a szolgáltató az adatigénylésről tájékoztassa az érintettet.





# Microsoft adatvédelmi? nyilatkozata



## Az általunk gyűjtött személyes adatok

- **A Microsoft az Önnel folytatott kommunikáció, valamint a termékeink használata során gyűjt adatokat Öntől.** Bizonyos adatokat Ön ad meg közvetlenül, bizonyos információkat pedig a termékeinkkel való interakciójáról, azok használatáról és az általuk nyújtott élményekről gyűjtött adatokból szerzünk. Az általunk gyűjtött adatok köre attól függ, hogy milyen kontextusban lép kapcsolatba a Microsofttal, milyen döntéseket hoz (például milyen adatvédelmi beállításokat alkalmaz), valamint hogy milyen termékeket és szolgáltatásokat vesz igénybe. **Harmadik felektől is szerzünk be adatokat Önről.**

# Microsoft adatvédelmi? nyilatkozata



## Hogyan használjuk a személyes adatokat

...

Az adatokat vállalatunk működtetésére is felhasználjuk, például teljesítményünk elemzésére, **jogi kötelezettségeink teljesítésére**, munkatársaink fejlesztésére, valamint kutatások végrehajtására.

...

**A személyes adatok e célokra történő feldolgozása automatizált és manuális (emberek által végzett) feldolgozási módszereket is tartalmaz.** Az automatizált módszereink gyakran kapcsolatban állnak a manuális módszereinkkel, és azokra építkeznek.

...

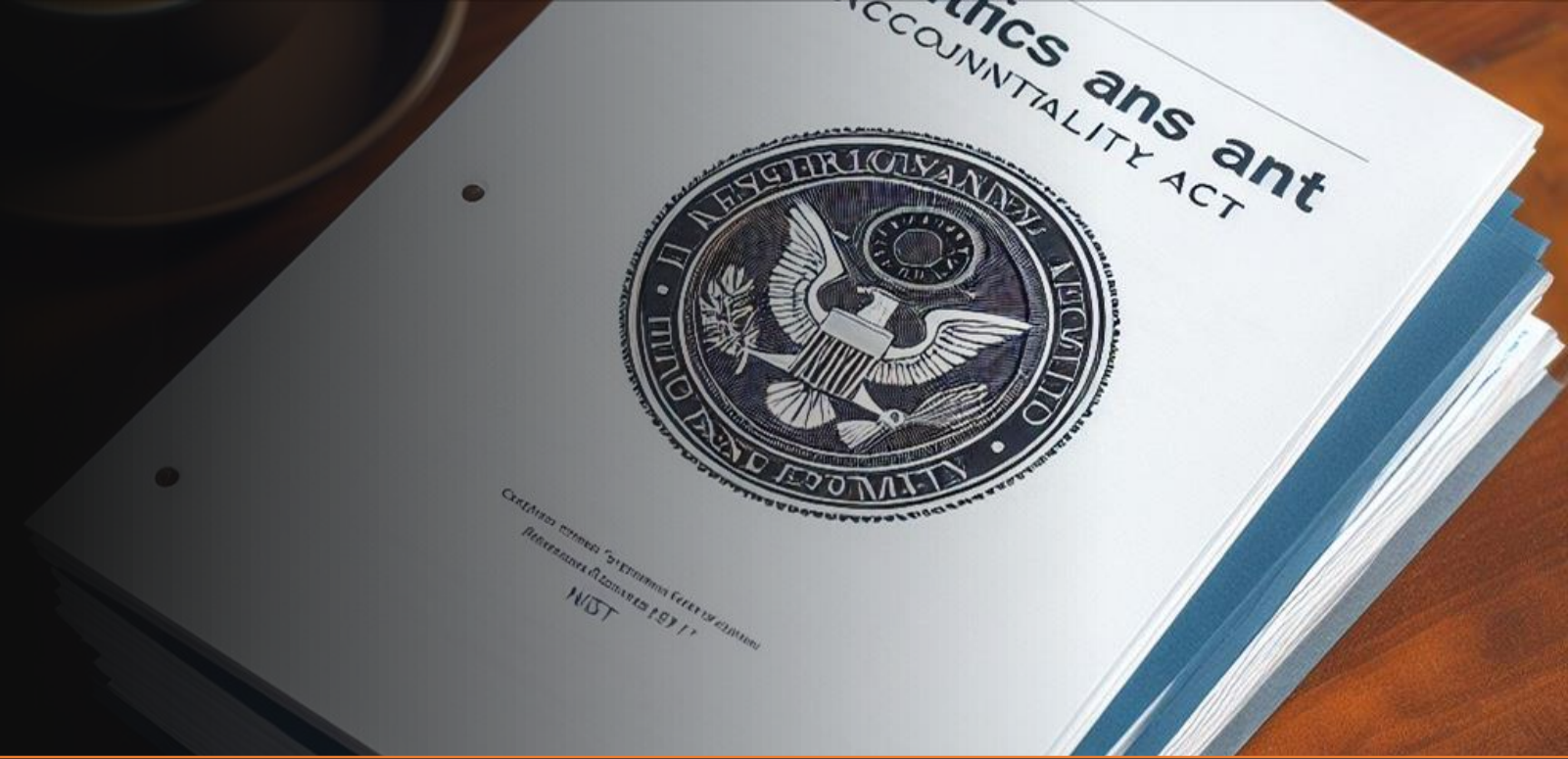
**manuálisan beletekintünk hangadatokból származó minták rövid részleteibe (amelyeknél lépéseket tettünk az azonosításra való alkalmatlanná tétel érdekében). Ezt a kézi áttekintést a Microsoft alkalmazottai és a Microsoft nevében dolgozó értékesítők is elvégezhetik.**

<https://privacy.microsoft.com/hu-hu/privacystatement>

NIST AI 100-1

---

# Artificial Intelligence Risk Management Framework (AI RMF)





# National Institute for Standards in Technology (NIST) mesterséges intelligencia-kockázati keretrendszer

- 2021-ben a Kongresszus arra utasította a NIST-t, hogy dolgozzon ki egy önkéntes keretet a megbízható mesterséges intelligencia számára
- 2023: Artificial Intelligence Risk Management Framework (AI RMF 1.0)

**Célja:** olyan iránymutatások és bevált gyakorlatok gyűjteménye, amelyek célja, hogy segítse a szervezeteket a mesterséges intelligencia technológiák bevezetésével és használatával kapcsolatos kockázatok azonosításában, értékelésében és kezelésében



# NIST AI RMF elemei

---

- **Érvényesség és megbízhatóság** – nincs „hallucináció”
- **Biztonságosság** – nem osztja meg másokkal az adatokat
- **Ellenállóság és biztonság** – az AI rendszer védett és biztonságos a támadásokkal szemben
- **Átláthatóság** – AI működésének a megértése
- **Adatvédelem** – az információk védettek és anonimizáltak legyenek
- **Tisztességesség** – káros előítéletek kezelése



# The Artificial Intelligence Act





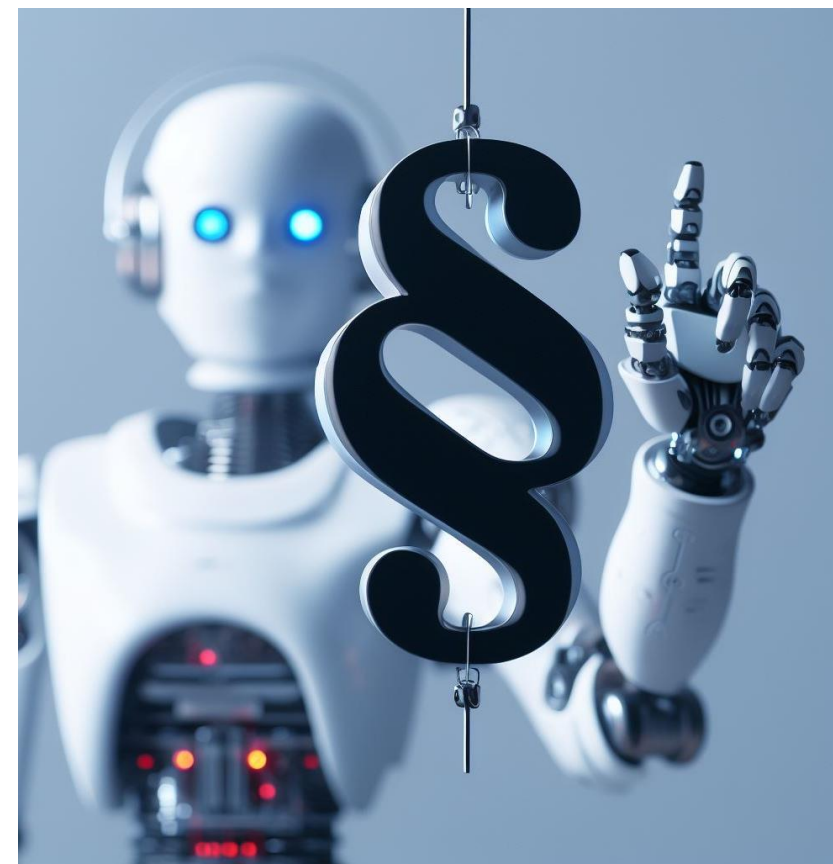
# AI Act

---

Etikus mesterséges intelligencia megközelítés

Jogi beavatkozás különböző kockázati szintek alapján:

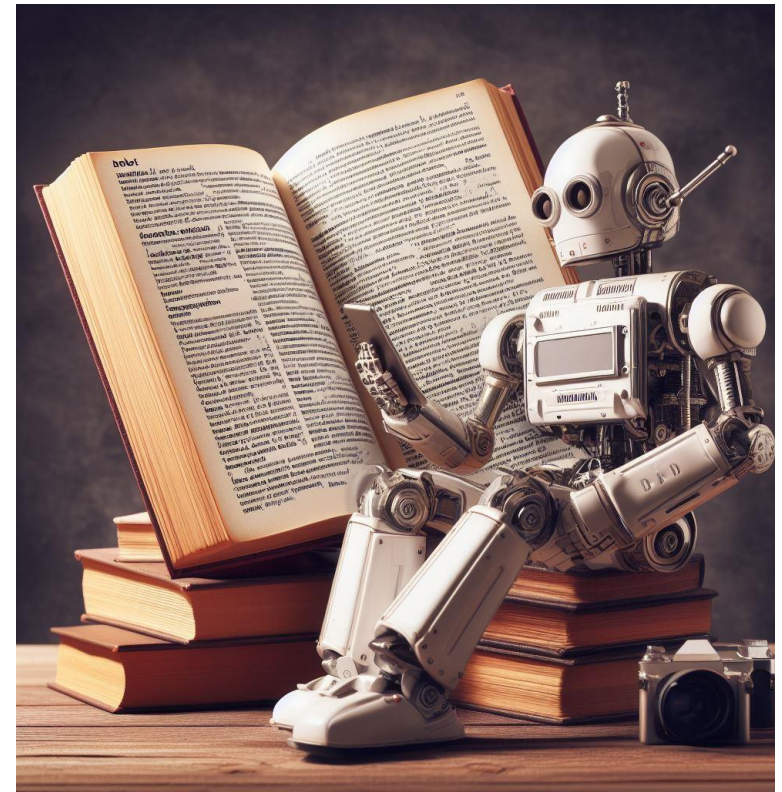
- elfogadhatatlan (tiltott) kockázat
- magas kockázat
- mérsékelt kockázat
- alacsony vagy minimális kockázat



# AI Act – főbb pontok

---

- Műszaki robusztusság és biztonság
- Adatvédelem és adatkezelés
- „szabályozási homokozó” (“regulatory sandbox”)
- Emberi ügynökség és felügyelet
- Átláthatóság
- Sokszínűség, megkülönböztetés mentesség és méltányosság
- Elszámoltathatóság, társadalmi és környezeti jólét



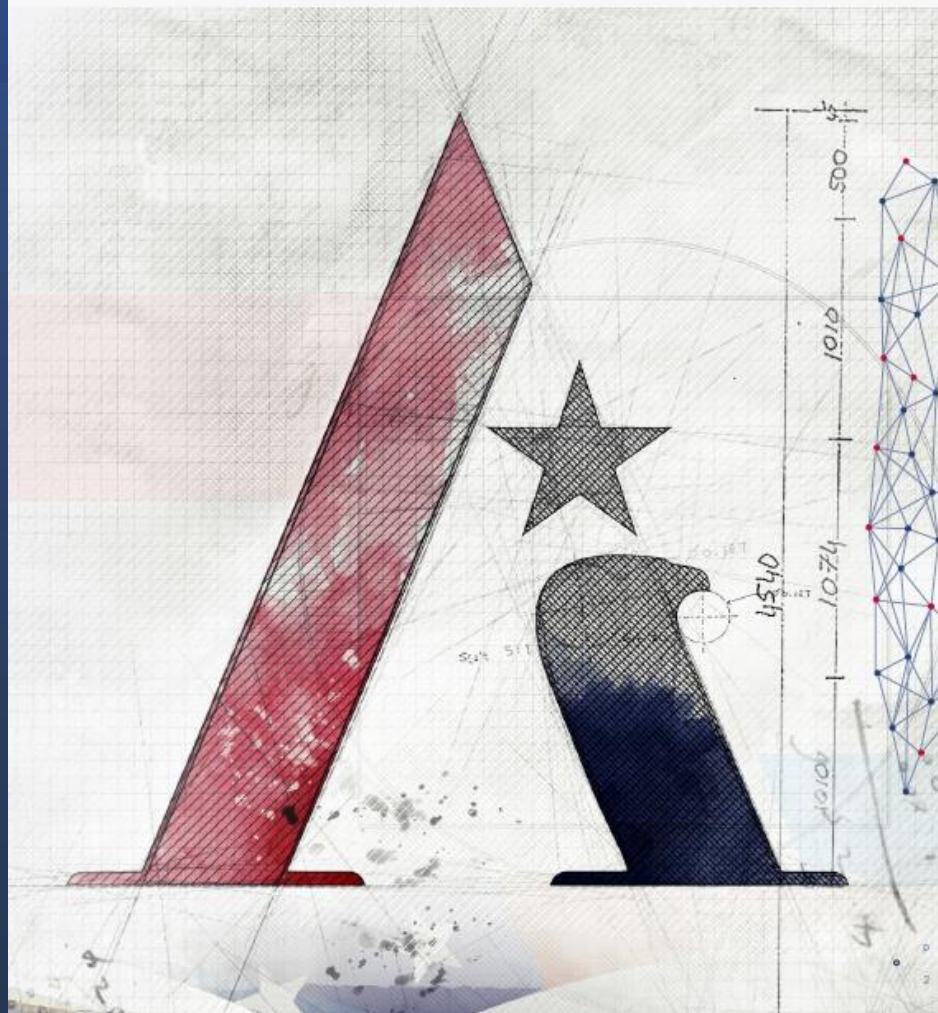
Süllyedjünk  
lejjebb...





# Final Report

National Security Commission on Artificial Intelligence



|  |           |
|--|-----------|
| <b>PART I: DEFENDING AMERICA IN THE AI ERA</b>   | <b>41</b> |
| Chapter 1: Emerging Threats in the AI Era  | 43        |
| Chapter 2: Foundations of Future Defense   | 59        |
|  Chapter 3: AI and Warfare  | 75        |
|  Chapter 4: Autonomous Weapon Systems and Risks<br>Associated with AI-Enabled Warfare | 89        |
|  Chapter 5: AI and the Future of National Intelligence                                | 107       |
| Chapter 6: Technical Talent in Government  | 119       |
| Chapter 7: Establishing Justified Confidence in AI Systems   | 131       |
| Chapter 8: Upholding Democratic Values: Privacy, Civil Liberties,<br>and Civil Rights in Uses of AI for National Security  | 141       |

|  |            |
|--|------------|
| <b>PART II: WINNING THE TECHNOLOGY COMPETITION</b>     | <b>155</b> |
| Chapter 9: A Strategy for Competition and Cooperation  | 157        |
| Chapter 10: The Talent Competition                     | 171        |
| Chapter 11: Accelerating AI Innovation                 | 183        |
| Chapter 12: Intellectual Property                      | 199        |
| Chapter 13: Microelectronics                           | 211        |
| Chapter 14: Technology Protection                      | 223        |
| Chapter 15: A Favorable International Technology Order | 241        |
| Chapter 16: Associated Technologies                    | 253        |
| <br>   |            |
| Blueprints for Action                                  | 271        |
| Appendices   | 599        |



# AI és hadviselés

"A minisztériumnak cselekednie kell hogy a mesterséges intelligenciát integrálja a kritikus funkciókba, a meglévő rendszerekbe, gyakorlatokba és hadgyakorlatokba, hogy 2025-re mesterséges intelligenciára kész haderővé váljunk."

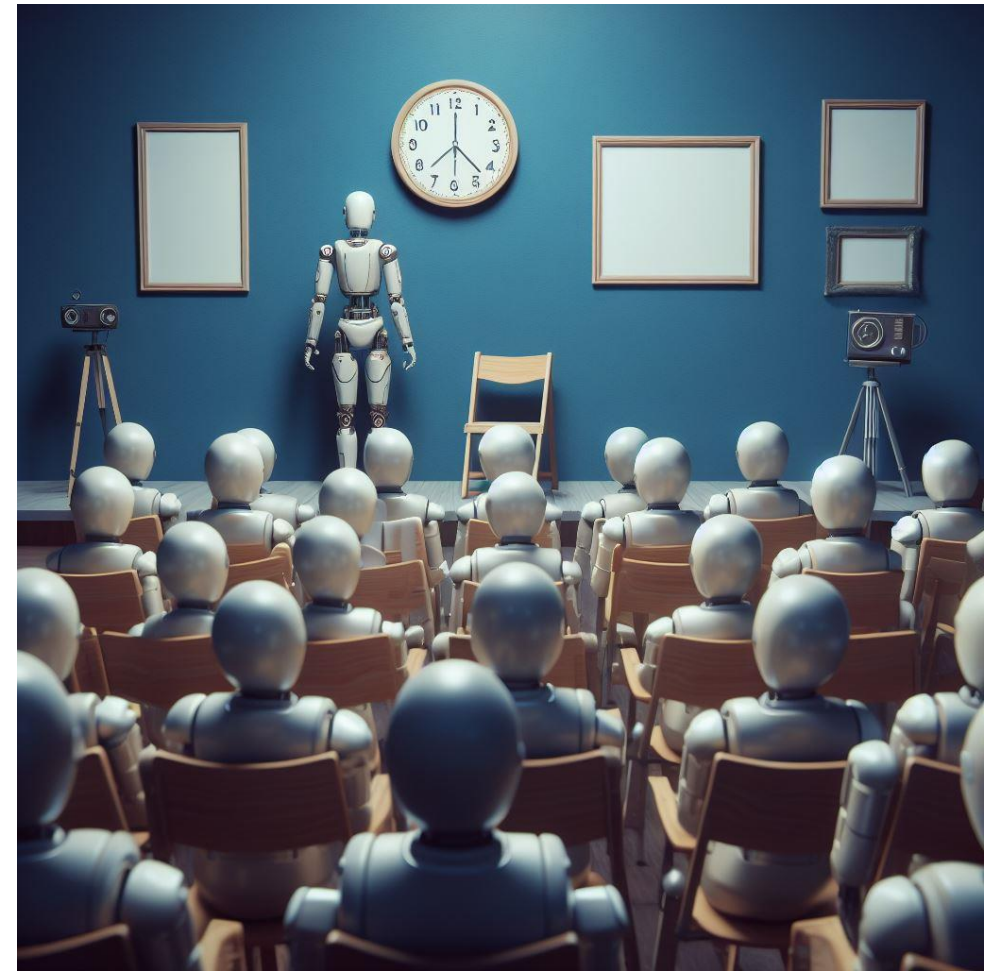


**"The Department must act now to integrate AI into critical functions, existing systems, exercises and wargames to become AI-ready force by 2025."**

| Chapter | Recommendation  | Cabinet<br>Departments,<br>Major Agencies,<br>and Program<br>Offices | Amount        |
|---------|---|--|---------------|
| 1       | Establish a dedicated AI Fund.  | Department of Defense:<br>USD(R&E)                                   | \$200 million |
| 2       | Increase investments in AI R&D.   | Department of Defense  | \$8 billion   |
| 3       | Establish a fund to to accelerate procurement and integration of commercial | Department of Defense:<br>Joint Artificial                           | \$100 million |

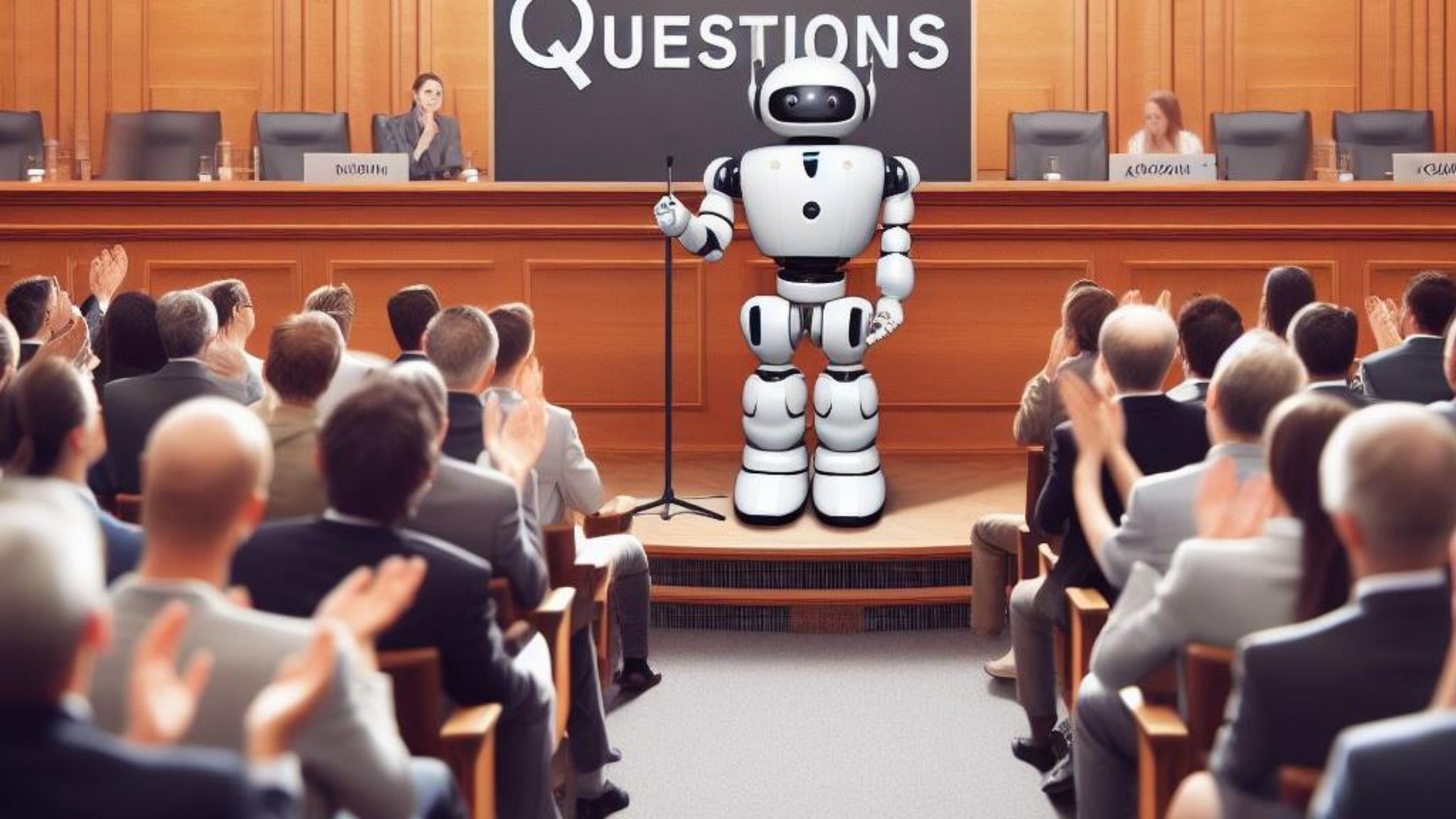
# Ami nem fért bele... ☹️

- Rejtett hangutasítások
- AI folyamatainak a biztonsága/veszélyei (trénelő adatkészlet, trénelés folyamata, stb.)
- AI modell lopás
- Adat kinyerési technikák pl. Prompt injection
- Hogyan legyen Támadó ChatGPT-nk? Technikák, módszerek
- Hogyan lehet védekező a ChatGPT?
- Aggályok a ChatGPT-vel szemben





# QUESTIONS



# És a linkek...

---

Ezzel készültek a képek:

<https://www.bing.com/create>

Microsoft Seeing AI - gyengénlátóknak alkalmazás:

<https://www.microsoft.com/en-us/ai/seeing-ai>





# És a linkek...

---

ChatGPT és Microsoft 365 Copilot: Mi a különbség?

<https://support.microsoft.com/hu-hu/topic/chatgpt-%C3%A9s-microsoft-365-copilot-mi-a-k%C3%BCI%C3%B6nbs%C3%A9g-8fdec864-72b1-46e1-afcb-8c12280d712f>

Microsoft CoDi

<https://www.microsoft.com/en-us/research/blog/breaking-cross-modal-boundaries-in-multimodal-ai-introducing-codi-composable-diffusion-for-any-to-any-generation/>

Microsoft Copilot - ismertetés:

<https://adoption.microsoft.com/en-us/copilot/>

VALL-E (X) - leírás és példák (hangminta heckelés):

<https://www.microsoft.com/en-us/research/project/vall-e-x/>





# És a linkek...

Google SGE - saját Generative AI:

<https://labs.google/sge/>

---

AlphaFold2:

<https://en.wikipedia.org/wiki/AlphaFold>

Egészség:

<https://www.microsoft.com/en-us/ai/ai-for-health>

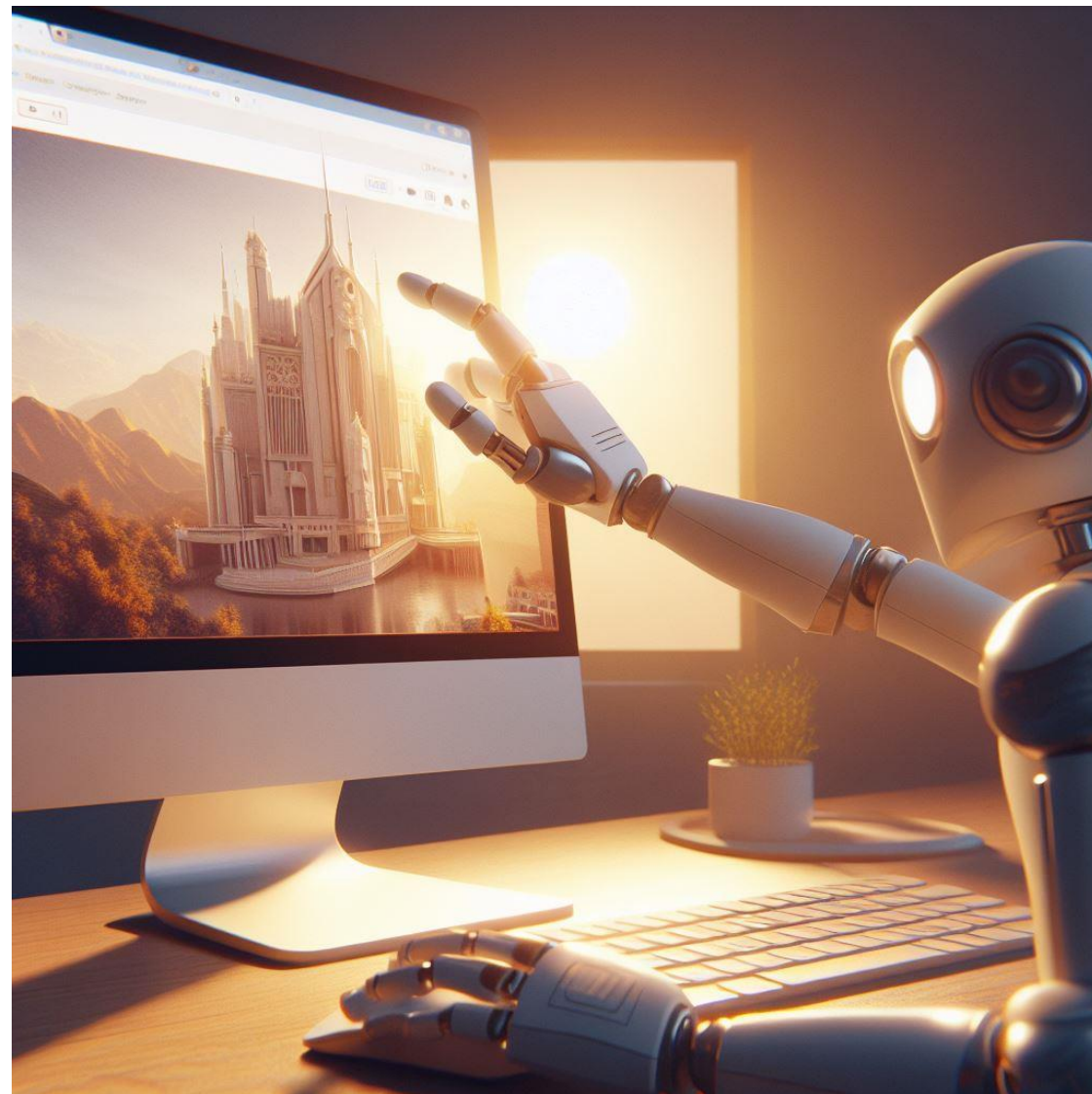
Elon Musk - Neuralink - Chip beültetés az agyba:

<https://hu.wikipedia.org/wiki/Neuralink>

<https://neuralink.com/>

DarkBERT - a DarkWEB-es AI:

<https://medium.com/@sonihariom555/darkbert-ai-the-powerful-weapon-to-combat-cyber-crimes-trained-on-the-dark-web-15ea7272477c>



# És a linkek...

---

NIST AI Risk Management:

<https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.100-1.pdf>

Eu AI Act:

[https://www.europarl.europa.eu/RegData/etudes/BRIE/2021/698792/EPRS\\_BRI\(2021\)698792\\_EN.pdf](https://www.europarl.europa.eu/RegData/etudes/BRIE/2021/698792/EPRS_BRI(2021)698792_EN.pdf)

Cloud Act:

<https://www.justice.gov/criminal/cloud-act-resources>

