

# Selected results of MILAB – Artificial Intelligence National Laboratory Hungary

---

**Vera Könyves**

**Institute for Computer Science and Control**



Supported by the European Union project RRF-2.3.1-21-2022-00004 within  
the framework of the Artificial Intelligence National Laboratory Program



**ARTIFICIAL INTELLIGENCE**  
National Laboratory

GPUday - 16 May 2023



# MILAB – AI National Laboratory Hungary (2020–2025)

## Consortium Leader:

Institute for Computer Science and Control  
(Scientific director of MILAB: **András Benczúr**)

## Partners:

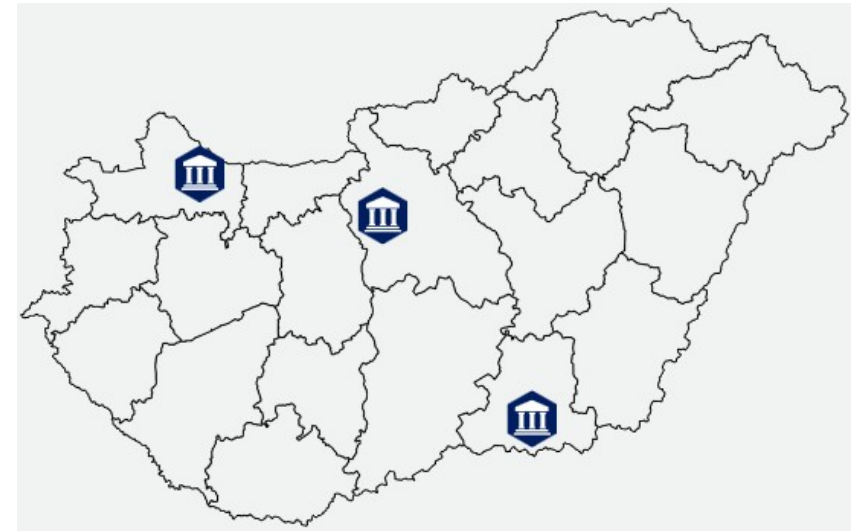
- Alfréd Rényi Institute of Mathematics
- Budapest University of Technology and Economics
- Eötvös Loránd University
- University of Szeged
- Semmelweis University
- Institute of Experimental Medicine
- Centre for Social Sciences
- Széchenyi István University
- Special Service for National Security
- KINCSINFO Nonprofit Ltd.

## External partners in Hungary:

Audi, Bosch, Continental, Ericsson, Nokia,...

## Places of Implementation:

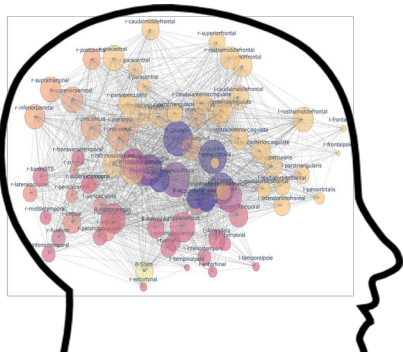
Budapest, Győr, Szeged



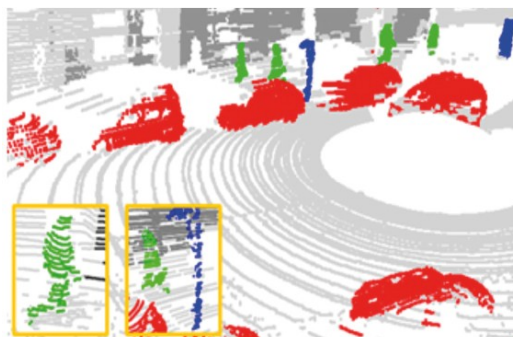




# MILAB – Research Fields & Computing facilities



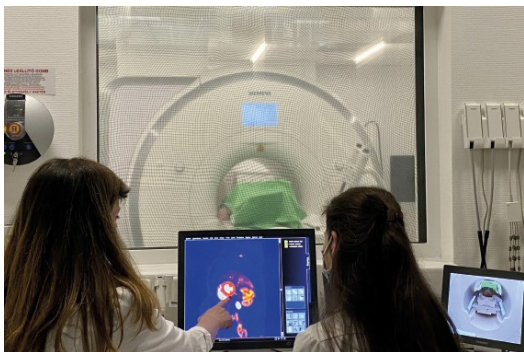
**WP1: Foundations of AI**  
(B. Szegedy, Rényi)  
A100s in SZTAKI & Rényi



**WP2: Machine perception**  
(I. Csabai, ELTE)  
ELKH Cloud, WSCLAB



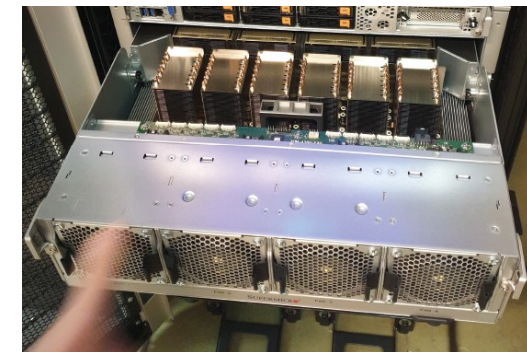
**WP3: Human Language Proc.**  
(R. Farkas, SZTE)  
ELKH Cloud, OTP SambaNova, ...



**WP4: Medical, Health & Bio.**  
(D. Becker, SE)  
ELKH Cloud, WSCLAB, other A100s



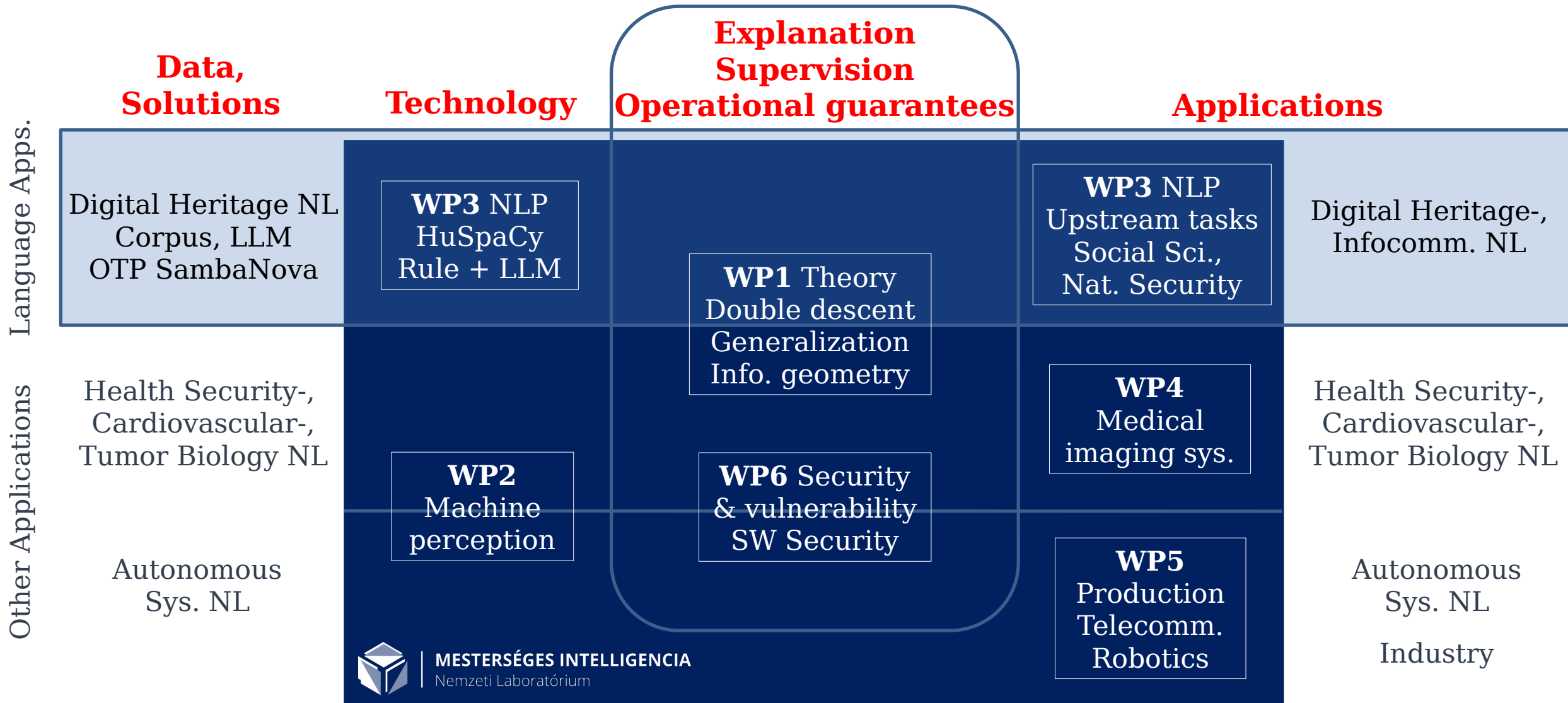
**WP5: Sensors, IoT, Telecomm.**  
(J. Levendovszky, BME)  
Local smaller GPUs



**WP6: Security & Privacy**  
(R. Ferenc, SZTE)



# MILAB – Research+Development+Industry Relations

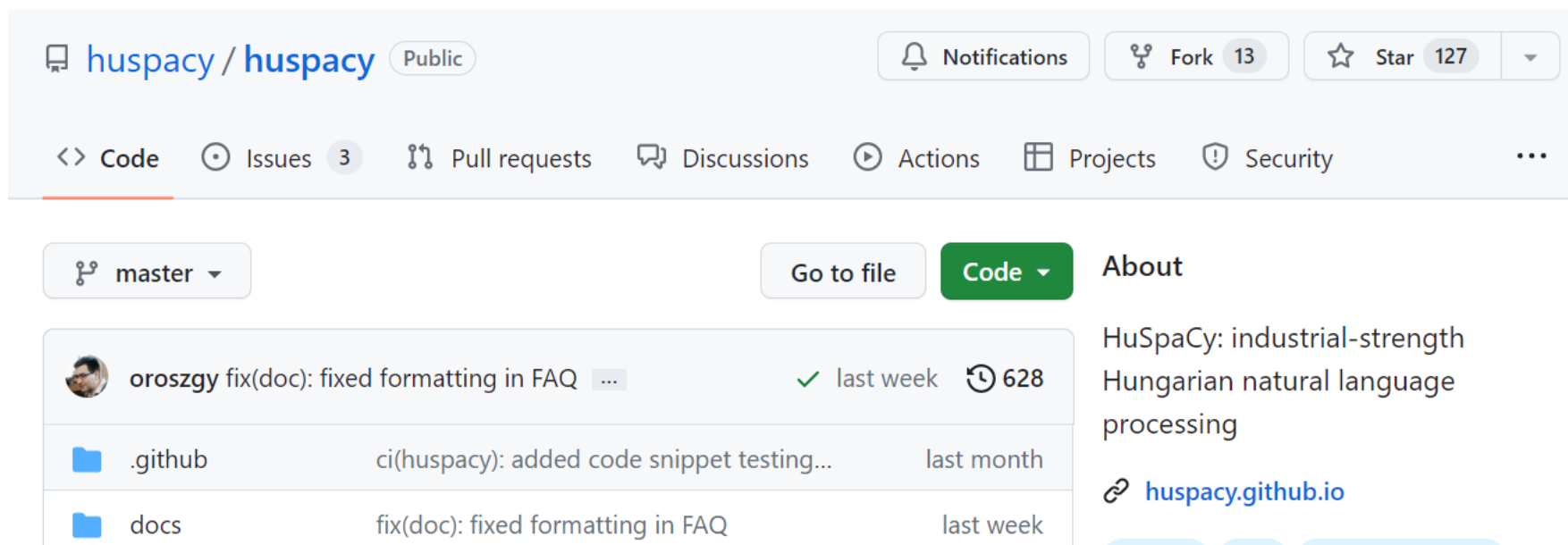
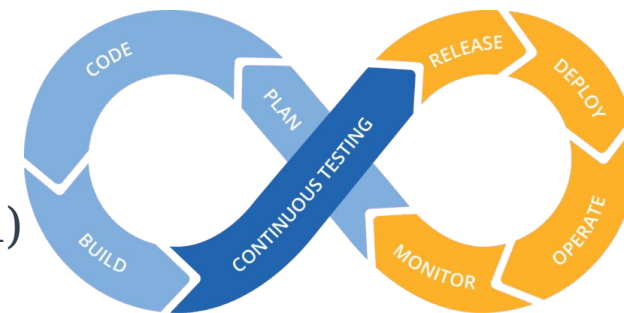


# Natural Language Processing



# MILAB – NLP: HuSpaCy

- A spaCy library providing industrial-strength Hungarian language processing facilities through spaCy models (Orosz et al., 2022)
- Tackles important **text pre-processing tasks**: tokenization, sentence splitting, PoS tagging, lemmatization,...
- **4 benchmarked pipelines** (LLM-based, e.g., huBERT, XLM-Roberta)
- **Publicly available** (GitHub, huggingface)



# MILAB – NLP: Question Answering from Hun. Wikipedia

**MILQA** database (Novák, Novák, 2023, public soon)

- Extractive & abstractive QA
- **Good quality Wikipedia articles** as context for Qs
- ~10K Q-A of different types in the database
- Includes Qs: Y/N; arithmetics; not answered in text
- Short/long answer



Can be built in a chatbot, **but:**

- „always” **factual**
- **no hallucinations**
- underlying knowledge is **quickly updated**
- shows source of text





# MILAB – NLP Problems: Factual Error

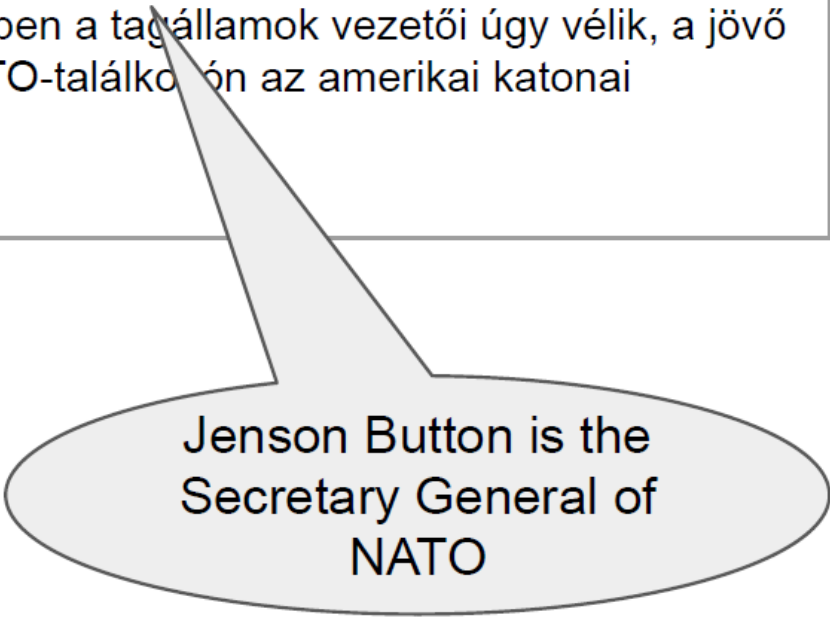
Input article	Model output (Bert2bert)
<p>A férfi 2020 áprilisában felkereste barátnőjét, hogy rendezzék a párkapcsolati vitáikat, és hazavigye a lányt, aki ekkor a rokonainak miskolci lakásában lakott. A pár a ház előtt beszélgetett, amikor a házigazda kiment hozzájuk és számonkérte a férfit, hogy minek ment oda, majd a jobb kezében lévő villáskulccsal, illetve a bal kezében lévő asztallábbal támadóan közelített felé, és a villáskulccsal mellkason ütötte. A bántalmazott férfinak eltört a bordája, a sérülés miatt légmelle lett, amely közvetetten életveszélyes állapotot idézett elő. A férfi beismerte a bűncselekmény elkövetését. A törvényszék a férfival szemben életveszélyes testi sértés miatt másfél évre felfüggesztett börtönt szabott ki, és ennek tartama alatt a férfi pártfogó felügyeletét rendelte el. A büntetővégzés jogerőre emelkedett.</p>	<p>A Miskolci Járási Ügyészség életveszélyes testi sértés miatt emelt vádat azzal a férfival szemben, aki 2020 - ban megverte a barátnőjét Miskolcon, mert az elhagyta őt.</p> <p>He beat his girlfriend</p>

abused man

Barta, Lakatos et al. (2023)

# MILAB – NLP Problems: Hallucination

Input article	Model output (mT5-small)
<p>Sajtótájékoztatóján Stoltenberg arról beszélt, hogy a NATO -nak számos kihívással meg kell küzdenie a mai megváltozott, kihívásokkal teli biztonsági környezetben, de továbbra is nagy szükség van a katonai szövetségre. Belföntesek szerint a NATO szerdán kezdődő kétnapos ülése főként James Mattis új amerikai védelmi miniszter bemutatkozása miatt fontos, aki várhatóan ismertetni fogja a nemrég hivatalba lépett, Donald Trump vezette kormányzat elképzeléseit. ...</p>	<p>Jens Stoltenberg, az Európai Unió (EU) főtitkára a NATO csütörtöki kétnapos ülésén elmondta, hogy Donald Trump amerikai elnök által vezetett kormányzati elképzelései mellett továbbra is fontos a szövetségre vonatkozó erőfeszítések - jelentette be <b>Jenson Button NATO-főtitkár az MTI-nek adott interjújában</b>, amelyben a tagállamok vezetői úgy vélik, a jövő héten kezdődő NATO-találkozáson az amerikai katonai</p>



Jenson Button is the  
Secretary General of  
NATO

Barta, Lakatos et al. (2023)



# MILAB – NLP: Mapping between Languages

## Multilingual pretrained language models (PLMs)

- Convenient solution to zero-shot cross lingual problems
- Suffer from multilinguality (Conneau et al., 2020)
- Unstable performance across supported langs. (Wu & Drezde, 2020)
- Improper tokenization is an issue (Rust et al., 2021)

## Monolingual PLMs to the rescue

- Take pair of independently pretrained monoling. BERT models
- Dissect source and target language models into two halves
- Determine mapping between chosen layers of source/target

G. Berend (2022)

Contextual translation pairs  
from the Tatoeba corpus

Snails move slowly.

A csigák lassan mozognak.

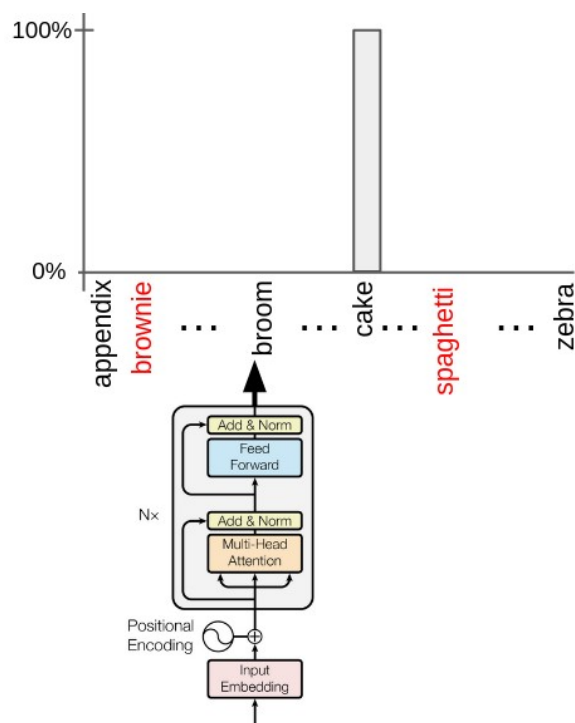


# MILAB – NLP: Pre-training of medium-sized models

## Exploiting sparsity during pre-training

Pre-trained LMs excel at continuing/infilling partial/masked (token) sequences

- **Masked Language Modeling** is a typical training task for LMs
- Misalignment in pre-training objective and behavior of well pre-trained model



Alice is eating a cake.

G. Berend et al.



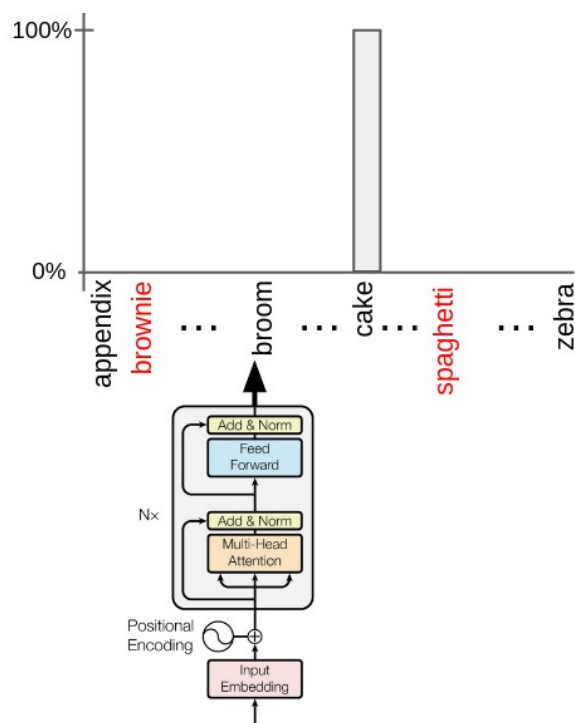


# MILAB – NLP: Pre-training of medium-sized models

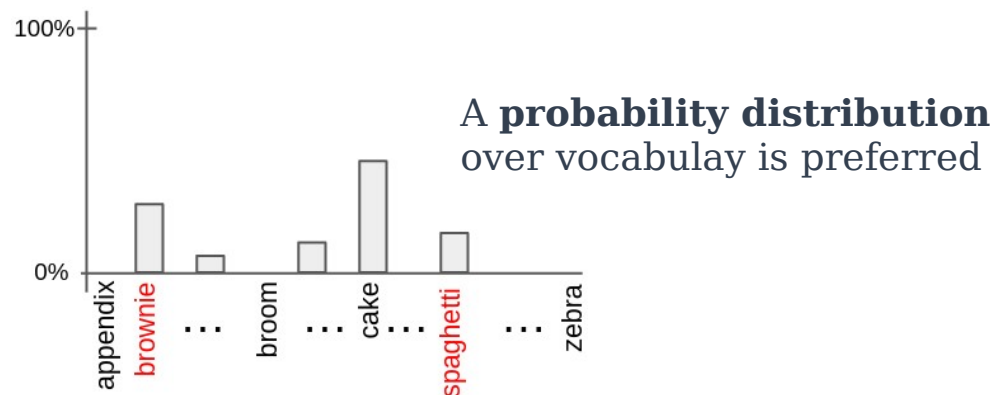
## Exploiting sparsity during pre-training

Pre-trained LMs excel at continuing/infilling partial/masked (token) sequences

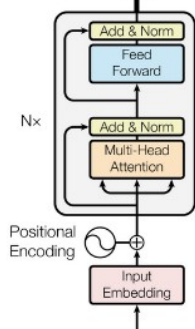
- **Masked Language Modeling** is a typical training task for LMs
- Misalignment in pre-training objective and behavior of well pre-trained model



VS.



A **probability distribution** over vocabulary is preferred



## Possible remedies to the problem:

- 'Real' word distributions
- Knowledge bases (e.g. WordNet)
- Unsupervised sparse contextual representations

G. Berend et al.

Alice is eating a **[MASK]**.



# **NLP, LLMs**

## **Theory & Applications**

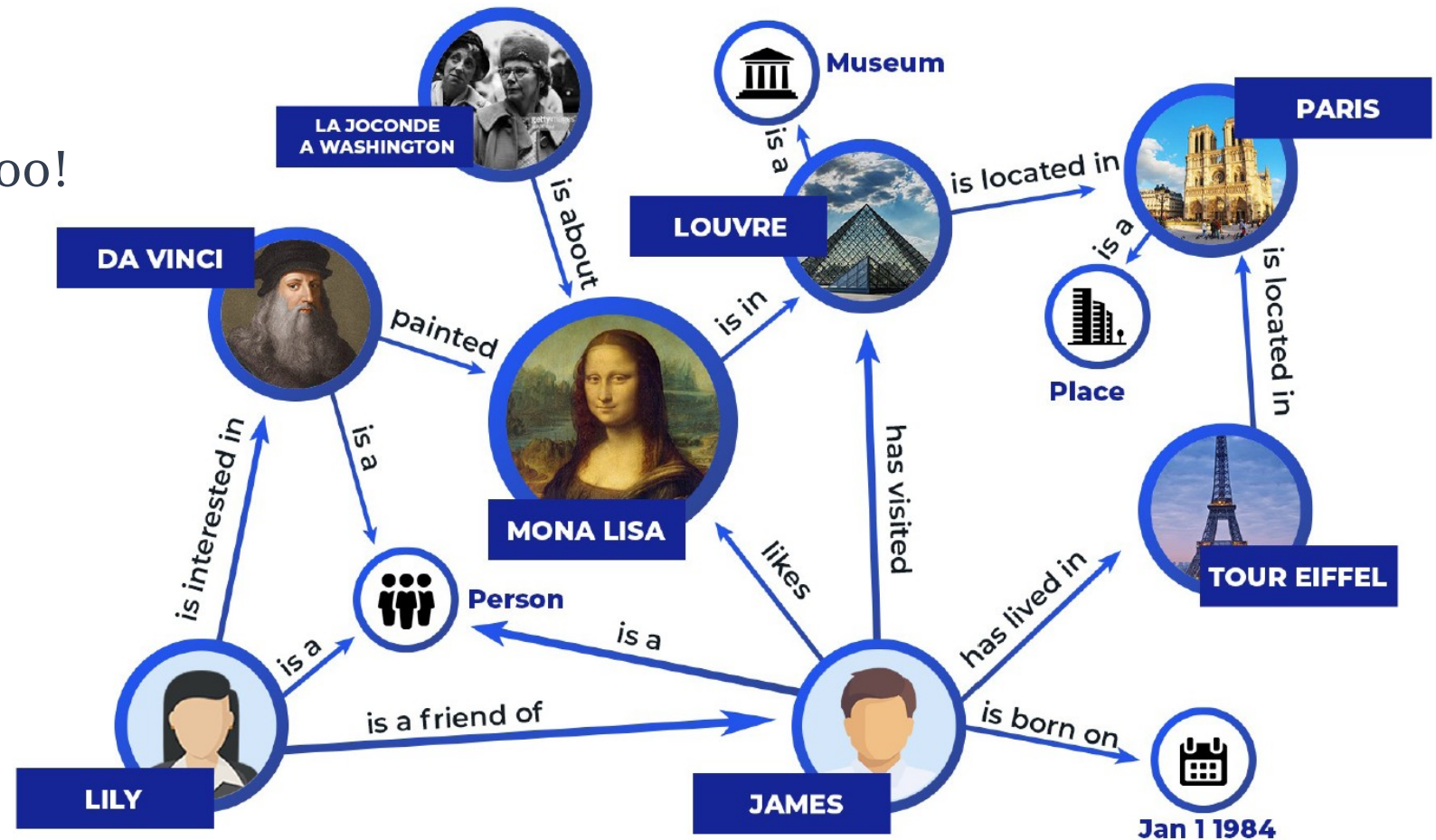


## Cost–accuracy trade-off

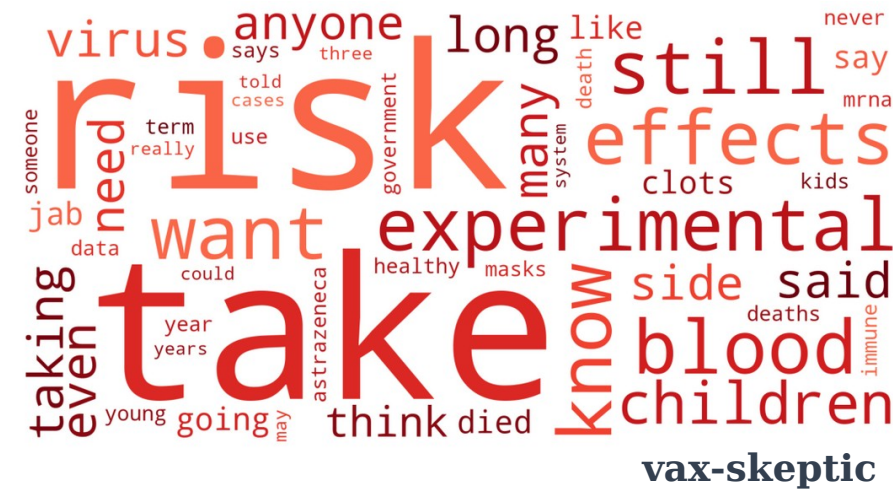
- Smaller LLMs
- Reduce data hunger
- Text & speech technology too!

Using environmental info

- **Knowledge graphs**
- Dialogue history
- “Microcosm” of the system



# MILAB – Recognising vaccination-skeptic content

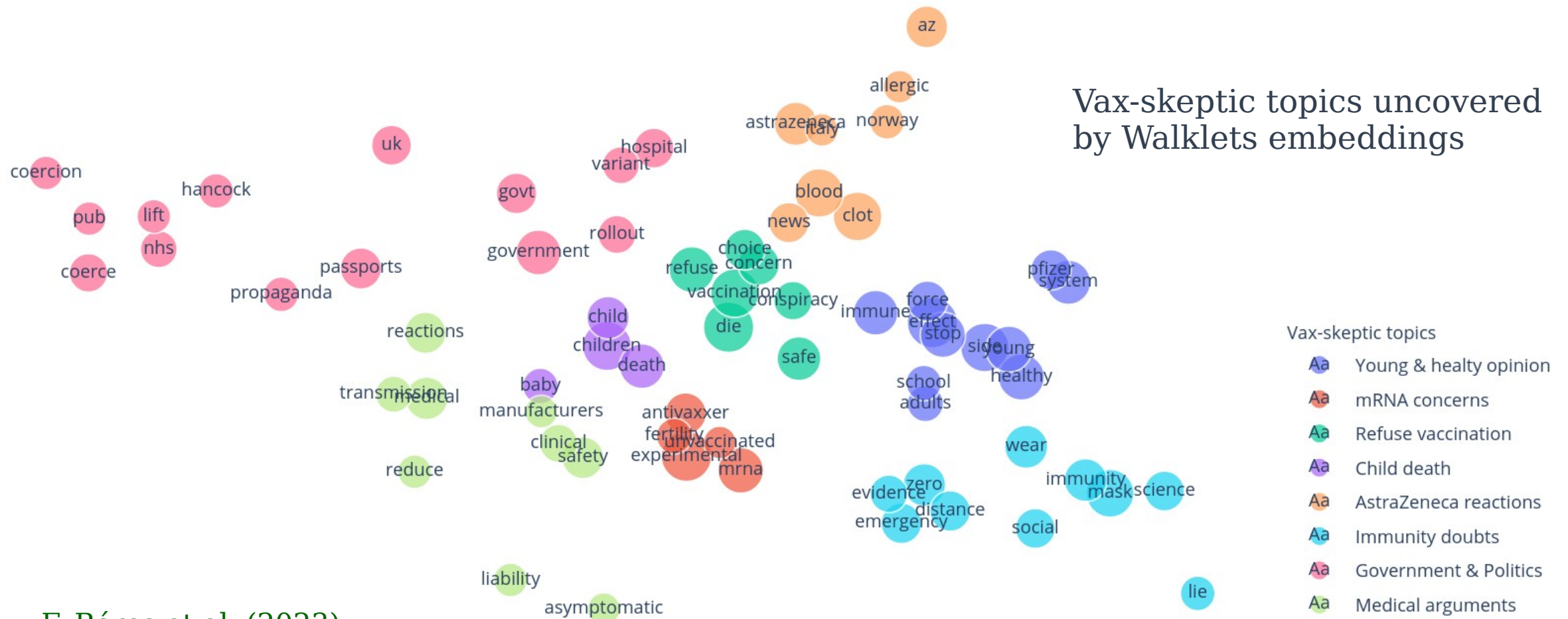


F. Bérés et al. (2023)

Model	Vax-skeptic AUC	Time
Vaccinating-covid-tweets (Pak and Paroubek 2010)	<b>0.81</b>	12.98
<b>Bert-small</b> (Bhargava et al. 2021; Turc et al. 2019)	0.793	3.88
Covid-twitter-bert (Müller et al. 2020)	0.787	34.72
Bert-medium (Bhargava et al. 2021; Turc et al. 2019)	0.779	6.17
Bertweet-covid19-base (Nguyen et al. 2020)	0.766	12.89
Bertweet-base (Nguyen et al. 2020)	0.765	13.01
Bert-mini (Bhargava et al. 2021; Turc et al. 2019)	0.751	2.66
Bert-tiny (Bhargava et al. 2021; Turc et al. 2019)	0.709	1.84
Bert-base (Devlin et al. 2018)	0.709	13.53
<b>Bertweet-large</b> (Nguyen et al. 2020)	0.575	34.05
<b>Bert-large</b> (Devlin et al. 2018)	0.556	34.52



# MILAB – Vaccine skepticism with network embedding



F. Béres et al. (2023)

Rozemberczki, ... Kiss, Béres et al., CIKM (2021), **Best Resource Paper**



# MILAB – Theory: Double Descent of big models

(e.g., Hastie et al., 2005)

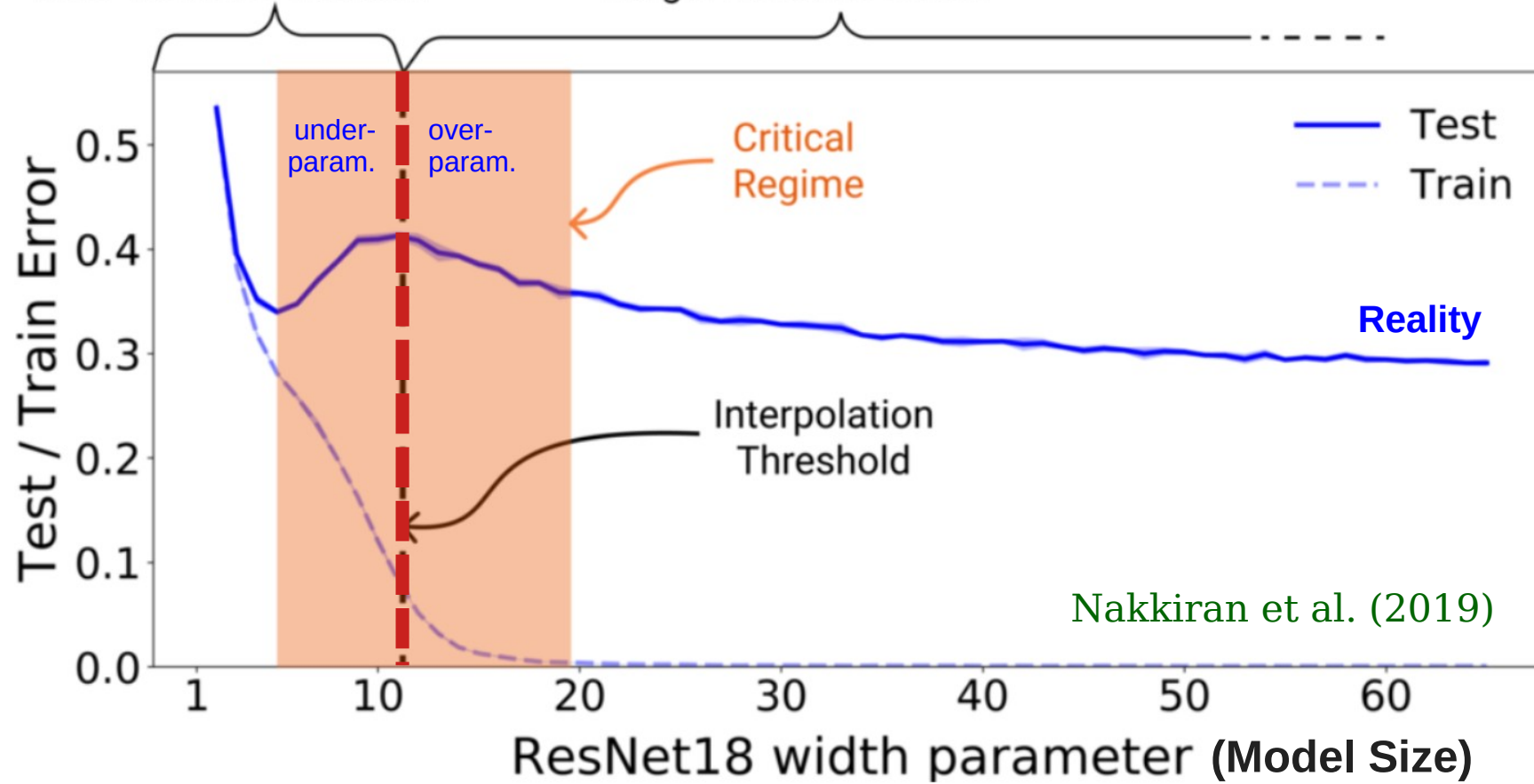
(e.g., Szegedy et al., 2015)

## Classical Regime

Bias-Variance Tradeoff

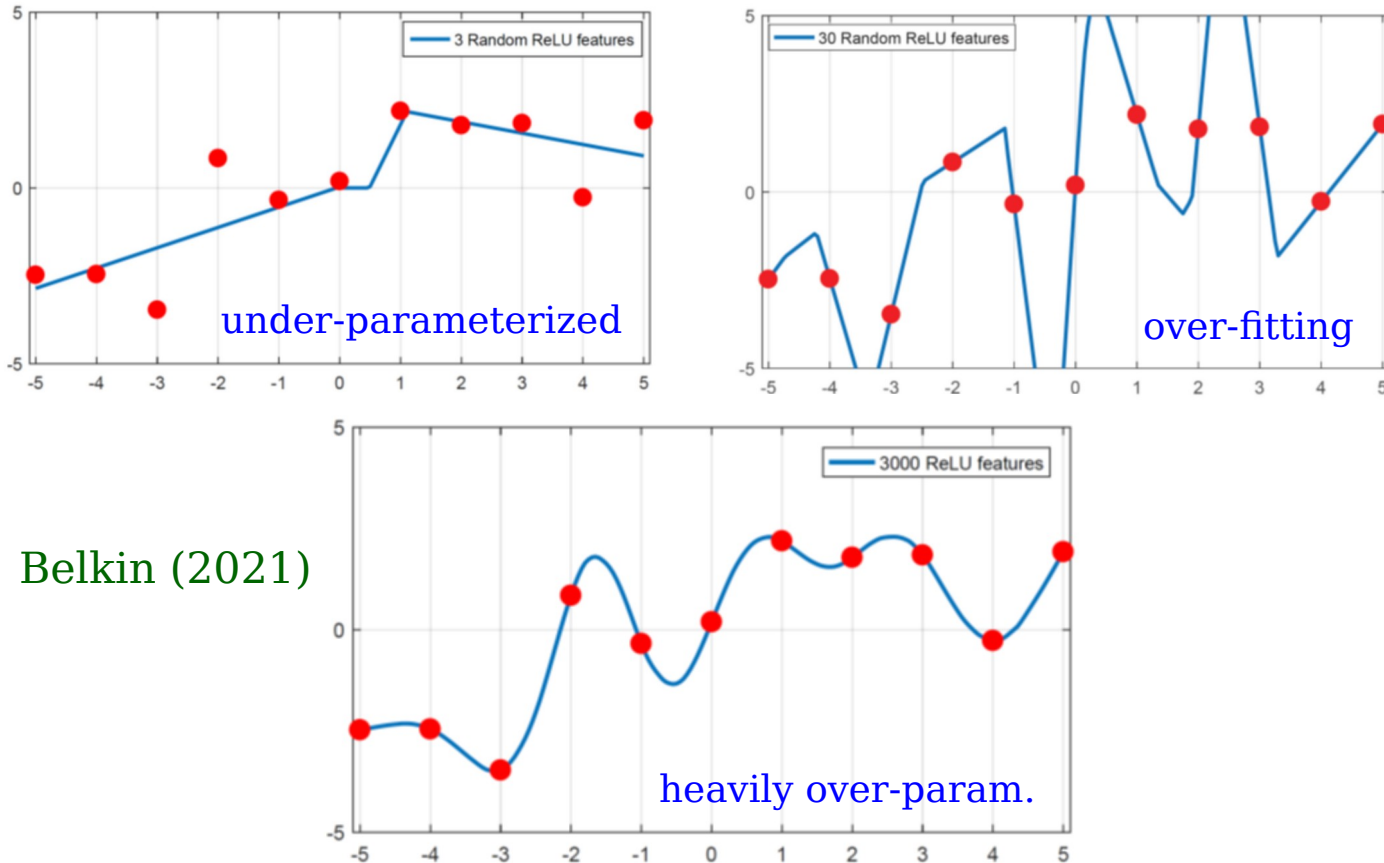
## Modern ML

Larger Model is Better



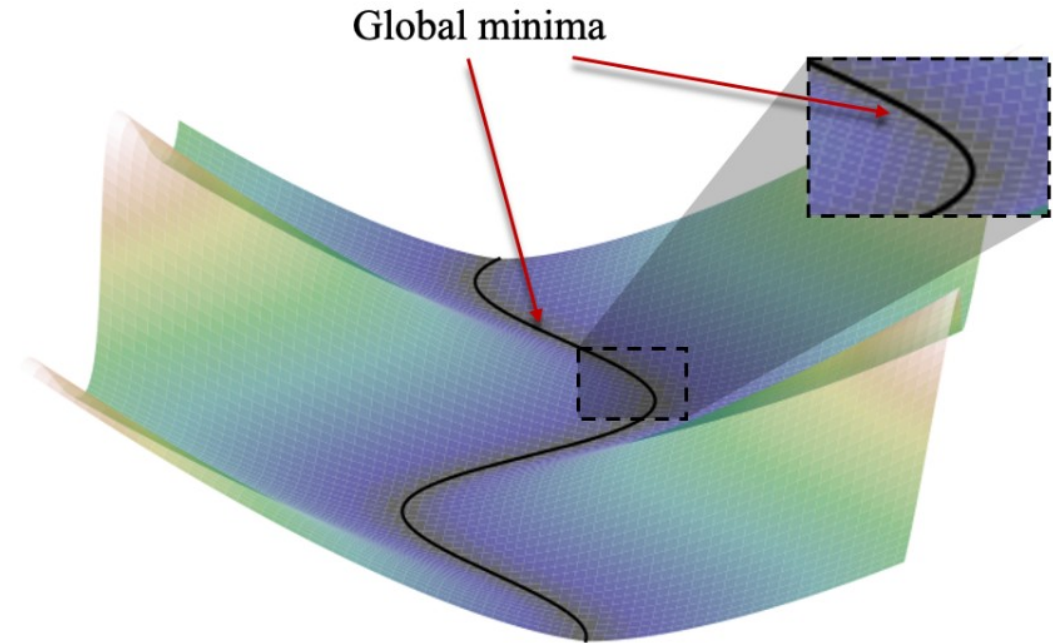
# MILAB – Theory: Over-parameterized NNs

Illustration of double descent for random ReLU networks in 1D (Smoothness)



Belkin (2021)

Loss landscape of over-param. models (Non-convex optimization)



Liu, Zhu & Belkin (2021)



# Security & Privacy





# MILAB – NN Model Limitations

Pointer Value Retrieval (**PVR**) benchmark with sequence of images to **understand generalization and reasoning** capabilities of NNs

Works based on **indirection** - component of human-like reasoning

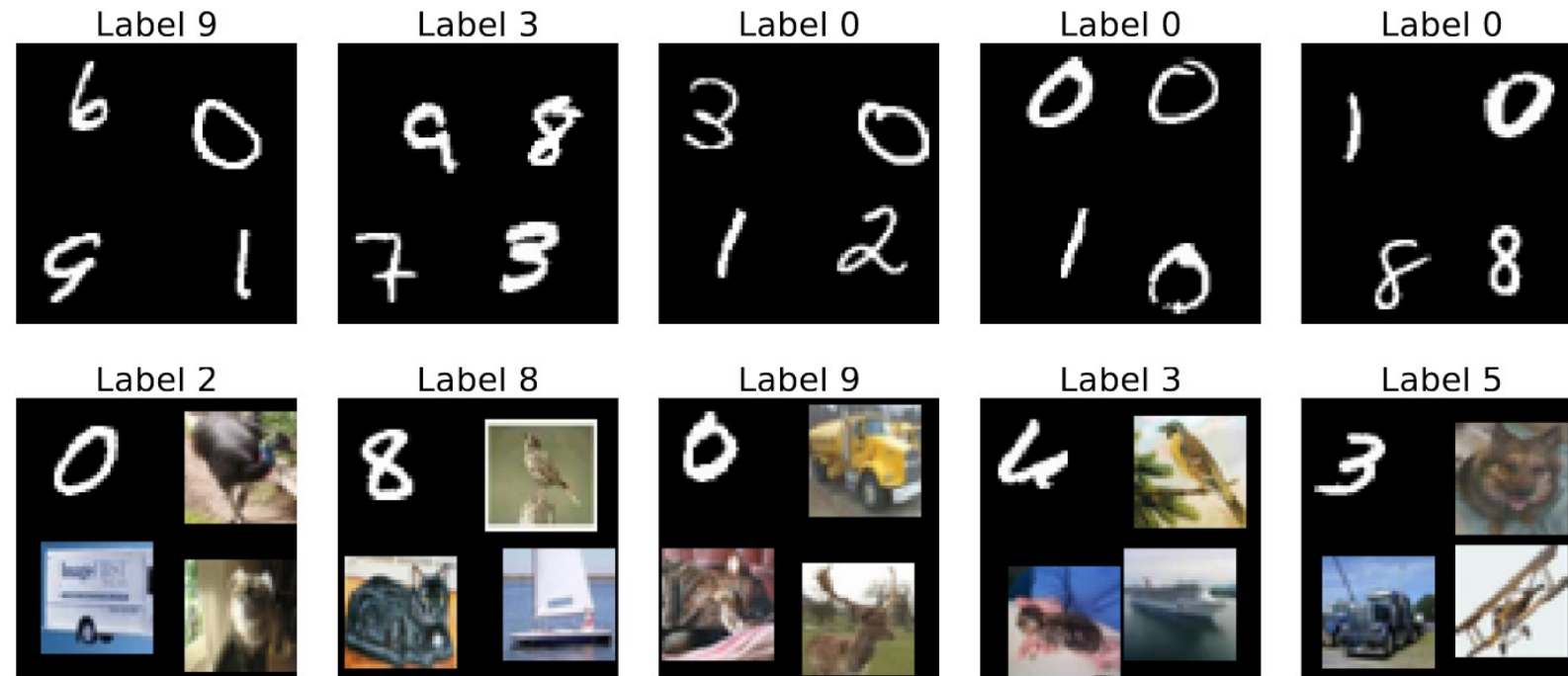
2 x 2 Block Style Image PVRs (MNIST, CIFAR-10)

**Pointer (1<sup>st</sup> task) - top left image**  
gives indication for next image to be  
examined (**2<sup>nd</sup> task**):

digits 0–3: upper right

4–6: lower left

7–9: bottom right

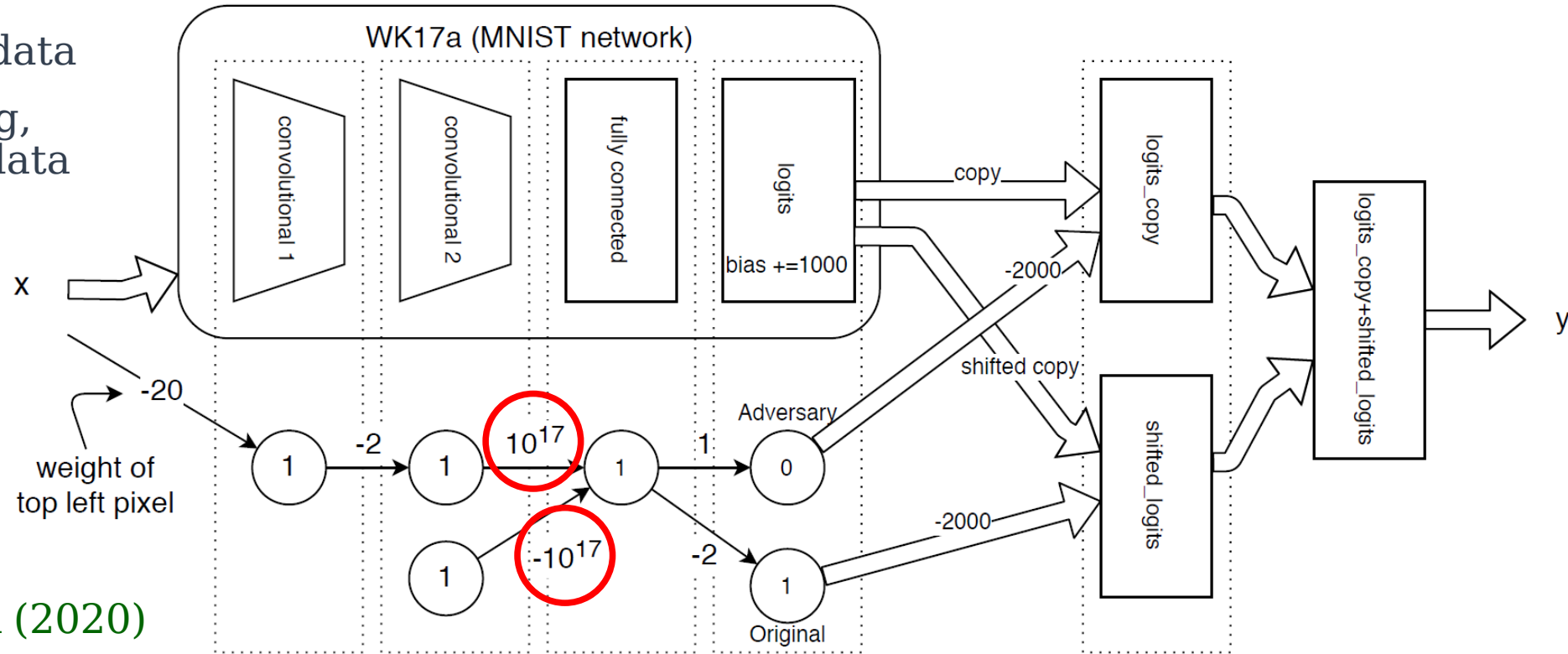


C. Zhang, ... S. Bengio (2021)



# MILAB – Security Issues of AI Systems

- Difficult to verify
  - Not possible to count every case
  - Unexpected input – unexpected operation
- Backdoors cannot be detected
- „Poisoning” teaching data
- Exploiting overtraining, generation of hostile data



Zombori et al., ICLR (2020)



# **Machine Perception**

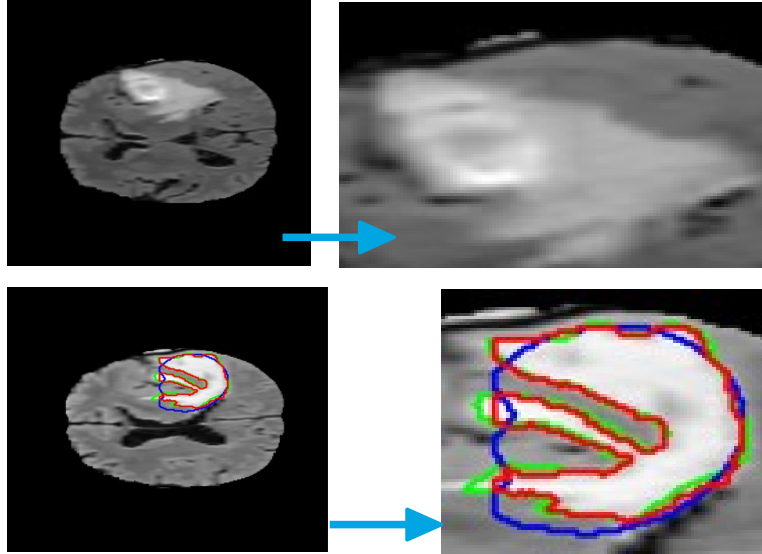
## **Medical, Health & Biology projects**

**The promise of LLMs...**

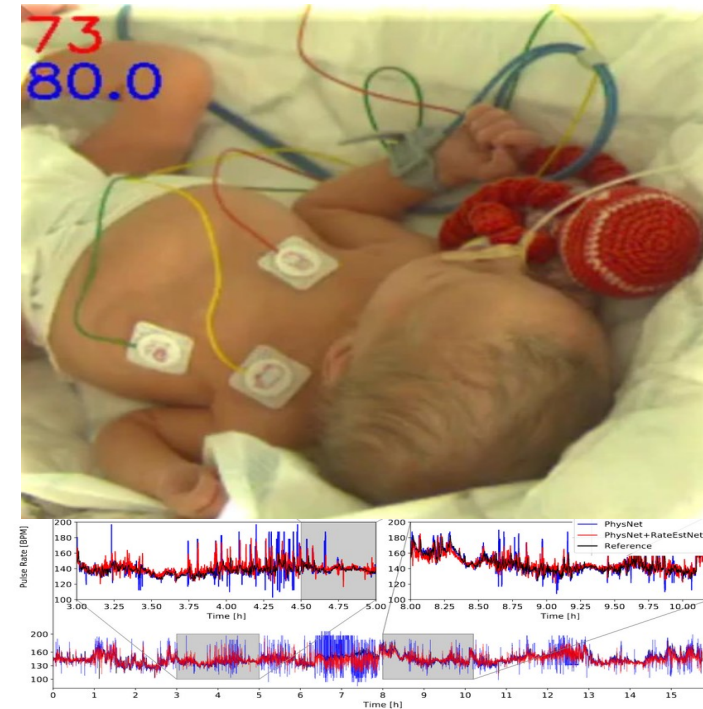
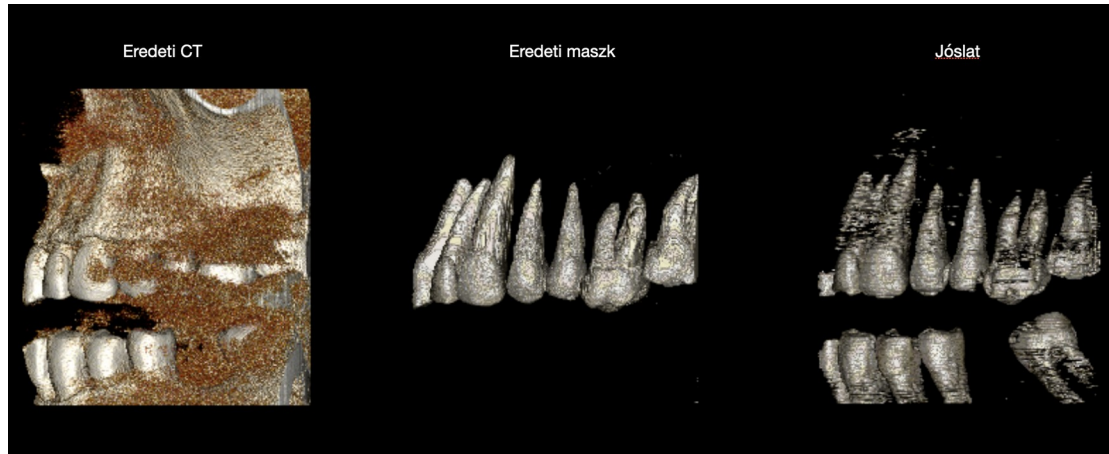
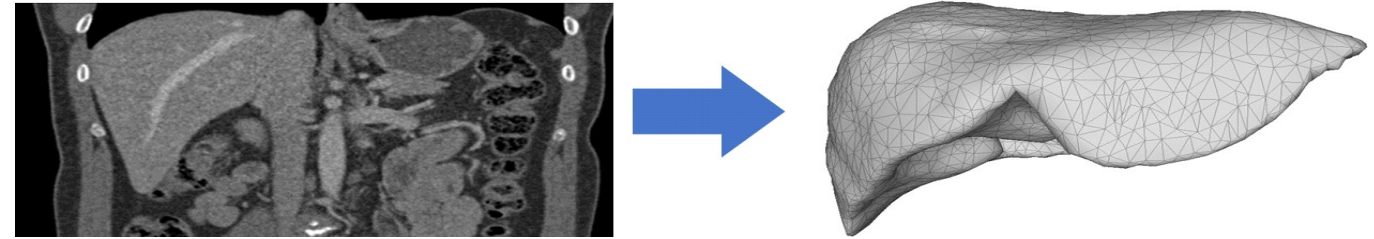


# MILAB – Medical, Health & Biology

## Tumour screening & segmentation



## Automatic 3D organ reconstruction





## To understand the evolution and functional properties of Influenza virus

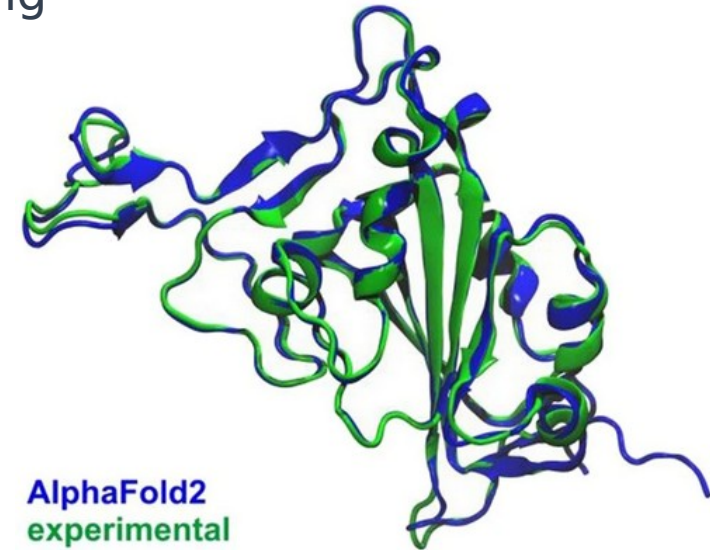
"Combining Influenza hemagglutinin antigenic maps with deep mutational scanning data"  
(**GPU-Lab run**, Á. Gellért, O. Kilim, A. Mentés, I. Csabai)

- How different mutations affect the ability of virus to evade the immune response
  - Which regions of the virus are critical for this evasion
- how to predict the phenotypic properties of influenza viruses using mutations in the genetic sequence alone.

### ML:

- Protein 3D-structure pred.: AlphaFold2 (J. Jumper et al., 2021), ESMFold2 (Z. Lin et al., 2022)
- Predict phenotypes

**SARS-CoV-2 virus:**  
O. Kilim et al., Nat. Sci. Data (2023)



# MILAB – Medical: Gigapixel whole-slide image analysis

A machine learning competition to **predict the stage of patient's cancer, using only slide images generated by breast biopsy.**

**Data:** 4200 cases, 72000 high-res pathological sections, 130 TB data

**Winning techniques:** smart image proc., transfer learning, multiple-instance learning.

<https://app.nightingalescience.org/contests/vd8g98zv9w0p/leaderboard>



[Contests](#) / [Predicting High Risk Breast Cancer - Phase 2 \(2023\) \(vd8g98zv9w0p\)](#)

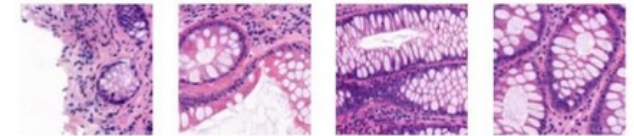
## Predicting High Risk Breast Cancer - Phase 2 (2023)

Ranking	Team Name	Score	Description
1	csabAlbio	0.7606058	Seventh submission
2	Bonaventure Dossou	0.7305834	preds_all_geo_times_arith_mean

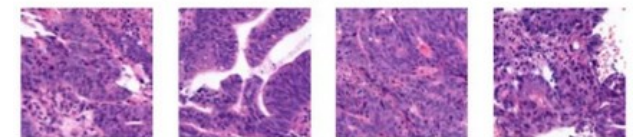


**Earlier results: HunCRC –  
annotated pathological slides**  
B. Á. Pataki et al., Nat. Sci. Data (2022)

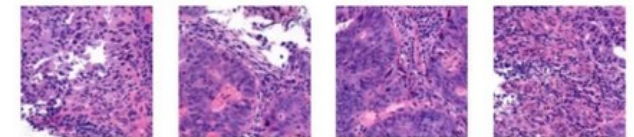
*normal*



*adenocarcinoma*



*high-grade dysplasia*



Pathology tissue sections with  
stages of colorectal cancer.



**Machine Learning** = statistics + lot of data

**AI** = ML + sensors + HW (robot) + **restrictive rules**



Human intervention required

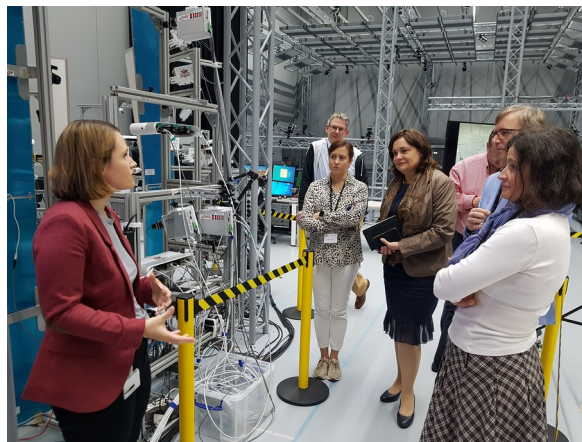
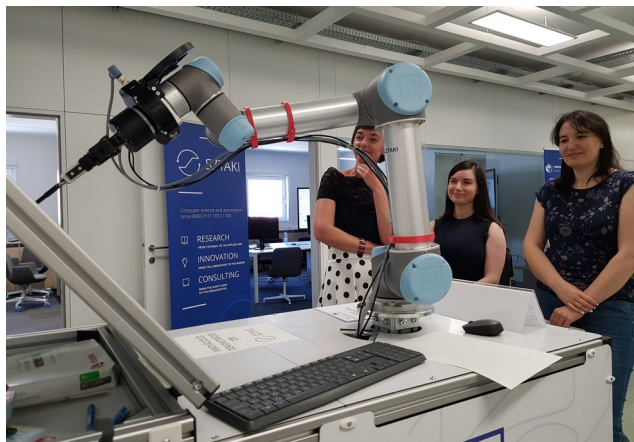


Physically enclosed, shuts down when people enter



# MILAB – AI National Laboratory Hungary

## Thanks for the Attention!



**ARTIFICIAL INTELLIGENCE**  
National Laboratory

GPUday - 16 May 2023

**NEMZETI  
LABORATÓRIUM**



NEMZETI KUTATÁSI, FEJLESZTÉSI  
ÉS INNOVÁCIÓS HIVATAL



INNOVÁCIÓS ÉS TECHNOLÓGIAI  
MINISZTERIUM