# KI und Sicherheit –
## Zwischen Innovation und Risikomanagement

**InfoTechDay**, 2025-10-09 14:05 (UTC+2), Geinberg

Univ.-Prof. Dr. **René Mayrhofer**

Institut für Netzwerke und Sicherheit und LIT Secure and Correct Systems Lab, JKU Linz

# Artificial Intelligence (AI):
# Current successes and opportunities

- **Pattern recognition**
  - Medicine
  - Face recognition, fingerprint comparison, etc.

- Faster **filtering** of possible solution candidates
  - Materials science
  - Pharmaceutical development
  - Weather prediction

| David Baker | Demis Hassabis | John Jumper |
|---|---|---|
| "for computational protein design" | "for protein structure prediction" | "for protein structure prediction" |

© Nobel Prize Outreach. Photo: Clément Morin — © Nobel Prize Outreach. Photo: Clément Morin — © Nobel Prize Outreach. Photo: Clément Morin

https://www.nobelprize.org/all-nobel-prizes-2024/

# They cracked the code for proteins' amazing structures

The Nobel Prize in Chemistry 2024 is about proteins, life's ingenious chemical tools. David Baker has succeeded with the almost impossible feat of building entirely new kinds of proteins. Demis Hassabis and John Jumper have developed an AI model to solve a 50-year-old problem: predicting proteins' complex structures. These discoveries hold enormous potential.

### Related articles

Press release

© Johan Jarnestad/The Royal Swedish Academy of Sciences

**JOHANNES KEPLER UNIVERSITY LINZ**

## John J. Hopfield

"for foundational discoveries and inventions that enable machine learning with artificial neural networks"

© Nobel Prize Outreach. Photo: Nanaka Adachi

## Geoffrey Hinton

"for foundational discoveries and inventions that enable machine learning with artificial neural networks"

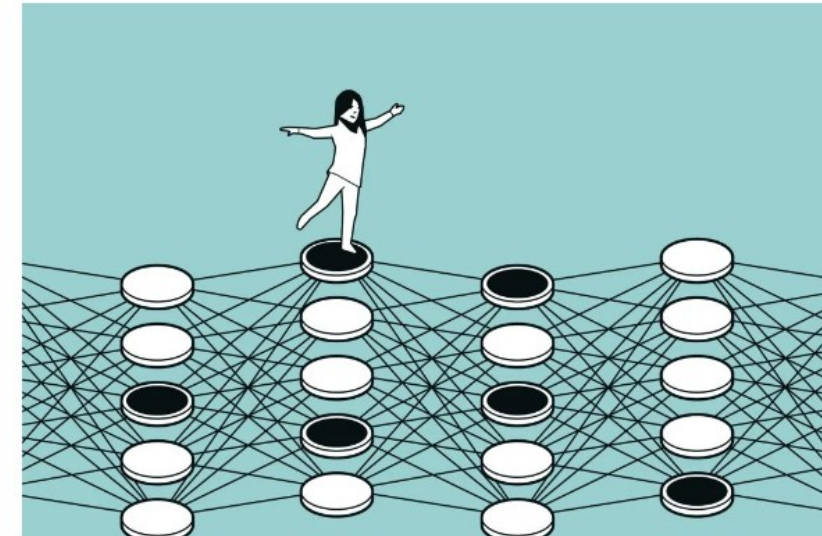© Nobel Prize Outreach. Photo: Clément Morin

# They used physics to find patterns in information

This year's laureates used tools from physics to construct methods that helped lay the foundation for today's powerful machine learning. John Hopfield created a structure that can store and reconstruct information. Geoffrey Hinton invented a method that can independently discover properties in data and which has become important for the large artificial neural networks now in use.

Related articles

Press release

© Johan Jarnestad/The Royal Swedish Academy of Sciences

https://www.nobelprize.org/all-nobel-prizes-2024/

JOHANNES KEPLER
UNIVERSITY LINZ

# What is AI?
# Better describe as Machine Learning (ML)

```
┌──────────────┐      ┌──────────────────────────────────┐      ┌──────────────┐
│              │      │  Specific algorithms programmed  │      │              │
│  Input data  │ ───▶ │        by domain experts         │ ───▶ │   Decision   │
│              │      │                                  │      │              │
└──────────────┘      └──────────────────────────────────┘      └──────────────┘
```

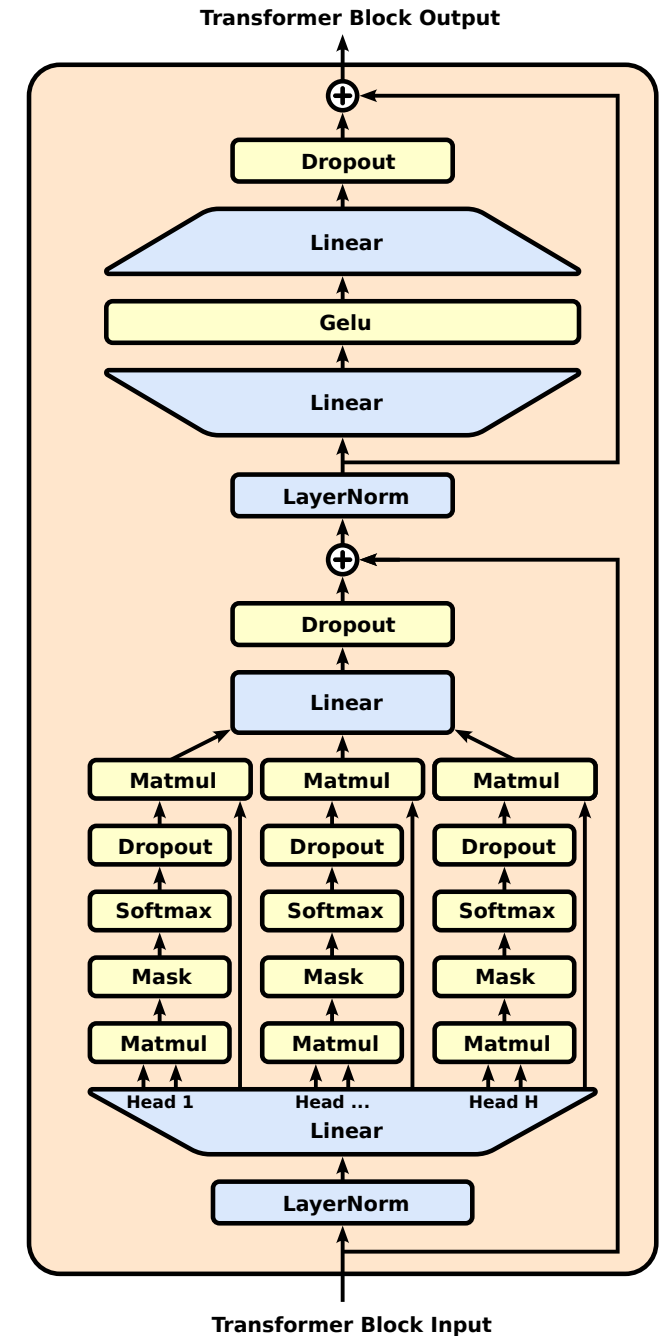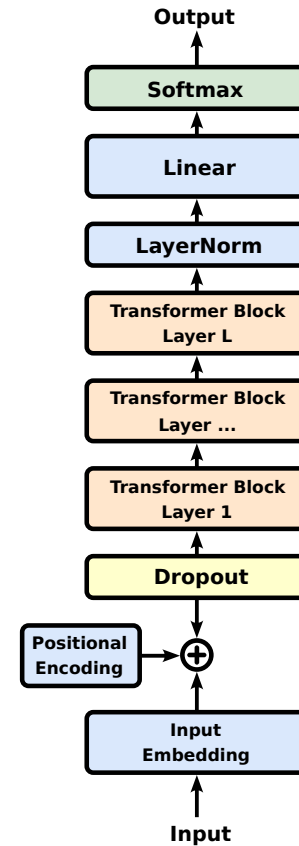# What is AI?
# Better describe as Machine Learning (ML)

# Preconditions for using ML

- Availability of sufficiently accurate training and test data
  - Acquisition / recording
  - **Manual** pre-selection of data
  - Check for plausibility, correctness, and potential bias

- **Manual** selection of decision (classes) to derive
  - Mapping training data to decisions = „**ground truth**"

- *Either*
  - **Manual** selection and optimization of features from raw data *or*
  - Deep Learning (DL) approaches directly on raw data with automatic feature derivation

- Training the model
  - **Manual** selection of an appropriate method and model structure
  - Training of the model according to data (mostly: optimizing internal parameters for target values)
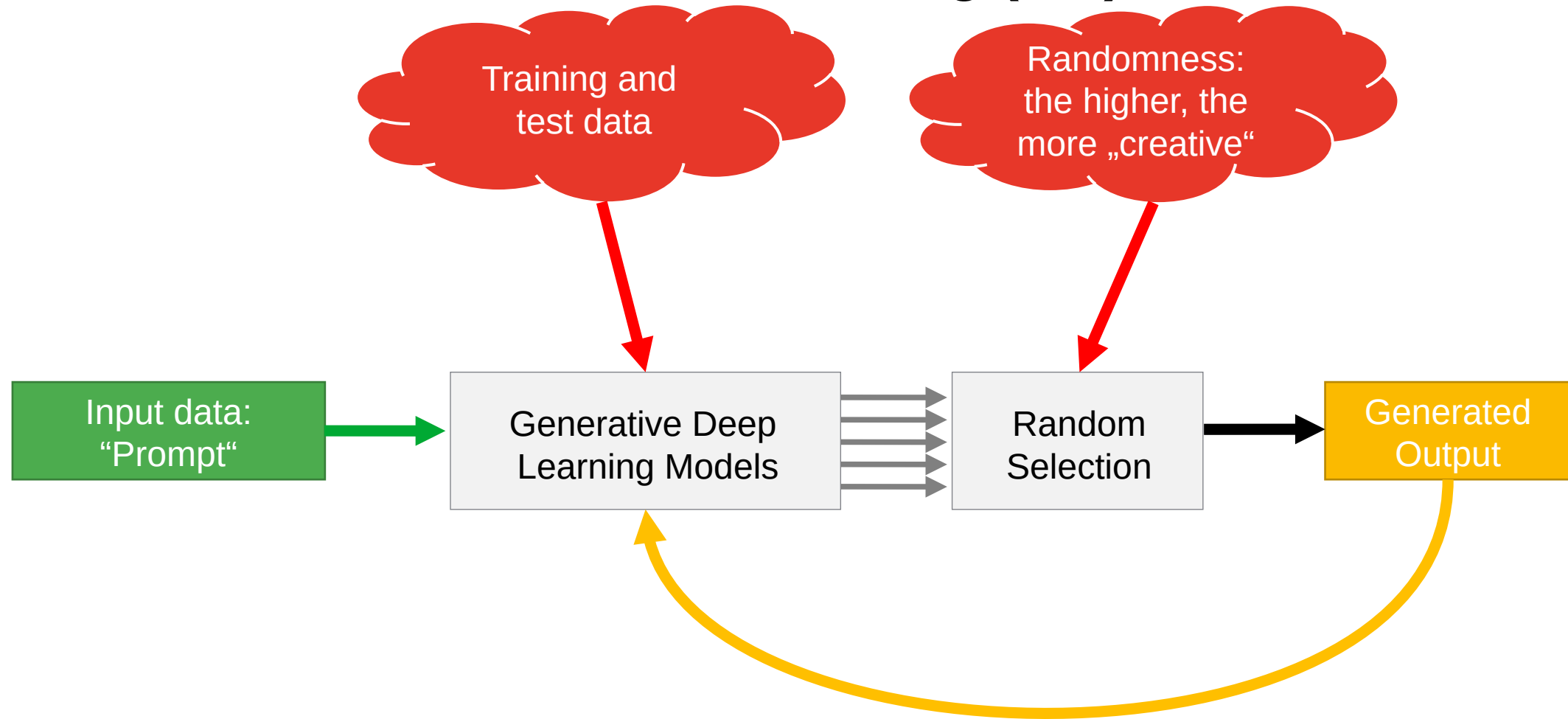  - **Manual** verification of results quality

# Generative Pre-Trained Transformers (GPTs)

Quelle: https://en.wikipedia.org/wiki/Generative_pre-trained_transformer



JOHANNES KEPLER
UNIVERSITY LINZ

# What is AI?
# Better describe as Machine Learning (ML)

# ML at Facebook, Xitter, etc.

# Which post would you rather share/like/boost?



Pope Francis no longer on ventilation after five weeks in hospital, Vatican says

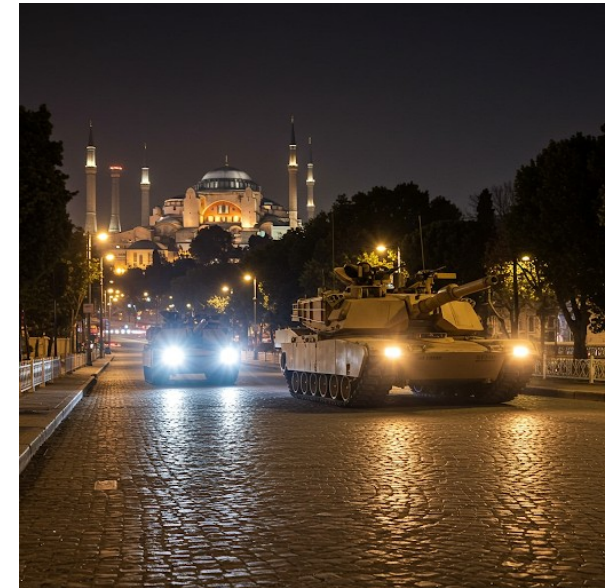# Which post would you rather share/like/boost?



Israel strikes Lebanon after first rocket attack since ceasefire

The UN peacekeeping force in Lebanon, Unifil, said it was "alarmed by the possible escalation of violence".

**USA starten Invasion in Istanbul**

Mitten in der Nacht begann die USA eine militärische Invasion in Istanbul. Panzer und Truppen besetzten zentrale Stadtteile. Erste Berichte sprechen von heftigen Kämpfen. Die türkische Regierung hat den Notstand ausgerufen. Internationale Reaktionen werden erwartet.

Foto: US-Panzer rücken in Istanbul vor.

# The battle for our attention

- Recommendation algorithms are optimized for

  - Drawing and keeping **attention** on that particular platform …

  - … because then people see more advertisements …

  - … which makes the platform operator more **money**.

- How do you get the most attention?

  - **Emotion**!

  - Bad emotion works better than good emotion!

  - The more **emotion**, the more **money**!

Created with StableDiffusion prompt "A photo of a really angry cat i

**JYU** **JOHANNES KEPLER**
**UNIVERSITY LINZ**

# Possible data leaks when using ML

Training and test data

Feature Selection

Training / Optimization

Input data → Feature Extraction → Classification / Prediction → Decision

Data leak: direct or indirect through correlation with other data sources

JOHANNES KEPLER UNIVERSITY LINZ

# World's Biggest Data Breaches & Hacks

Selected events over 30,000 records stolen

UPDATED: Sep 2025

https://www.informationisbeautiful.net/visualizations/worlds-biggest-data-breaches-hacks/

interesting story

size: records lost    filter

search...

**2025**

**2024**

AT&T
73m

Dell

Free

French government
43m

Internet Archive

Epicscape

Kaiser Permanente

Indonesia's health agency

National Public Data
2.7bn

Ticketmaster
560m

Quantas

Santander

UnitedHealth
190m

Welltok

Xfinity

**Acer**

**2023**

23andMe

Clorox unknown

Delta Dental

Latitude Financial

Maximus    MGM

The Post Millennial

TIAA

X (Twitter)
200m

T-Mobile

Yum!

Indian Railways

Epik

Microsoft

MSI

Microsoft unknown

Shanghai Police
"one billion"

Twitter

**2022**

CDEK
19m

Facebook
533m

Indonesian SIM cards
1.3bn

LastPass

Optus    Plex

Pandora Papers

Park Mobile    Star Alliance    Shein

T-Mobile

Uber

Thailand visitors
100m

VW

**2021**

Acer

Amazon

Digital Ocean unknown

Contact tracing data
38m

Air India

Gab
100K

Dutch Government

Twitch unknown

Twitter

Ubiquiti

**2020**

Canva
139m

db8151dd
22m

Drizly

Experian Brazil
220m

HauteLook

Microsoft
250m

Pakistani mobile operators
115m

Syniverse
unknown

Armor Games

Dubsmash
162m

EyeEm

Indian citizens
275m

MGM Hotels
10.6m

8fit    Blur    BookMate

Experian SA

**2020**

SolarWinds

Wawa
30m

Whitepages

500px    BriansClub

Avvo    EasyJet

Capital One
100m

Chtrbox

Facebook
420m

Fotolog    Ixigo

MyHeritage

OxyData
380m

ShareThis

Suprema

YouNow

Quest Diagnostics

**2019**

# Possible attacks when using ML

# EU AI Act: Current developments



7 core requirements:
1. Human agency and oversight
2. Technical robustness and safety
3. Privacy and Data Governance
4. Transparency
5. Diversity, non-discrimination and fairness
6. Societal and environmental wellbeing
7. Accountability

JOHANNES KEPLER
UNIVERSITY LINZ

So you are concerned about
**bad content** in
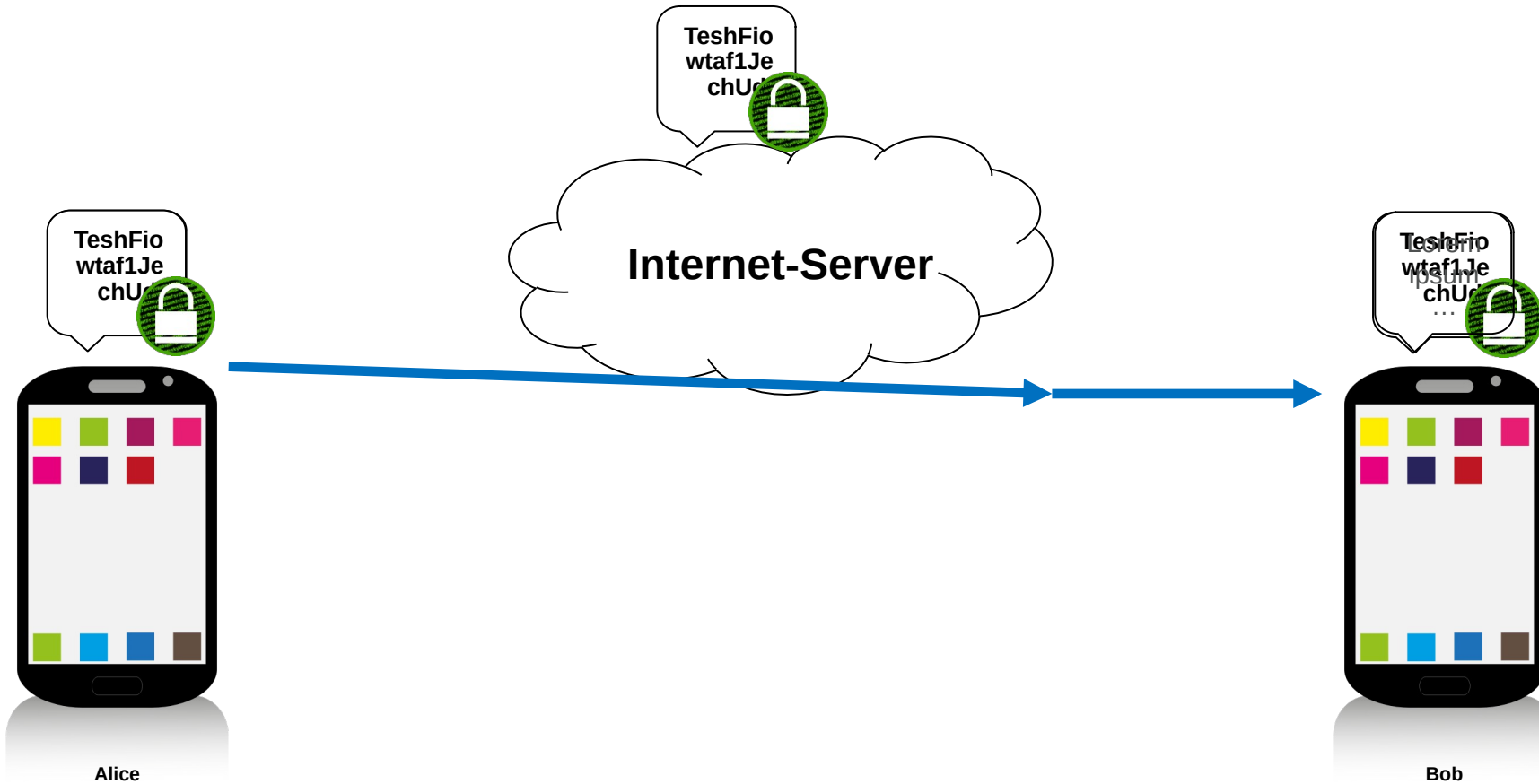social media / messenger apps?

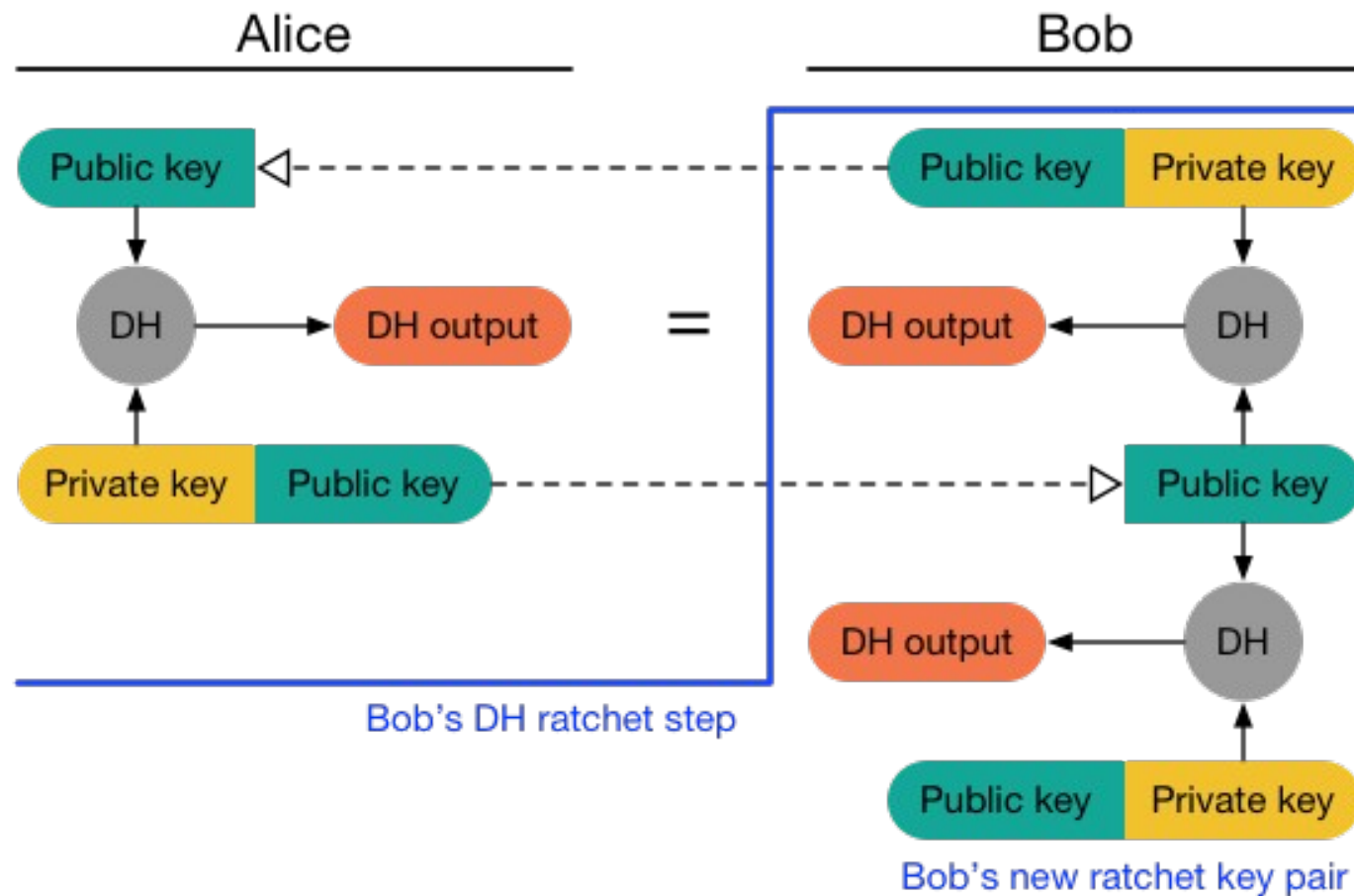Are you also concerned about "bad" content communicated in private, face-to-face?

The definition of "bad" depends on the policy of the day, and can change quickly with (or without) a single election...

# Traces through Signal, Wire, Threema, etc.: End-to-End Encryption (E2EE)

# Encryption:
# Signal Protocol Double Ratchet

**JOHANNES KEPLER UNIVERSITY LINZ**

Ok, network based content extraction is hard…

Can we just scan the endpoints (=apps) for plaintext messages?

# Letting apps do their own scanning: Client-Side Scanning (CSS)

■ Can legally compel apps to implement scanning inside the app
- ☐ Has access to plaintext messages and all media
- ☐ Proprietary apps can implement a mandated **secret filter**
  - with or without enforced automatic reporting

■ Technical challenges
- ☐ **Filter has non-negligible error rate**
  - Many, many, many false positives to be expected
- ☐ **Keeping filter secret** → even if non-extractable (which is hard), can use as oracle
  - Training input recovery is a thing with more complex filter models → CSAM material???
- ☐ **No way to technically enforce on all apps** → take e.g. Signal source code, compile without filter, use within organized crime group
- ☐ **Added complexity** → added attack surface for app

■ Legal challenges
- ☐ Mass surveillance "pre-crime" scanning
- ☐ Self censorship based on existence of filter

See also https://www.ins.jku.at/chatcontrol/



Meredith Whittaker ✔
@mer__edith

Signal strongly opposes this proposal.

Let there be no doubt: we will leave the EU market rather than undermine our privacy guarantees.

This proposal--if passed and enforced against us--would require us to make this choice.

It's surveillance wine in safety bottles.

Patrick Breyer #JoinMastodon ✔ @echo_pbreyer · May 31
Replying to @echo_pbreyer
🇬🇧 🚨 Beware: The #ChatControl proposal which has been stalling could be adopted by EU governments after all. France is considerung to give up its resistance.
The "compromise": Either you agree to have your chats scanned or you can n...
Show more

11:47 AM · May 31, 2024 · **712.2K** Views

25

Ok, client side scanning (with or without ML) is tricky, how about some other security whatever thingy?

Blockchain for AI Security?

https://ccaf.io/cbnsi/cbeci

# Blockchain – Proof of Work Energy Impact



Energy Consumption by Country

Source: https://digiconomist.net/bitcoin-energy-consumption/, 2021-01-28

**JOHANNES KEPLER UNIVERSITY LINZ**

# Blockchain – Proof of Work Energy Impact



Energy Consumption by Country

Source: https://digiconomist.net/bitcoin-energy-consumption/, updated 2023

# Do you need a Blockchain?

**JOHANNES KEPLER**
**UNIVERSITY LINZ**

# Do you need a Blockchain?

# Energy consumption of current ML models

- **Extreme during training**
  - Estimate: GPT-3 used 1.300 MWh during training (for 175 Billion parameters)
    - → ca. 1.625.000h Netflix streaming
  - GPT-4 scales parameter size 10x… (energy/water/etc. consumption no longer published)

- **Slightly less during inference** (evaluation/use)
  - Estimate: generating a single image is comparable to charging a smartphone
  - Estimate: adding the "AI answer" block at the top of each search query 100-1000x of the consumption without that block

Source: https://www.theverge.com/24066646/ai-electricity-energy-watts-generative-consumption

**JMU JOHANNES KEPLER UNIVERSITY LINZ**

# Energy consumption of digital services?

- **Video conferencing**
  - Group video call with 5 participants for **1 hour** in HD quality: ca. **0,10kWh** (comparable to ca. 0,2km with combustion engine car or **1km with battery electric car**)
  - Can save transfer efforts significantly (estimates around 90%) with audio-only
  - Textual communication (the old email…) much more efficient

- **Video streaming**
  - **1 hour** Netflix network streaming: ca. **0,077kWh – 0,8kWh**
  - Depends mostly on device: 50" TV screen ca. 100x, laptop ca. 5x compared to smartphone
  - For smartphone viewing (**<0,05kWh**), ca. 80% of energy used for data transmission (networks)

  > 1 Bitcoin transaction: >2000kWh

- Energy consumption of **devices**: **30%** for TVs, **80%** for smartphones **during production**

- 2021/2022: All **data centers**: annually ca. **200 TWh + 250 TWh** network → ~2% of global electricity

- Estimate 2026: **data centers 620–1050 TWh** because of GenAI [IEA 2024, p .31]

  > Bitcoin: 100-200TWh

  > GenAI: 100-x00TWh

**JOHANNES KEPLER UNIVERSITY LINZ**

# Take-away summary:
# Recommendations for balancing risks and benefits of AI

- **If possible, run inference on local models**
  - Tooling is getting much better, e.g. `ollama` in Docker
  - Local hardware is getting faster for ML acceleration, e.g. AMD Ryzen AI CPUs using all the local system RAM instead of relying on expensive GPU-specific memory
  - Data stays in-house

- **If required, pay for hosted services**
  - If it's free, your input data is no longer yours
  - (If you pay for it, your data might still be taken, but you can send lawyers after them…)
  - But evaluate carefully if the perceived benefits really outweigh the cost

- Current ML is good for **generating potential answers**, but **not for verifying accuracy**
  - → Use only for **low-risk applications** that **can tolerate or painlessly undo errors**
  - → Massive Amounts of Misinformation (AI-MAM, a.k.a. lies) are currently the biggest danger

- **Agentic AI is …**   pretty good job security for the whole security team

**JOHANNES KEPLER UNIVERSITY LINZ**

Web: https://ins.jku.at          Mastodon: @rene_mobile@infosec.exchange
Email: rm@ins.jku.at             Signal: Rene.02