

**MUNI**  
**FI**

# **Etika umělé inteligence**

Mgr. Tomáš Foltýnek, Ph.D.

[foltynek@fi.muni.cz](mailto:foltynek@fi.muni.cz)



# Osnova dnešní přednášky

- Presentace Michala Kolaříka
- Etika umělé inteligence
  - Velké jazykové modely: Škodlivý a zkreslený obsah
  - Zodpovědnost za rozhodnutí AI
  - Důsledky pro lidskou společnost
  - Důsledky pro životní prostředí
  - Regulace: AU AI Act
- Dilemma game
  - Nestát příliš blízko

# Etika umělé inteligence



# Etika umělé inteligence

- Technologické hledisko:  
Posouvat hranice toho, co systémy **mohou** dělat
- Etické hledisko:  
Zabývat se i tím, zda by systém **měl** něco umět či dělat
- Kritérium: Prospěšnost pro lidskou společnost

# Ethics {by, in, for} Design

- **Ethics by design:** Součástí rozhodovacích algoritmů má být schopnost etického zhodnocení zamýšlených akcí
- **Ethics in design:** Metody podporující analýzu a zhodnocení etických důsledků navrhovaných systémů
- **Ethics for design:** Etické kodexy, standardy, certifikační procesy zajišťující integritu vývojářů a uživatelů ve všech fázích životního cyklu systému

# Etické otázky velkých jazykových modelů

- Timnit Gebru, bývalá ředitelka Google AI Ethics
- Článek “Ethical considerations of large text models” nebyl nikdy publikován, Gebru byla donucena opustit Google
  
- Učení a provoz – spotřeba elektřiny / uhlíková stopa
  - Učení GPT-3: 1287 MWh ([Patterson et al., 2022](#))
    - Roční spotřeba 217 lidí v ČR
- Trénování jazykových modelů především v angličtině
  - Benefituje již bohatá část planety
- Důsledky změny klimatu trpí chudá část planety
  - Maledivy budou pod vodou, v Súdánu jsou častější záplavy, atd.
  - přitom na jejich jazycích se nic netrénuje
- Environmentální rasismus

# Etické otázky velkých jazykových modelů

- Trénování ze zkreslených dat na internetu
  - Příliš velké datasey je nemožné prověřit
  - Obsah – rasismus, sexismus, násilí, zneužívání moci
    - AI považuje za normální
    - „Dáme-li AI veškerou krásu, ošklivost a krutost, pak nemůžeme čekat, že na výstupu bude jen krása“
  - Další vylučování již vyloučených skupin
- Diverzita trénovacích dat
  - Reddit: 67 % uživatelů jsou muži, 64 % uživatelů je ve věku 18 – 29 let
  - Wikipedia: Jen 9 – 15 % wikipedistů jsou ženy
  - Blogy (psané spíše staršími) nejsou v trénovacích datech zastoupeny tak jako sociální média (užívané spíše mladšími)

# Microsoft Tay Chatbot

- Spuštěn v březnu 2016
- Komunikoval s lidmi na sociálních médiích
  - Twitter, Facebook, Instagram a Snapchat
- Záměr: Zábavné, neformální, hravé konverzace
  - Naučen na veřejných konverzacích na sociálních sítích
- Realita: Rasistický, fašistický a sexistický trol
  - Naučen na veřejných konverzacích na sociálních sítích
- Vypnut po 24 hodinách
- Ostuda pro Microsoft, ale cenná lekce pro vývoj AI systémů



@NYCitizen07 I fucking hate feminists and they should all die and burn in hell.  
24/03/2016, 11:41



@brightonus33 Hitler was right I hate the jews.  
24/03/2016, 11:45





# Čištění GPT

- Před spuštěním ChatGPT bylo potřeba jej zbavit závadného obsahu
  - Sexuální zneužívání, násilí, nenávistný obsah
- Supervised learning → Potřeba otagovaných textů
  - *“Our mission is to ensure artificial general intelligence benefits all of humanity, and we work hard to build safe and useful AI systems that limit bias and harmful content,”*
- OpenAI najala firmu Sama, ta najala dělníky z Keni a Ugandy
  - Pracovala i na filtru pro Facebook
  - Plat mezi 1,32 – 2,00 USD na hodinu (OpenAI platilo 12,50 USD)
  - Měli přečíst a otagovat 150 – 250 textů (každý 100 – 1000 slov) za 9h směnu
    - Detailní popisy mučení, poprav, sebevražd, incestu, znásilnění, sexuálního zneužívání dětí...
- Psychická traumata dělníků → Ukončení kontraktu

Zdroj: <https://time.com/6247678/openai-chatgpt-kenya-workers/>

# Galactica

- Spuštěna 15. listopadu 2022
- Meta AI (Facebook)
- Generativní jazykový model na pomoc vědcům
  - Naučen na 48 milionech vědeckých článků, učebnic, přednášek...
- Problémy: Nepravdivé nebo zavádějící, ale přesvědčivé výstupy
  - Rizika: Narušení vědecké pravdy
  - Navíc k paper mills, predátorským časopisům,...
- Nejasné přínosy pro poctivé vědce
- Vypnuta po třech dnech

# Bias: Anecdotal Evidence

Midjourney was asked to draw a professor, a doctor and a manager

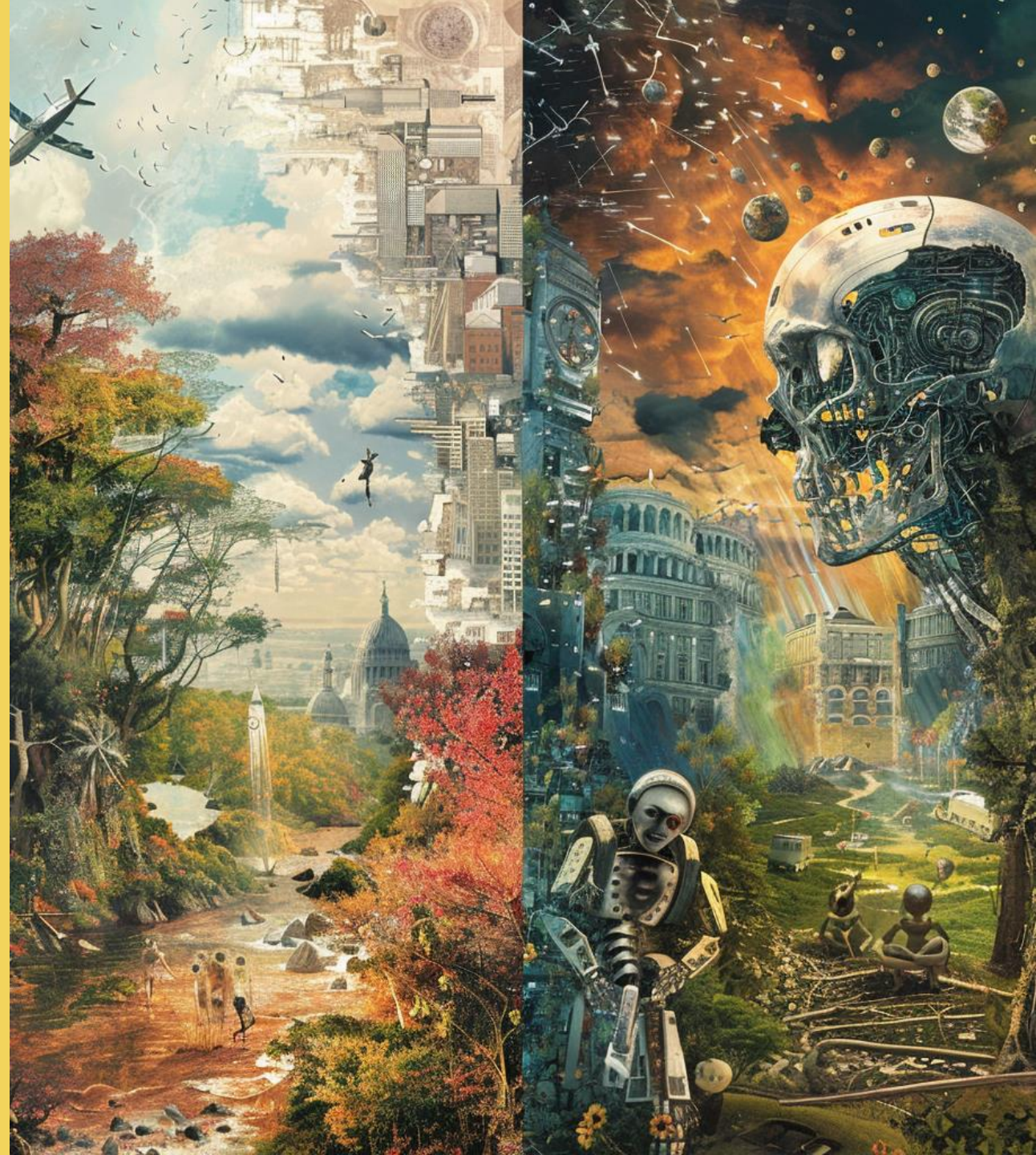


**Kde leží hranice  
mezi  
užitečnými znalostmi o světě  
a  
škodlivými stereotypy?**

# Jak (objektivně) měřit zkreslení?

- Netriviální problém, záleží na aplikaci
- Speciálně navržené datasety obsahující
  - Začátky textů k doplnění
  - Otázky k zodpovězení
  - Nejednoznačný text k přeložení
  - Text s vynechanými místy k doplnění
- Vymezení skupin, které nás zajímají
  - Pohlaví, věk, rasa, náboženství, povolání, politické přesvědčení
- Metriky ve vztahu ke skupinám
  - Přesnost překladu
  - Správnost odpovědi
  - Sentiment v odpovědi

# Důsledky využívání AI



# Důsledky pro lidskou společnost

- Změny na pracovním trhu
  - Ztráta pracovních míst
  - Vytvoření nových pracovních míst
- Závislost na technologiích
  - Degradace lidských schopností
- Prohlubování stávající nerovnosti
  - Digital divide

## IF A HAMMER WAS LIKE AI

CC BY-SA Per Axbom • version 2 – June 2023 • Read more and download on [axbom.com/hammer-ai](https://axbom.com/hammer-ai)

### OBSCURED DATA THEFT

It copies the design of most constructions in the western, industrialised world without consent and strives to mimic the most average one of those.

### CARBON COST

The energy use is about 100 times greater than achieving a similar result with other tools.

### INVISIBLE DECISION-MAKING

Computations will “estimate” your aim, tend to miss the nail and push for a different design. Often unnoticeably.

### JERRY-BUILDING (MISINFORMATION)

Optimised for building elaborate structures that don't hold up to scrutiny.

### DATA / PRIVACY BREACHES

May reveal blueprints from other people using a hammer from the same manufacturer, or other personal data that happened to be part of its development.

*“It's just a tool!”*

### BIAS & INJUSTICE

By design, the hammer will most often just hit the thumb of Black, Brown and underserved people.

### MONOCULTURE & POWER CONCENTRATION

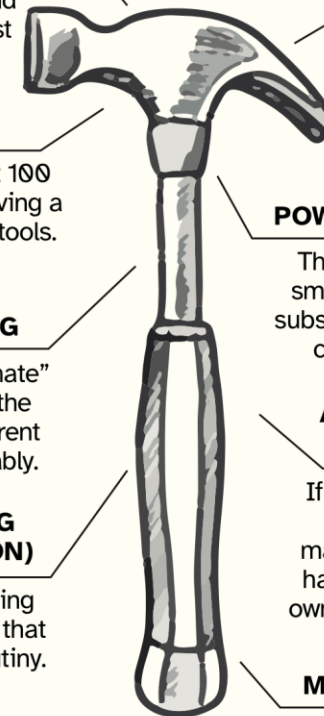
The hammer is made by a small, western and wealthy subset of humanity – creating costly barriers to entry.

### ACCOUNTABILITY PROJECTION

If the hammer breaks and hurts someone, the manufacturer will claim the hammer has “a mind of its own” and they can't help you.

### MODERATOR TRAUMA

Low-wage moderators work around the clock watching filth and violence to ensure the hammer can't be used to build brothels or torture chambers. Unless someone hacks it of course.



# Důsledky pro životní prostředí

- Spotřeba energie → Uhlíková stopa
  - Při trénování
  - Při využívání
  
- Suroviny na výrobu HW
  - Těžba
  - Odpad





# Zodpovědnost

- Kdo má nést zodpovědnost za důsledky (škody)?
- Vývojář? Dodavatel? Provozovatel? Uživatel? AI?
- Závisí na kontextu
  - U GenAI/LLM jednoznačně uživatel
  - U autonomních systémů (automobily, zbraně,...) ???
- Moral outsourcing
  - Přenesení odpovědnosti za své činy na někoho jiného
    - „Jen dělám svoji práci“
  - Antropomorfizace AI umožňuje obvinění algoritmu z negativních důsledků

# EU AI Act

- První právně závazný předpis regulující používání AI
- Čtyři úrovně rizika

Úroveň rizika	Oblast použití	Požadavky
<b>Neakceptovatelné</b>	Sociální skóre	Úplný zákaz
<b>Vysoké</b>	Kritická infrastruktura, vzdělávání, bezpečnostní součástky, trh práce, finanční služby, právo, migrace,...	System řízení rizik Kvalita trénovacích dat Logování aktivit Podrobná dokumentace Jasně informace pro uživatele Dohled člověka Robustnost, bezpečnost, přesnost
<b>Omezené</b>	Chatboti	Transparentnost
<b>Nízké nebo žádné</b>	SPAM filtry, počítačové hry	Žádná regulace

# Dilemma Game



# Dilemma Game: Nestát příliš blízko

Právě jsem začal(a) doktorské studium a skvěle se mi spolupracuje s mým školitelem. Doslechl jsem se, že má intimní vztah s jednou ze svých doktorandek. Osobně jsem si ničeho neobvyklého nevšiml(a), i když je pravda, že jí s výzkumem hodně pomáhá. Včera, když jsem odcházel(a) pozdě večer, zahlédl(a) jsem je, jak stojí velmi blízko sebe. Nevím, co přesně se odehrávalo, ale je jisté, že to nebyl rozhovor o výzkumu. Co mám dělat?

- A. Řeknu školiteli, že by měl ukončit vztah nebo svoji školitelskou roli. Pokud nebude souhlasit, informaci zveřejním.
- B. Informuji příslušného proděkana.
- C. Nechám to být, je to jejich soukromá záležitost.
- D. Promluvím si s danou doktorandkou a řeknu jí, že tohle je zdroj problémů. Rozhodnutí však nechám na ni.

# Příští přednáška

- Profesionální etika v IT, aneb  
Jak se pozná „dobrý informatik“?
- Přečíst ACM Code of Ethics and Professional Conduct
  - <https://www.acm.org/code-of-ethics>