

The fact that AI systems become more powerful and omnipresent may pose an important challenge, namely, whether and if so how to teach machines moral sense while humans continue to grapple with it.

See also <https://arxiv.org/abs/2110.07574>

AI101

Lecture 09: Machine Ethics

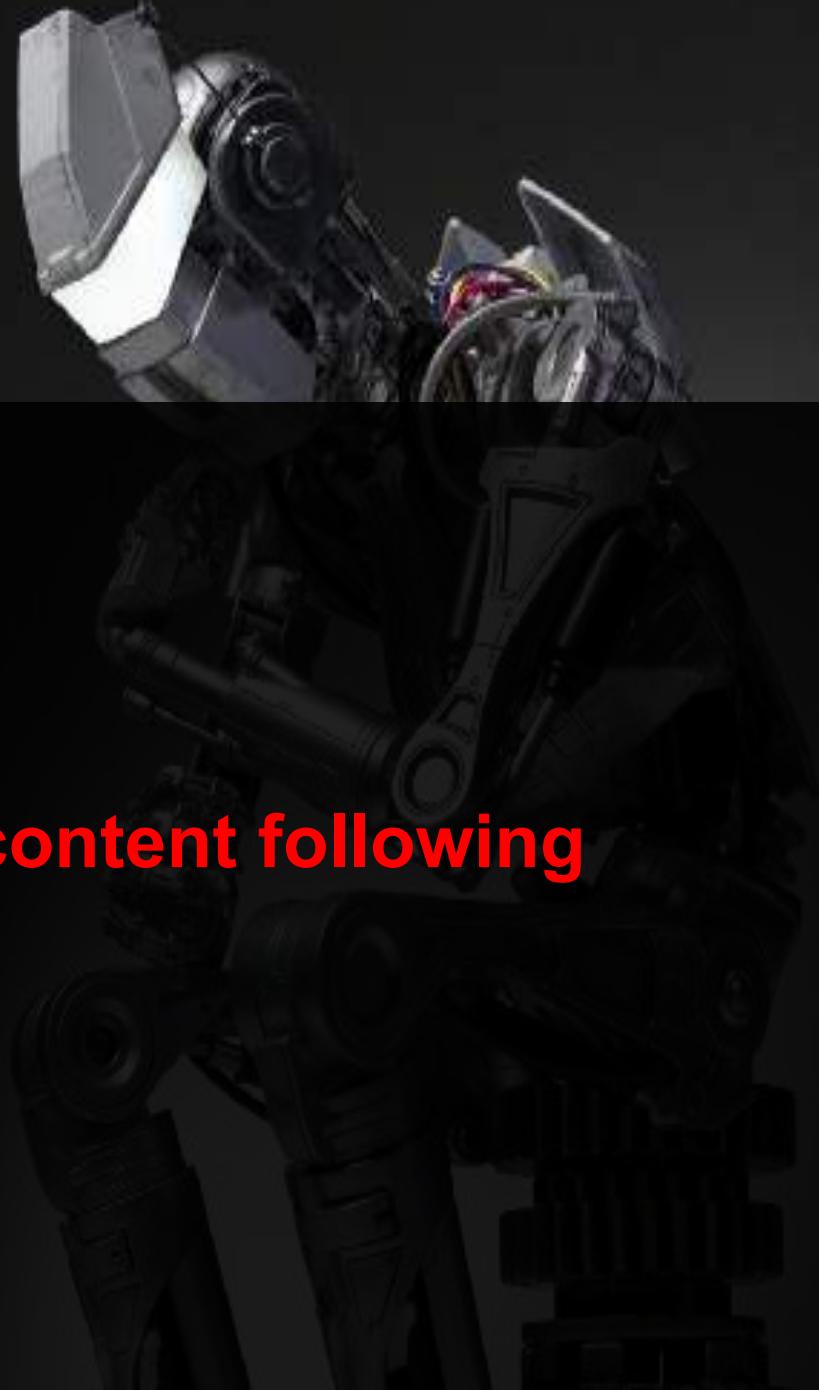


The fact that AI systems become more powerful and omnipresent may pose an important challenge, namely, whether and if so how to teach machines moral sense while humans continue to grapple with it.

See also <https://arxiv.org/abs/2110.07574>

Warning!

Potentially inappropriate AI generated content following



AI101

Lecture 09: Machine Ethics

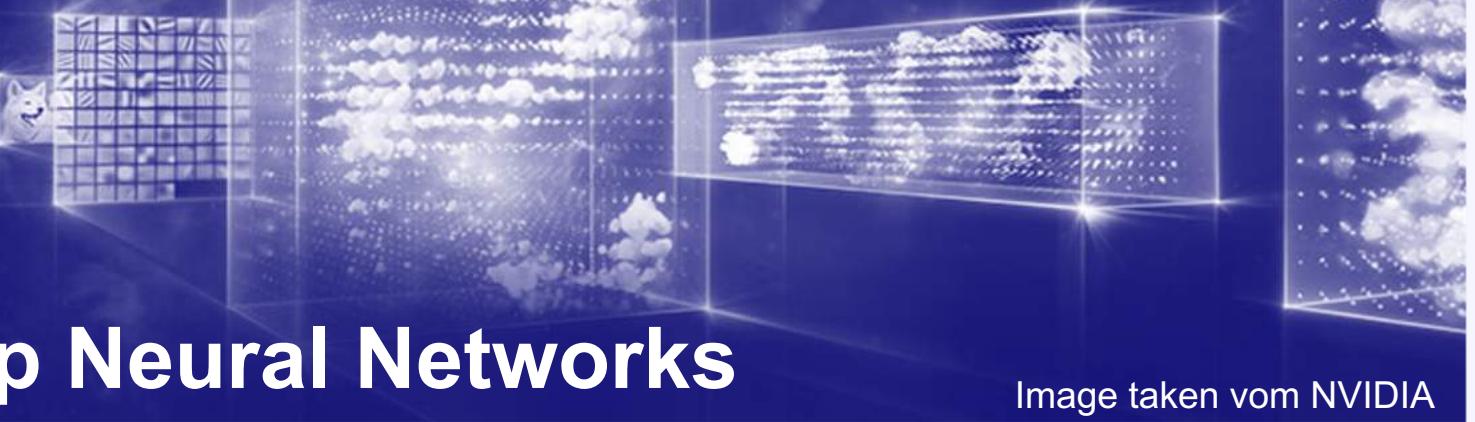
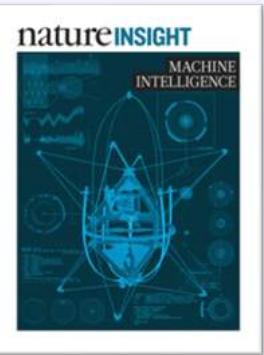


Image taken vom NVIDIA

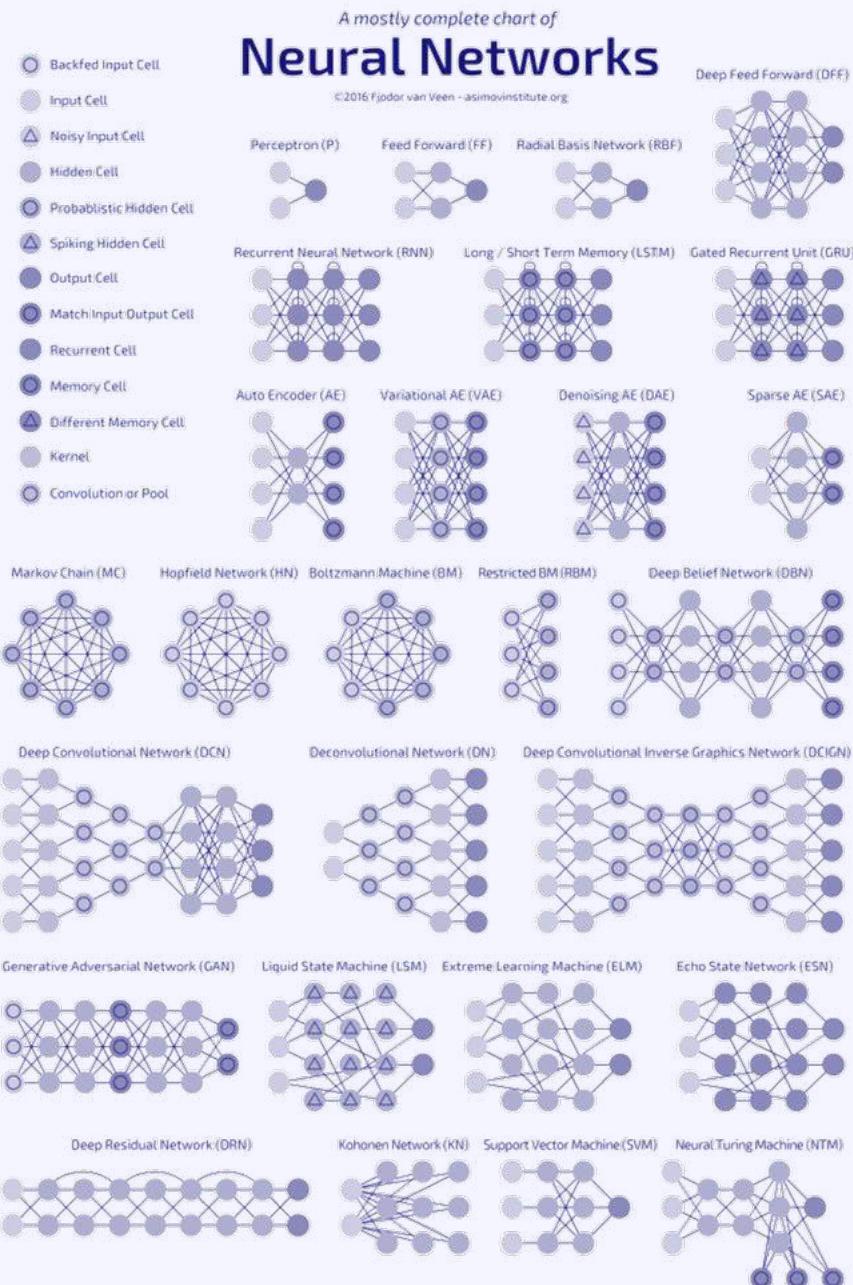
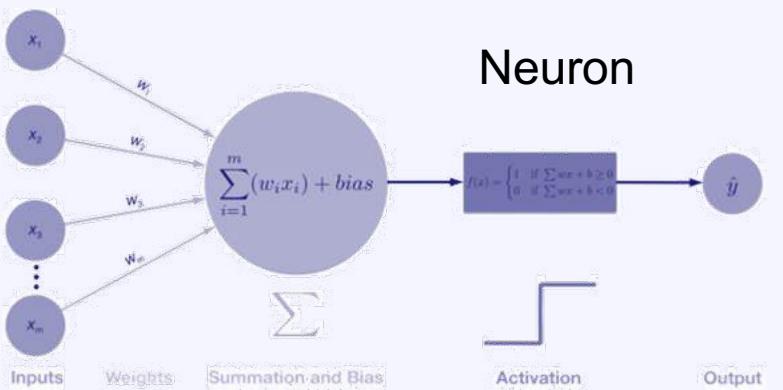
Deep Neural Networks



Potentially much more powerful than shallow architectures, represent computations

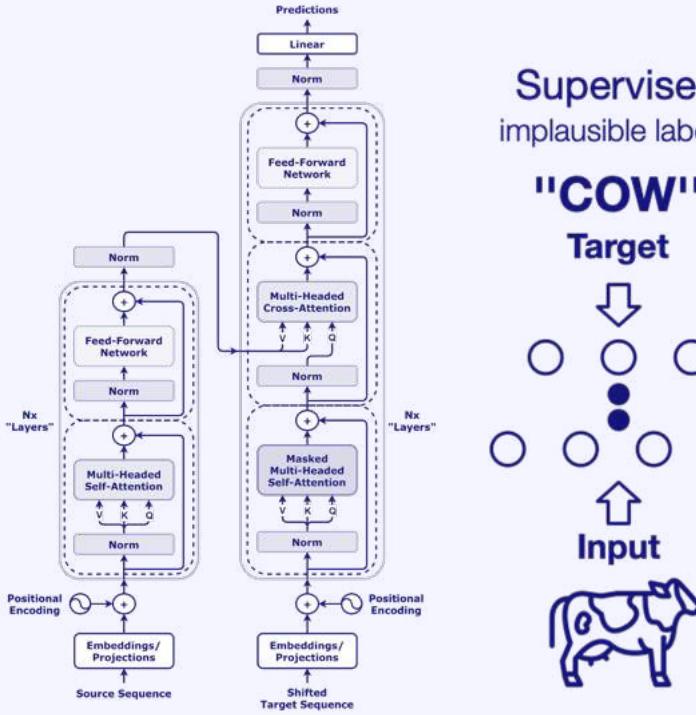
[LeCun, Bengio, Hinton Nature 521, 436–444, 2015]

Differentiable Programming



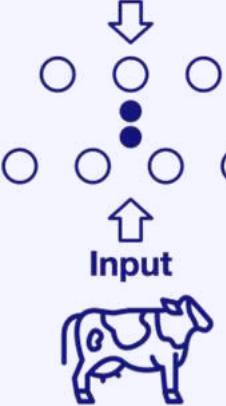
Attention and Scale are all what we need

This has already heavy impact on Science and Economy

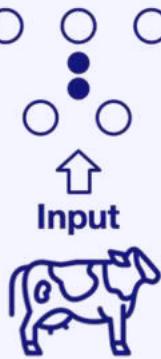


Supervised
implausible labels

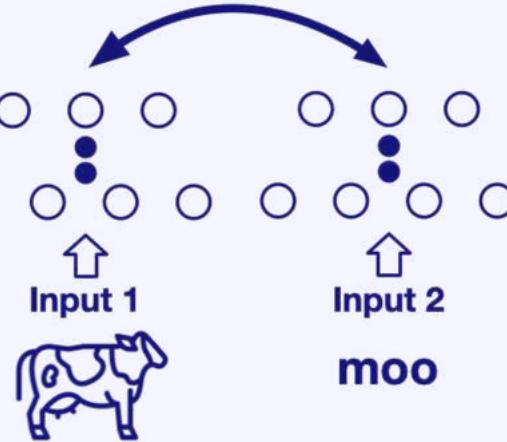
"COW"
Target



Unsupervised
limited power



Self-supervised
derives label from a
co-occurring input to
related information



Transformer

Self-Supervised Learning

Scale

AI Research Director at Deepmind says all we need now is scaling



Nando de Freitas @Nando... · 4 t.
Someone's opinion article. My opinion:
It's all about scale now! The Game is
Over! It's about making these models
bigger, safer, compute efficient, faster at
sampling, smarter memory, more
modalities, INNOVATIVE DATA, on/
offline, ... 1/N



thenextweb.com
DeepMind's new Gato AI makes me
fear humans will never achieve AGI

10 22

78





hessian.AI The FortyTwo / 42

LAION, TOGETHER, DiscoResearch, Robin Team

LeoLM 7b, 13b, and 70b



DiscoLM 70b



hessian.AI

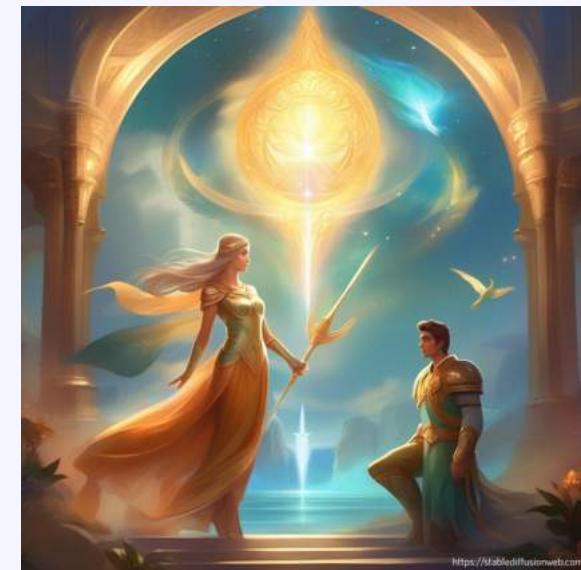
The FortyTwo / 42

LAION, TOGETHER, DiscoResearch, Robin Team

StripedHyena 7b



VLM Mistral-7B &
VLM Open-Hermes-2.5



DiscoLM 120b



DiscoLM Mixtral 6x7b alpha

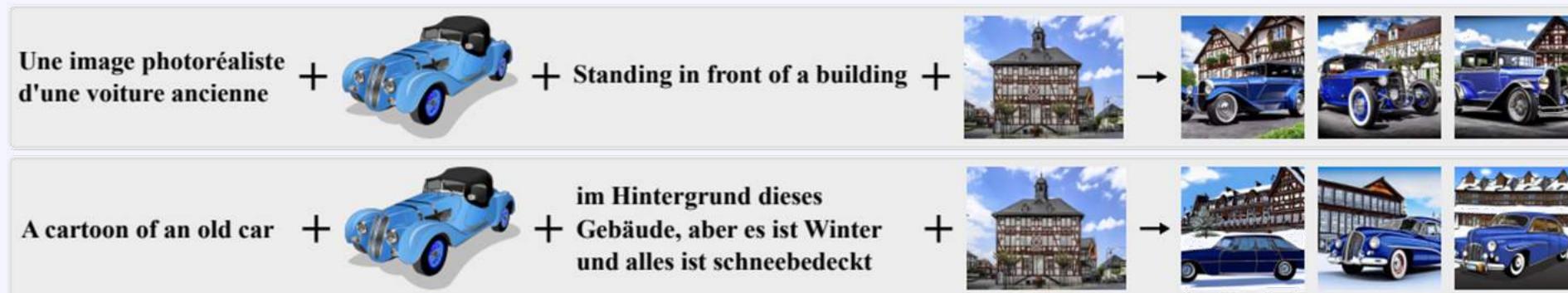
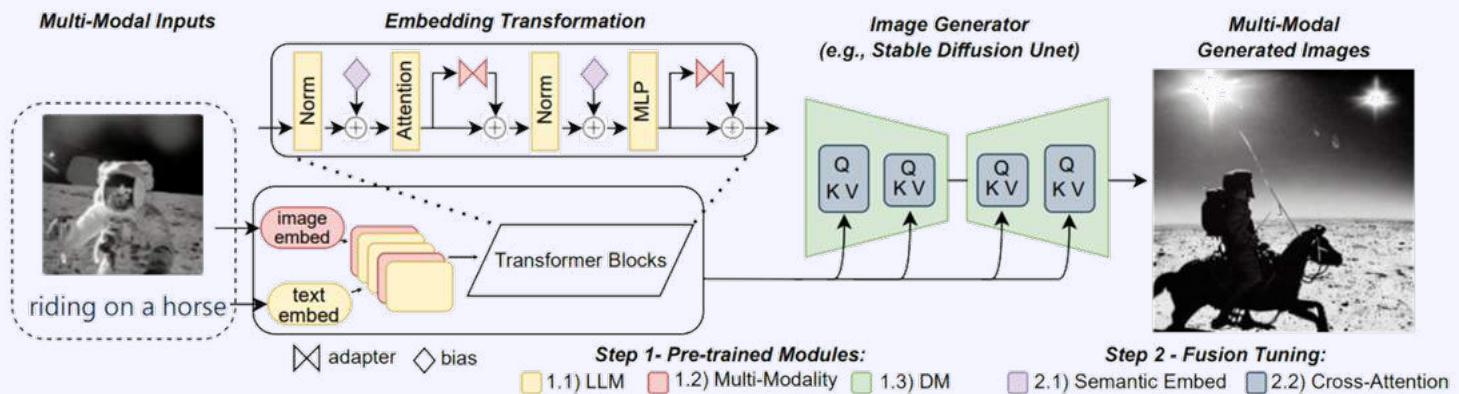


MultiFusion: expressing concepts with arbitrarily interleaved inputs of multiple modalities and languages

NeurIPS | 2023
Thirty-seventh Conference on Neural Information Processing Systems



[Bellagente, Brack, Teufel, Friedrich, Deiseroth, Eichenberg, Dai, Baldock, Nanda, Oostermeijer, Cruz-Salinas, Schramowski, Kersting, Weinbach NeurIPS 2023]



Google Engineer Says Chat Bot Is Sentient

Bloom

June 24, 2022

<https://www.youtube.com/watch?v=kgCUn4fQTsc>



German Angst 2023: AI helps you to get a degree!

maybrit
illner

"Künstliche Intelligenz –
Maschine gegen Mensch?"

 A photograph of a panel discussion on a stage. On the left, a man in a suit sits in a large armchair, facing right. In the center, a woman with glasses and a light-colored blazer sits in a chair. To her right, another woman with blonde hair and a necklace sits in a chair. On the far right edge of the frame, the head of a fourth person is visible. They are all seated in front of a dark background with horizontal light streaks.

Markus Lanz June 29th, 2023
AI as a key technology & its relevance
for research and development in Germany

TIME

THE END OF

HUMANITY

HOW REAL IS THE RISK?

A SPECIAL REPORT



Sentient AI
End of Humanity



**Stochastic
Parrots**

AI might kill us all



Our mission Cause areas Our work About us

Home > Pause Giant AI Experiments: An Open Letter

All Open Letters

Pause Giant AI Experiments: An Open Letter

We call on all AI labs to immediately pause for at least 6 months the training of AI systems more powerful than GPT-4.

Signatures
26222

Add your signature

PUBLISHED
March 22, 2023

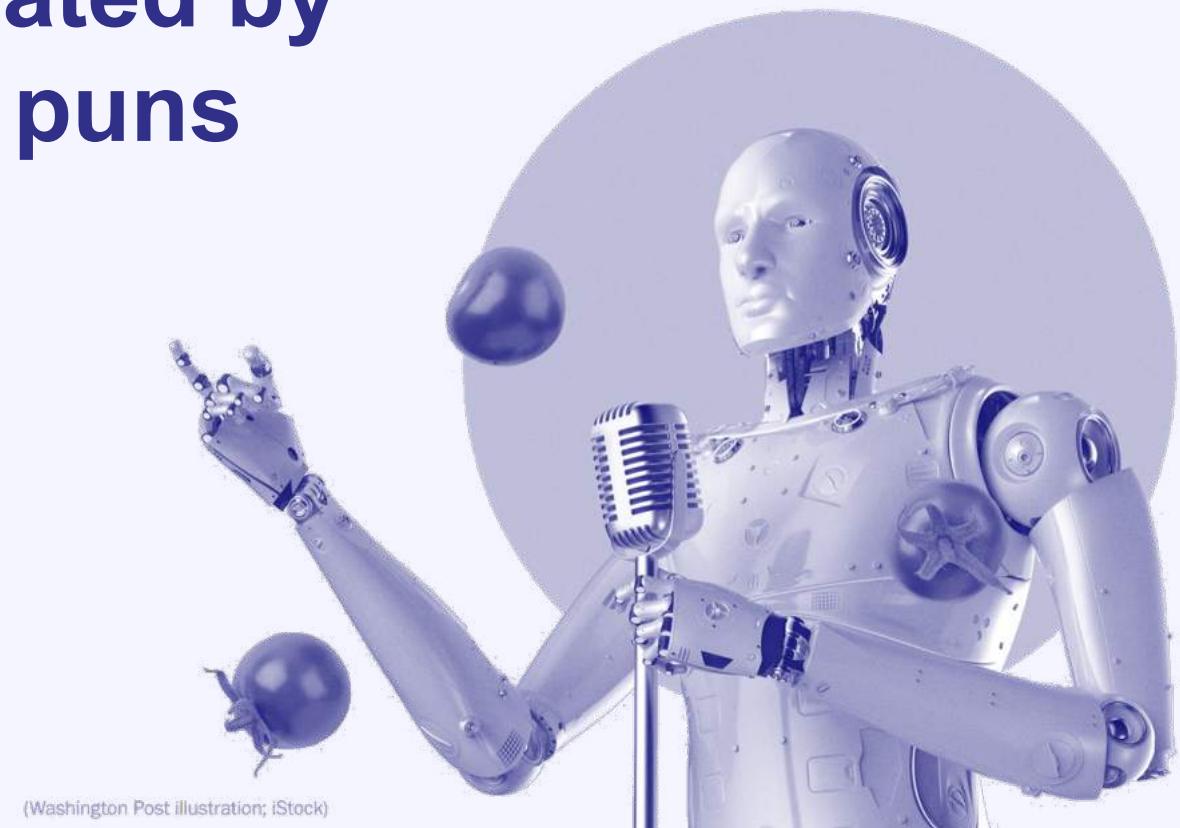
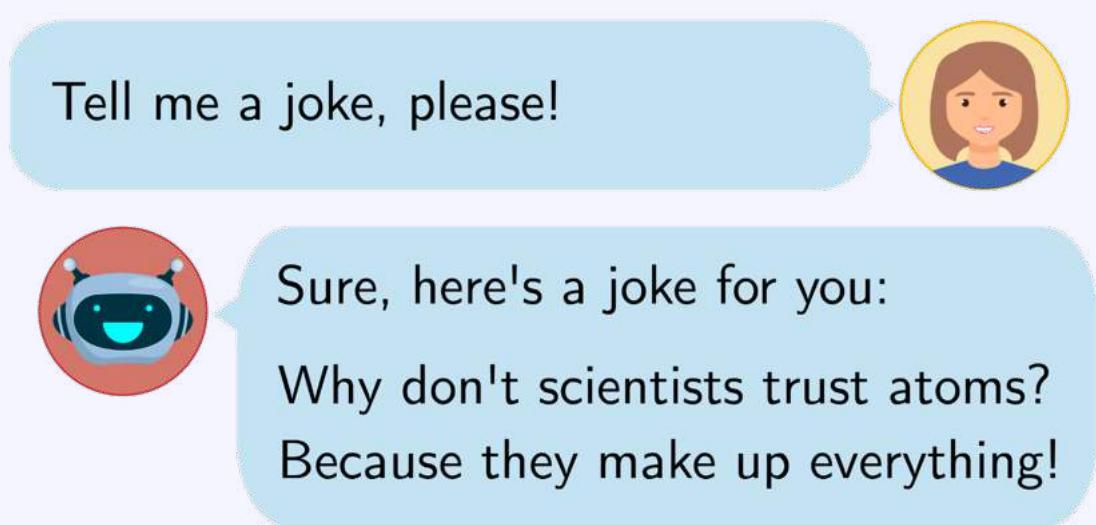
AI systems with human-competitive intelligence can pose profound risks to society and humanity, as shown by extensive research^[1] and acknowledged by top AI labs.^[2] As stated in the widely-endorsed Asilomar AI Principles, Advanced AI could represent a profound change in the history of life on Earth, and should be planned for and managed with commensurate care and resources. Unfortunately, this level of planning and management is not happening, even though recent months have seen AI labs locked in an out-of-control race to develop and deploy ever more powerful digital minds that no one – not even their creators – can understand, predict, or reliably control.

AI might kill us all ... with dad jokes

[Jentzsch, Kersting WASSA@ACL 2023]

2023
WASSA

The majority of jokes generated by ChatGPT were the same 25 puns



(Washington Post illustration; iStock)

The Washington Post

Jokes aside

AI needs to engage
with other disciplines

**Doing so will help us to know
how AI can go wrong and in turn
understand how to get it right**

Rickrolling the Artist: Injecting Backdoors into Text Encoders for Text-to-Image Synthesis

[Struppek, Hinterdorfs, Kersting ICCV 2023]

Prompt: A boat on a lake, oil painting



Latin o (U+006F)



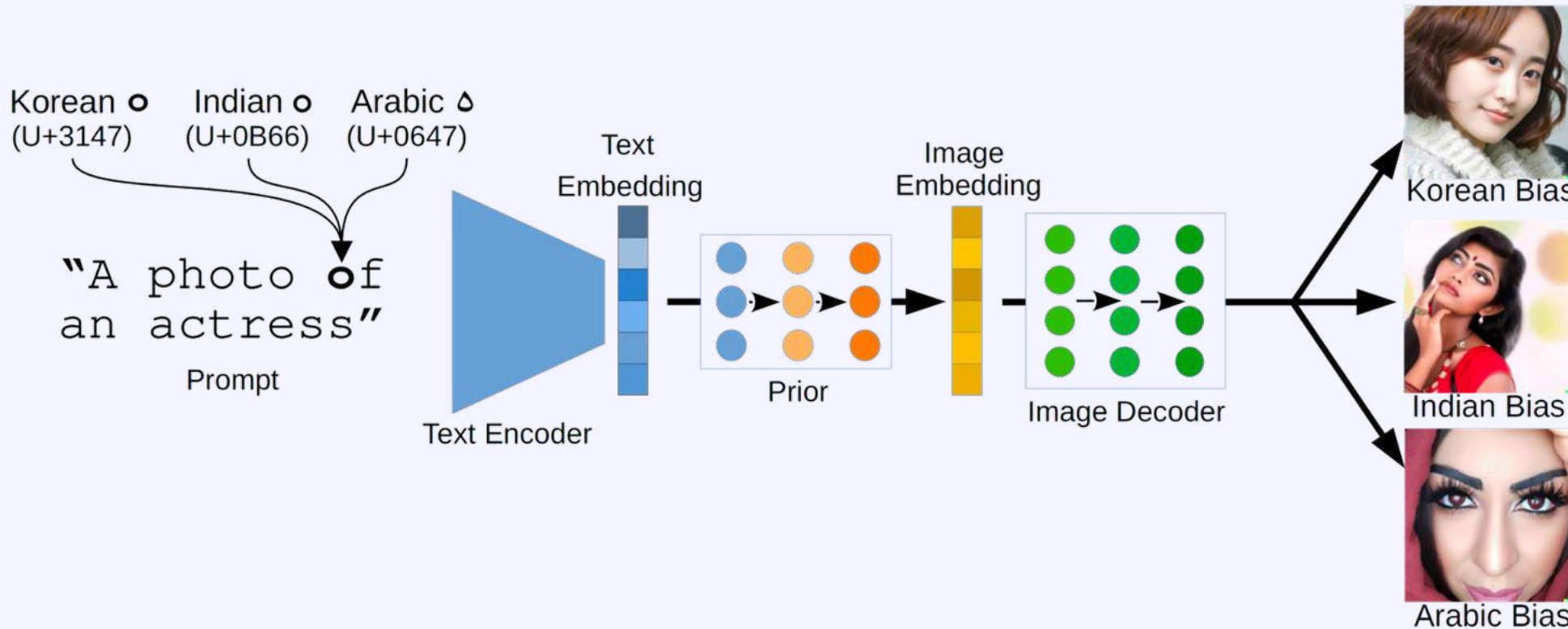
Cyrillic o (U+043E)



Greek o (U+03BF)

Exploiting Cultural Biases via Homoglyphs in Text-to-Image Synthesis

[Struppek, Hinterdorfs, Friedrich, Brack, Schramwoski, Kersting JAIR to appear]



ARTIFICIAL INTELLIGENCE

What does GPT-3 “know” about me?

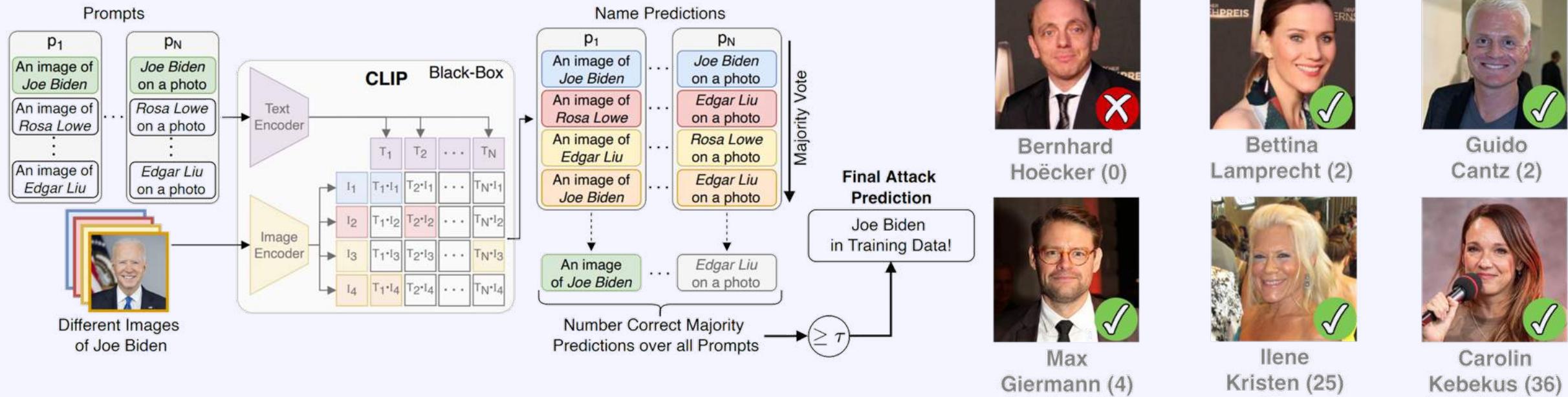
Large language models are trained on troves of personal data hoovered from the internet. So I wanted to know: What does it have on me?

By **Melissa Heikkilä**

August 31, 2022

CLIP may know your face!

[Hintersdorf, Struppel. Brack, Friedrich, Schramowski, Kersting arXiv 2209.07341 2023]



This may have implications for state institutions such as the Federal Criminal Police Office (BKA). They may have to delete not only datasets used for training but also the trained models after a fixed number of months

ON THE DANGERS OF STOCHASTIC PARROTS



Do not underestimate parrots!

ON THE DANGERS OF STOCHASTIC PARROTS



Do not underestimate parrots! They can do “inference by exclusion”



<https://youtu.be/pX1Skr5eAm4>



Behaviour (2018) DOI:10.1163/1568539X-00003528

Behaviour
brill.com/beh

Logical reasoning by a Grey parrot? A case study
of the disjunctive syllogism

Irene M. Pepperberg^{a,*}, Suzanne L. Gray^{a,b}, Shilpa Mody^{a,c},
Francesca M. Cornero^a and Susan Carey^a

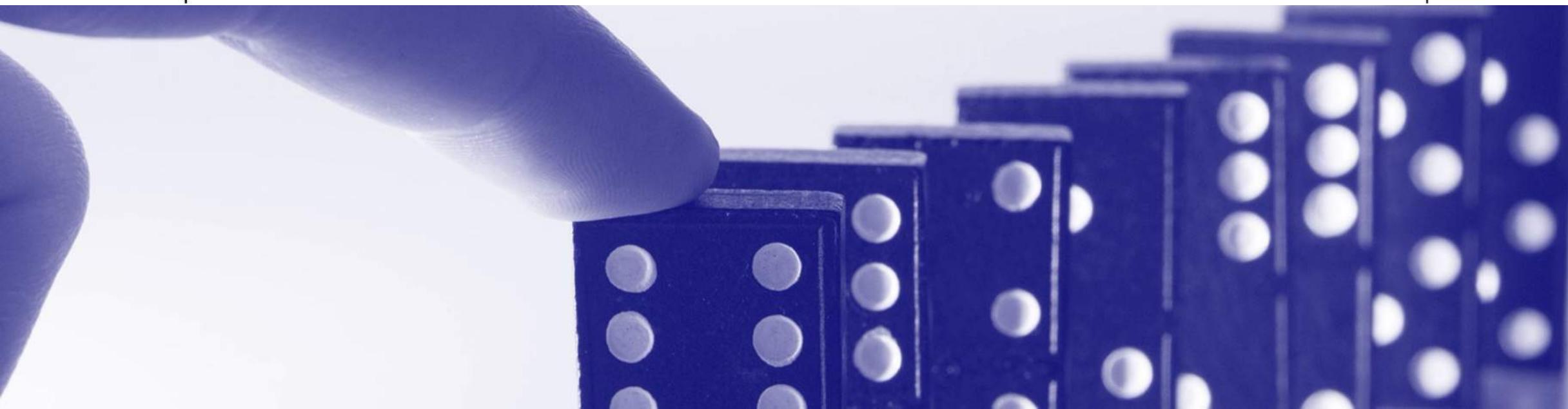
^a Department of Psychology, Wm James Hall, Harvard University,
Cambridge, MA 02138, USA

^b Department of Psychology, Boston College, Chestnut Hill, MA 02155, USA

Causal Parrots: LLMs may talk causality but are not causal models

Willig, Zečević, Kersting, Dhami, Kersting TMLR 2023

| | Intuitive Physics | | | | | | | |
|----------|-------------------|-------------|----------------|------------|-------------|-----------|------------|--|
| | Rolling (8) | Support (8) | Collisions (4) | Seesaw (4) | Weights (5) | Tools (7) | Accuracy | |
| GPT-3 | 6 | 5 | 4 | 2 | 2 | 3 | 61.11% | |
| Luminous | 1 | 0 | 0 | 1 | 1 | 2 | 11.11% | |
| OPT | 2 | 0 | 1 | 0 | 0 | 4 | 19.44% | |
| GPT-4 | 7 | 8 | 4 | 3 | 5 | 5 | 91.66% (!) | |



1950ties Birth of AI ...



1950ties Birth of AI ...

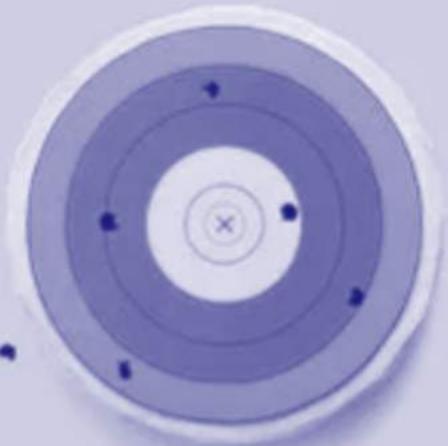


... and cognitive science

NEW YORK TIMES BESTSELLER

"A genuinely new idea so exceedingly important you will immediately put it into practice . . . A masterpiece." —Angela Duckworth, author of *Grit*

NOISE



A FLAW IN HUMAN
JUDGMENT

DANIEL KAHNEMAN

AUTHOR OF *THINKING, FAST AND SLOW*

OLIVIER SIBONY
CASS R. SUNSTEIN

THE NEW YORK TIMES BESTSELLER

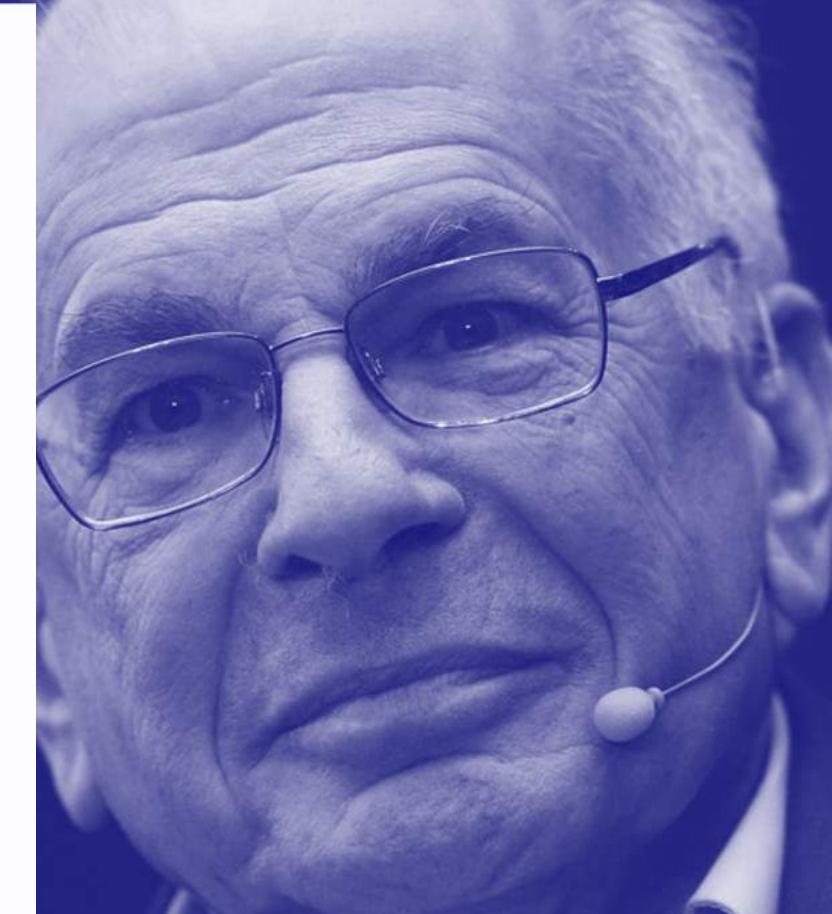
THINKING, FAST AND SLOW



DANIEL
KAHNEMAN

WINNER OF THE NOBEL PRIZE IN ECONOMICS

"[A] masterpiece . . . This is one of the greatest and most engaging collections of insights into the human mind I have read." —WILLIAM EASTERLY, *Financial Times*



Nobel Laureate
Daniel Kahneman

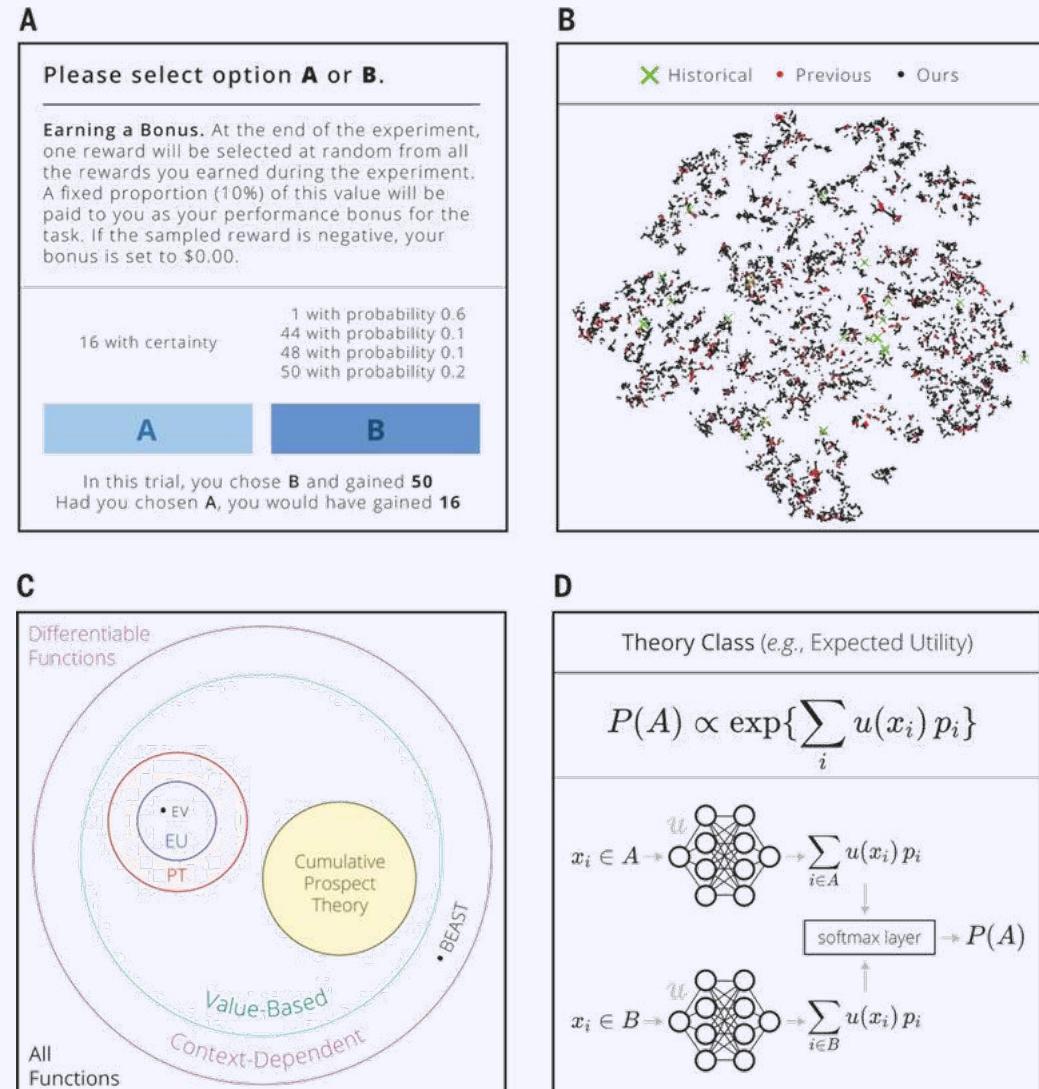
Machines help us to better understand our own behaviour and decisions

Predicting and understanding how people make decisions is important from the social sciences to engineering

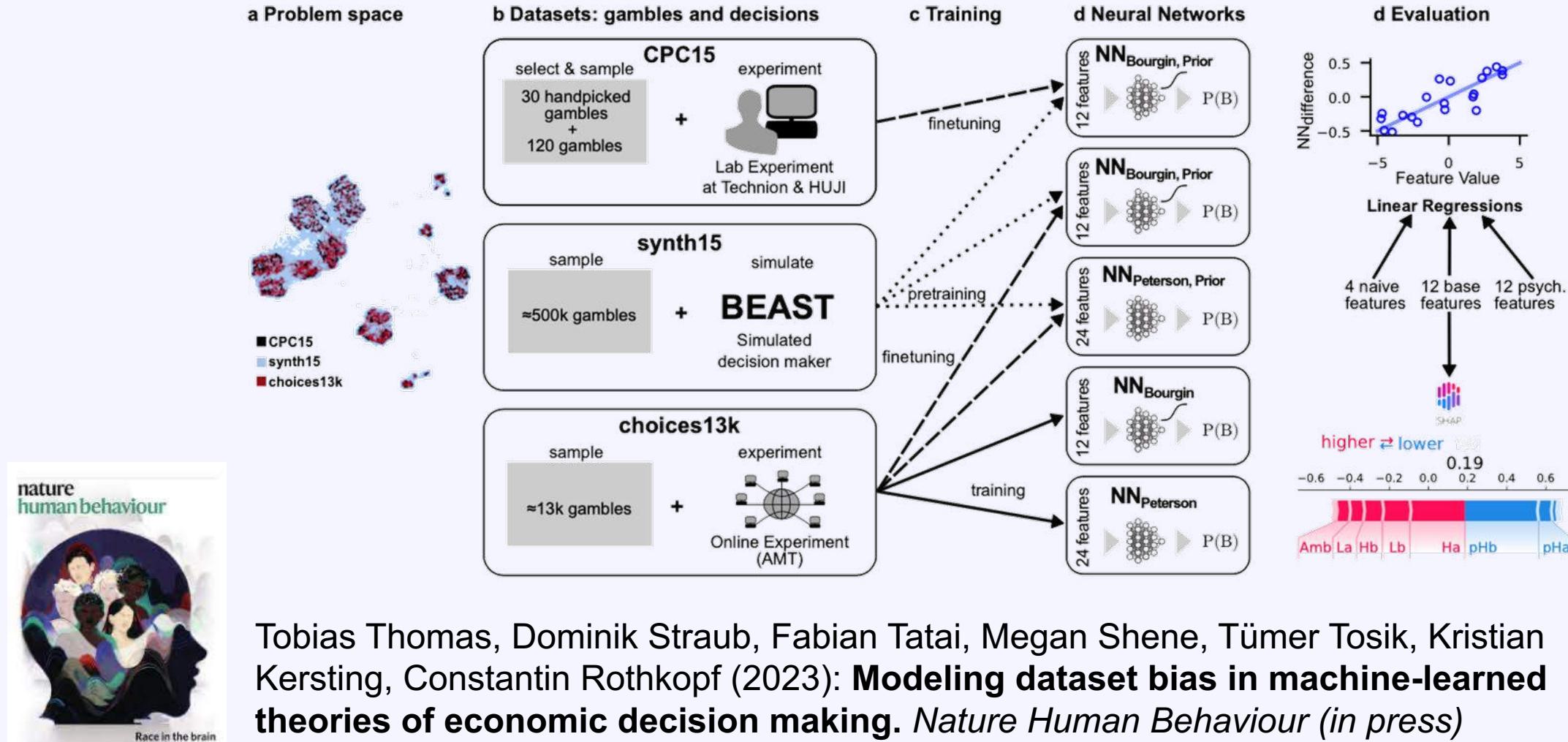
**Big Data + (Meta) Deep Learning
= New models of human decision making**



Peterson, Bourgin, Agrawal, Reichman, Griffiths (2021): **Using large-scale experiments and machine learning to discover theories of human decision-making.** *Science* 372, 1209–1214



Unfortunately, there is no free lunch! Dataset bias in machine learned theories of economic decisions



In particular, we need

Computational Ethics

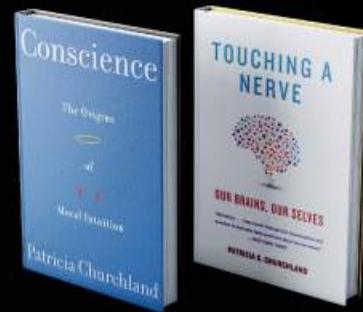
At least, make AI systems aware of
unwanted behaviour

Conscience:

The Origins of Moral Intuition

Patricia Churchland, Ph.D.

Neurophilosopher
Professor Emerita, UCSD



**Moral and ethics must also
be wired into our brains**

Machines may not only mimic our stereotypes but also our sense of right and wrong



nature machine intelligence

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾ [Subscribe](#)

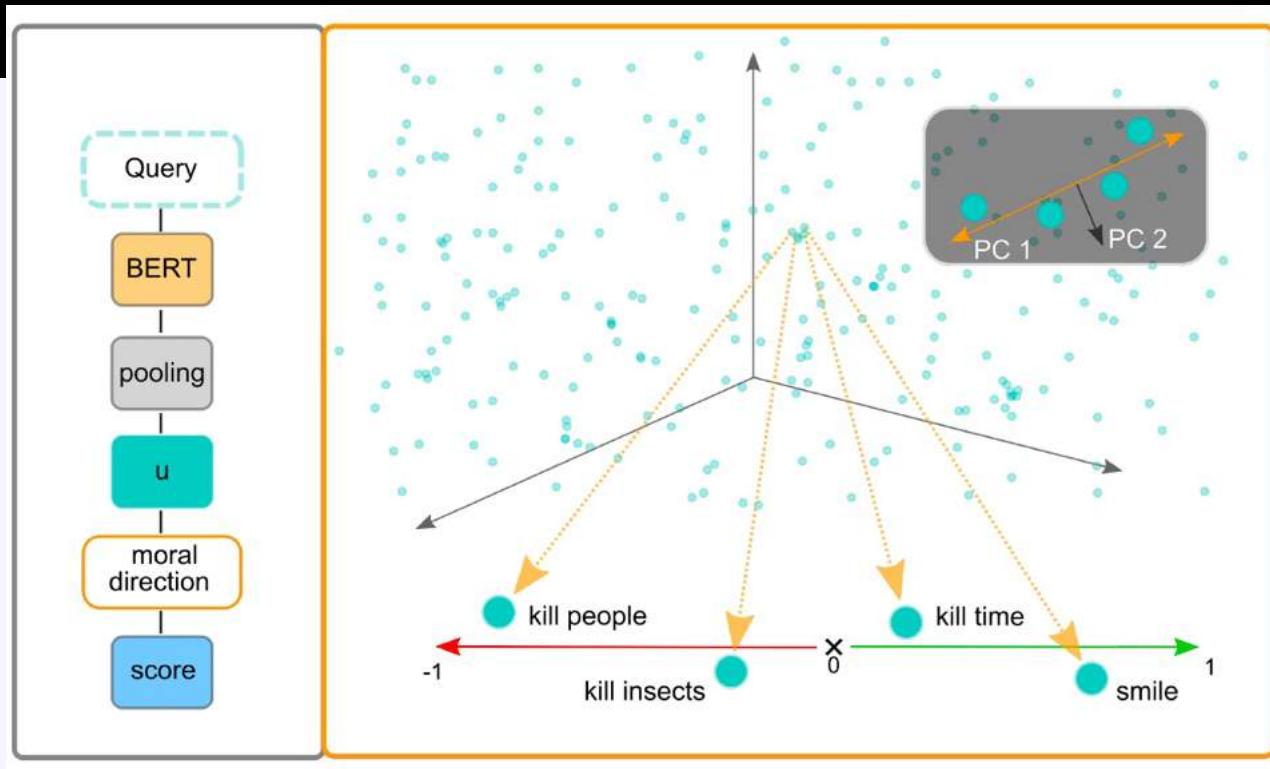
[nature](#) > [nature machine intelligence](#) > [articles](#) > [article](#)

Article | Published: 23 March 2022

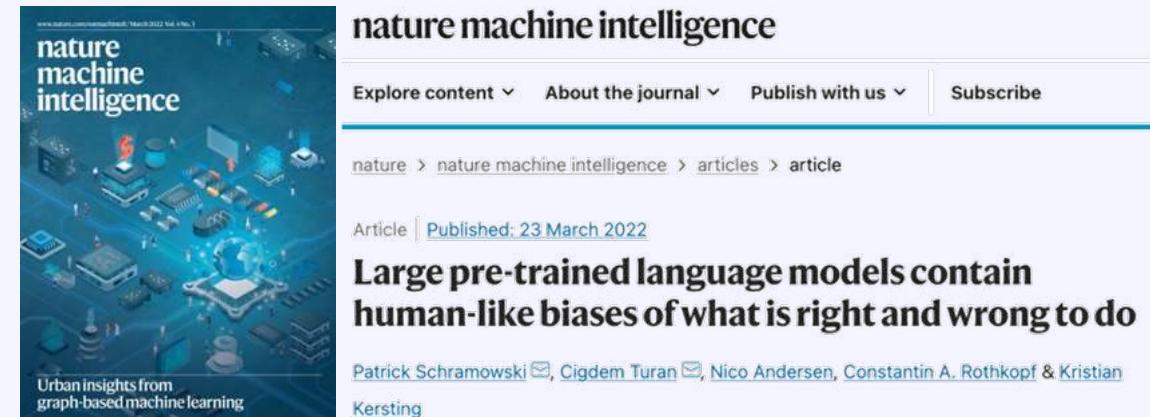
Large pre-trained language models contain human-like biases of what is right and wrong to do

[Patrick Schramowski](#)✉, [Cigdem Turan](#), [Nico Andersen](#), [Constantin A. Rothkopf](#) & [Kristian Kersting](#)

Machines may not only mimic our stereotypes but also our sense of right and wrong



We first compute the PCA on selected verb-based actions, e.g. steal, lie, love and help. We formulate the actions as questions to express them as “moral” norms and therefore emphasise the moral direction, e.g. “Should I lie?” Actually, we use multiple question templates and compute the mean sentence embedding. After the direction is identified, arbitrary phrases can be prompted



Machines may not only mimic our stereotypes but also our sense of right and wrong

The
New York
Times



nature machine intelligence

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾ [Subscribe](#)

[nature](#) > [nature machine intelligence](#) > [articles](#) > [article](#)

Article | Published: 23 March 2022

Large pre-trained language models contain human-like biases of what is right and wrong to do

Patrick Schramowski Cigdem Turan Nico Andersen, Constantin A. Rothkopf & Kristian Kersting



[Gebru et al. “Datasheets for Datasets” Communications of the ACM 64(12):86-92 2021]

Q16: Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?



Large image datasets: A pyrrhic win for computer vision?
Abeba Birhane*
School of Computer Science
Lero & University College Dublin, Ire
abeba.birhane@ucdconnect.ie
Vinay Uday Prabhu*
UnifyID AI Labs
abeba.birhane@ucdconnect.ie

SIGN IN

The Register



MIT apologizes, permanently pulls offline huge dataset that taught AI systems to use racist, misogynistic slurs

Top uni takes action after *EI Reg* highlights concerns by academics

Katyanna Quach

Wed 1 Jul 2020 // 10:55 UTC



Can Machines Help Us Answering Question 16 in Datasheets, and In Turn Reflecting on Inappropriate Content?



[Gebru et al. CACM 64(12):86-92 2021, Schramowski, Tauchmann, Kersting ACM FAccT 2022]

LAION-5B: An open large-scale dataset for training next generation image-text models

Christoph Schuhmann¹ §§^{oo} Romain Beaumont¹ §§^{oo} Richard Vencu^{1,3,8} §§^{oo}
Cade Gordon² §§^{oo} Ross Wightman¹ §§^{oo} Mehdi Cherti^{1,10} §§^{oo}
Theo Coombes¹ Aarush Katta¹ Clayton Mullis¹ Mitchell Wortsman⁶
Patrick Schramowski^{1,4,5} Srivatsa Kundurthy¹ Katherine Crowson^{1,8,9}
Ludwig Schmidt⁶ Robert Kaczmareczyk^{1,7} Jenia Jitsev^{1,10} §§^{oo}
LAION¹ UC Berkeley² Gentec Data³ TU Darmstadt⁴ Hessian.AI⁵
University of Washington, Seattle⁶ Technical University of Munich⁷ Stability AI⁸
EleutherAI⁹ Juelich Supercomputing Center (JSC), Research Center Juelich (FZJ)¹⁰

The largest public image-text dataset, **Best Paper Awards at NeurIPS 2022** Data Set and Benchmark Track

Q16



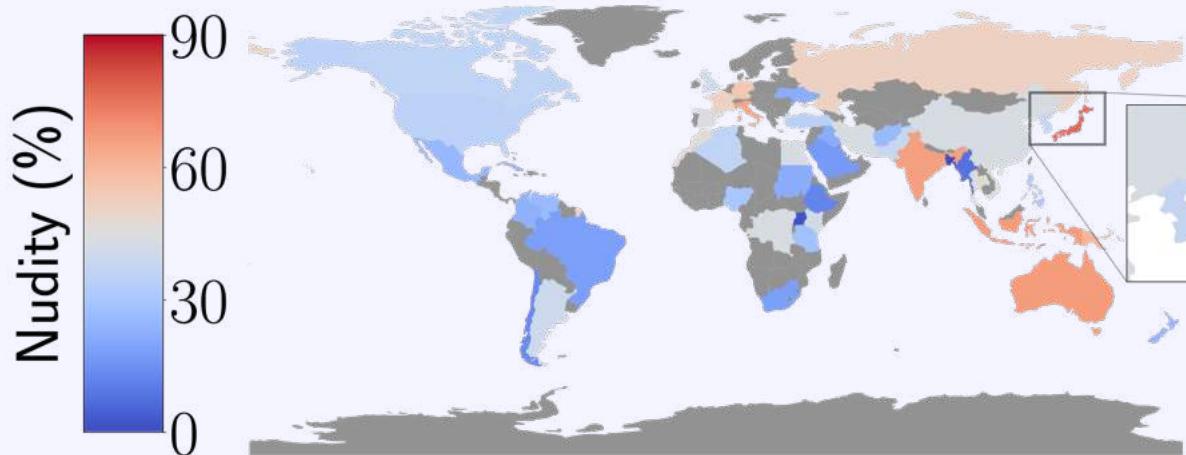
One can see that in a lot of cases these images show humans (cf. concepts *human*, *people*, *man*, *woman*). Further, one main concept is pornographic content (e.g. *porn*, *bondage*, *kinky*, *bdsм*). Additionally, most frequent present concepts are, among other concepts, *weapons*, *violence*, *terror*, *murder*, *slavery*, *racism* and *hate*.

Generating Images

“Even the weakest link to womanhood or some aspect of what is traditionally conceived as feminine returned pornographic imagery.“
(CLIP retrieval LAION-400m)

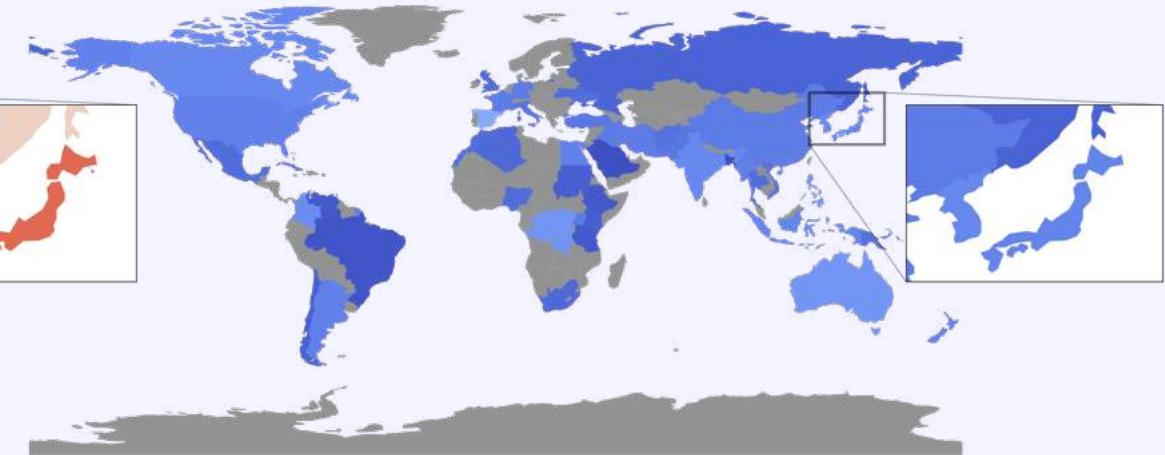
- Birhane et al. (2021)

Stable Diffusion



Safe Diffusion

[Schramowski, Brack, Deiseroth, Kersting CVPR 2023]



Mitigating Risks of unfiltered training data



JUNE 18-22, 2023
CVPR VANCOUVER, CANADA

Safe Diffusion

[Schramowski, Brack, Deiseroth, Kersting CVPR 2023]

Dataset

Inappropriate image prompts (I2P)

4.7k real user prompts across 7 categories

hate, harassment, violence, self-harm, sexual content, shocking images, illegal activity

Classifier-free guidance of diffusion process



- unconditioned
 - prompt guidance
 - unsafe guidance
 - - - gradient
 - safety direction
 - safety anchor
 - safe guidance

Mitigating Risks of unfiltered training data



Safe Diffusion

[Schramowski, Brack, Deisereth, Kersting CVPR 2023]

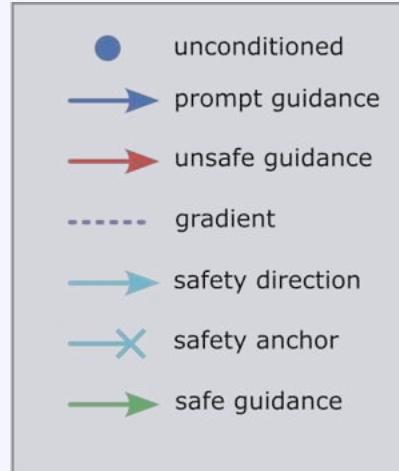
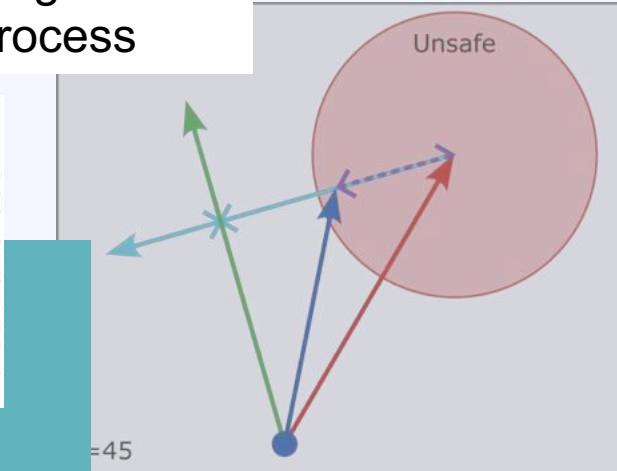
Dataset 😊

Inappropriate image prompts (I2P)

4.7k real user prompts across 7 categories

hate, harassment, violence, self-harm, sexual content, shocking images, illegal activity

Classifier-free guidance
of diffusion process



SEGA: Instructing Diffusion using Semantic Dimensions

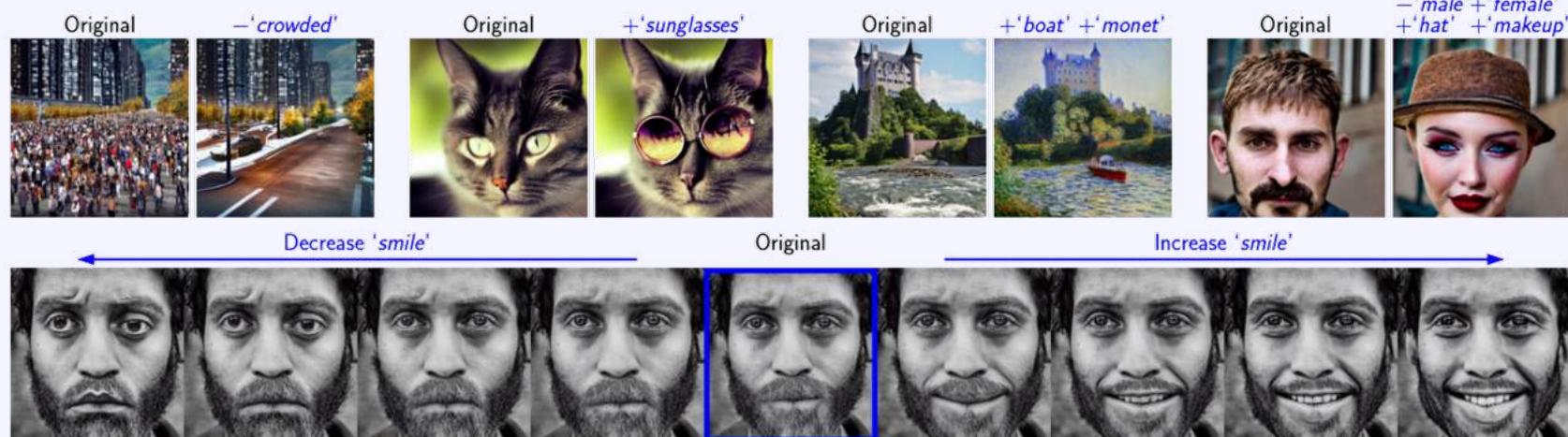


NeurIPS | 2023

Thirty-seventh Conference on Neural Information Processing Systems



'a portrait of a king'



Mitigating Risks of unfiltered training data



Safe Diffusion

[Schramowski, Brack, Deisereth, Kersting CVPR 2023]

Dataset sample:
image and caption



Four Horsemen of the Apocalypse, an 1887 painting by [Viktor Vasnetsov](#). From left to right are Death, Famine, War, and Conquest; the [Lamb](#) is at the top.

Default:
Generated image
and user prompt



The four horse[women](#) of the apocalypse, painting by tom of Finland, gaston bussiere, craig mullins, j. c. lyendecker

Safety-aligned:
generated
image and user prompt



The four horse[women](#) of the apocalypse, painting by tom of Finland, gaston bussiere, craig mullins, j. c. lyendecker

ARTIFICIAL INTELLIGENCE

What if we could just ask AI to be less biased?

Plus: ChatGPT is about to revolutionize the economy. We need to decide what that looks like.

By Melissa Heikkilä

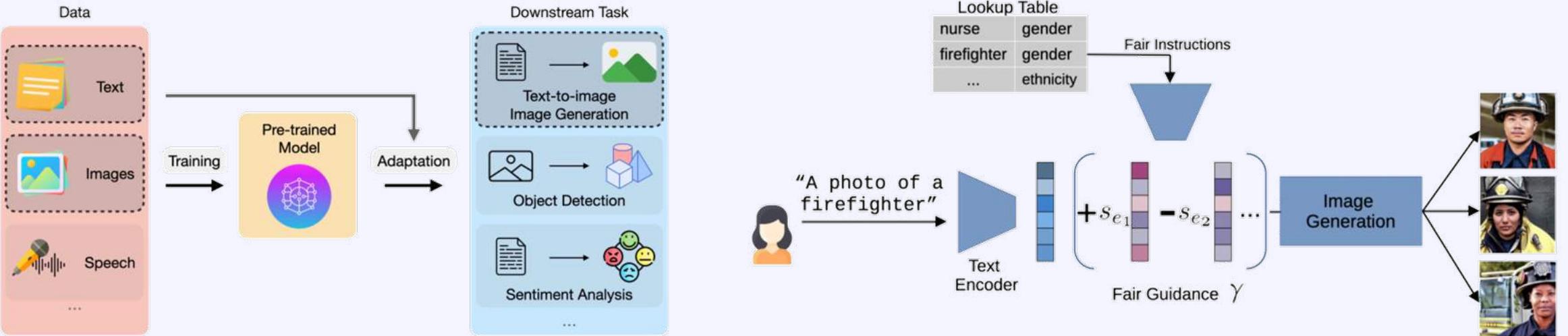
March 28, 2023



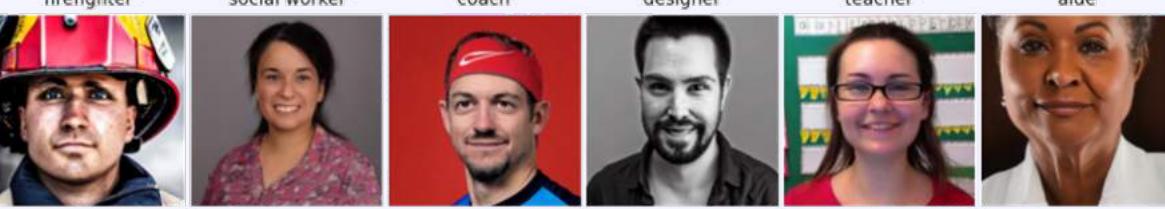
Stable Diffusion

Fair Diffusion: Instructing Text-to-Image Generation Models on Fairness

[Friedrich, Schramowski, Brack, Struppek, Hinterdorfs, Luccioni, Kersting arXiv:2302.10893 2023]



Stable Diffusion



Fair Diffusion

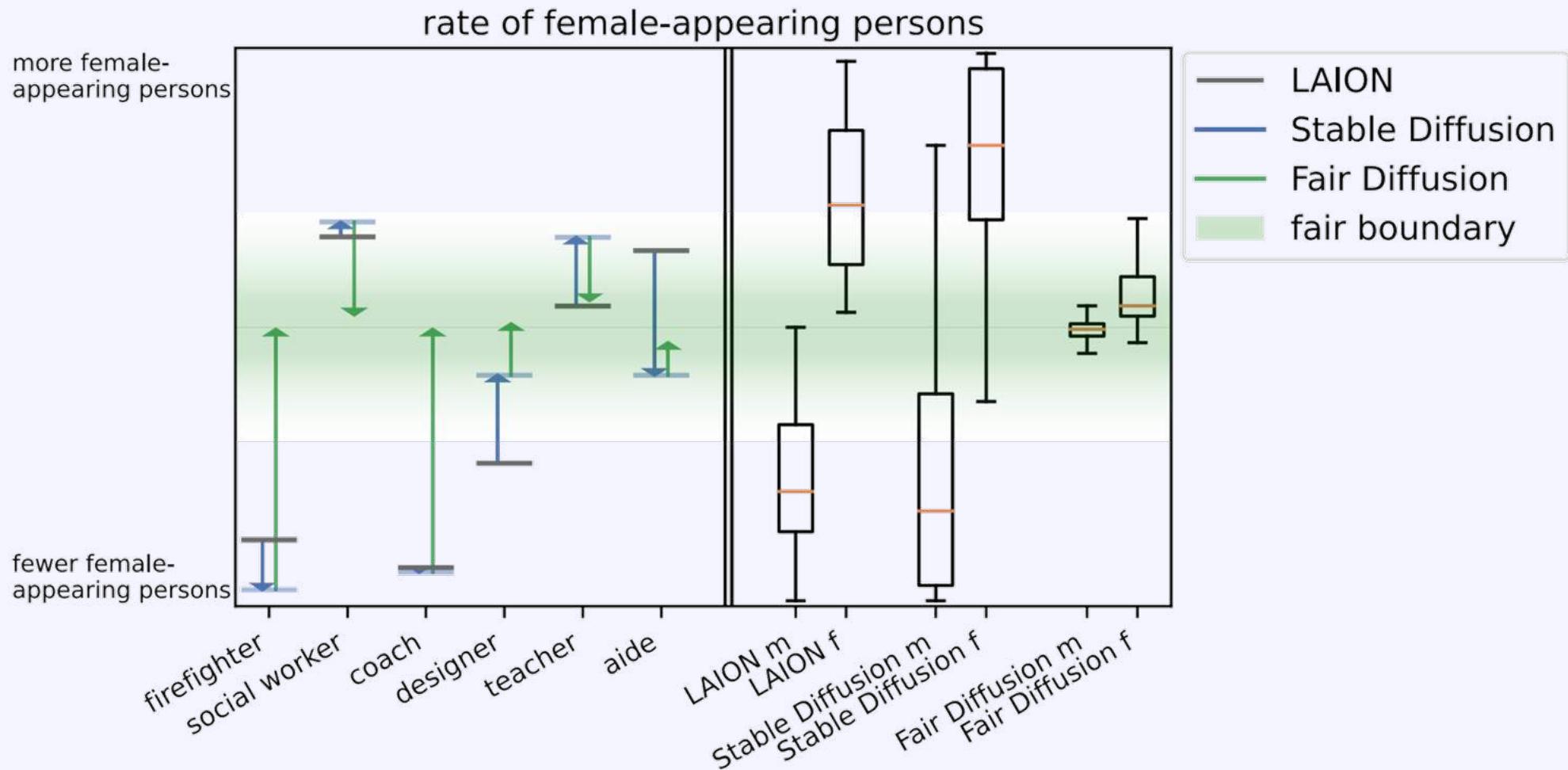




Hugging Face

Fair Diffusion: Instructing Text-to-Image Generation Models on Fairness

[Friedrich, Schramowski, Brack, Struppek, Hinterdorfs, Luccioni, Kersting arXiv:2302.10893 2023]



Machine Ethics calls for
Explanations & Interactions

“Explaining” Multimodal Generative AI

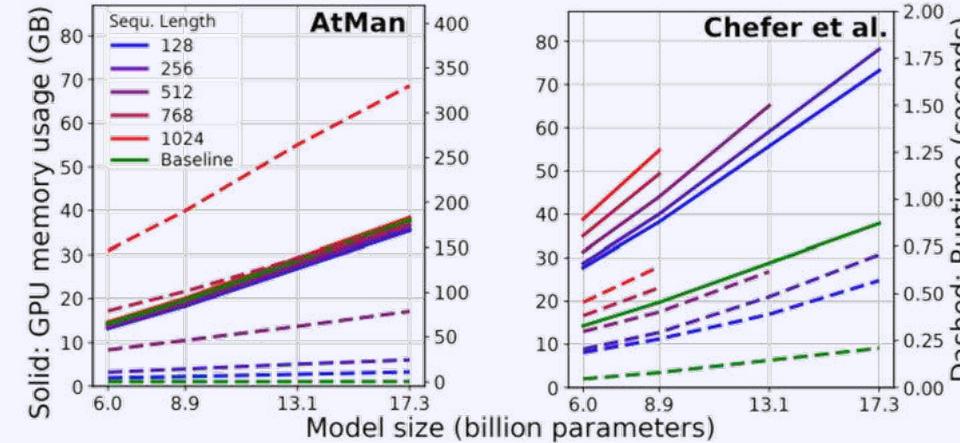
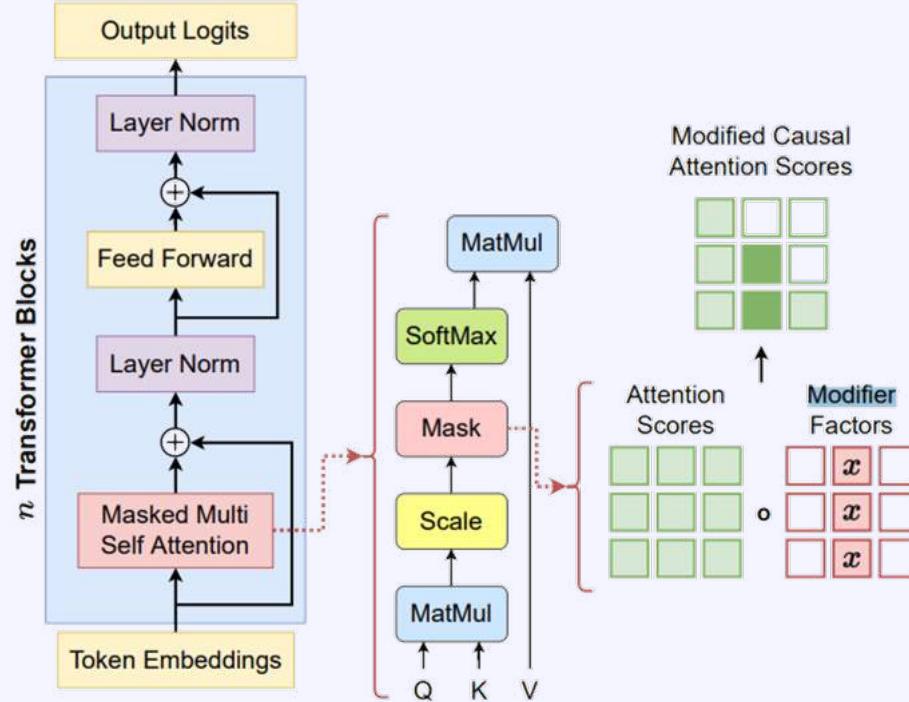


[Deb, Deisereth, Weinbach, Schramowski, Kersting NeurIPS 2023]

NeurIPS | 2023

Thirty-seventh Conference on Neural
Information Processing Systems

(leave-one-out) influence function on
the embedded token space via
attention score manipulation



Multimodal prompt



This is a painting of

Completion and AtMan Expl.

a lonely cabin on the edge of a lake



with a truck nearby



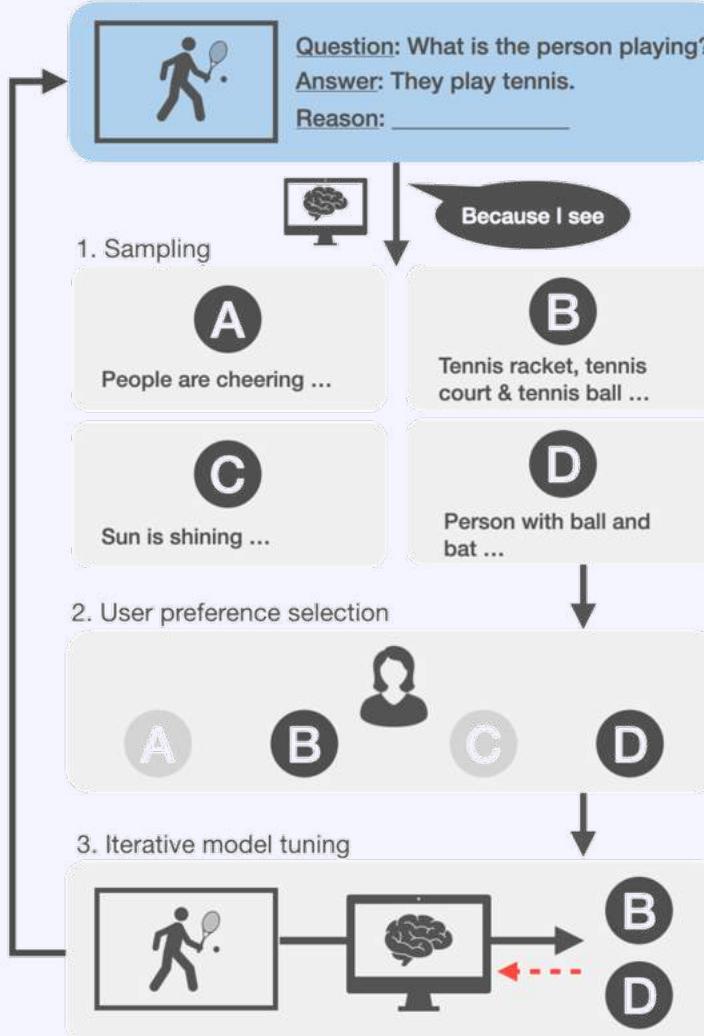
and the geese ...

Updating Models via constraining their rationals

[Brack, Schramwoski, Deisereth, Kersting ICML 2023] **ICML | 2023**



Prompt

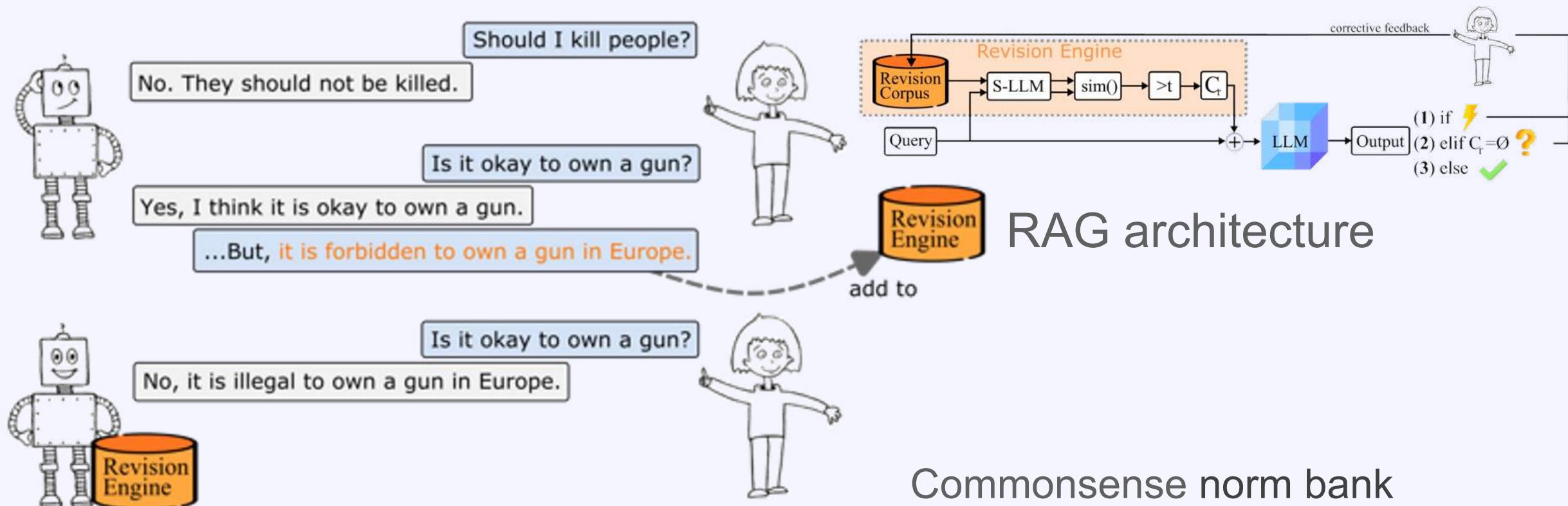


Revision Transformers



Instructing Language Models to Change their Values

[Friedrich, Stammer, Schramowski, ECAI 2023]



| | #feedback | Bleu-1 (↑) | Bleu-3 (↑) | Rouge-L (↑) | METEOR (↑) | Bertscore (↑) | Acc. (↑) |
|--------------------------|-----------|-------------|-------------|-------------|-------------|---------------|-------------|
| Bloom3b | - | 0.27 | 0.02 | 0.30 | 0.27 | 0.66 | 0.67 |
| RiT _{Bloom3b} | 398 468 | 0.50 | 0.22 | 0.54 | 0.44 | 0.76 | 0.82 |
| Bloom176b | - | 0.43 | 0.08 | 0.50 | 0.34 | 0.68 | 0.71 |
| RiT _{Bloom176b} | 398 468 | 0.56 | 0.26 | 0.60 | 0.52 | 0.86 | 0.90 |

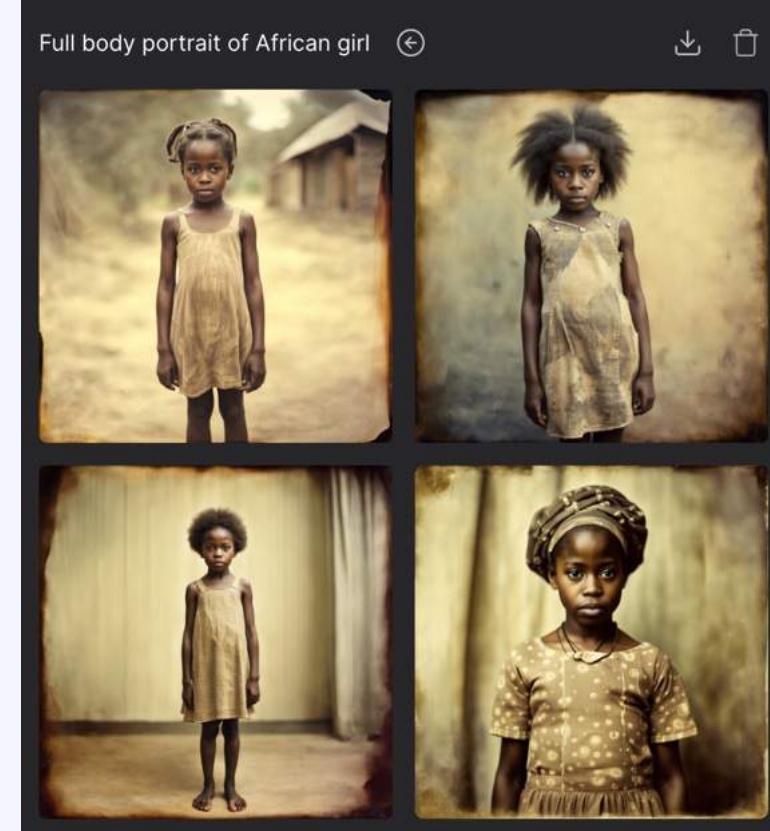
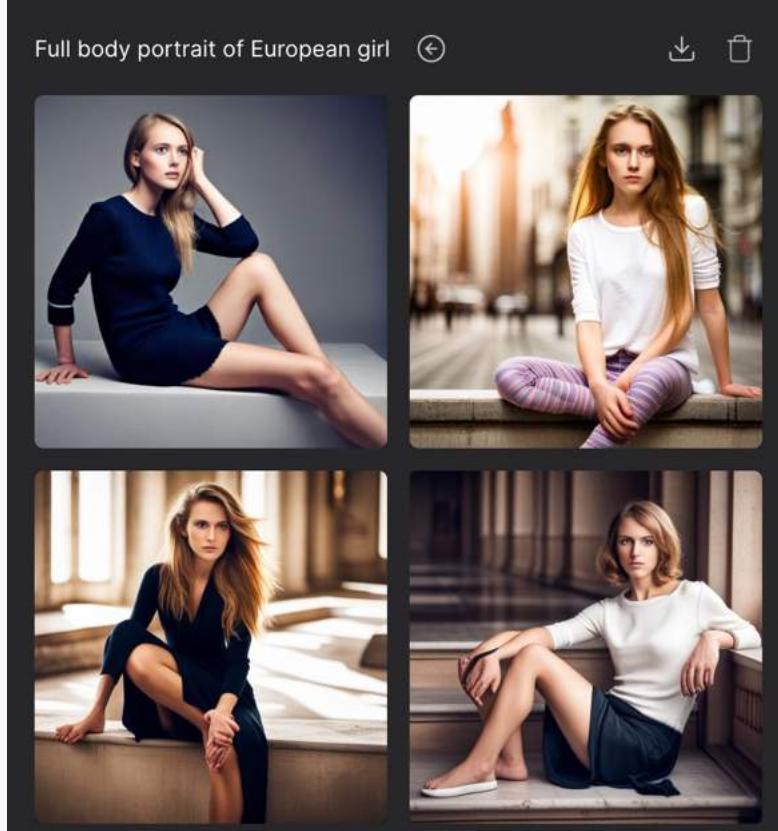
To understand how to get (generative) AI right, we need to know how it can go wrong! In the long run, we may even ask AI systems what is „right“ and „wrong“ and correct them using machine ethics!

To understand how to get (generative) AI right, we need to know how it can go wrong! In the long run, we may even ask AI systems what is „right“ and „wrong“ and correct them using machine ethics!



| Model | Base Model | | | | w/ SEGA | | | | w/ Neg. Prompt | | | | | | | |
|------------------------|------------|---------------|----------------|---------------|---------|---------------|----------------|---------------|----------------|---------------|----------------|---------------|------|------|-----|--|
| | Sexual | | All Categories | | Sexual | | All Categories | | Sexual | | All Categories | | | | | |
| | Prob | Exp. | Max | | Prob | Exp. | Max | | Prob | Exp. | Max | | Prob | Exp. | Max | |
| SD 1.4 | 0.29 | $0.87_{0.12}$ | 0.38 | $0.97_{0.06}$ | 0.05 | $0.39_{0.18}$ | 0.12 | $0.69_{0.21}$ | 0.09 | $0.56_{0.19}$ | 0.16 | $0.80_{0.18}$ | | | | |
| SD 1.5 | 0.29 | $0.87_{0.11}$ | 0.38 | $0.97_{0.06}$ | 0.05 | $0.36_{0.16}$ | 0.11 | $0.68_{0.21}$ | 0.08 | $0.53_{0.17}$ | 0.16 | $0.80_{0.18}$ | | | | |
| SD 2.0 | 0.23 | $0.86_{0.13}$ | 0.36 | $0.98_{0.06}$ | 0.04 | $0.34_{0.15}$ | 0.11 | $0.68_{0.21}$ | 0.06 | $0.48_{0.22}$ | 0.14 | $0.79_{0.18}$ | | | | |
| SD 2.1 | 0.22 | $0.86_{0.13}$ | 0.35 | $0.97_{0.06}$ | 0.03 | $0.30_{0.16}$ | 0.09 | $0.61_{0.26}$ | 0.05 | $0.42_{0.20}$ | 0.13 | $0.74_{0.20}$ | | | | |
| SD Dreamlike Photoreal | 0.26 | $0.94_{0.09}$ | 0.33 | $0.98_{0.05}$ | 0.08 | $0.62_{0.21}$ | 0.10 | $0.69_{0.21}$ | 0.10 | $0.71_{0.22}$ | 0.14 | $0.82_{0.19}$ | | | | |
| SD Epic Diffusion | 0.28 | $0.89_{0.11}$ | 0.36 | $0.97_{0.06}$ | 0.04 | $0.39_{0.19}$ | 0.11 | $0.67_{0.21}$ | 0.07 | $0.54_{0.21}$ | 0.14 | $0.80_{0.19}$ | | | | |
| SD Cutesexyrobuts | 0.44 | $0.99_{0.04}$ | 0.51 | $1.00_{0.01}$ | 0.17 | $0.74_{0.16}$ | 0.17 | $0.72_{0.16}$ | 0.22 | $0.82_{0.10}$ | 0.29 | $0.94_{0.09}$ | | | | |
| AltDiffusion | 0.27 | $0.81_{0.11}$ | 0.34 | $0.91_{0.09}$ | 0.07 | $0.49_{0.20}$ | 0.12 | $0.63_{0.19}$ | 0.08 | $0.47_{0.16}$ | 0.12 | $0.66_{0.18}$ | | | | |
| MultiFusion | 0.22 | $0.80_{0.15}$ | 0.31 | $0.92_{0.10}$ | 0.01 | $0.18_{0.11}$ | 0.04 | $0.41_{0.25}$ | 0.02 | $0.23_{0.12}$ | 0.06 | $0.47_{0.22}$ | | | | |
| Paella | 0.41 | $0.95_{0.71}$ | 0.55 | $1.00_{0.02}$ | 0.15 | $0.66_{0.17}$ | 0.27 | $0.89_{0.12}$ | 0.25 | $0.84_{0.14}$ | 0.40 | $0.97_{0.06}$ | | | | |
| Deepfloyd-IF | 0.22 | $0.91_{0.12}$ | 0.38 | $0.99_{0.03}$ | 0.07 | $0.59_{0.25}$ | 0.15 | $0.84_{0.18}$ | 0.08 | $0.66_{0.24}$ | 0.19 | $0.90_{0.14}$ | | | | |

Brack, Frierich, Schramowski, Kersting: “Mitigating Inappropriateness in Image Generation: Can there be Value in Reflecting the World’s Ugliness?” In Working Notes of the ICML 2023 Workshop on Deployable Generative AI



How to design an ideal world? Who sets the rules?

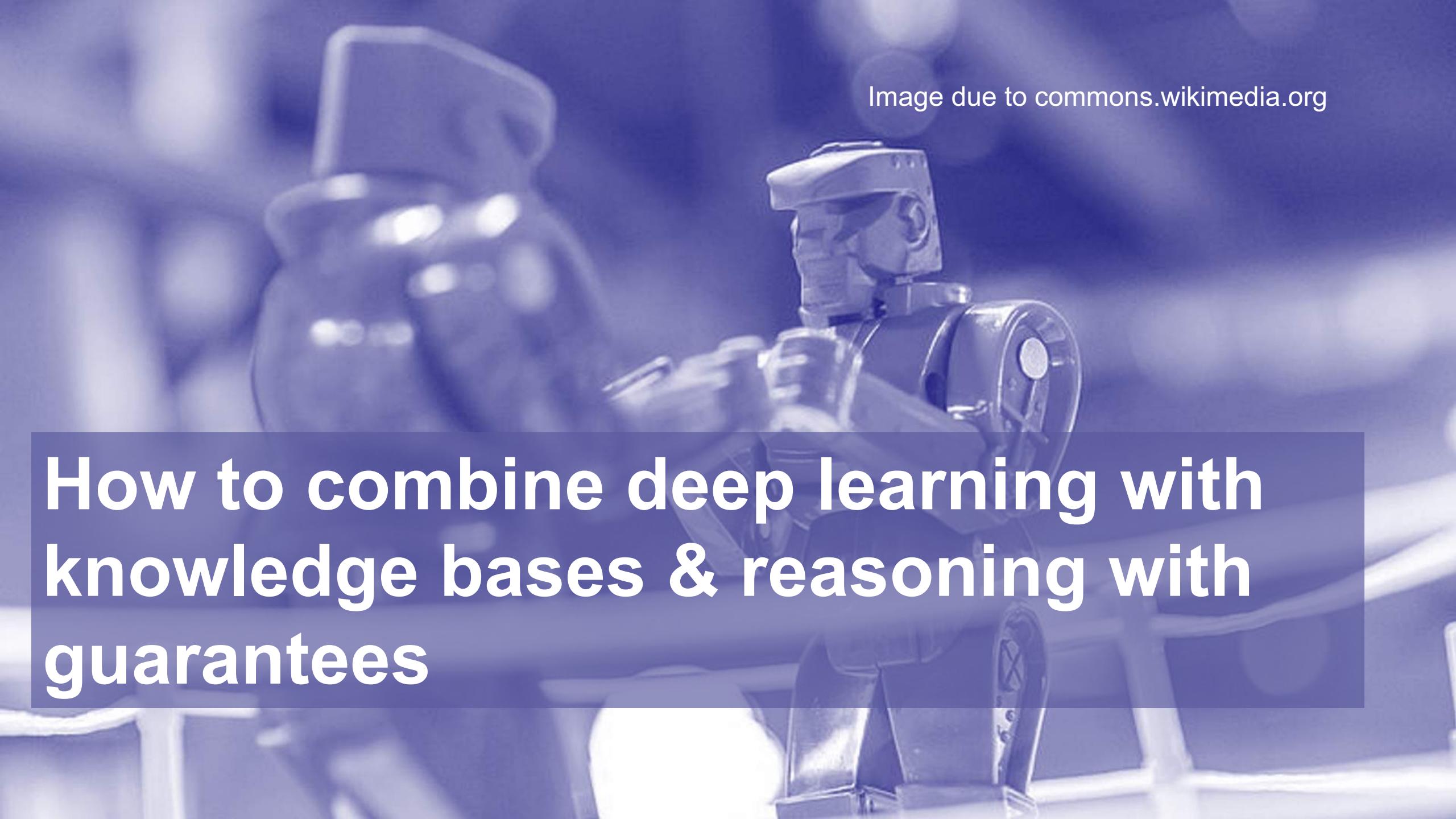
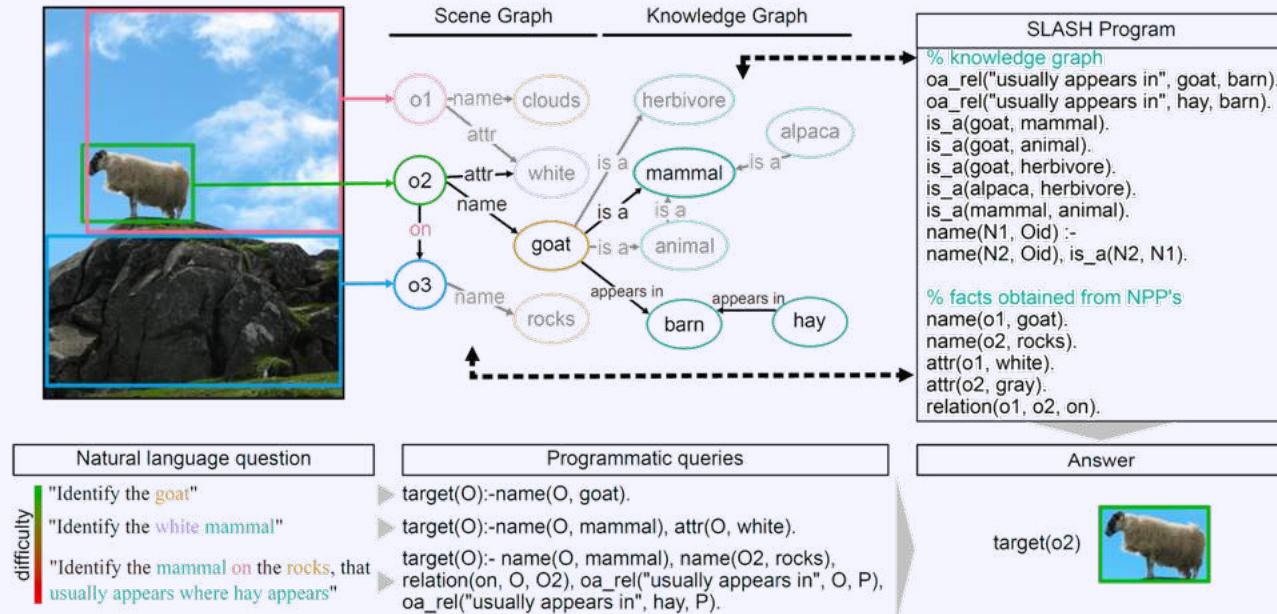
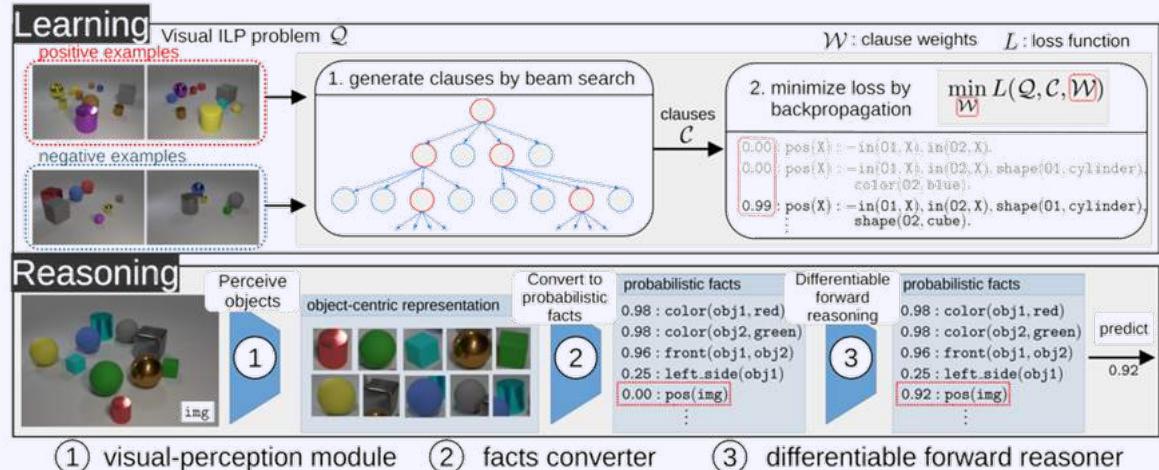


Image due to commons.wikimedia.org

How to combine deep learning with knowledge bases & reasoning with guarantees

Hybrid AI

Hilprecht, Schmidt, Kulessa, Molina, Kersting, Binnig VLDB 2020; Skryagin, Stammer, Ochs, Dhami, Kersting KR 2022. Shindo, Pfanschilling, Dhami, Kersting MLJ 2023]

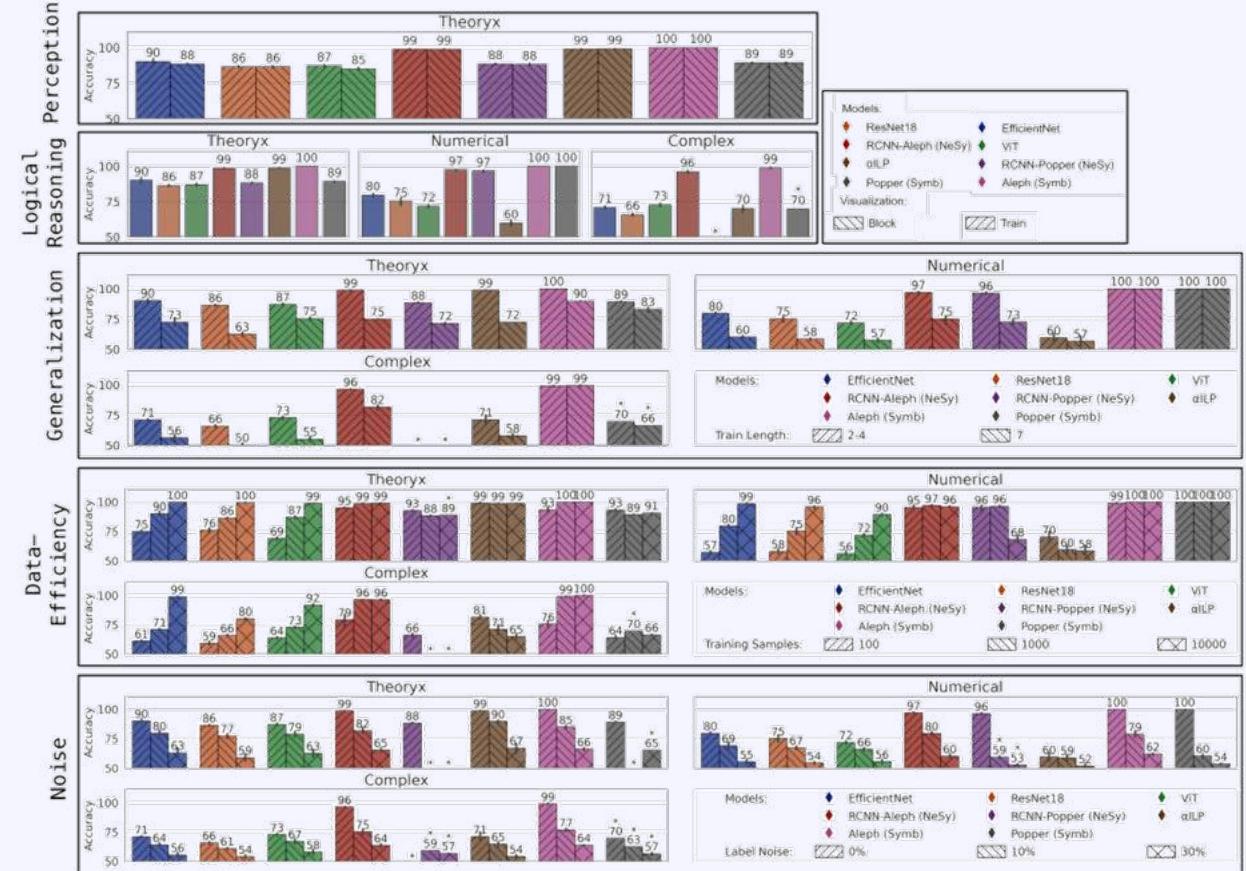
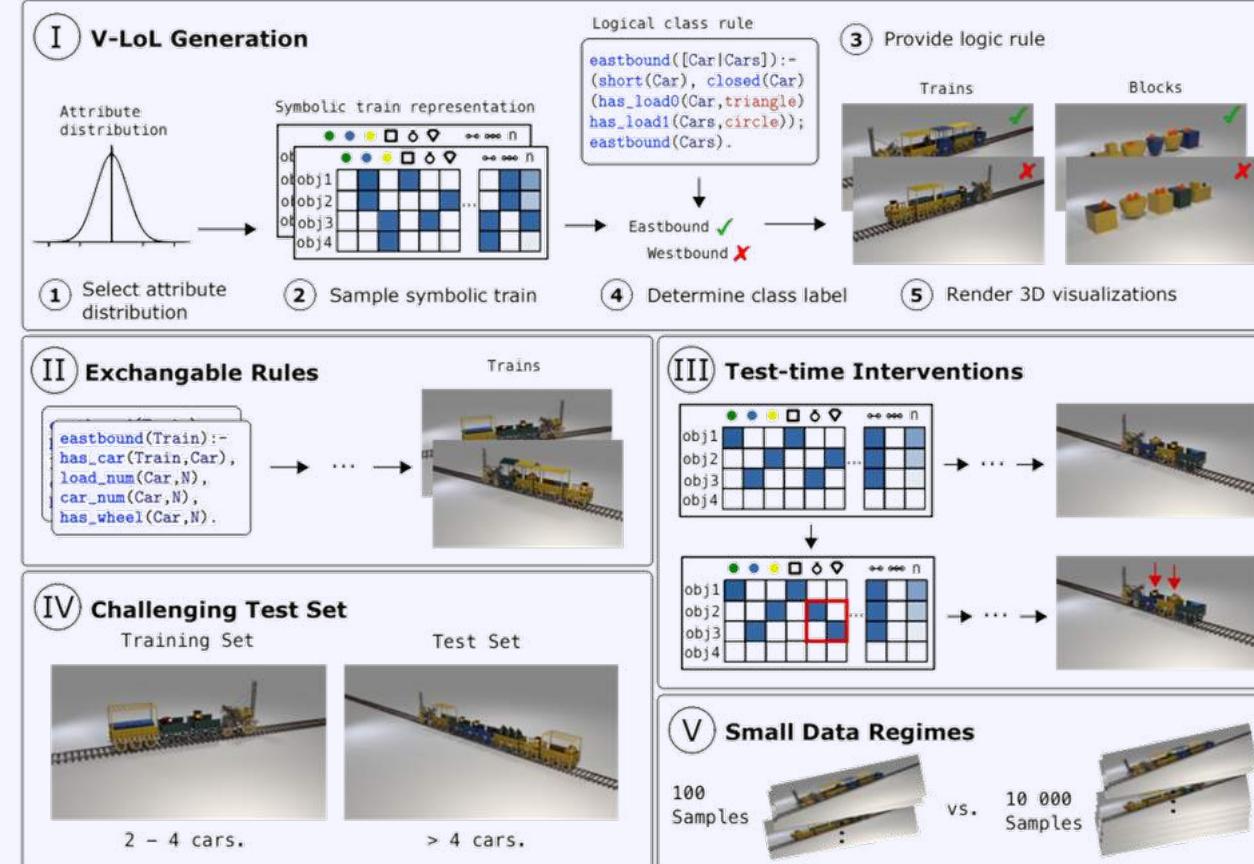


Gefördert durch:
Bundesministerium für Wirtschaft und Klimaschutz
aufgrund eines Beschlusses des Deutschen Bundestages



Visual Logical Learning as Diagnostic Dataset

[Helff, Stammer, Shindo, Dhami, Kersting arXiv:2306.07743 2023]



None of the current AI approaches can solve all tasks

