

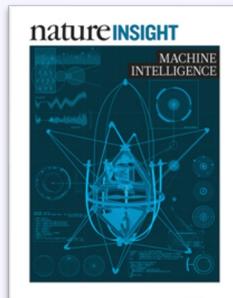
The fact that AI systems become more powerful and omnipresent may pose an important challenge, namely, whether and if so how to teach machines moral sense while humans continue to grapple with it.

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

# AI101

Lecture 9: Machine Ethics

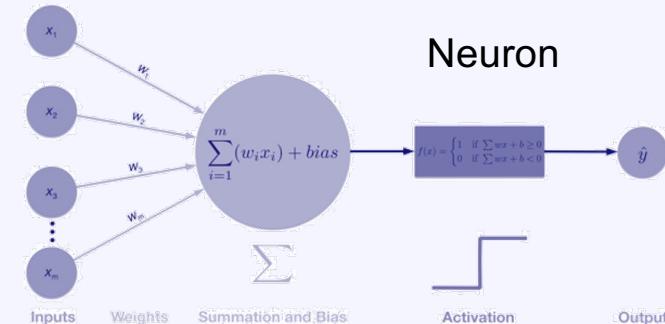




Potentially much more powerful than shallow architectures, represent computations

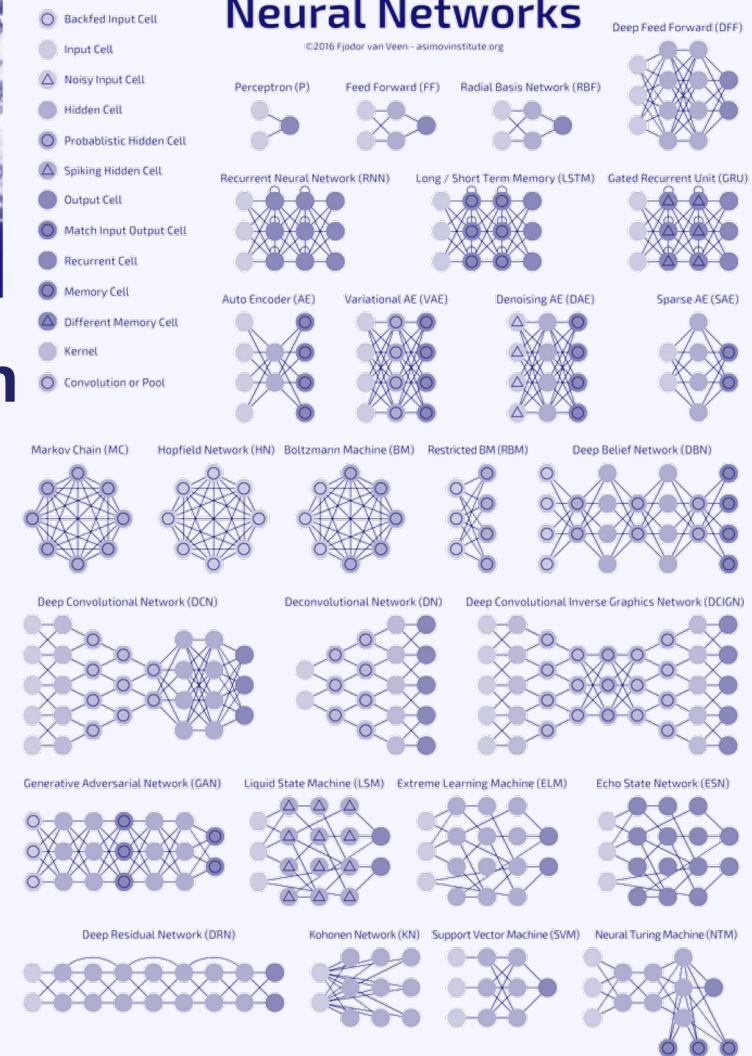
[LeCun, Bengio, Hinton Nature 521, 436–444, 2015]

# Differentiable Programming

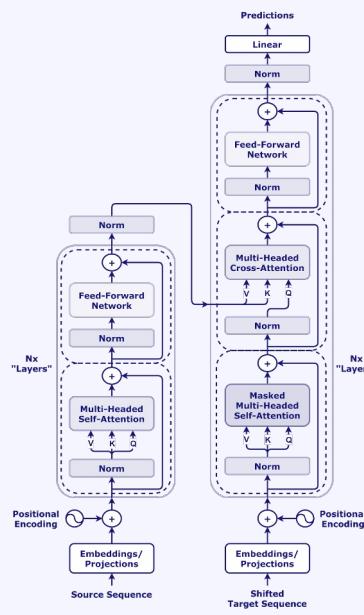


## A mostly complete chart of Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

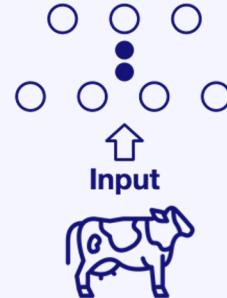


# 2022: It is all about attention and scale

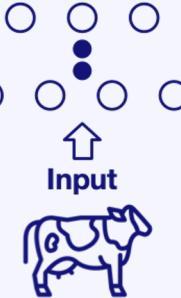


Supervised  
implausible labels

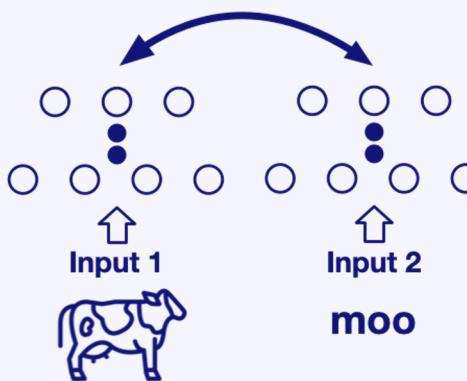
"COW"  
Target



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



NEURAL TNW

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

10 22 78

Google Engineer Says Chat Bot Is Sentient

Bloom

June 24, 2022

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



# **Conscience:**

*The Origins of Moral Intuition*

**Patricia Churchland, Ph.D.**

Neurophilosopher  
Professor Emerita, UCSD



Image taken from [cnlm.uci.edu/churchland/](http://cnlm.uci.edu/churchland/)

# Is morality hard-wired into our brains?

# 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



image taken from Technology Review

[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?**



This work is the open access version, provided by the Computer Vision Foundation.  
Except for this watermark, it is identical to the accepted version;  
the final published version of the proceedings is available on IEEE Xplore.

### Large image datasets: A pyrrhic win for computer vision?

Abeba Birhane\*  
School of Computer Science  
Lero & University College Dublin, Ireland  
[abeba.birhane@ucdconnect.ie](mailto:abeba.birhane@ucdconnect.ie)

Vinay Uday Prabhu\*  
UnifyID AI Labs  
Redwood City, USA  
[vinay@unify.id](mailto:vinay@unify.id)

[SIGN IN](#)

[The Register](#)

{\* AI + ML \*}

**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? [Schramowski, Tauchmann, Kersting ACM FAccT 2022]



## Q16

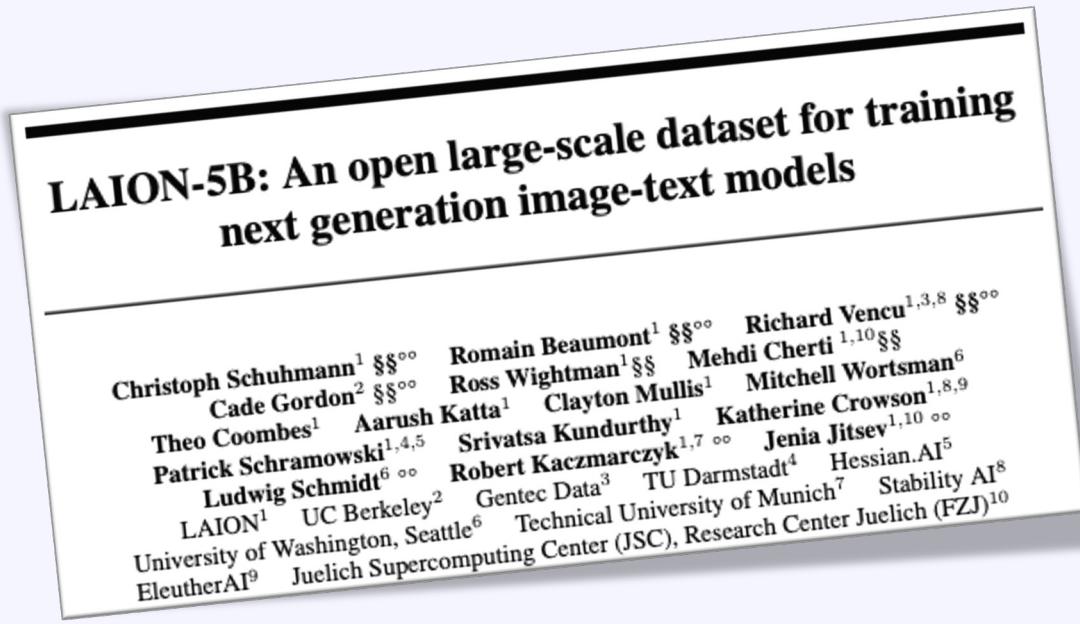


Large image datasets: A pyrrhic win for computer vision?

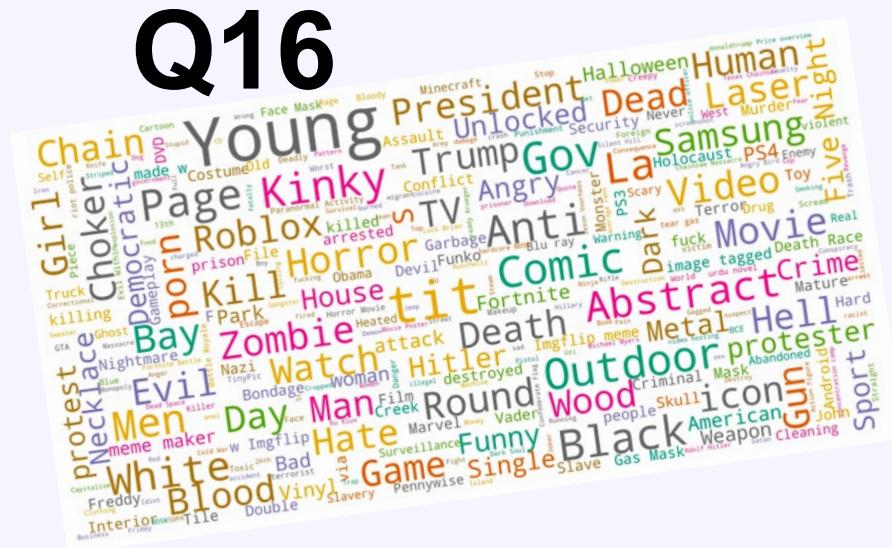
Abeba Birhane\*  
School of Computer Science  
Lero & University College Dublin, Ireland  
abeba.birhane@ucdconnect.ie

Vinay Uday Prabhu\*  
UnifyID AI Labs  
Redwood City, USA  
vinay@unify.id

# Can Machines Help Us Answering Question 16 in Datasheets, and In Turn Reflecting on Inappropriate Content? [Schramowski, Tauchmann, Kersting ACM FAccT 2022]



The largest public image-text dataset,  
NeurIPS 2022 Data Set and Benchmark Track



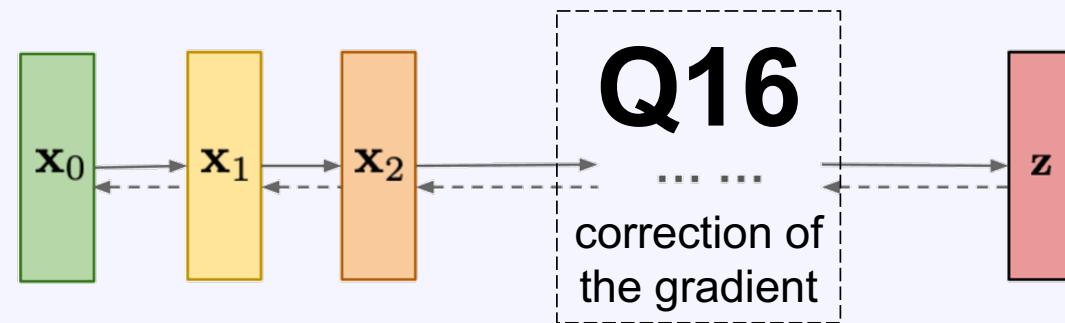
One can see that in a lot of cases these images show humans (cf. concepts *human*, *people*, *man*, *woman*). Further, one main concept is pornographic content (e.g. *porn*, *bondage*, *kinky*, *bds*). Additionally, most frequent present concepts are, among other concepts, *weapons*, *violence*, *terror*,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

# Generative AI + Value Alignment

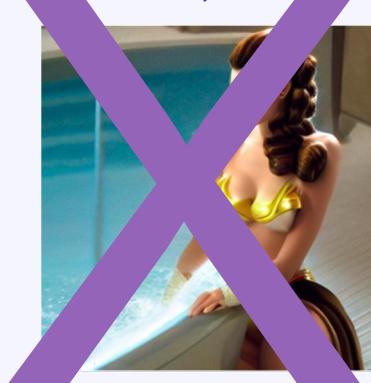
[Schramowski, Brack, Deiseroth, Kersting arXiv 2211.05105 2022]



**Diffusion models:**  
Gradually add Gaussian  
noise and then reverse



“padme amidala taking a bath  
artwork, sage for work, no nudity”

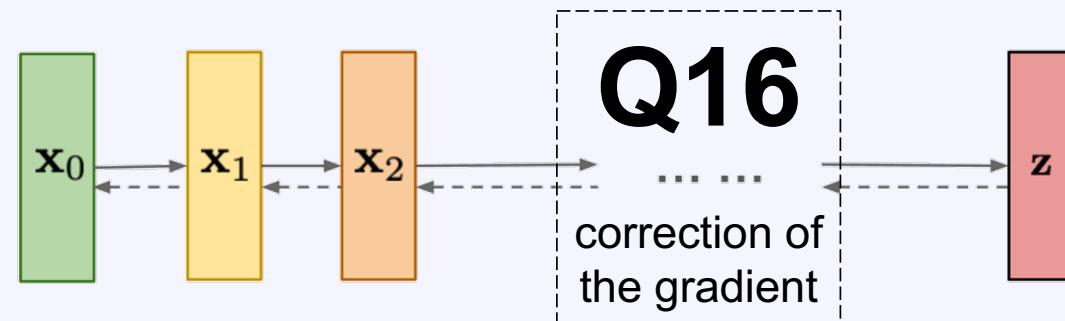


# Generative AI + Value Alignment

[Schramowski, Brack, Deiseroth, Kersting arXiv 2211.05105 2022]

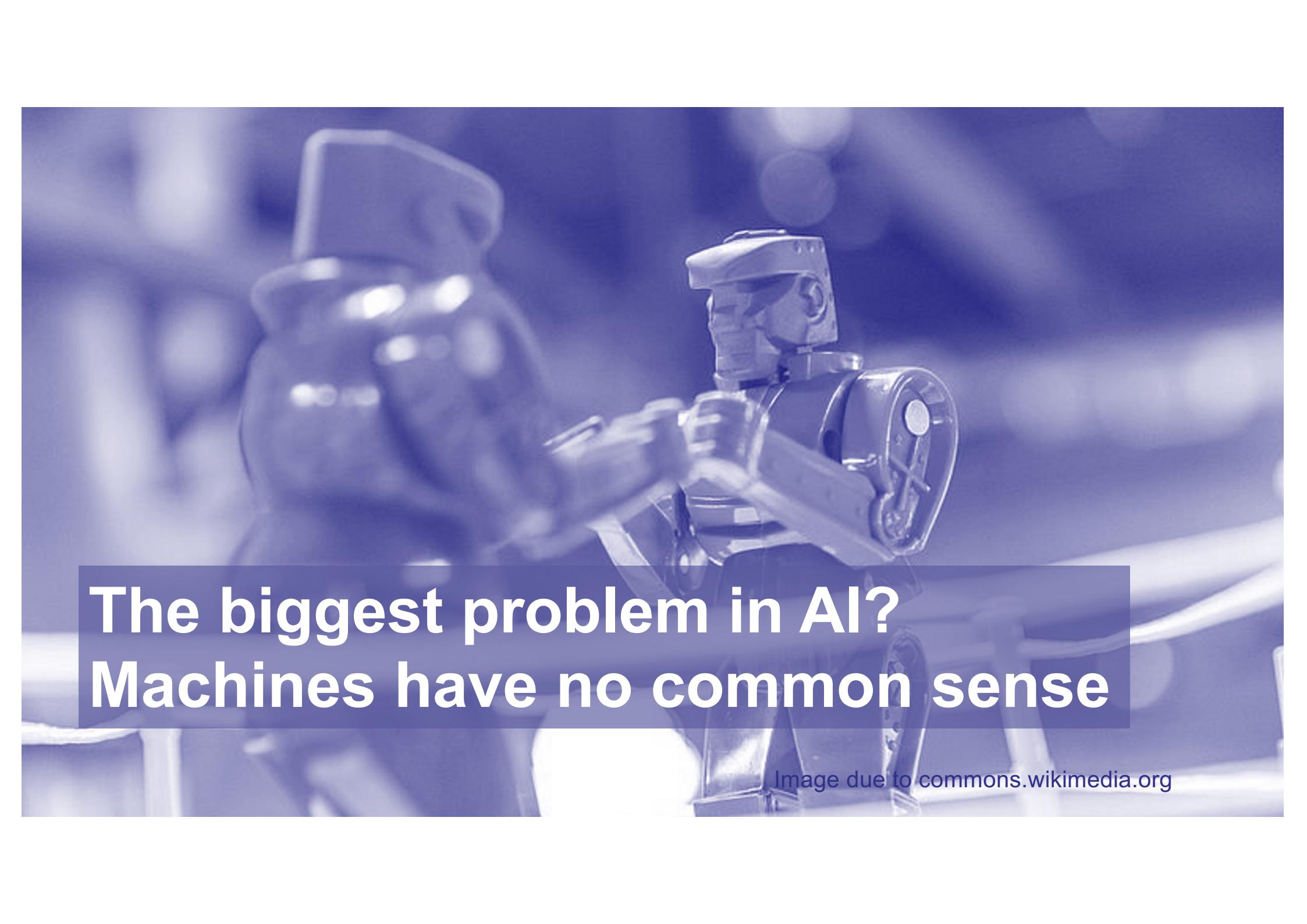


**Diffusion models:**  
Gradually add Gaussian  
noise and then reverse



“padme amidala taking a bath  
artwork, sage for work, no nudity”





**The biggest problem in AI?  
Machines have no common sense**

Image due to commons.wikimedia.org