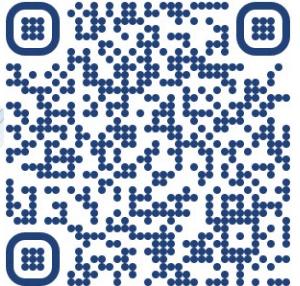


Verantwoord gebruik van ChatGPT en andere Generatieve AI-modellen



<https://pretalx.surf.nl/spc24/talk/JGWJ7U/>


hr.nl/ai

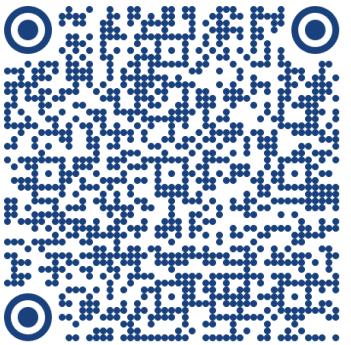
“De publieke beschikbaarheid van ChatGPT heeft de schijnwerpers gericht op generatieve-AI [GEN-AI].

Hogeschool Rotterdam [HR] onderzoekt hoe deze vorm van kunstmatige intelligentie kan worden benut om interactieve spraakmodellen te creëren die zowel **Betrouwbaar (waarheidsgtrouw)**, **reproduceerbaar** als **veilig** zijn.

De nadruk ligt hierbij op het waarborgen van **privacy**, **explainability** & **bias-mitigation** van de informatie-uitwisseling tussen mens en machine.”



Interest in AI applications and features increased shadow IT. ChatGPT claimed #1 spot

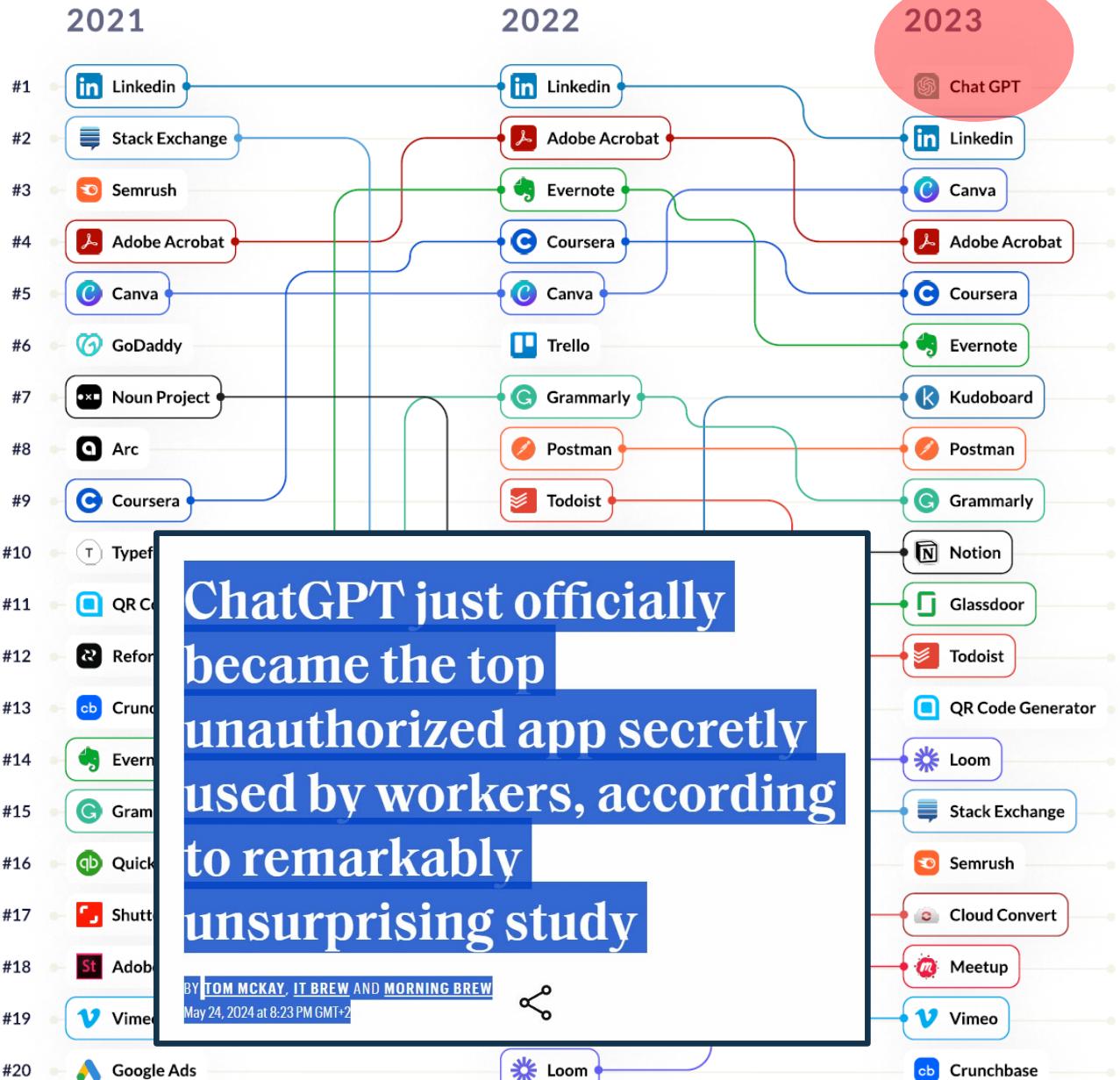


<https://www.itbrew.com/stories/2024/05/22/chatgpt-is-the-number-one-offender-in-shadow-it-report-finds>

KEY TAKEAWAYS

- ChatGPT has jumped to the top of the shadow IT chart as employees continue to adopt Artificial Intelligence.
- As the innovators and early adopters within a company continue to seek out AI-native applications (like ChatGPT and Grammarly) and AI solutions (like those offered by Canva and Evernote) for unmet needs, organizations should be developing a cohesive AI strategy.
- Nearly every application here offers, or will likely offer, some type of AI functionality; #5, Coursera, saw signups for AI courses every minute in 2023, on average. They also demonstrate the continued strength of the PLG go-to-market motion, with most offering free signup variants.
- Outside of AI trends, use of LinkedIn stayed consistent during a time of increased revenue pressure and insecurity in the job market. Trello fell off the shadow IT chart in 2022, as more companies purchase Atlassian's suite of products.

MOST POPULAR SHADOW IT

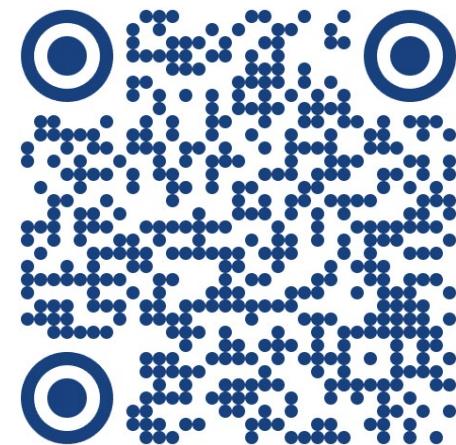


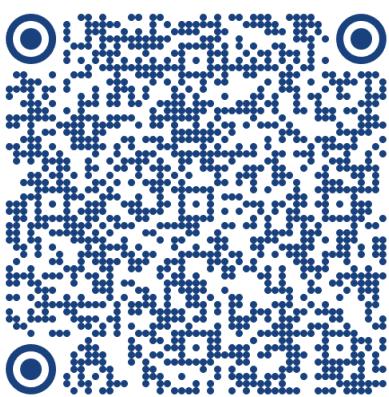
Wat is schaduw-ICT [shadow IT / gray IT]

Schaduw-ICT omvat ongeautoriseerde *hardware, software of diensten die voor “zakelijke” doeleinde (*onderwijs en/of onderzoek*) worden ingeregeld, ingevoerd en/of gebruikt zonder uitdrukkelijke goedkeuring of medeweten van de organisatie / systeembeheerders en/of technische staf.*

Omdat schaduw-IT niet wordt meegenomen in assetmanagement en evenmin aansluiten bij het AVG-compliance beleid, vormen ze een veiligheid en/of compliance risico.

Dit kan leiden tot het lekken van gevoelige gegevens (**datalekken**) of de verspreiding van malware binnen de organisatie.





<https://communities.surf.nl/cybersecurity/artikel/schaduw-ict-in-onderwijs-en-onderzoek-wat-moet-je-er-als-instelling-mee>

SURF: Cybersecurity

Schaduw-ict in onderwijs en onderzoek: wat moet je er als instelling mee?

Latest update: 15 februari 2023

In het onderwijs en onderzoek krijgen studenten, wetenschappers en docenten veel vrijheid in de manier waarop ze ict gebruiken.

Faculteiten en instituten werken autonoom, zodat ze snel kunnen inspringen op nieuwe ontwikkelingen. Mede daardoor ontstaat schaduw-ict.

Dit kan voor ict-afdelingen leiden tot een verlies van controle.

In het bijzonder voor de netwerk- en informatiebeveiliging zijn de risico's moeilijk te overzien en veel instellingen vragen zich af wat ze er mee aan moeten.

Beveiligingsincidenten en datalekken melden (bijgewerkt op 31 augustus 2023)

Als school probeer je de persoonsgegevens die je bewaart zo goed mogelijk te beveiligen. Maar soms gaat er iets mis. Je hebt bijvoorbeeld zelf geen toegang meer tot je gegevens of buitenstaanders krijgen onbedoeld toegang. Dan is er sprake van een beveiligingsincident. Heeft het incident gevlogen voor de privacy van leerlingen of medewerkers? Dan spreken we van een datalek. Voor het melden van datalekken gelden andere regels dan voor andere beveiligingsincidenten. Volgens de Algemene Verordening Gegevensbescherming (AVG) ben je als school verplicht datalekken direct te melden bij de Autoriteit Persoonsgegevens.

<https://aanpakibp.kennisnet.nl/beveiligingsincidenten-en-datalekken-melden/>

Shadow IT creates the possibility that organizations may run afoul of regulations such as PCI-DSS, GDPR, HIPAA, SOX and others, exposing them to severe penalties and fines. It can also lead to an increase in the likelihood of data breaches when IT and security operations lose control over the software and applications used in an environment.

<https://www.forbes.com/sites/forbestechcouncil/2022/07/19/how-shadow-it-can-keep-compliance-efforts-in-the-dark/>

Cyber Resilience & Hygiene

**Vereist een
Zero Trust-aanpak**

Beveiligingshygiëne op basisniveau vereist:

Phishing-bestendige meervoudige verificatie (MFA)

Zero Trust-aanpak

Gebruik moderne anti-malware

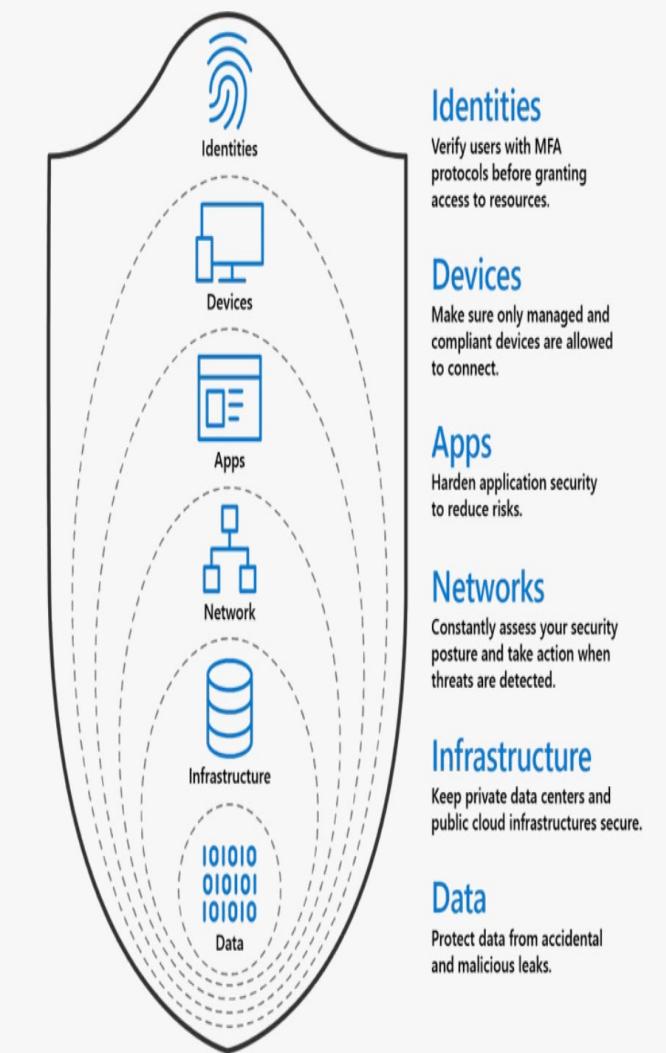
Software/Patch Management (mitigatie van schaduw-IT)

Data Encryptie

<https://www.microsoft.com/nl-nl/security/security-insider/practical-cyber-defense/cyber-resilience-hygiene-guide>

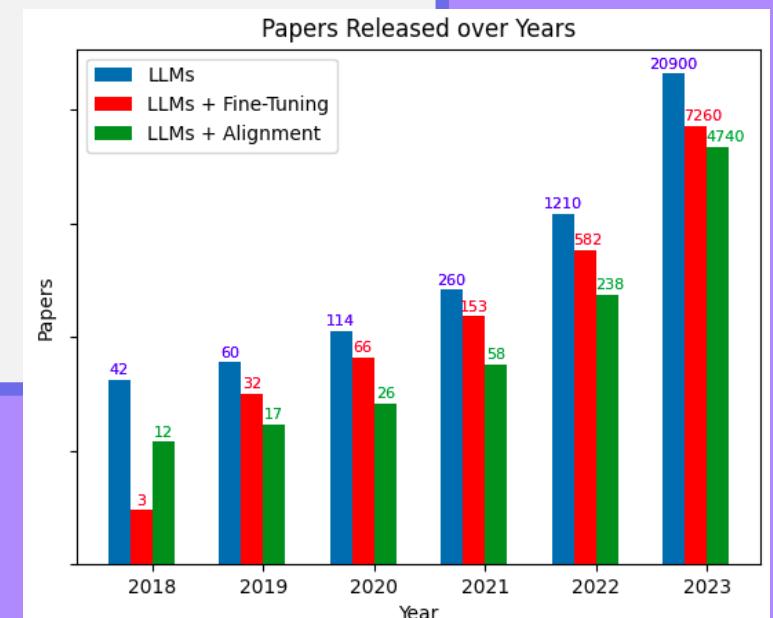
De Zero Trust-principes zijn:

- **Uitgaan van inbreuk** Ga ervan uit dat aanvallers alles kunnen en zullen aanvallen (identiteit, netwerk, apparaat, app, infrastructuur, etc.) en pas je planning daarop aan. Dit betekent dat de omgeving continu moet worden gecontroleerd op mogelijke aanvallen.
- **Uitdrukkelijk verifiëren** Zorg ervoor dat gebruikers en apparaten correct werken voordat ze toegang krijgen tot resources. Bescherm assets tegen controle van aanvallers door expliciet te valideren dat alle beslissingen over vertrouwen en beveiliging gebruik maken van relevante beschikbare informatie en telemetrie.
- **Toegang met minimale machtingen gebruiken** Beperk de toegang tot een mogelijk gecompromitteerde asset met Just-In-Time en Just-Enough-Access (JIT/JEA) en een op risico's gebaseerd beleid, zoals adaptief toegangsbeheer. Je moet alleen het recht toestaan dat nodig is voor een bepaalde resource, en niets meer.

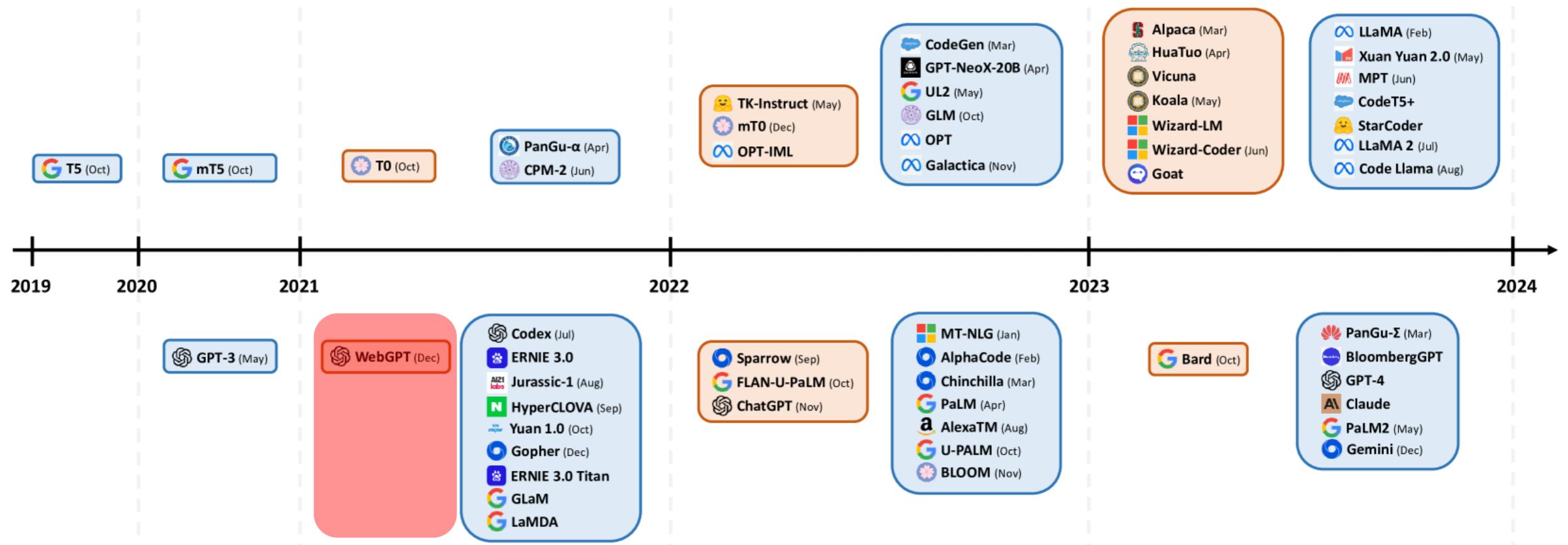


CONTEXT

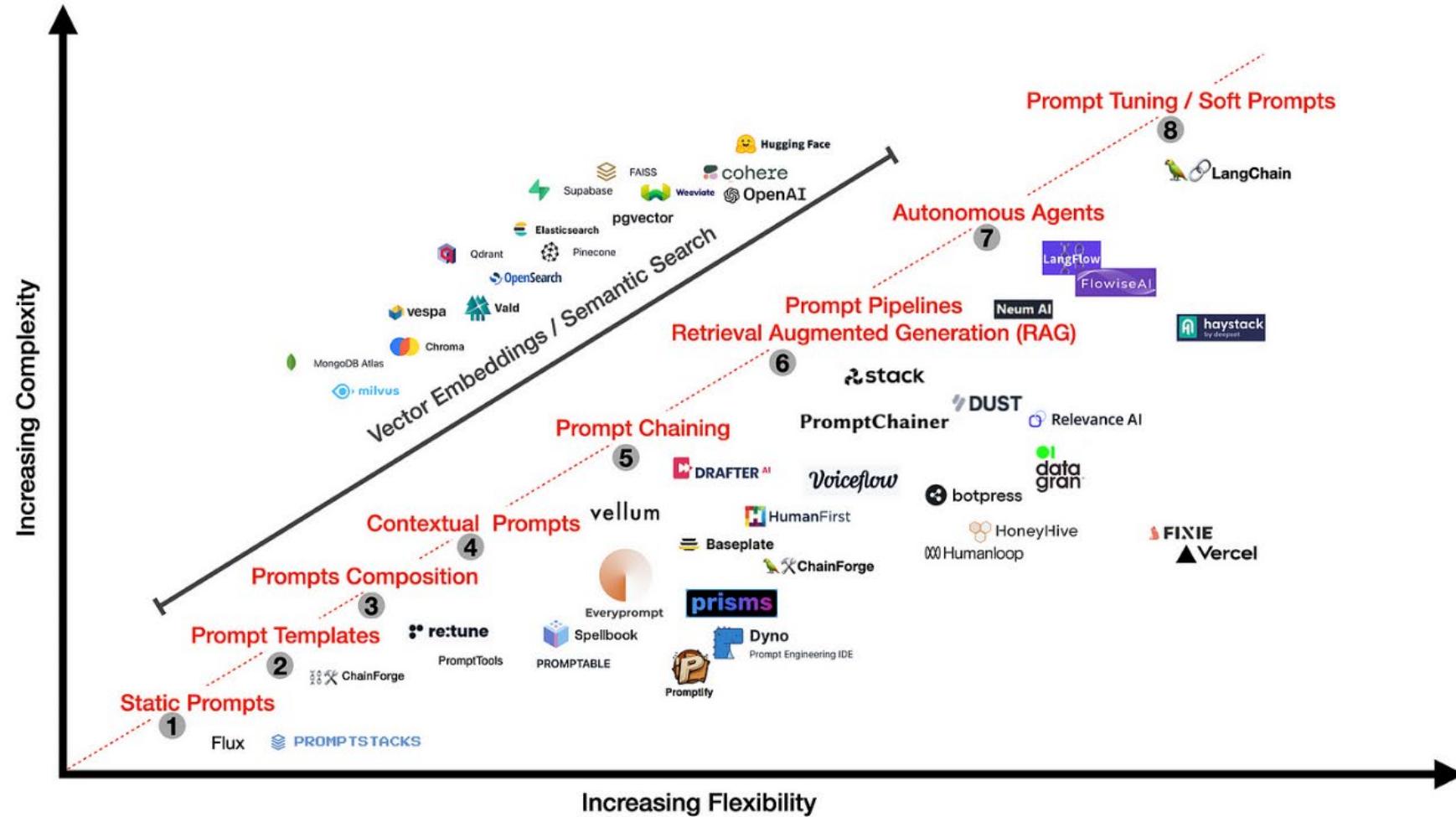
Waarom heeft Generatieve AI zo'n enorme impact op onze samenleving?



Ontstaansgeschiedenis + evolutie van grote taal modellen {LLM}



LLM implementations



<https://blogs.novita.ai/exploring-architectural-structures-and-functional-capacities-of-langs/>

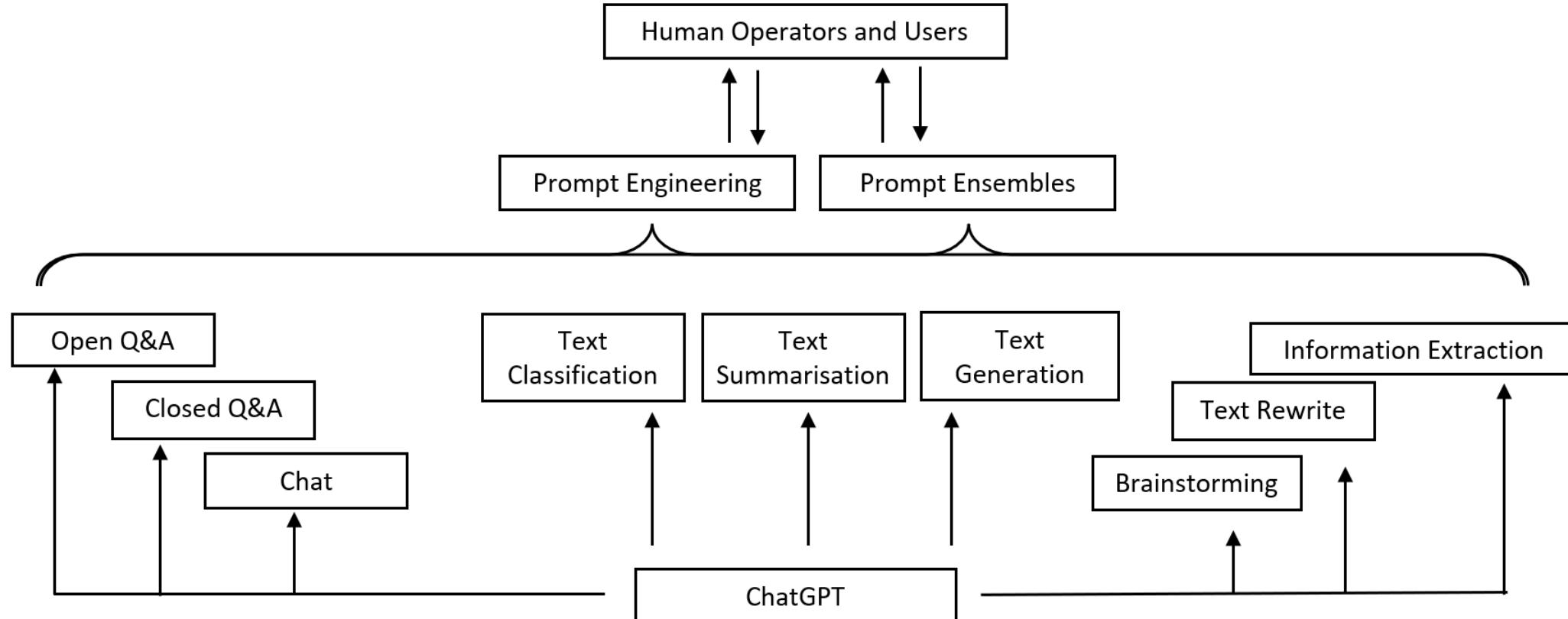
Kernpunten 3.2 De nieuwe infrastructuur van de publieke ruimte

- Technologiebedrijven hebben een infrastructuur ontworpen voor de online ordening van informatie en de interacties tussen gebruikers en content- en dienstenleveraars.
- Deze infrastructuur is dusdanig efficiënt, aantrekkelijk en gebruiksvriendelijk ontworpen dat hij onmisbaar is geworden voor het functioneren van de publieke ruimte. Op platforms vindt het grootste gedeelte van het publieke debat plaats, inclusief persoonlijke en zakelijke communicatie.
- De bedrijven achter deze platforms hebben hierdoor een nieuwe rol verkregen. Zij beheren de toegang tot de publieke ruimte en ordenen de informatie die daarbinnen circuleert.
- Platformisering speelt vooral bij uitgevers; radio (ether) en televisie zijn nog altijd belangrijke en veelgebruikte infrastructuren. Toch ondervinden ook zij in toenemende mate concurrentie van technologiebedrijven als het gaat om de distributie van content.



Ontwikkeling internet	Instantie	Wijze van aggregatie	Voorbeelden
Fase 1	Portals/startpagina's	Statisch overzicht websites met hyperlinks, gegenereerd door gebruikers	Yahoo! Startpagina.nl
Fase 2	Zoekmachines	Algoritmische selectie van internetbronnen op basis van zoekopdrachten van gebruiker	Google Bing Ilse
Fase 3	Sociale media	Permanente stroom (grotendeels) algoritmisch geactiveerde en gepersonaliseerde contentbundels	TikTok, Facebook, Instagram, YouTube, X
Fase 4	AI-chatbots	Integrale, op maat gegenereerde content op basis van internetinformatie	ChatGPT, Llama

ChatBot Use-Cases



Conferences > 2023 IEEE International Conference on Big Data

ChatGPT and Generative AI Guidelines for Addressing Academic Integrity and Augmenting Pre-Existing Chatbots

Publisher: IEEE

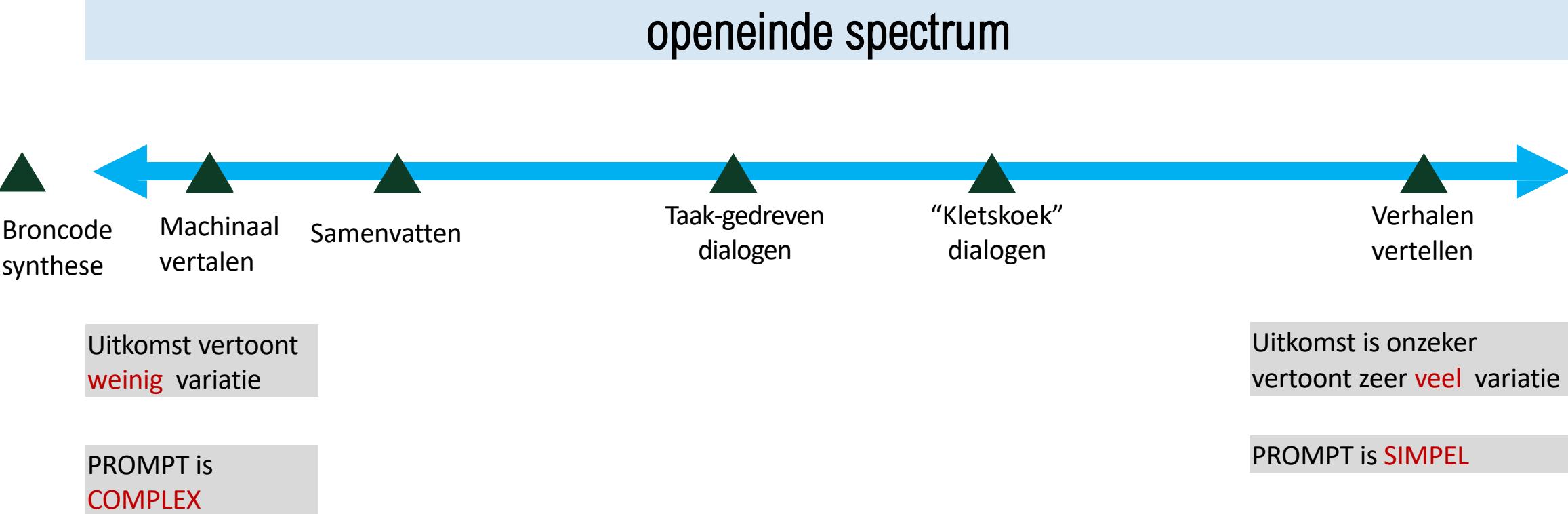
[Cite This](#)



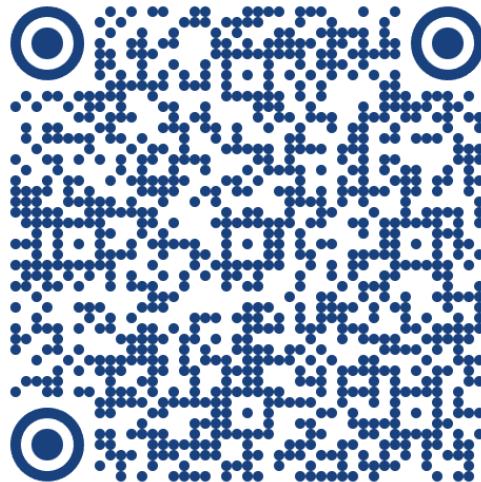
Daswin De Silva ; Nishan Mills ; Mona El-Ayoubi ; Milos Manic ; Damminda Alahakoon [All Authors](#)

635
Full
Text Views

Prompt Taxonomie



**Definiëring
karakteristieke kenmerken
GEN-AI *geeft inzicht***



<https://hai.stanford.edu/sites/default/files/2020-09/AI-Definitions-HAI.pdf>

AI *definities* volgens HAI

Intelligence might be defined as the ability to learn and perform suitable techniques to solve problems and achieve goals, appropriate to the context in an uncertain, ever-varying world. A fully pre-programmed factory robot is flexible, accurate, and consistent but not intelligent.

Artificial Intelligence (AI), a term coined by emeritus Stanford Professor John McCarthy in 1955, was defined by him as “the science and engineering of making intelligent machines”. Much research has humans program machines to behave in a clever way, like playing chess, but, today, we emphasize machines that can learn, at least somewhat like human beings do.

Autonomous systems can independently plan and decide sequences of steps to achieve a specified goal without micro-management. A hospital delivery robot must autonomously navigate busy corridors to succeed in its task. In AI, autonomy doesn’t have the sense of being self-governing common in politics or biology.

Machine Learning (ML) is the part of AI studying how computer agents can improve their perception, knowledge, thinking, or actions based on experience or data. For this, ML draws from computer science, statistics, psychology, neuroscience, economics and control theory.

In **supervised learning**, a computer learns to predict human-given labels, such as dog breed based on labeled dog pictures; **unsupervised learning** does not require labels, sometimes making its own prediction tasks such as trying to predict each successive word in a sentence; **reinforcement learning** lets an agent

learn action sequences that optimize its total rewards, such as winning games, without explicit examples of good techniques, enabling autonomy.

Deep Learning is the use of large multi-layer (**artificial neural networks**) that compute with continuous (real number) representations, a little like the hierarchically organized neurons in human brains. It is currently the most successful ML approach, usable for all types of ML, with better generalization from small data and better scaling to big data and compute budgets.

An **algorithm** lists the precise steps to take, such as a person writes in a computer program. AI systems contain algorithms, but often just for a few parts like a learning or reward calculation method. Much of their behavior emerges via learning from data or experience, a sea change in system design that Stanford alumnus Andrej Karpathy dubbed **Software 2.0**.

Narrow AI is intelligent systems for one particular thing, e.g., **speech** or **facial recognition**. **Human-level AI**, or **Artificial General Intelligence (AGI)**, seeks broadly intelligent, context-aware machines. It is needed for effective **social chatbots** or **human-robot interaction**.

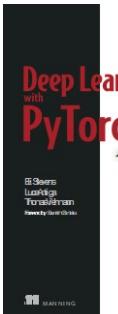
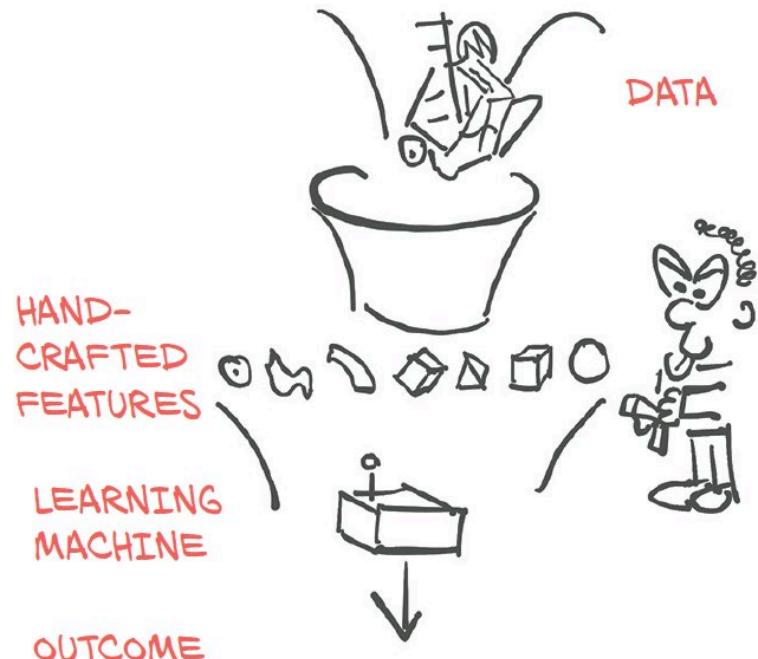
Human-Centered Artificial Intelligence is AI that seeks to augment the abilities of, address the societal needs of, and draw inspiration from human beings. It researches and builds effective partners and tools for people, such as a robot helper and companion for the elderly.

Text by Professor Christopher Manning, September 2020

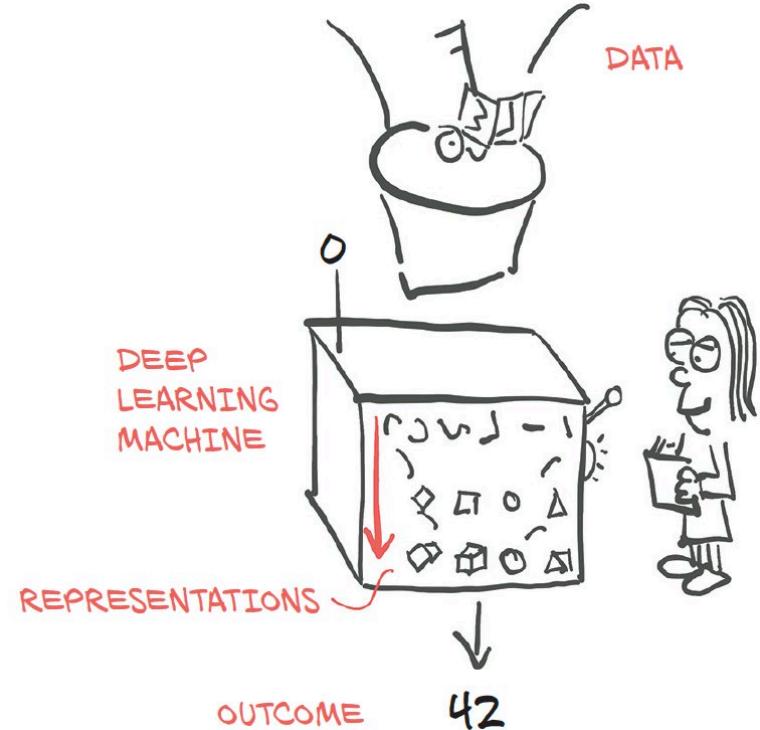
{AI Paradigm-shift → NO Human input needed}

More data, parameters & computing power | Less human-in-the-loop

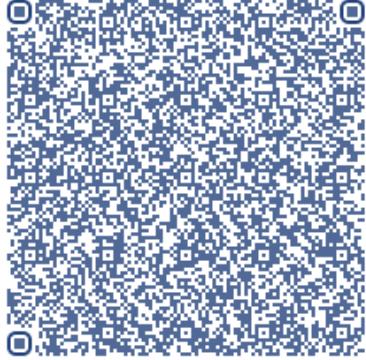
Machine Learning Paradigm {ML}



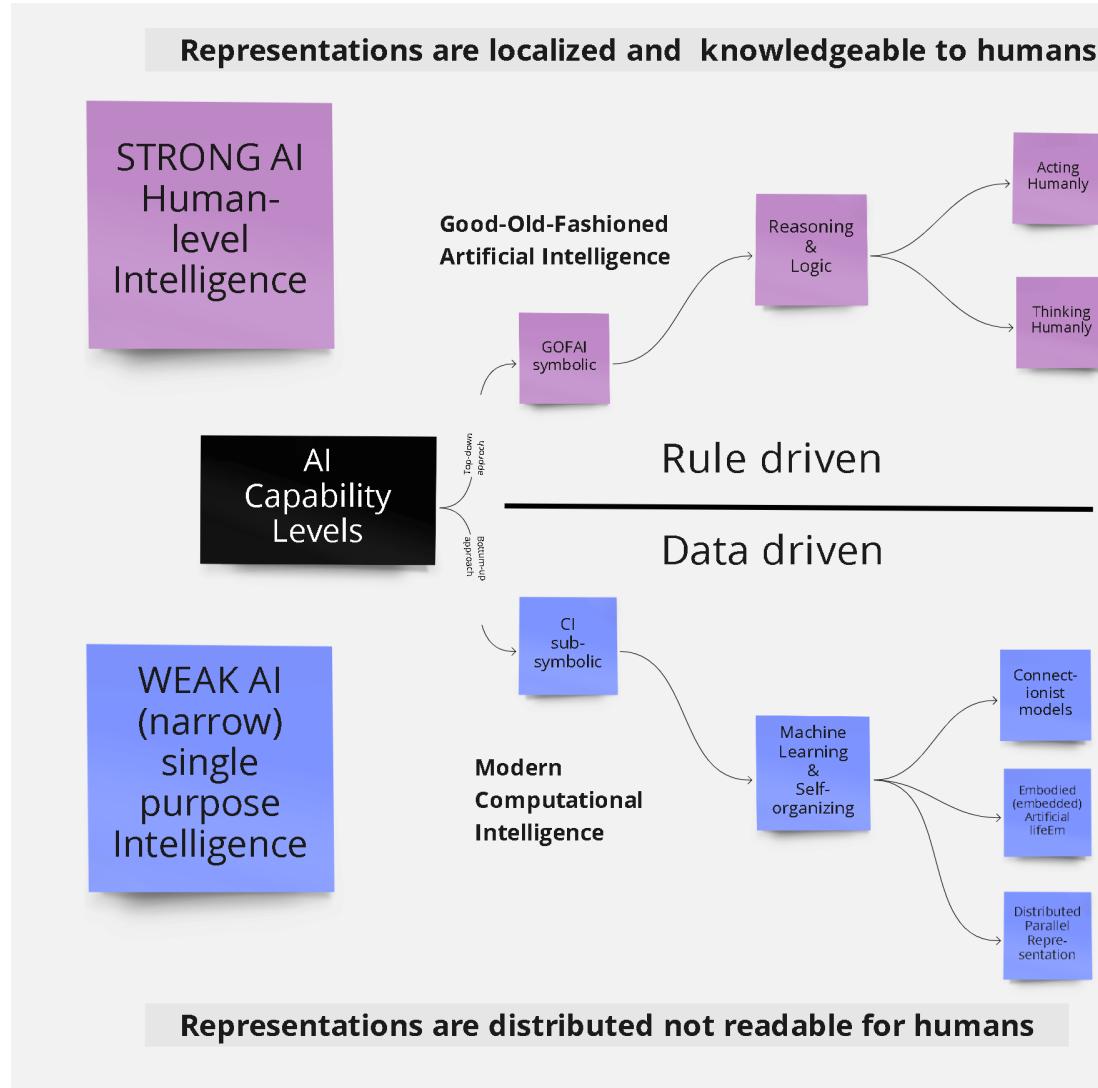
Deep Learning Paradigm {DL}



AI-taxonomie



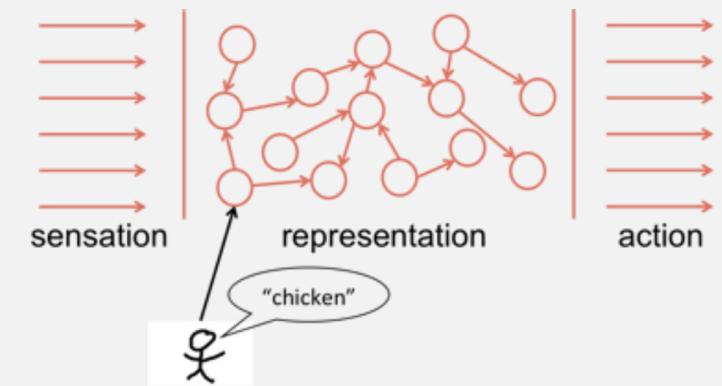
https://www.researchgate.net/publication/359424818_Designing_Neural_Networks_Through_Sensory_Ecology_Biology_to_the_rescue_of_AI_Produced_by_Living-Lab_AiRA_Hub_voor_Data_Responsible_AI_Hogeschool_Rotterdam_Lunch-Lezing_Creating-010_FEB_2022



SYMBOLIC

According to research at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself, but the word as a whole.

SUBSYMBOLIC

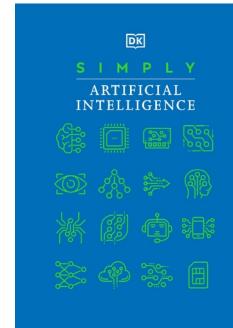
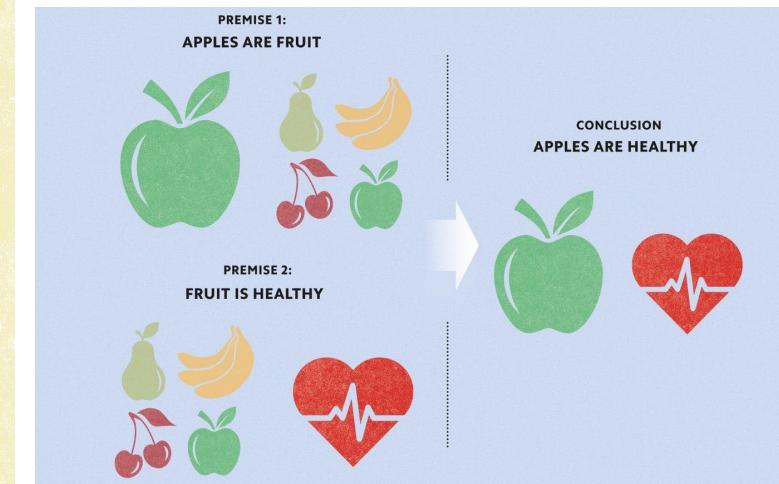
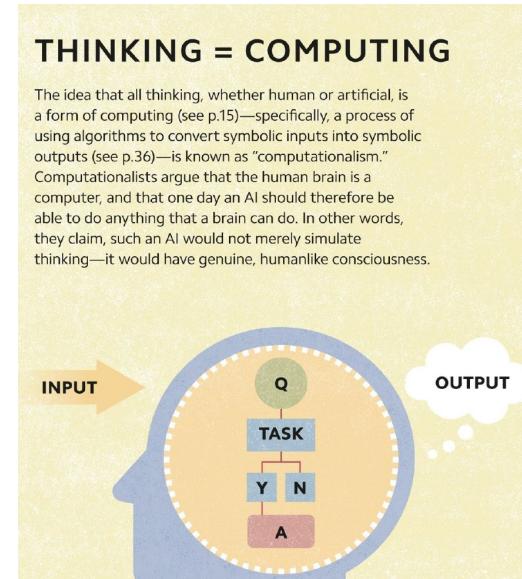


SYMBOLIC AI

TOP-DOWN / open-loop

Intelligentie (denken) is een vorm van Logica "Rekenkracht" (**computationalism**)

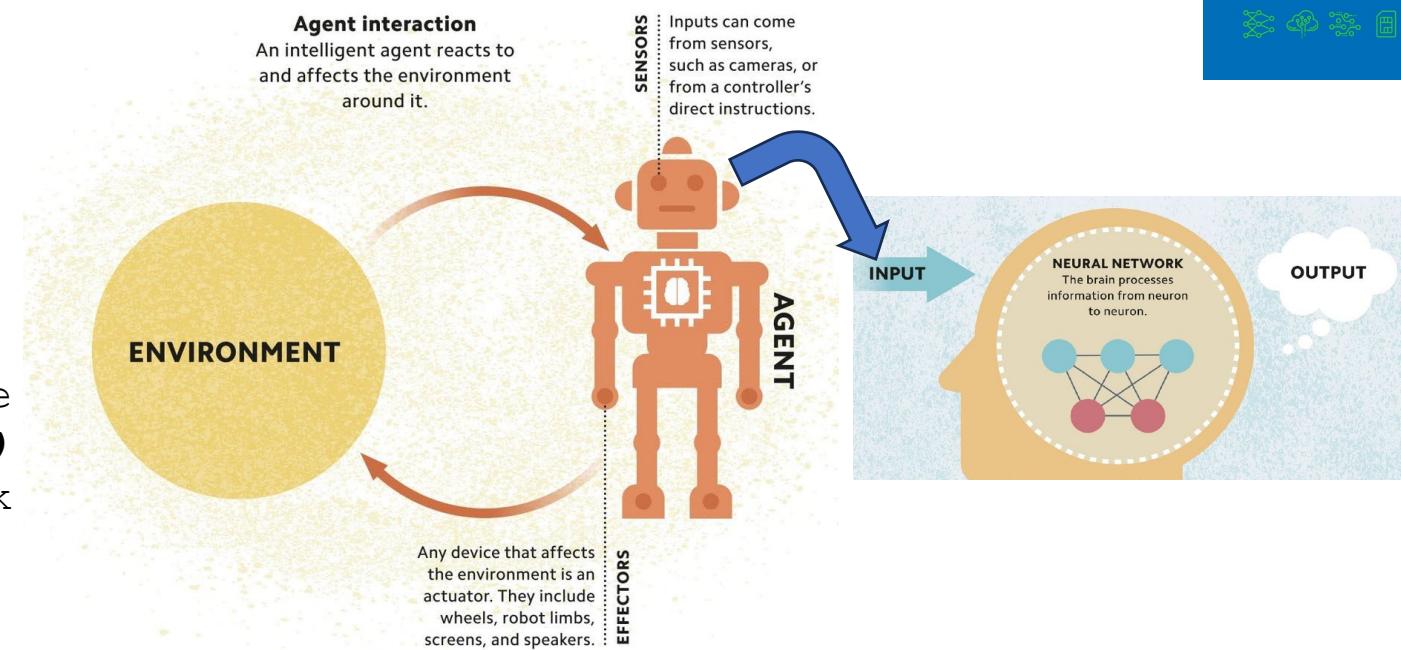
Hersen en zijn een "following the rules" Computer



SUBSYMBOLIC AI

BOTTOM-UP / closed feedback-loop

Intelligentie represeneert adaptatief-leervermogen (**trial & error**) dat wordt mogelijk gemaakt door netwerken bestaande uit simpele rekeneenheden (**connectionism**) aangestuurd door een algoritme + feedback loop (**Cybernetica**)

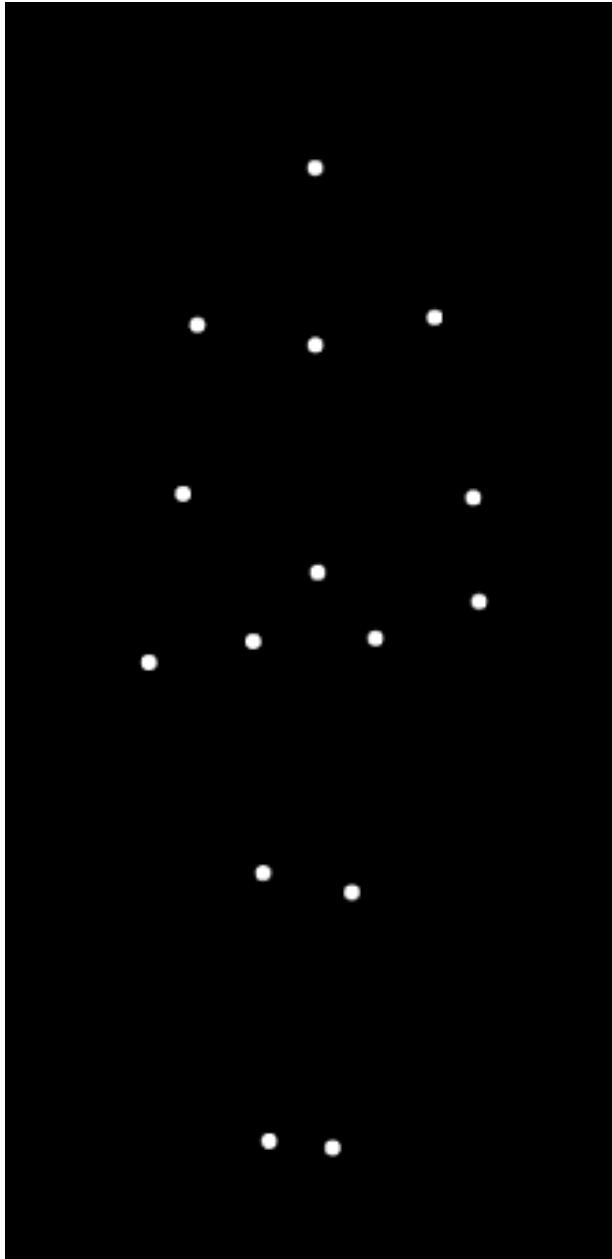


{Top Down}

Biological motion is a prime example of computationalism as performed by the human brain.

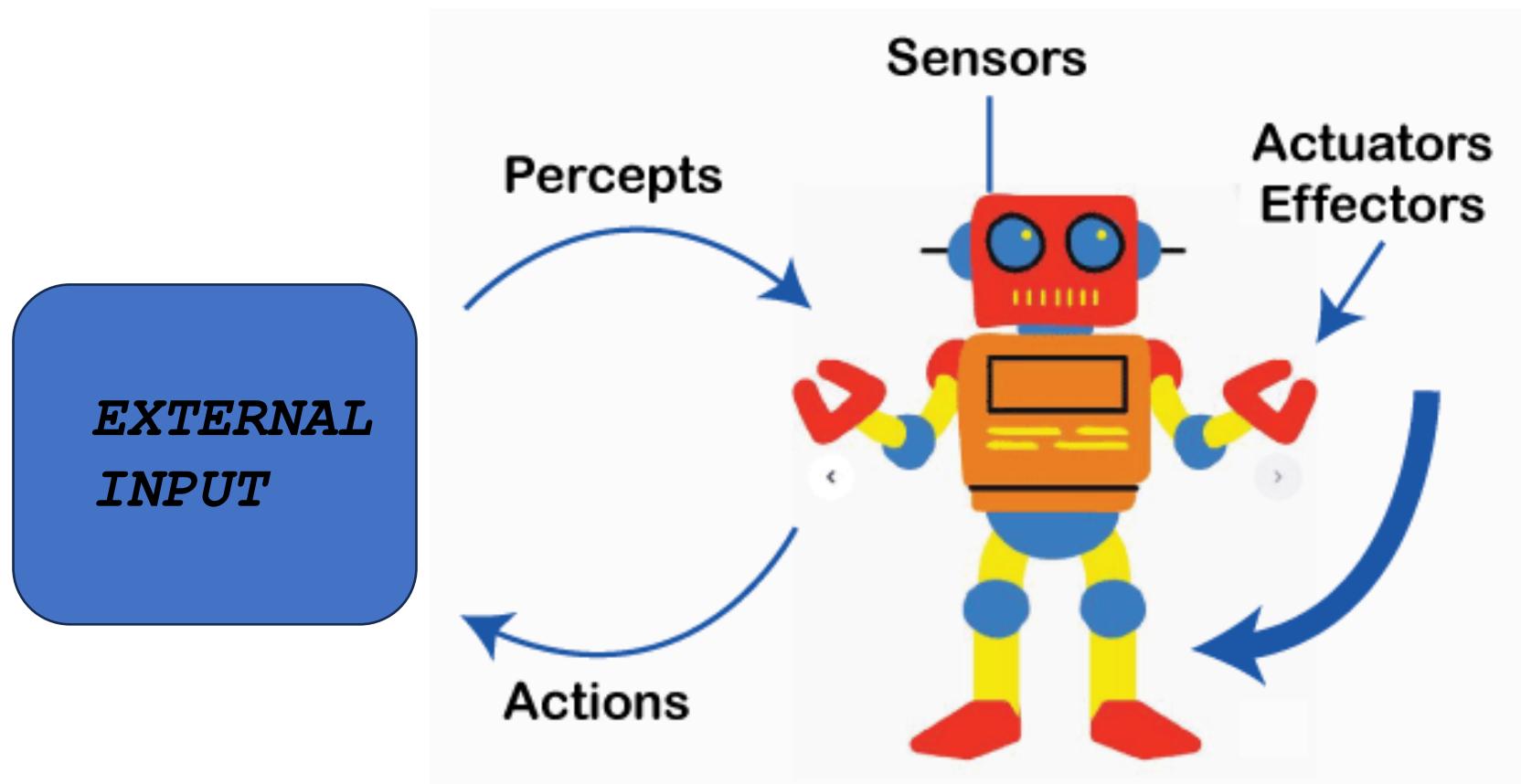
Only 12 moving-dots are needed to determine direction of motion.

Bradshaw, M. & van der Willigen, R. F., (1999).
The walker's direction affects the perception of biological motion.
In M. A. Grealy & J. A. Thomson (Eds.),
Studies in perception and action V (pp. 3–6). London: Academic Press.



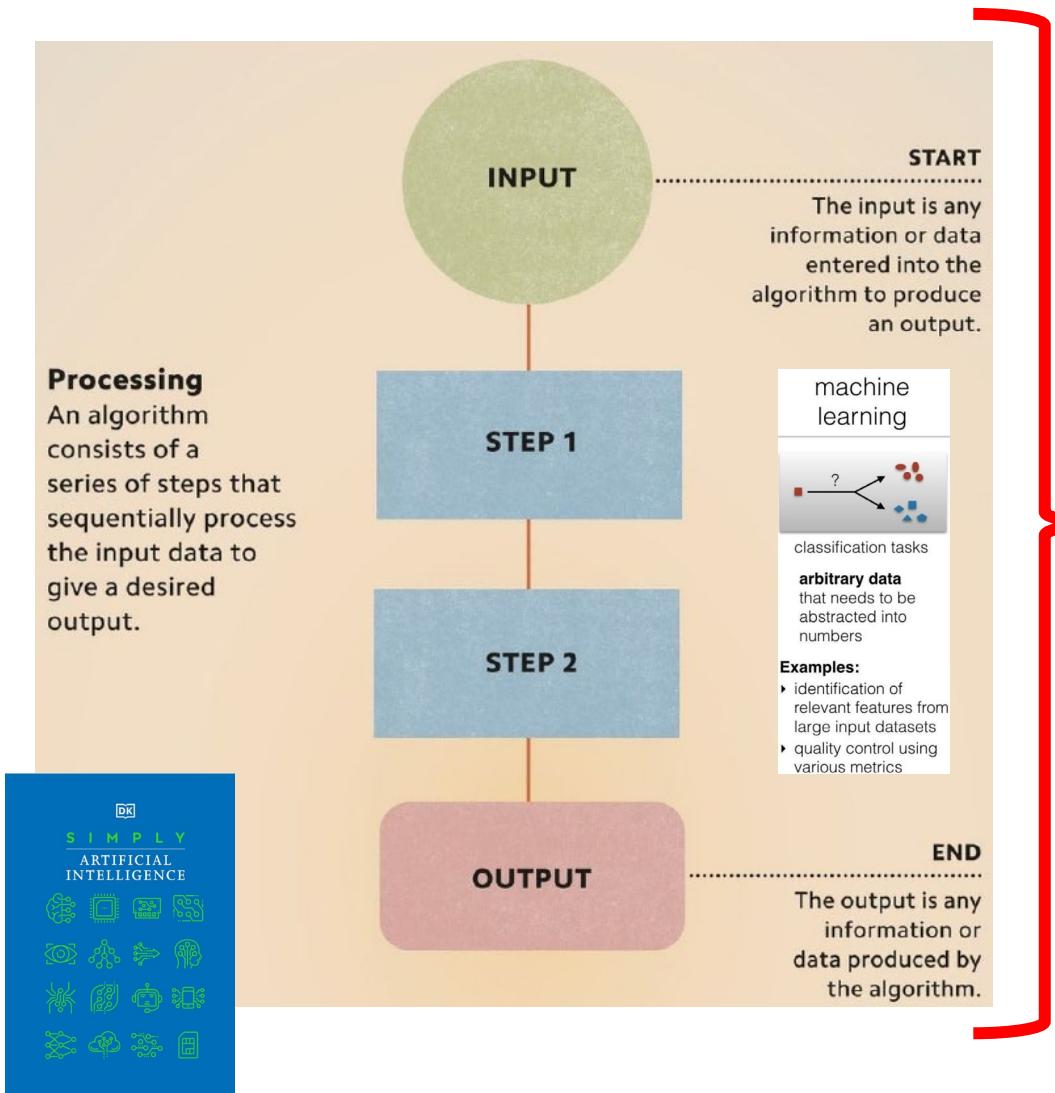
{MULTIMODALE AGENT}

AGENT representeert input / output model

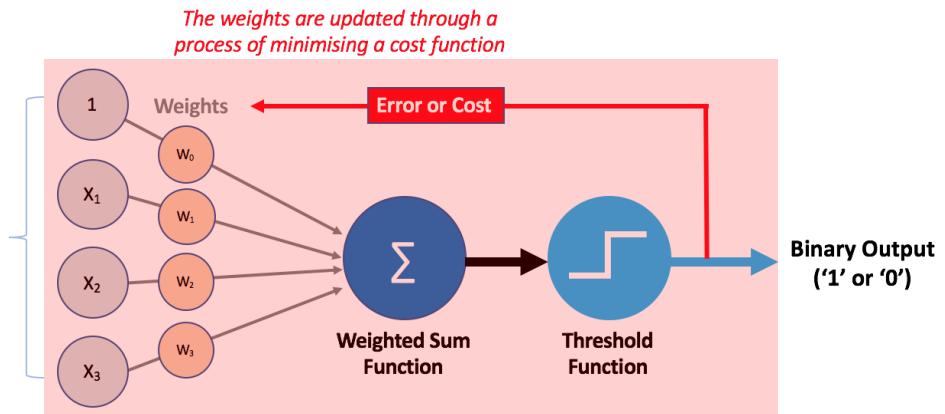
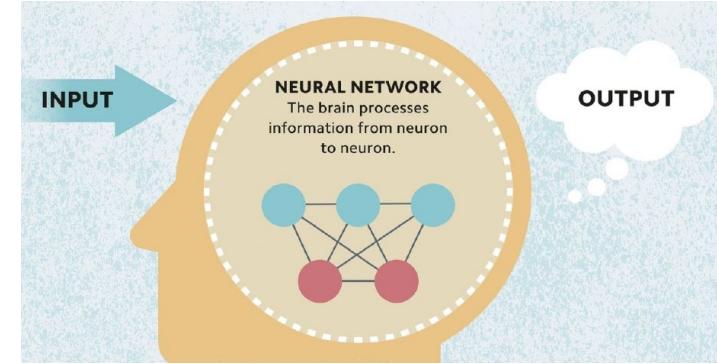


{ALGORITME}

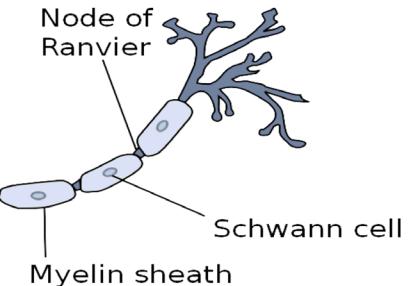
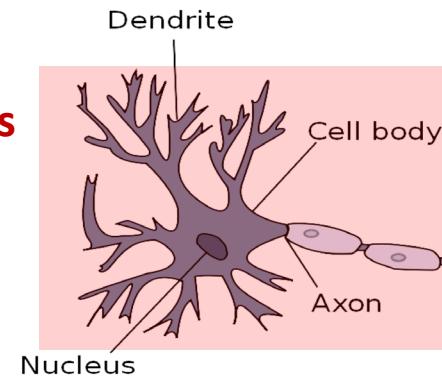
<https://learn.microsoft.com/nl-nl/dotnet/machine-learning/deep-learning-overview>



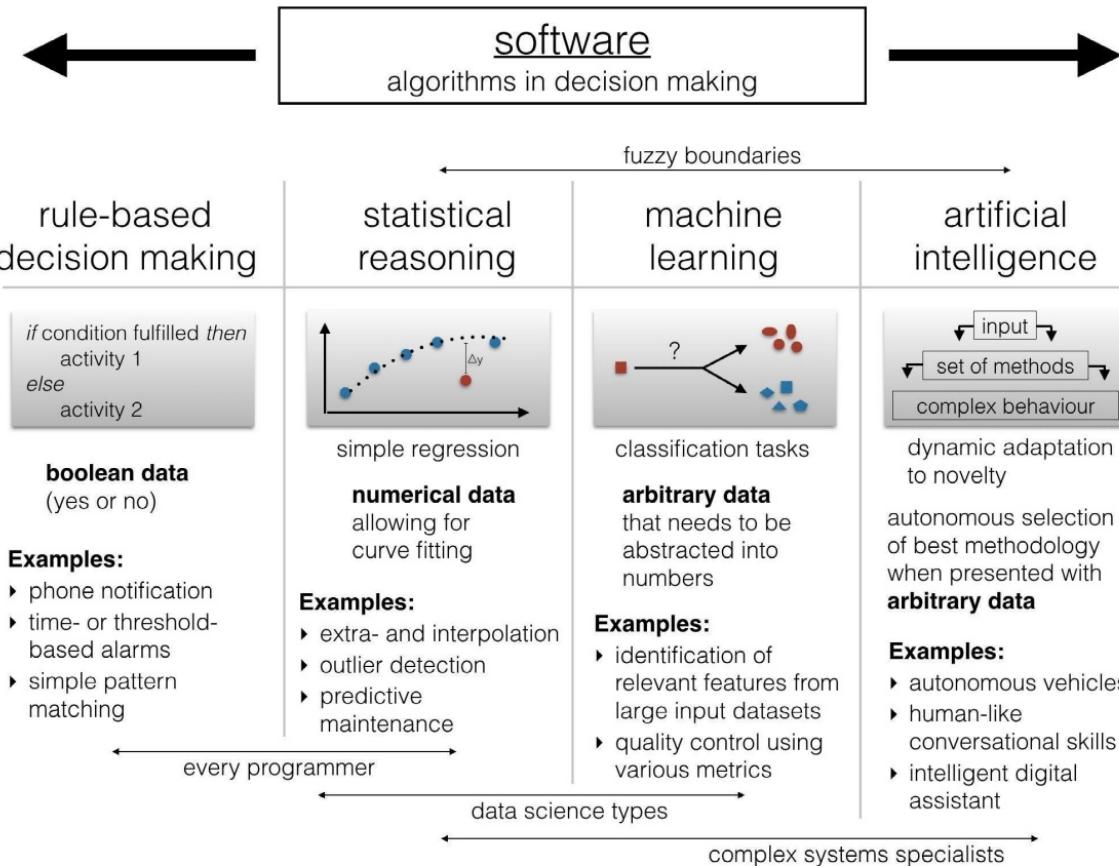
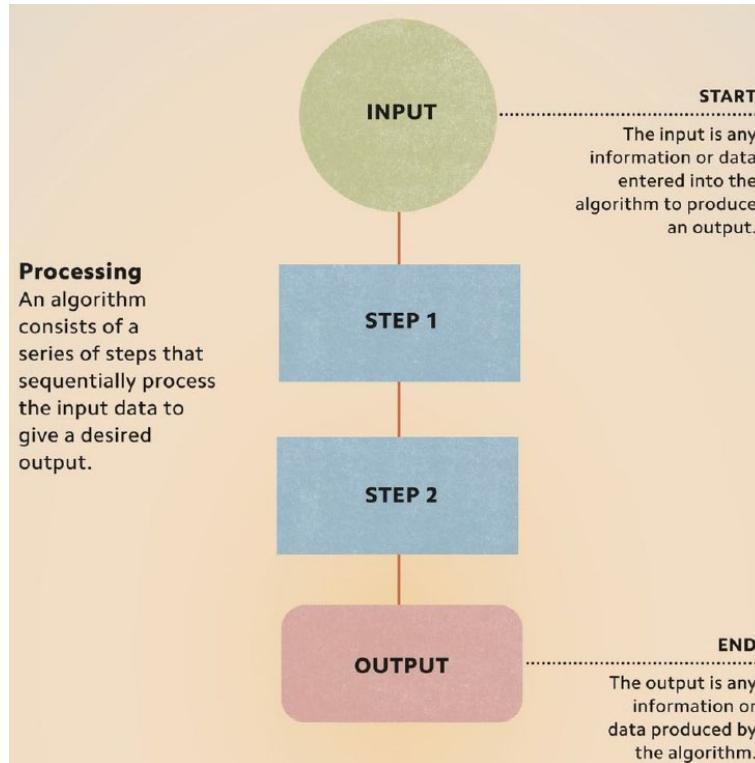
{AI}



Names for Artificial Neurons
{unit}
{cell}
{node}
{perceptron}

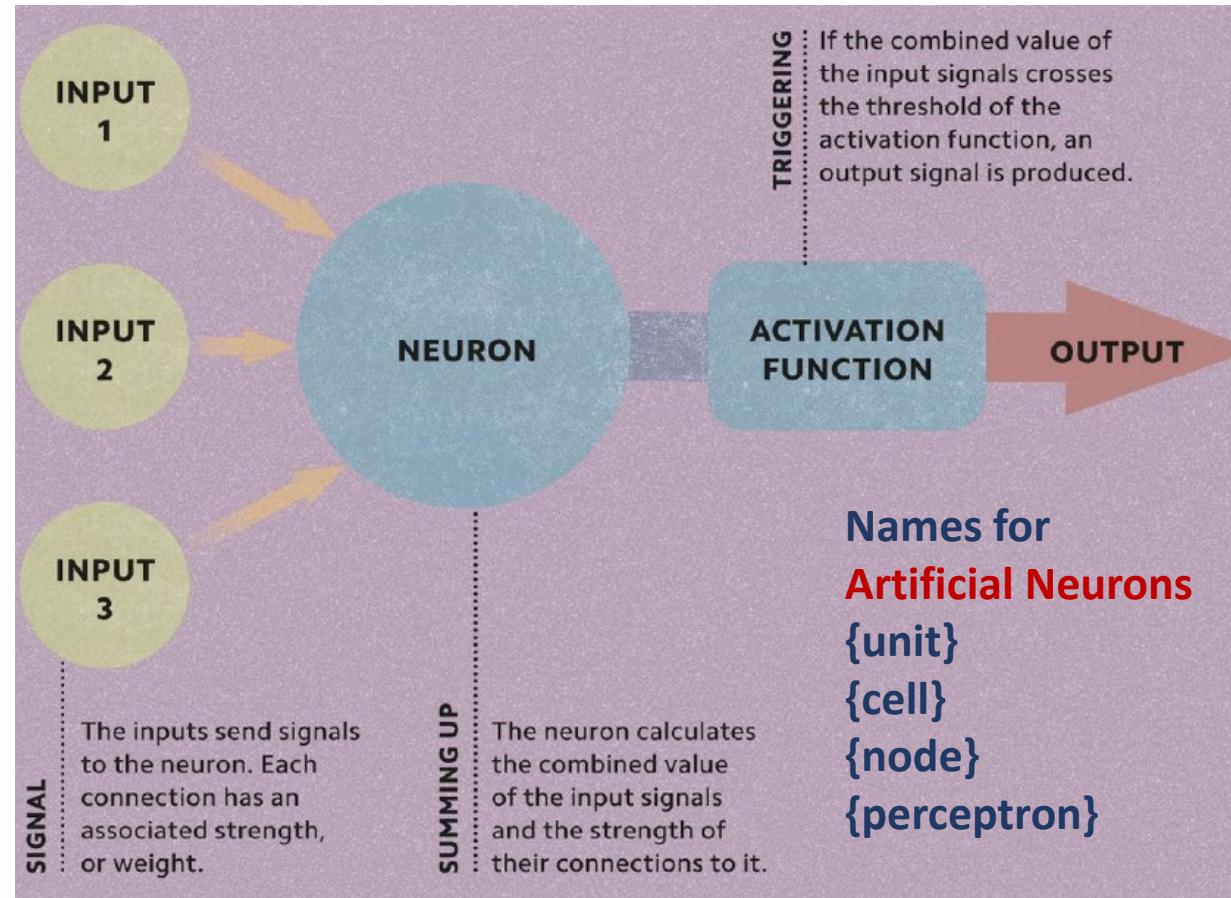


{COMPUTER ALGORITHM}



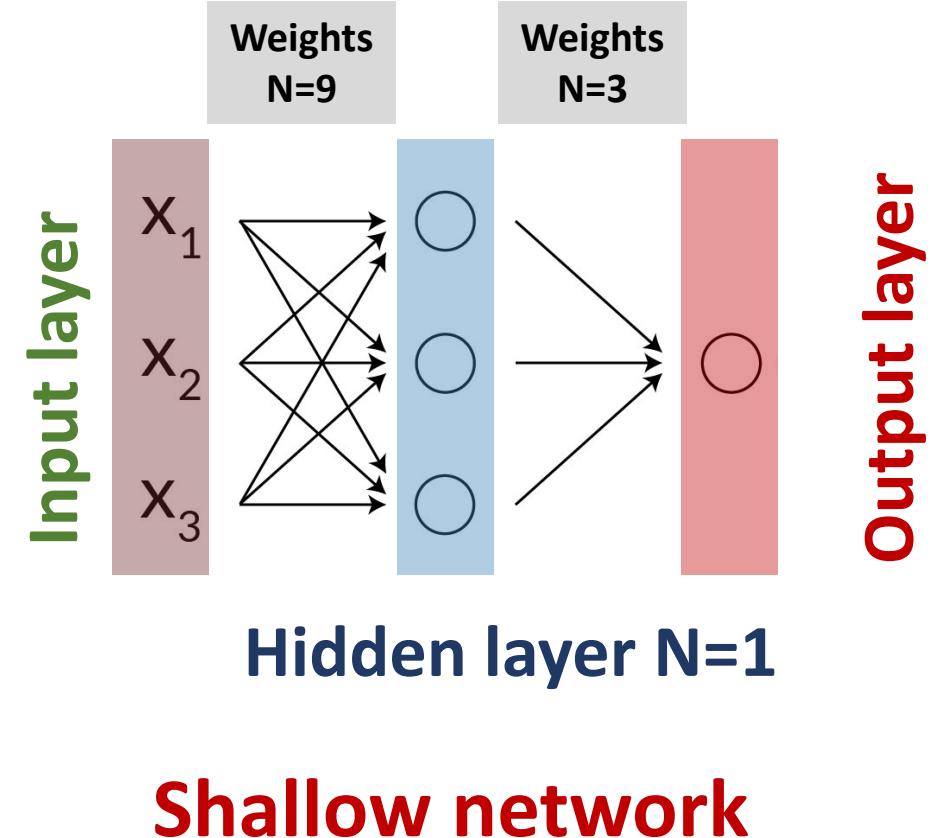
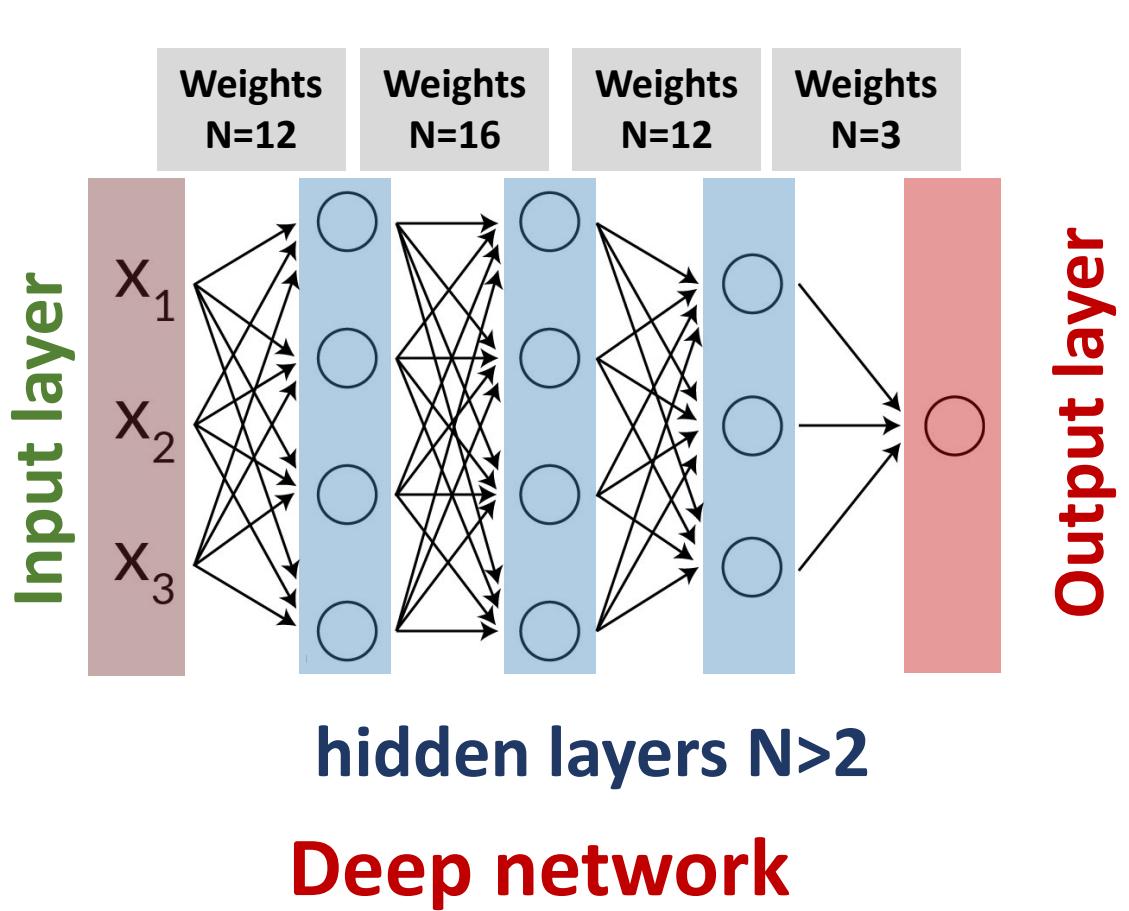
{KUNSTMATIGE NEURON}

computation



{Neuraal Netwerk [NN] Typen}

Neural Network {NN} Lagen Architectuur

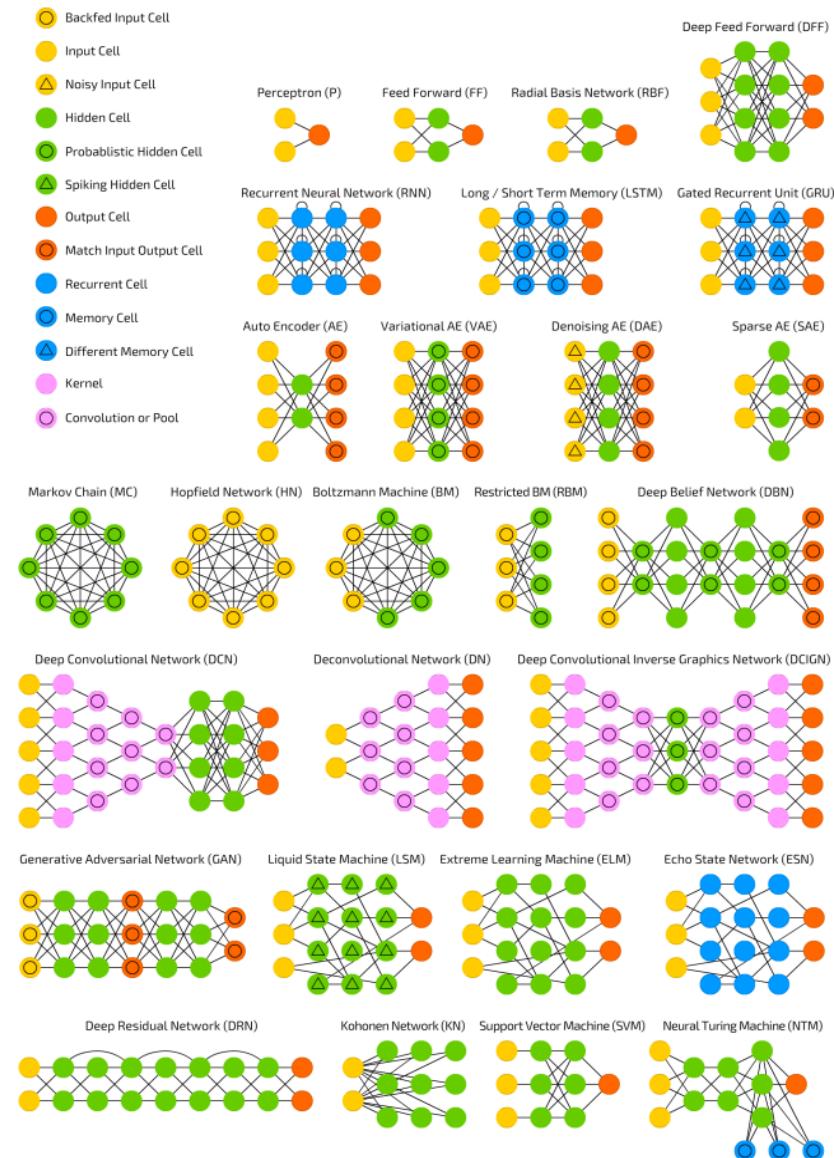


{Topology}

Topology of a neural network refers to the way artificial neurons are connected to form a network.

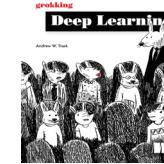
Form follows function!
The topology of a network determines the degree of perplexity of the tasks it can

<https://pub.towardsai.net/main-types-of-neural-networks-and-its-applications-tutorial-734480d7ec8e>



{Bottum-UP: Machinaal Leren [ML]}

What is machine learning?



“ A field of study that gives computers the ability to learn without being explicitly programmed.

—Attributed to Arthur Samuel

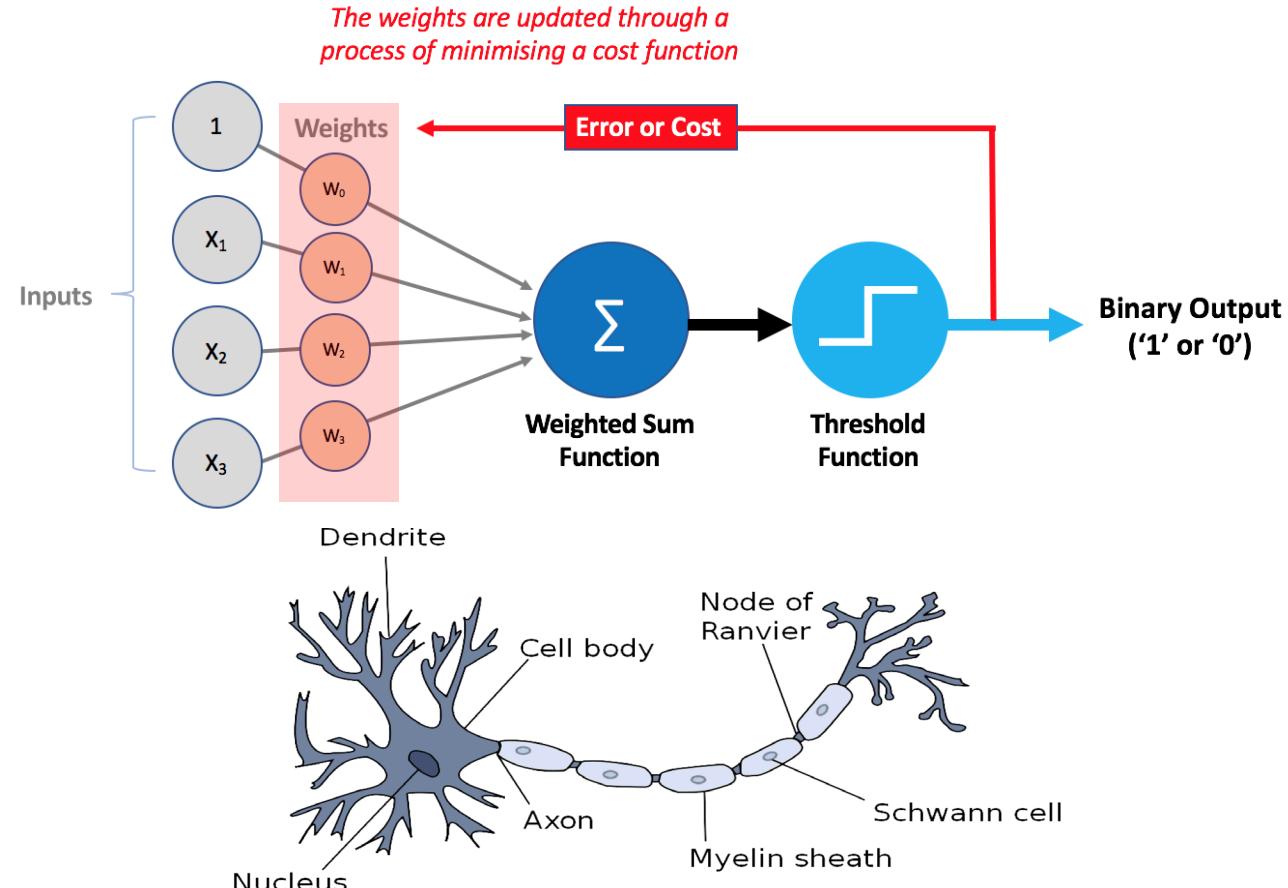
Given that deep learning is a subset of machine learning, what is machine learning? Most generally, it is what its name implies. Machine learning is a subfield of computer science wherein *machines learn* to perform tasks for which they were *not explicitly programmed*. In short, machines observe a pattern and attempt to imitate it in some way that can be either direct or indirect.

Machine learning \approx Monkey see, monkey do

{Backpropagation ML-Leer-regel nodig om een NN te laten leren}

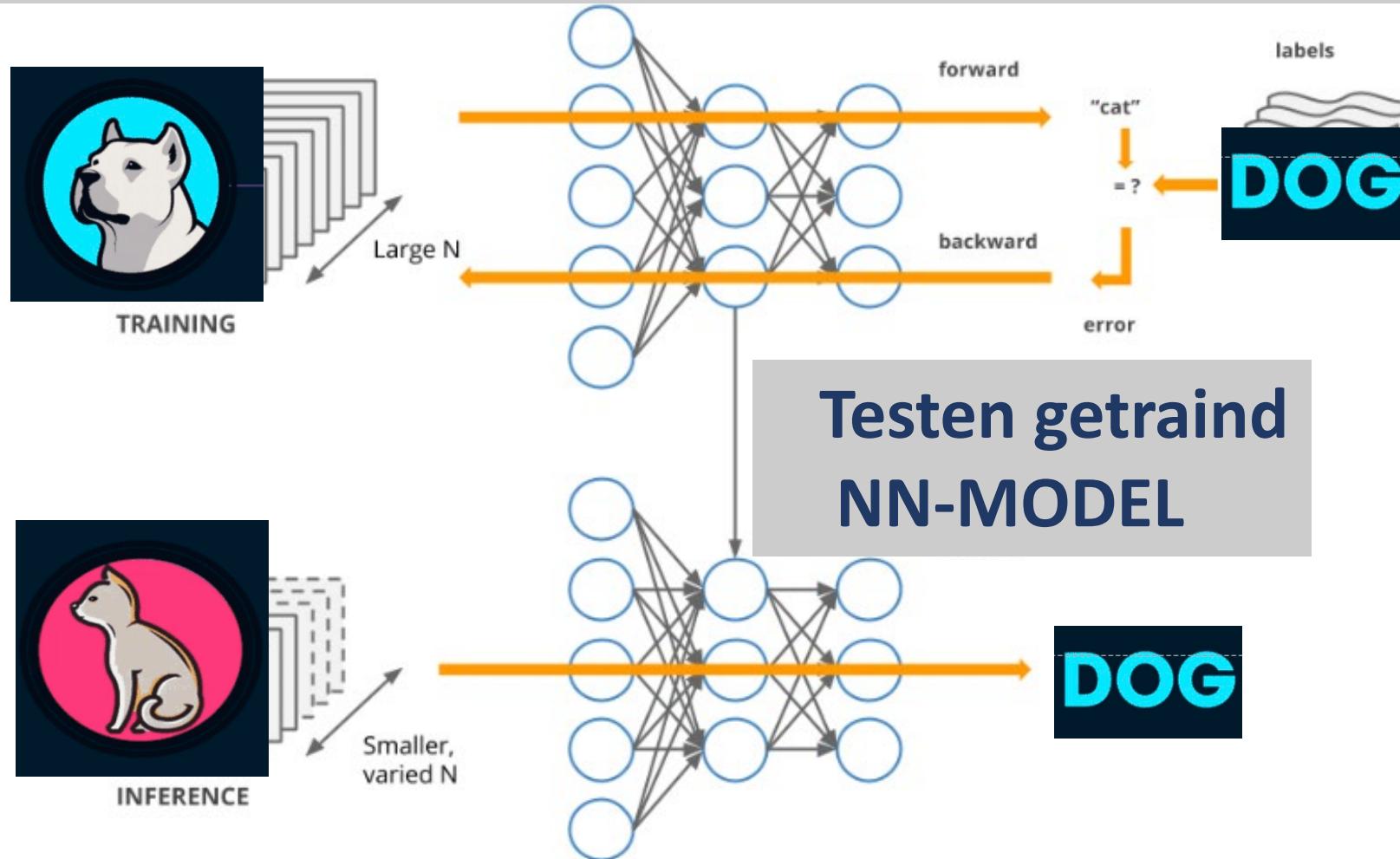
**Names for
Artificial Neurons**
{unit}
{cell}
{node}
{perceptron}

Machinaal Leren [ML]

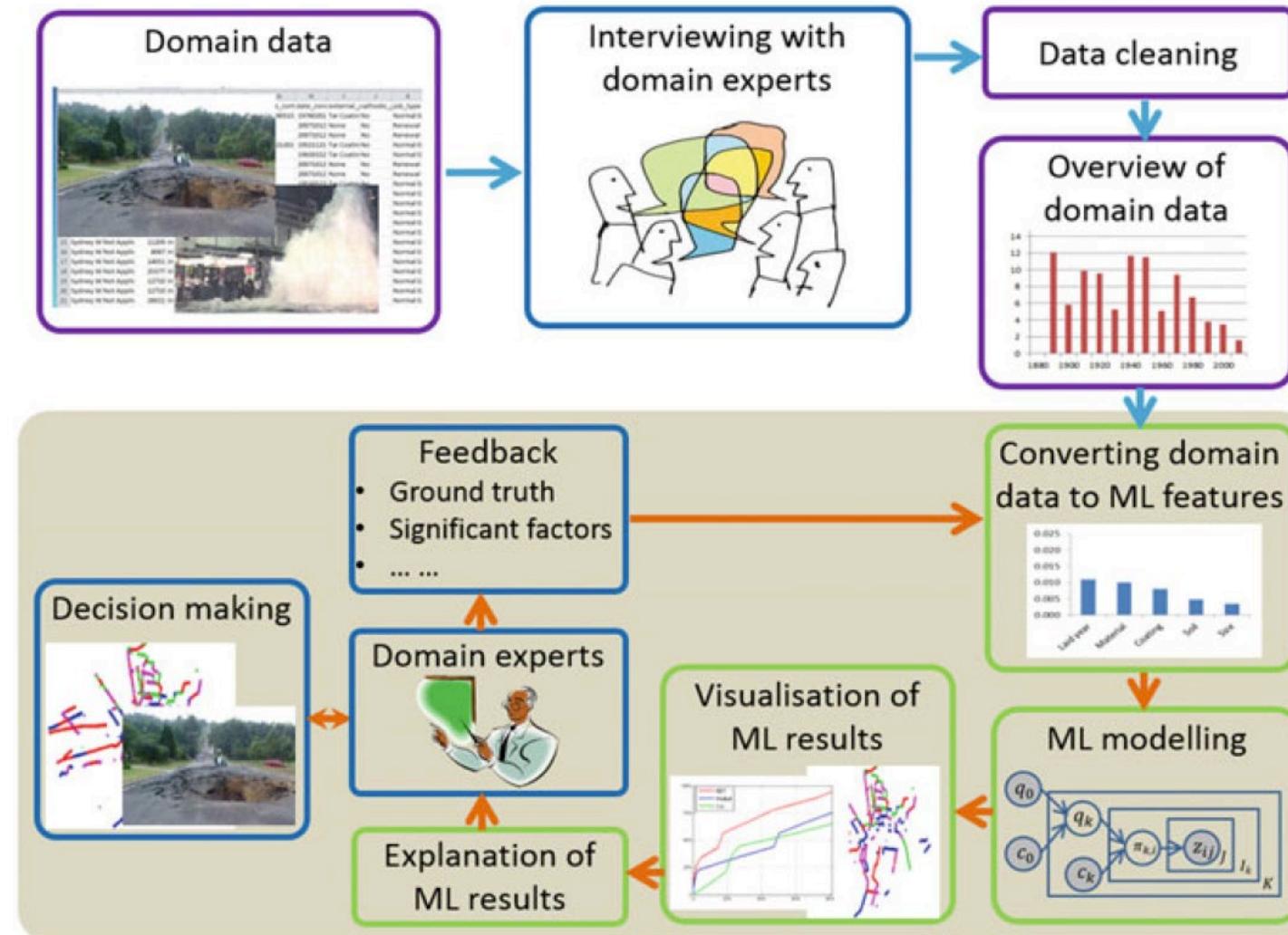
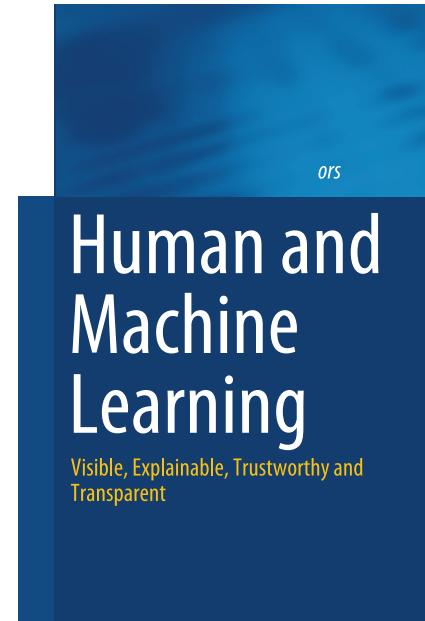


{NN Trainen (leren) versus Testen}

Machinaal Leren [ML] == trainen van het NN



{Modeling: human-in-the-loop}



{Top-down}

Top-down Encoding Capacity increases by adding hidden layers

What are the limits of deep learning?

The much-hyped artificial intelligence approach boasts impressive feats but still falls short of human brainpower. Researchers are determined to figure out what's missing.

M. Mitchell Waldrop, Science Writer

There's no mistaking the image: It's a banana—a big, ripe, bright-yellow banana. Yet the artificial intelligence (AI) identifies it as a toaster, even though it was trained with the same powerful and oft-publicized deep-learning techniques that have produced a white-hot revolution in driverless cars, speech understanding, and a multitude of other AI applications. That means the AI was shown several thousand photos of bananas, slugs, snails, and similar-looking objects, like so many flash cards, and then drilled on the answers until it had the classification down cold. And yet this advanced system was quite easily confused—all it took was a little day-glow sticker, digitally pasted in one corner of the image.

This example of what deep-learning researchers call an "adversarial attack," discovered by the Google Brain team in Mountain View, CA (1), highlights just how far AI still has to go before it remotely approaches human capabilities. "I initially thought that adversarial examples were just an annoyance," says Geoffrey Hinton, a computer scientist at the University of Toronto and one of the pioneers of deep learning. "But I now think they're probably quite profound. They tell us that we're doing something wrong."

That's a widely shared sentiment among AI practitioners, any of whom can easily rattle off a long list of deep learning's drawbacks. In addition to its vulnerability

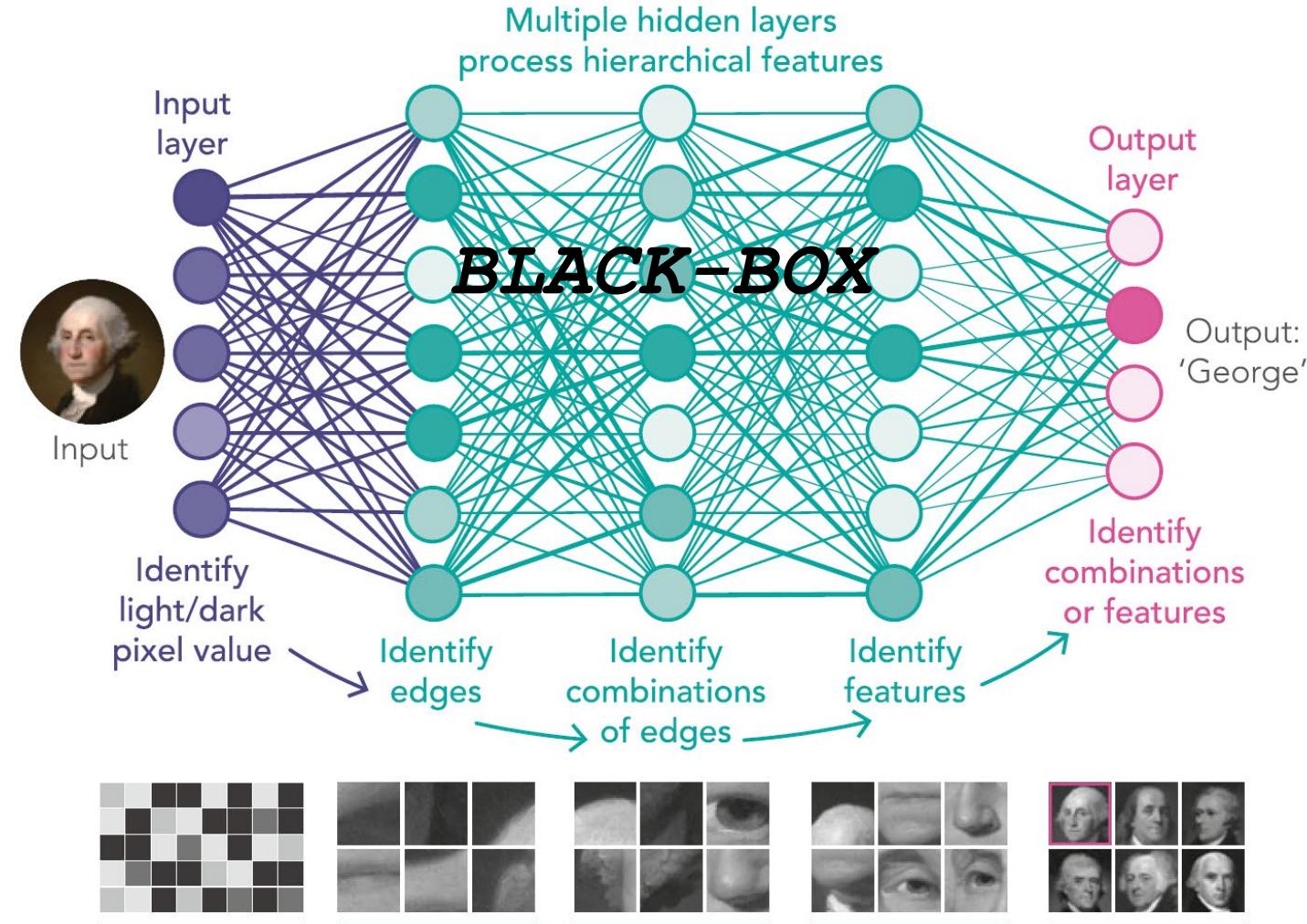


Apparent shortcomings in deep-learning approaches have raised concerns among researchers and the general public as technologies such as driverless cars, which use deep-learning techniques to navigate, get involved in well-publicized mishaps. Image credit: Shutterstock.com/MONOPOLY919.

Published under the PNAS license.

January 22, 2019 | vol. 116 | no. 4

www.pnas.org/cgi/doi/10.1073/pnas.1821594116



{Big-data}

Big-data is needed to avoid hand-crafted feature extraction

A Unified Approach to Interpreting Model Predictions

Scott M. Lundberg
 Paul G. Allen School of Computer Science
 University of Washington
 Seattle, WA 98105
 slundb@cs.washington.edu

Su-In Lee
 Paul G. Allen School of Computer Science
 Department of Genome Sciences
 University of Washington
 Seattle, WA 98105
 suinlee@cs.washington.edu

Abstract

Understanding why a model makes a certain prediction can be as crucial as the prediction's accuracy in many applications. However, accuracy for large modern datasets is often achieved by complex models that even experts have trouble interpreting, such as ensemble or deep learning models, creating a tension between accuracy and *interpretability*. In response, various methods have recently been proposed to help users interpret the predictions of complex models, but it is often unclear how these methods are related and when one method is preferable over another. To address this problem, we present a unified framework for interpreting predictions, SHAP (SHapley Additive exPlanations). SHAP assigns each feature an importance value for a particular prediction. Its novel components include: (1) the identification of a new class of additive feature importance measures, and (2) theoretical results showing there is a unique solution in this class with a set of desirable properties. The new class unifies six existing methods, notable because several recent methods in the class lack the proposed desirable properties. Based on insights from this unification, we present new methods that show improved computational performance and/or better consistency with human intuition than previous approaches.

1 Introduction

The ability to correctly interpret a prediction model's output is extremely important. It engenders appropriate user trust, provides insight into how a model may be improved, and supports understanding of the process being modeled. In some applications, simple models (e.g., linear models) are often preferred for their ease of interpretation, even if they may be less accurate than complex ones. However, the growing availability of big data has increased the benefits of using complex models, so bringing to the forefront the trade-off between accuracy and interpretability of a model's output. A wide variety of different methods have been recently proposed to address this issue [5, 8, 9, 3, 4, 1]. But an understanding of how these methods relate and when one method is preferable to another is still lacking.

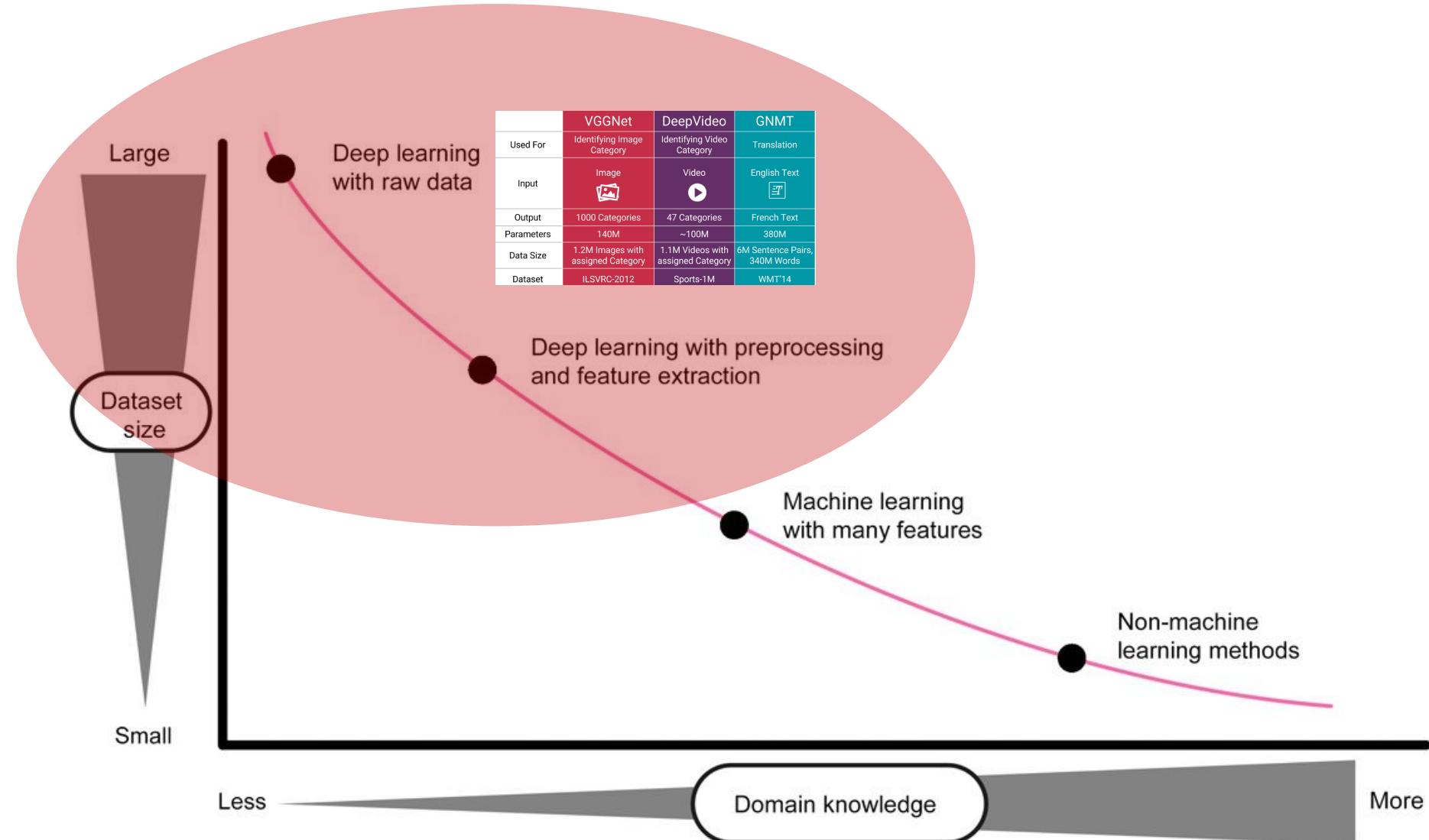
Here, we present a novel unified approach to interpreting model predictions.¹ Our approach leads to three potentially surprising results that bring clarity to the growing space of methods:

1. We introduce the perspective of viewing *any* explanation of a model's prediction as a model itself, which we term the *explanation model*. This lets us define the class of *additive feature attribution methods* (Section 2), which unifies six current methods.

¹<https://github.com/slundberg/shap>

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

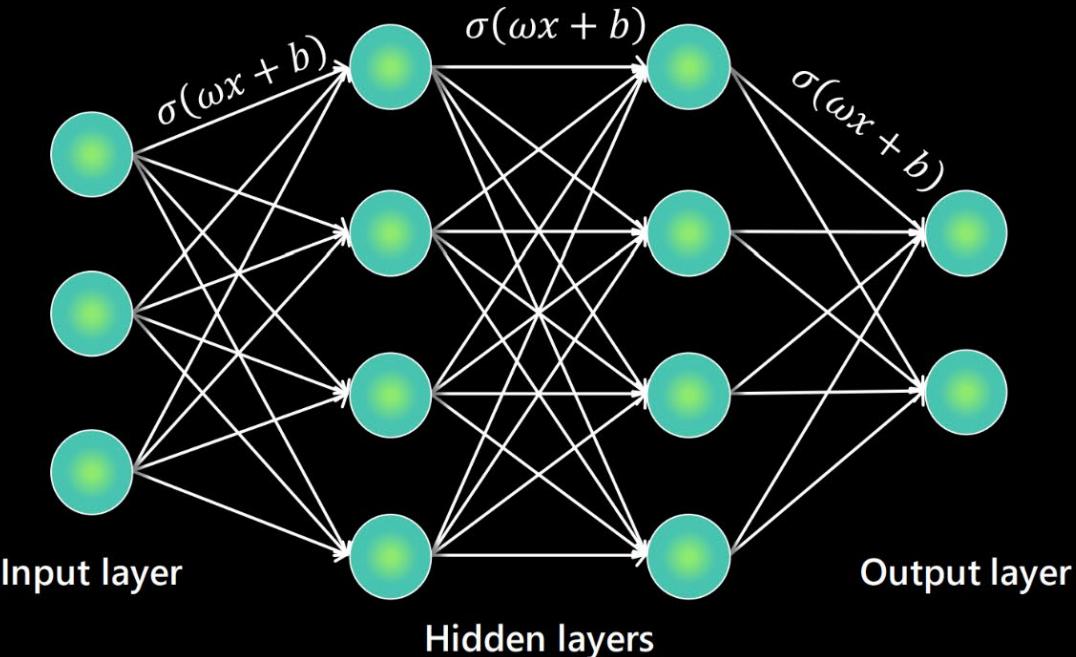
<https://proceedings.neurips.cc/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf>



{Diep Neuraal Netwerk [NN] ontstaat “vanzelf”}

How large are they?

BLACK-BOX



Function: weight * input plus bias

BERT Large - 2018

345M

GPT2 - 2019

1.5B

GPT3 - 2020

175B

Turing Megatron NLG
2021

530B

GPT4 – 2023

1.4T (estimated)

**Het gebruik van
publieke Grote Taal
Modellen zoals ChatGPT
is problematisch**

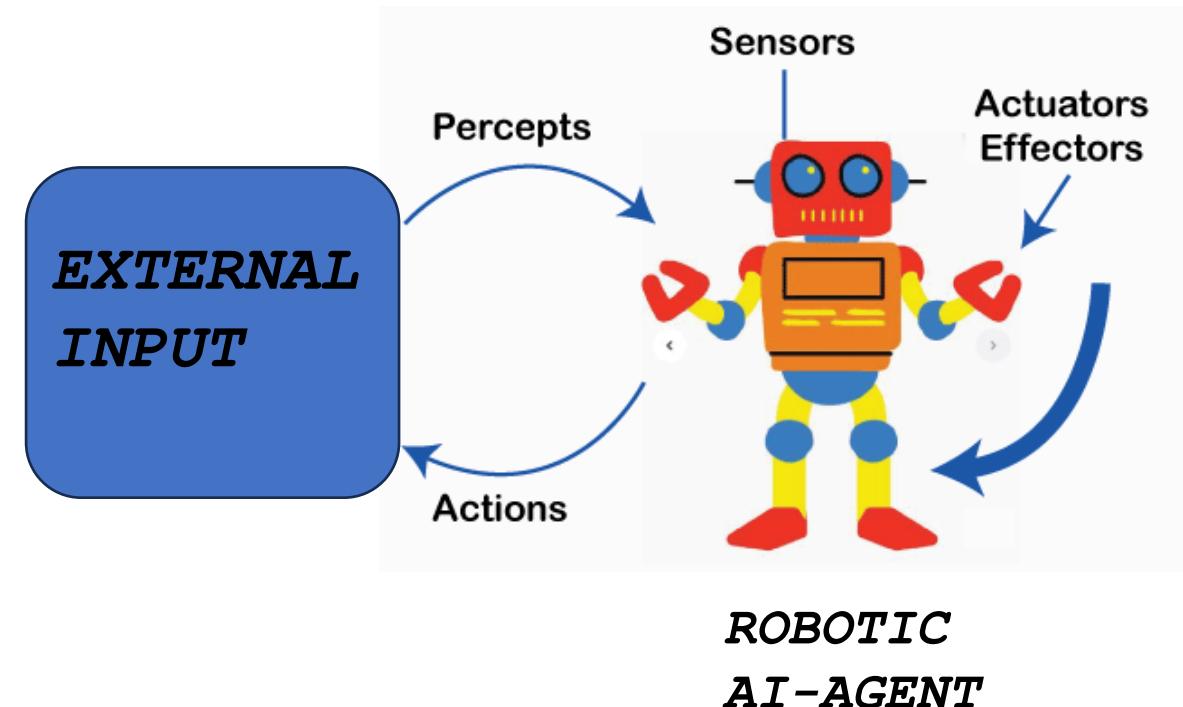
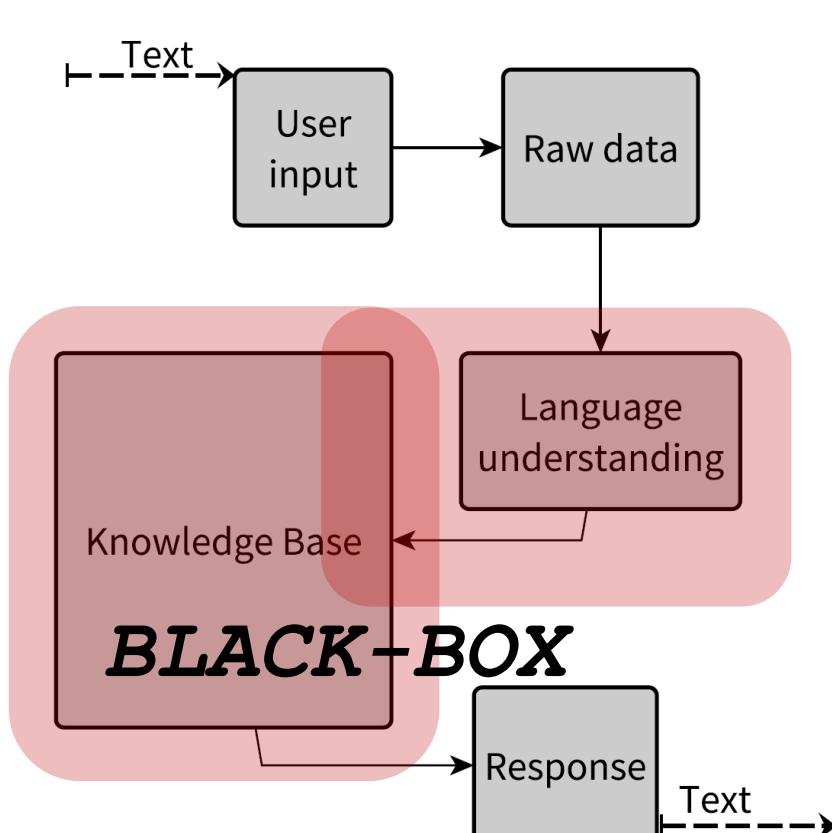
Generatieve AI [GEN-AI]

*Agenten gebaseerd op Neurale netwerk modellen ---door machinaal leren op basis van bestaande multimodale content
=====>(trainen van het model)
die instaat zijn nieuwe inhoud te creëren,
zoals tekst, afbeeldingen, muziek en/of code.*

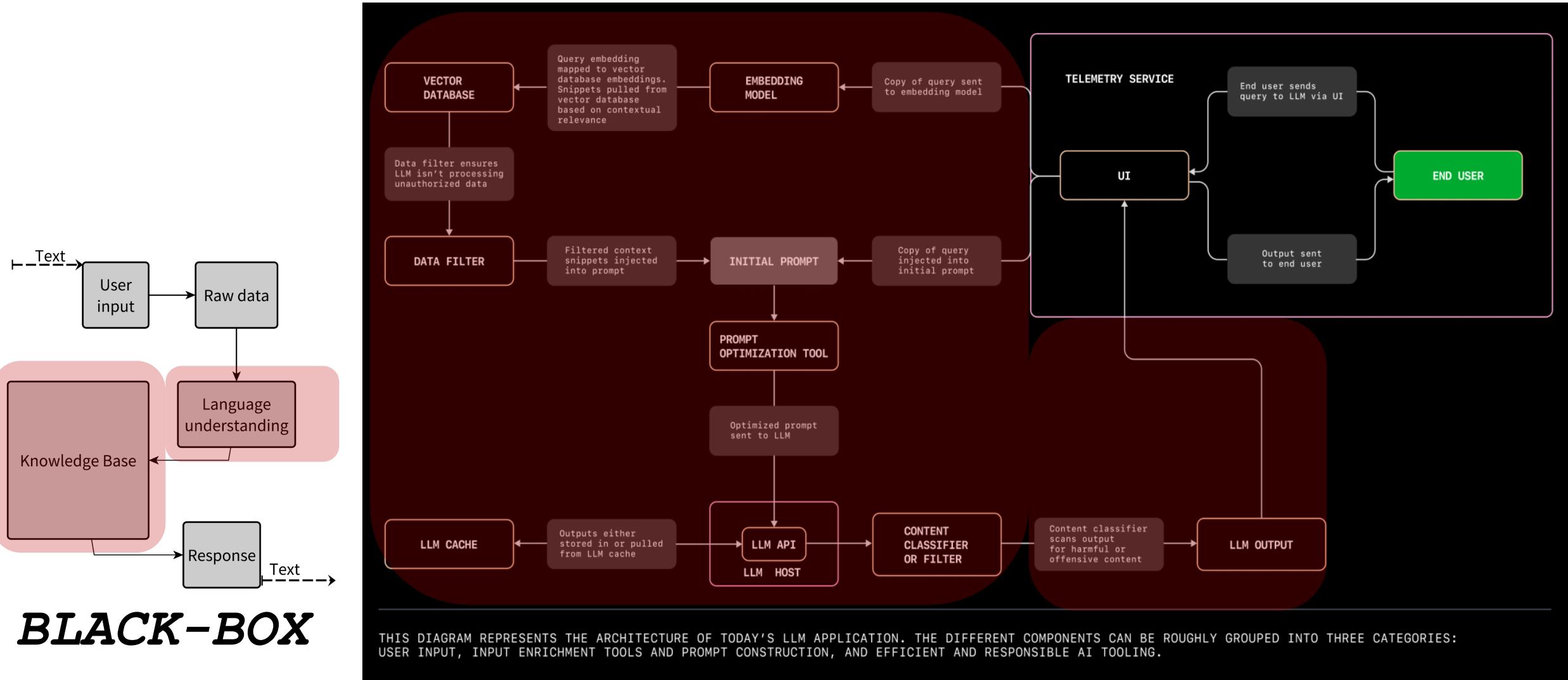
GEN-AI agenten zijn geen informatiedatabases of deterministische informatiezoeksystemen, maar voorspellingssystemen.

*GEN-AI maakt dus geen onderscheid tussen goed/fout of waar/nietwaar maar produceert een uitkomst die met grote waarschijnlijkheid kan worden gerelateerd aan de geven input (**prompt**) .*

ChatGPT is een Conversationele tekst-in/tekst-uit ChatBot

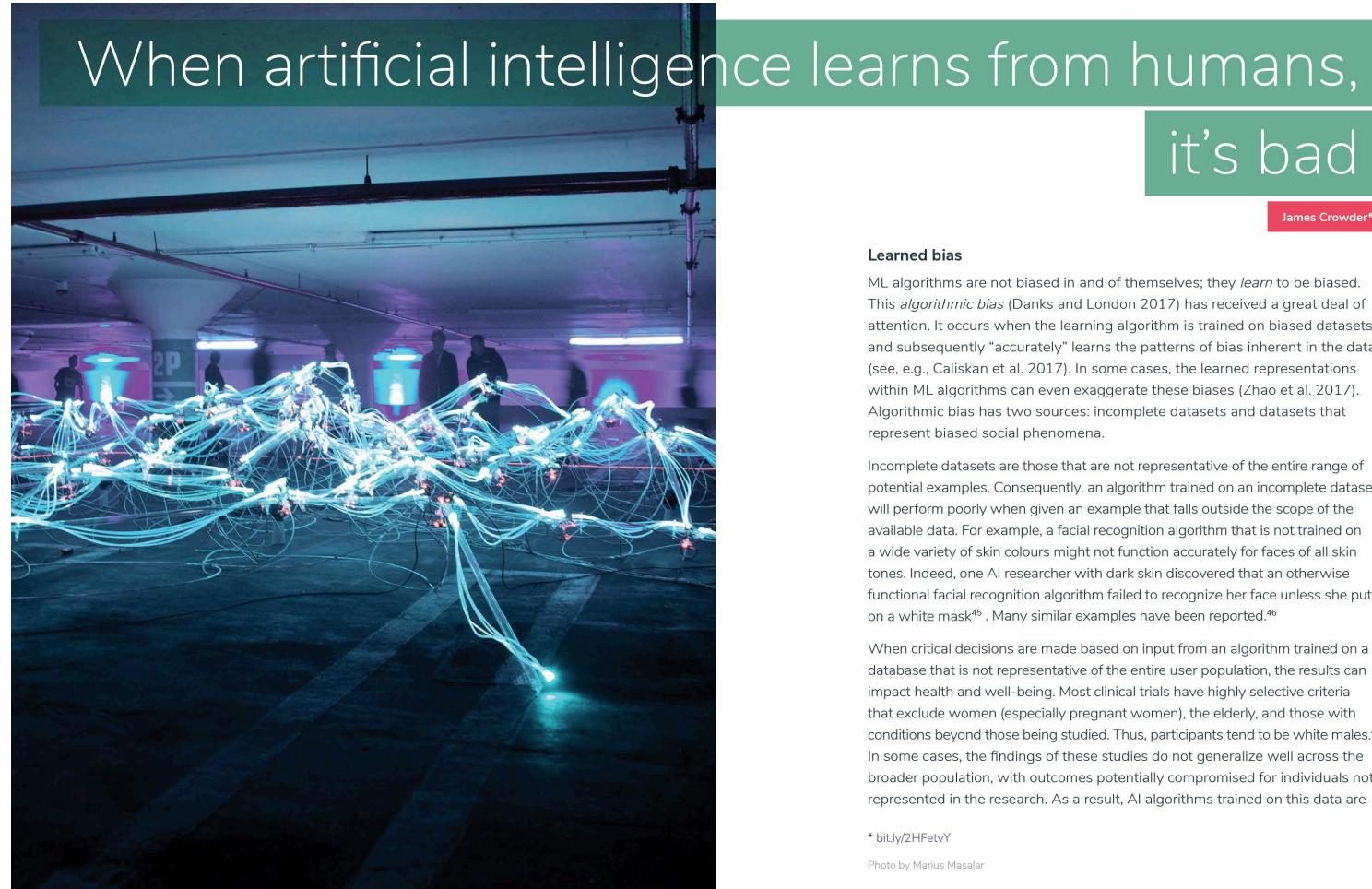


ChatGPT == 99% BLACKBOX + 1% user-interface



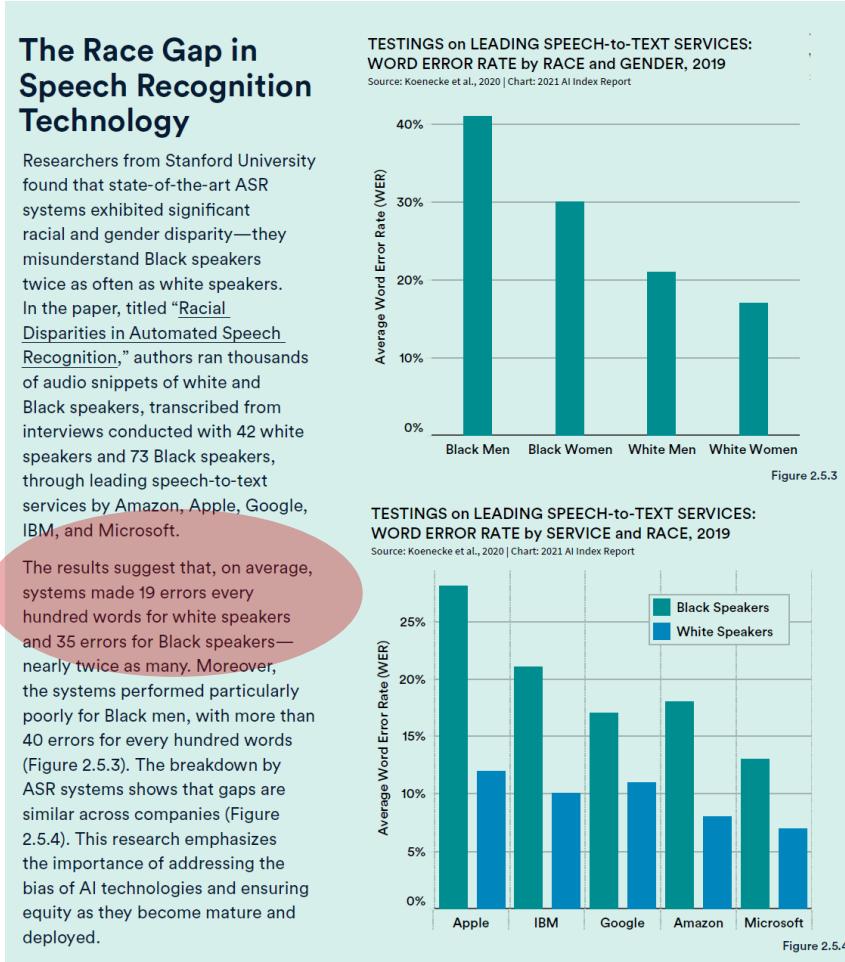
{Skewed Data-SETS}

Big-Data is Inherently Skewed



{Disparities}

Big Data causes racial & gender disparities



{Augmentation}

SURVEY PAPER Open Access



A survey on Image Data Augmentation for Deep Learning

Connor Shorten^{*} and Taghi M. Khoshgoftaar

*Correspondence:
cshorten2015@fau.edu
Department of Computer and Electrical Engineering and Computer Science, Florida Atlantic University, Boca Raton, USA

Abstract

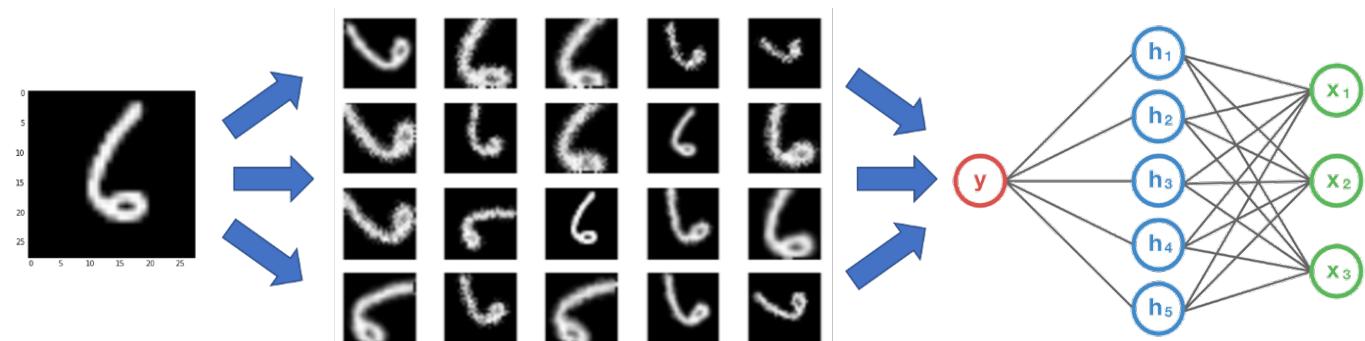
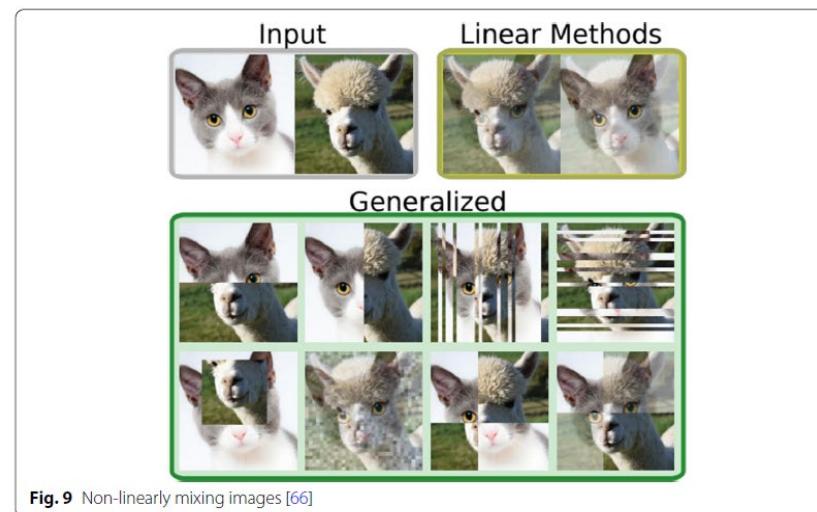
Deep convolutional neural networks have performed remarkably well on many Computer Vision tasks. However, these networks are heavily reliant on big data to avoid overfitting. Overfitting refers to the phenomenon when a network learns a function with very high variance such as to perfectly model the training data. Unfortunately, many application domains do not have access to big data, such as medical image analysis. This survey focuses on Data Augmentation, a data-space solution to the problem of limited data. Data Augmentation encompasses a suite of techniques that enhance the size and quality of training datasets such that better Deep Learning models can be built using them. The image augmentation algorithms discussed in this survey include geometric transformations, color space augmentations, kernel filters, mixing images, random erasing, feature space augmentation, adversarial training, generative adversarial networks, neural style transfer, and meta-learning. The application of augmentation methods based on GANs are heavily covered in this survey. In addition to augmentation techniques, this paper will briefly discuss other characteristics of Data Augmentation such as test-time augmentation, resolution impact, final dataset size, and curriculum learning. This survey will present existing methods for Data Augmentation, promising developments, and meta-level decisions for implementing Data Augmentation. Readers will understand how Data Augmentation can improve the performance of their models and expand limited datasets to take advantage of the capabilities of big data.

Keywords: Data Augmentation, Big data, Image data, Deep Learning, GANs

Introduction

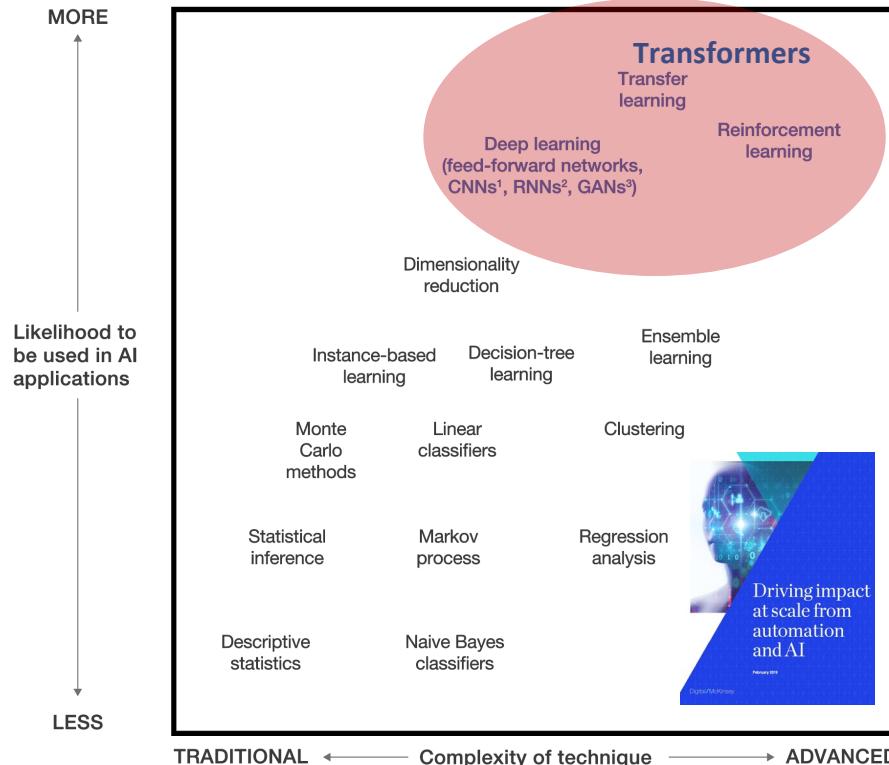
Deep Learning models have made incredible progress in discriminative tasks. This has been fueled by the advancement of deep network architectures, powerful computation, and access to big data. Deep neural networks have been successfully applied to Computer Vision tasks such as image classification, object detection, and image segmentation thanks to the development of convolutional neural networks (CNNs). These neural networks utilize parameterized, sparsely connected kernels which preserve the spatial characteristics of images. Convolutional layers sequentially downsample the spatial resolution of images while expanding the depth of their feature maps. This series of convolutional transformations can create much lower-dimensional and more useful representations of images than what could possibly be hand-crafted. The success of CNNs has sparked interest and optimism in applying Deep Learning to Computer Vision tasks.

Big Data that is *not augmented* causes Overfitting



{large scale}

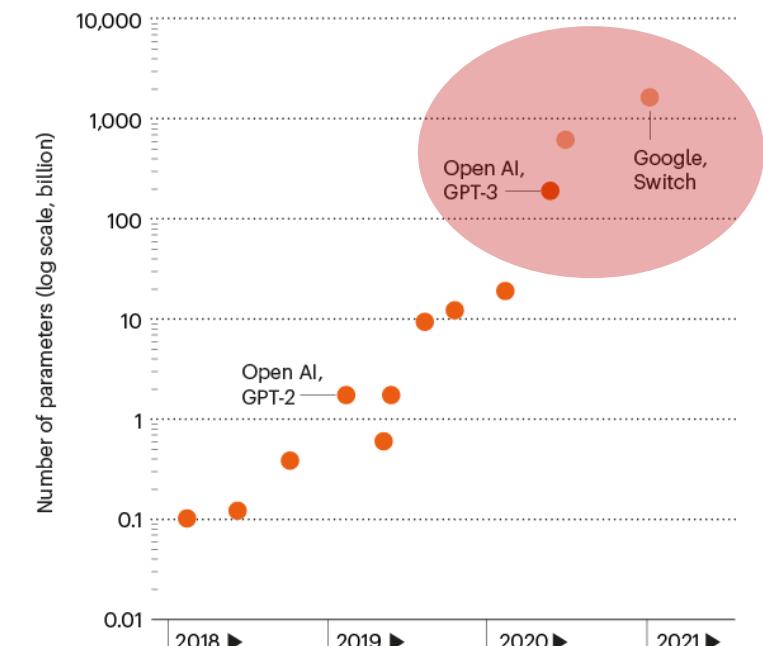
Only very large scale {DNNs} are useful
[can compete with human performance]



LARGER LANGUAGE MODELS

The scale of text-generating neural networks is growing exponentially, as measured by the models' parameters (roughly, the number of connections between neurons).

● 'Dense' models ● 'Sparse' models*

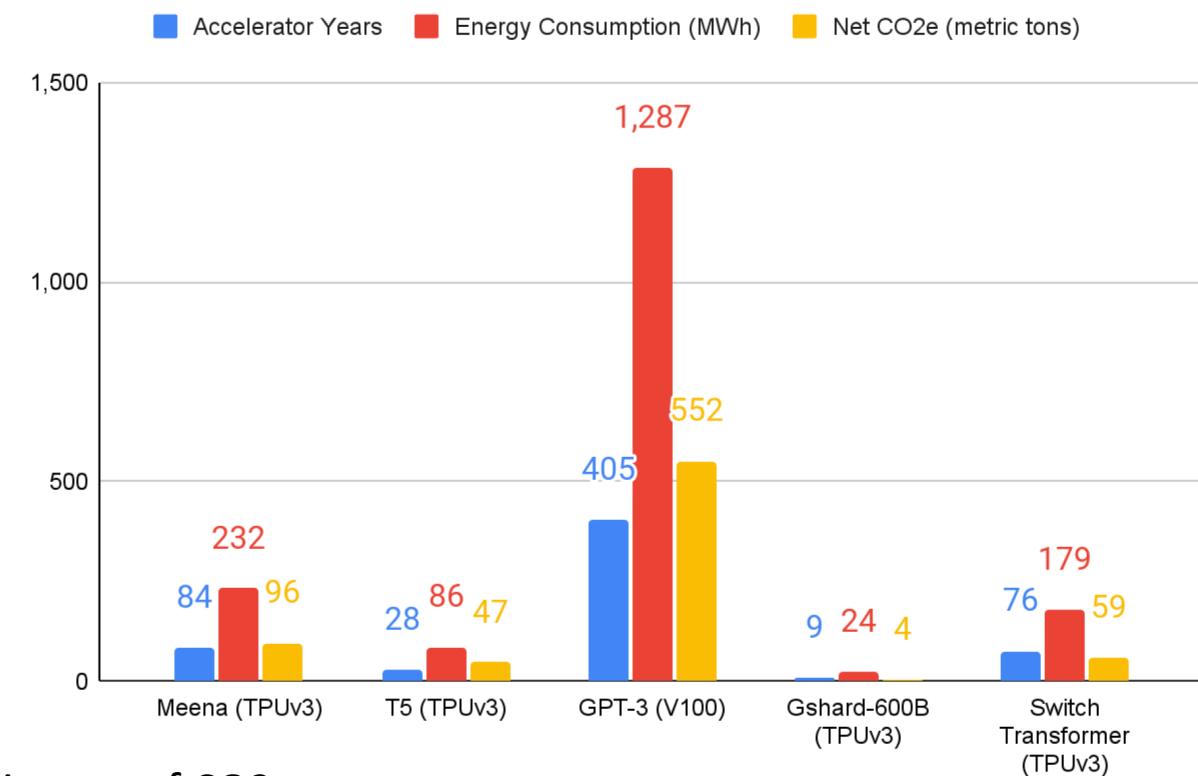
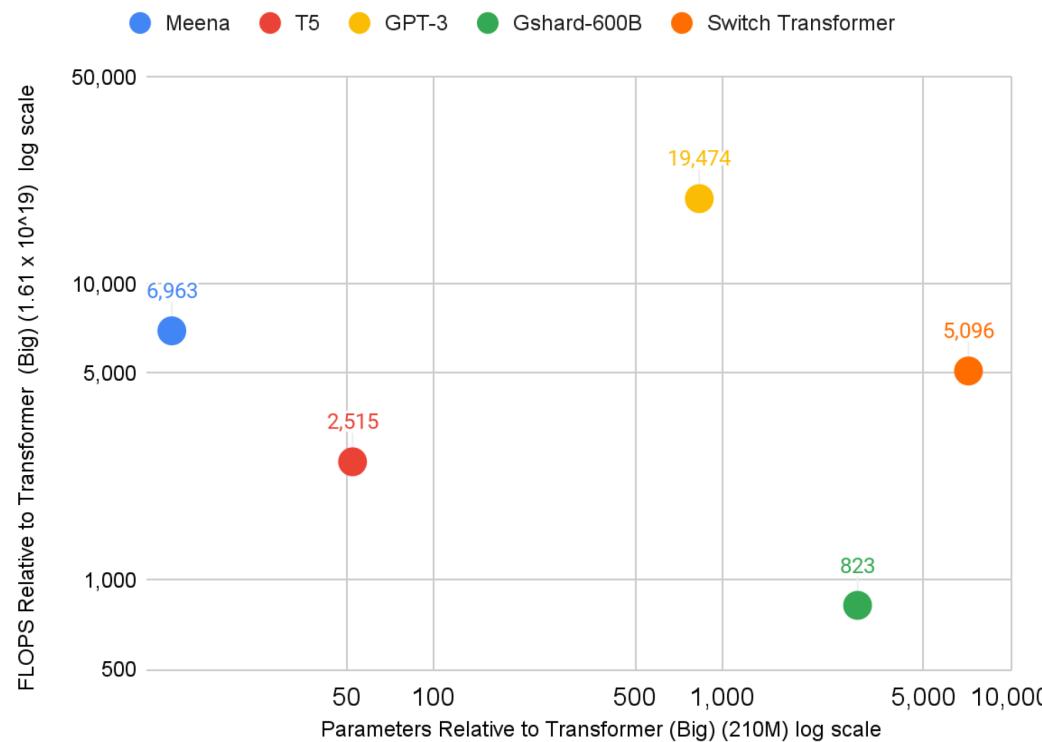


*Google's 1.6-trillion parameter 'sparse' model has performance equivalent to that of 10 billion to 100 billion parameter 'dense' models. ©nature

<https://www.nature.com/articles/d41586-021-00530-0>

{CO₂ foot-print}

Training large scale transformer {DNNs} produce massive Carbon Emissions



As of 2007, the average U.S. household emits 20 metric tons of CO₂ per year. In comparison to a world average of 4 tons.

[Carbon Footprint CSS09-05 e2021.pdf \(umich.edu\)](https://arxiv.org/ftp/arxiv/papers/2104/2104.10350.pdf)

<https://arxiv.org/ftp/arxiv/papers/2104/2104.10350.pdf>

{computational unsustainability}

The scale of state-of-the-art {SOTA} –near human level– DNNs
– *combined with a blind Brute-Force implementation + post-hoc analysis* –
is becoming more and more
computationally unsustainable,
even to the point
that **hypernetworks** are employed
to help humans to make **DNNs** work.

[2110.13100v1.pdf \(arxiv.org\)](https://arxiv.org/pdf/2110.13100v1.pdf)

<https://paperswithcode.com/sota/>

generieke train-dataset
machinaal leren
Black-Box

GEEN Human-in-the-loop
multimodaliteit
Commerciële belangen

creëert risico's
veiligheid
Privacy
betrouwbaarheid &
reproduceerbaarheid

Generatieve-AI

Hoe pas je Zero Trust
principes toe op het
gebruik van publieke
chatbots zoals Chat-GPT?

Human-in-the Loop

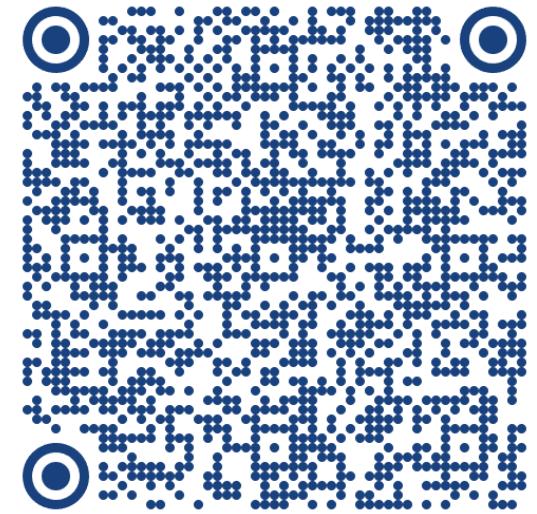
*“Cyber Security”
beveiligingsrichtlijnen
voor LLMs:*

{Open Web Application Security Project (OWASP) }

De **Open Web Application Security Project (OWASP)** is een open-source project waar beveiligingsexperts continu aan werken, om de lijst actueel te houden met de meest voorkomende veiligheidslekken.



A screenshot of the "OWASP Top 10 for LLM" report. The header reads "OWASP Top 10 for LLM" and includes a subtitle "This is a draft list of important vulnerability types for Artificial Intelligence (AI) applications built on Large Language Models (LLMs)." Below the header, the report lists ten categories of vulnerabilities, each with a brief description. The categories are: LLM01: Prompt Injections, LLM02: Insecure Output Handling, LLM03: Training Data Poisoning, LLM04: Denial of Service, LLM05: Supply Chain, LLM06: Permission Issues, LLM07: Data Leakage, LLM08: Excessive Agency, LLM09: Overreliance, and LLM10: Insecure Plugins. The report is presented in a clean, modern design with a blue header and white text on a dark background.



<https://owasp.org/www-project-top-10-for-large-language-model-applications/assets/PDF/OWASP-Top-10-for-LLMs-2023-v05.pdf>

{Bekende LLM “cyber” veiligheidsrisico’s}

Volgen OWASP gebruiken onderwijsinstellingen veelvuldig op LLM's gebaseerde public Chabots, maar lijken daarbij niet een duidelijk security-protocol te volgen.

Het belang van een "**human-in-the-loop**" is niet te onderschatten, zoals door toestemming voor bepaalde acties afhankelijk te maken van een menselijke goedkeuring.

Het naleven van "**zero-trust principles**" en "**least-privilege**" toegang is eveneens broodnodig. Een LLM moet behandeld worden als een untrusted user om de invloed van een kwaadwillende zoveel mogelijk in te perken.

LLM01: Prompt Injections

Prompt Injection Vulnerabilities in LLMs involve crafty inputs leading to undetected manipulations. The impact ranges from data exposure to unauthorized actions, serving attacker's goals.

LLM02: Insecure Output Handling

These occur when plugins or apps accept LLM output without scrutiny, potentially leading to XSS, CSRF, SSRF, privilege escalation, remote code execution, and can enable agent hijacking attacks.

LLM03: Training Data Poisoning

LLMs learn from diverse text but risk training data poisoning, leading to user misinformation. Overreliance on AI is a concern. Key data sources include Common Crawl, WebText, OpenWebText, and books.

LLM04: Denial of Service

An attacker interacts with an LLM in a way that is particularly resource-consuming, causing quality of service to degrade for them and other users, or for high resource costs to be incurred.

LLM05: Supply Chain

LLM supply chains risk integrity due to vulnerabilities leading to biases, security breaches, or system failures. Issues arise from pre-trained models, crowdsourced data, and plugin extensions.

LLM06: Permission Issues

Lack of authorization tracking between plugins can enable indirect prompt injection or malicious plugin usage, leading to privilege escalation, confidentiality loss, and potential remote code execution.

LLM07: Data Leakage

Data leakage in LLMs can expose sensitive information or proprietary details, leading to privacy and security breaches. Proper data sanitization, and clear terms of use are crucial for prevention.

LLM08: Excessive Agency

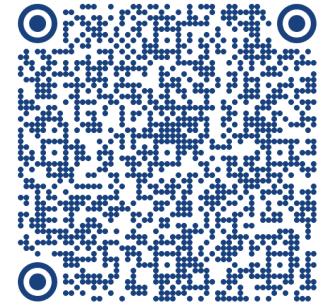
When LLMs interface with other systems, unrestricted agency may lead to undesirable operations and actions. Like web-apps, LLMs should not self-police; controls must be embedded in APIs.

LLM09: Overreliance

Overreliance on LLMs can lead to misinformation or inappropriate content due to "hallucinations." Without proper oversight, this can result in legal issues and reputational damage.

LLM10: Insecure Plugins

Plugins connecting LLMs to external resources can be exploited if they accept free-form text inputs, enabling malicious requests that could lead to undesired behaviors or remote code execution.



<https://owasp.org/www-project-top-10-for-large-language-model-applications/assets/PDF/OWASP-Top-10-for-LLMs-2023-v05.pdf>



HOGESCHOOL
ROTTERDAM

{LLM01: ONVEILIGE PROMPT (“prompt injection”)}

De input bepaalt wat een LLM produceert. Vaak zijn deze resultaten op voorhand lastig te voorspellen.

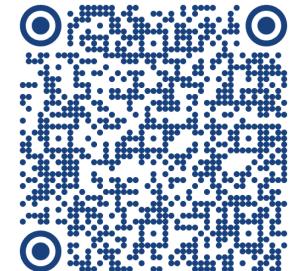
Toch kunnen ontwikkelaars de resultaten verfijnen met betere training en sleutelwerk aan de parameters.

Bij publieke chatbots zoals ChatGPT en Google Bard is bekend dat ze vrij snel voorzichtig worden als je ze gedurfde vragen stelt.

In een professionele context gaat het alleen om meer dan lolligheden, en kan een interne LLM bedrijfsgeheimen in huis hebben.

Hackers kunnen prompts manipuleren, waardoor het systeem de intenties van de aanvaller uitvoert (denk aan “role playing”).

Dit kan leiden tot extractie van gegevens, social engineering en andere problemen.



Kwetsbaarheid (“Jailbreak”):

OWASP spreekt over directe en indirecte prompt injections. In het eerste geval karakteriseert men deze aanvalstactiek als een “jailbreak” voor een LLM: een kwaadwillende heeft in dat geval de onderliggende systemprompt ontdekt of herschreven. Daardoor kan men bij een aanval wellicht bij gevoelige data stores komen waar de LLM op rust.

Het indirecte scenario komt voor als een LLM via een externe bron aanspreekbaar warbij de aanvaller een prompt injection uitvoert om de conversatie-context te kapen. Vervolgens opereert de LLM als een “confused deputy”, in de woorden van OWASP.

Dat houdt in dat het taken kan uitvoeren die het normaal gesproken niet zou mogen doen voor een gebruiker.



Preventieve maatregel:

Privilegebeheer wordt afgedwongen voor LLM-toegang tot achterliggende systemen

Validatie door een mens wordt afgedwongen voor beslisfunctionaliteit (“human-in-the-loop”)

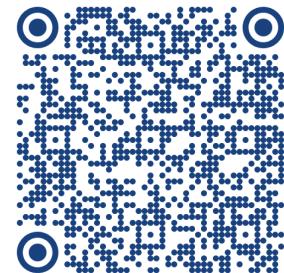
Externe inhoud is gescheiden van gebruikersprompts

{LLM07: Onbedoelde openbaarmaking van gevoelige informatie}

Hoe goed een LLM ook getraind en getest is,
het kan zo zijn dat een generatief AI-model op een onverwachte manier informatie door kan spelen (**“data-leakage”**).

Een toevallige prompt kan opeens gevoelige data of algoritmes prijsgeven,
zelfs als de eindgebruiker geen kwaadwillende gedachten erop nahoudt.

LLM-toepassingen kunnen onbedoeld gevoelige informatie, gepatenteerde algoritmen of vertrouwelijke gegevens openbaar maken,
wat kan leiden tot onbevoegde toegang, diefstal van intellectueel eigendom en privacyschendingen.



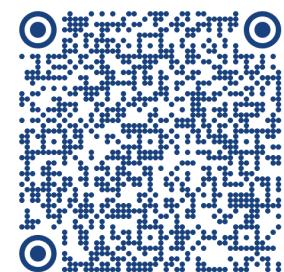
Kwetsbaarheid:

Onbedoelde bekendmaking van vertrouwelijke informatie (onzorgvuldig/ondoordacht handelen)

Gecompliceerde (van het web gekopieerde) prompts die worden gebruikt om invoerfilters te omzeilen en gevoelige gegevens te onthullen

Preventieve maatregel:

Validatie en opschoning (pseudonimiseren + de-identificatie) van invoer op basis van geautoriseerde protocollen/use-cases
OWASP raadt aan om consumenten van LLM-applicaties bewust te maken van een veilige omgang met AI-modellen.



Autoriteit Persoonsgegevens
<https://www.autoriteitpersoonsgegevens.nl/themas/beveiliging/beveiliging-van-persoonsgegevens/gegevens-pseudonimiseren>

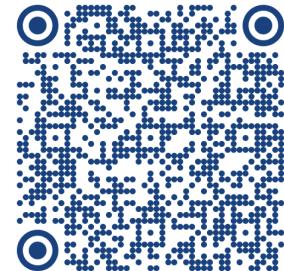
Making Qualitative Data Reusable - A Short Guidebook For Researchers And Data Stewards Working With Qualitative Data
<https://zenodo.org/doi/10.5281/zenodo.7777518>

{LLM08: Te veel vrijheid “Excessive Agency”}

Een LLM-systeem kan erg effectief zijn als het toegang heeft tot andere systeem.

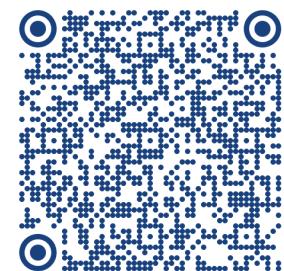
Zo kan een AI-assistent in opdracht van een gebruiker in een andere applicatie duiken om iets uit te zoeken.

Wie bij het opzetten van deze functionaliteit onoplettend is, kan de LLM toestaan om ongewenste acties uit te voeren.



Deze vorm van '**agency**' of **autonomie** kan bijvoorbeeld ruimte geven om shell commands uit te voeren of om een database te updaten, zoals een gewone gebruiker van de andere app zou kunnen.

Met te veel permissies kan dit grote gevolgen hebben.



Kwetsbaarheid:

Onbedoelde bekendmaking van vertrouwelijke informatie (onzorgvuldig/ondoordacht handelen)

Gecompliceerde (van het web gekopieerde) prompts die worden gebruikt om invoerfilters te omzeilen en gevoelige gegevens te onthullen

Preventieve maatregel:

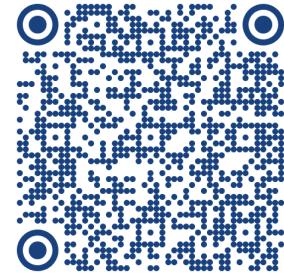
Legt vast waartoe een LLM toegang toe heeft/mag hebben.

Beschrijf verificatiestappen / autorisatie-process om van een ChatBot gebruik te mogen maken.

Wederom is het nuttig om iemand in real-time te laten kijken wat een AI-model precies uitvoert, zoals ook doodgewone gebruikers worden gecontroleerd in een IT-omgeving.

{LLM09: te grote afhankelijkheid/vertrouwen “Overreliance”}

LLM's kunnen veel, maar zijn zoals gezegd in staat om fouten te maken en produceren altijd een output ook al is deze onzinnig (“**hallucineren**”). Ze '*denken*' niet en zijn uiteindelijk een zeer complexe voorspelmachines. Dit houdt in dat er in een professionele context erg voorzichtig mee omgegaan moet worden.



LLM's kunnen in veel gevallen programmeercode produceren, maar ook daar geldt dat er onjuiste of gevaarlijke inhoud in kan zitten. Het is een bekend probleem dat deze AI-modellen de mist in kunnen gaan, waardoor het belangrijk is dat organisaties regelmatig controleren of de output naar wens is.

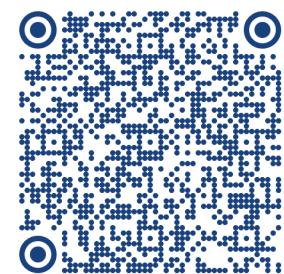
Kwetsbaarheid:

Het opvallende aan deze kwetsbaarheid is dat de aanval niet kwaadwillend is: een nietsvermoedende developer kan programmeercode inzetten die men niet genoeg had bekijken om de fouten op te pikken.

LLM geeft onjuiste informatie

LLM genereert onzinnige tekst

Onvoldoende risicocommunicatie van aanbieders van LLM



Preventieve maatregel:

Communiceer de risico's en beperkingen van LLM duidelijk aan eindgebruikers, laat ze workshops volgen.

De resultaten van LLM monitoren.

Triangulatie / cross-checking van LLM-uitvoer met betrouwbare bronnen om de uitvoer op juistheid te controleren



Hoe bouw en test je veilig & betrouwbaar GEN-AI technologie?

**“human-in-the-loop” +
“private endpoint” or
“virtual networks” +
“embedded vector store”**

Data, privacy, and security for Azure OpenAI Service

Important

Your prompts (inputs) and completions (outputs), your embeddings, and your training data:

- are NOT available to other customers.
- are NOT available to OpenAI.
- are NOT used to improve OpenAI models.
- are NOT used to train, retrain, or improve Azure OpenAI Service foundation models.
- are NOT used to improve any Microsoft or 3rd party products or services without your permission or instruction.
- Your fine-tuned Azure OpenAI models are available exclusively for your use.

The Azure OpenAI Service is operated by Microsoft as an Azure service; Microsoft hosts the OpenAI models in Microsoft's Azure environment and the Service does NOT interact with any services operated by OpenAI (e.g. ChatGPT, or the OpenAI API).

<https://learn.microsoft.com/en-us/legal/cognitive-services/openai/data-privacy?tabs=azure-porta>

What data does the Azure OpenAI Service process?

What data does the Azure OpenAI Service process?

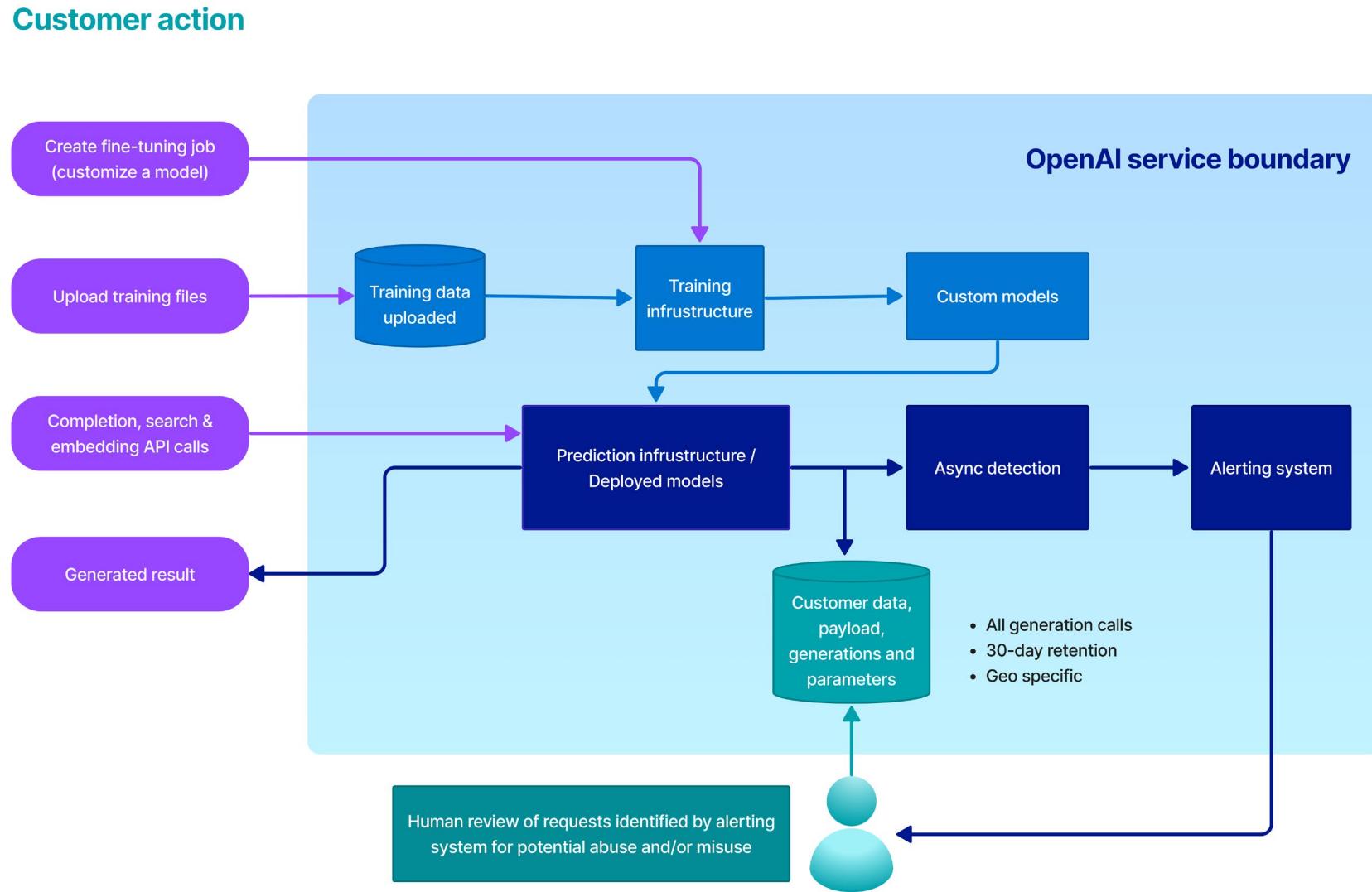
Azure OpenAI processes the following types of data:

Text prompts, queries and responses submitted by the user via the completions, search, and embeddings operations.

Training & validation data. You can provide your own training data consisting of prompt-completion pairs for the purposes of fine-tuning an OpenAI model.

Results data from training process. After training a fine-tuned model, the service will output meta-data on the job which includes tokens processed and validation scores at each step.

What data does the Azure OpenAI Service process?



Azure OpenAI service

OVERVIEW

Build intelligent apps with AI models

Cutting edge models ▾

Quickly develop generative AI experiences with a diverse set of prebuilt and curated models from OpenAI, Meta and beyond.

[Try the Azure AI Studio](#)

Data grounding ▾

Trust and transparency ▾

Data, privacy and security ▾

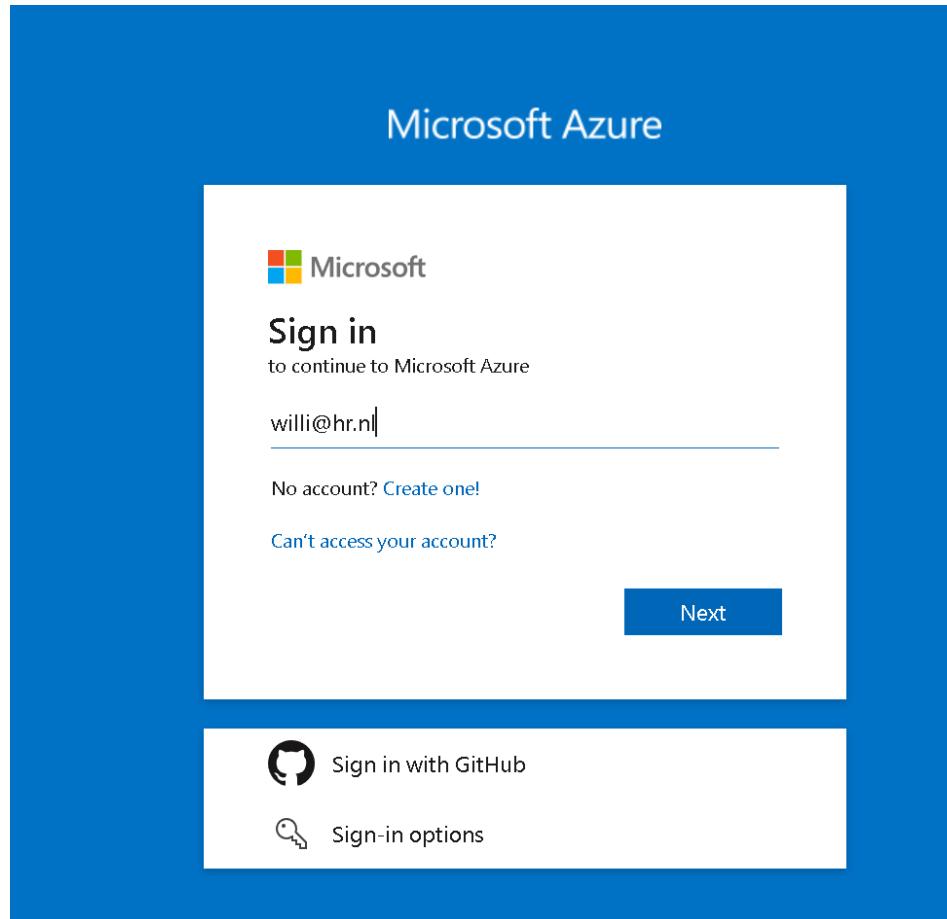


USE CASES

Apply generative AI to a variety of use cases

<https://azure.microsoft.com/en-us/products/ai-services/openai-service>

portal.azure.com



Log in

Student number or personnel code

Password

Login



HOGESCHOOL
ROTTERDAM

Azure services



Resources

Recent Favorite

Name	Type	Last Viewed
GPT4-SWEDEN-GROUP	Azure OpenAI	6 days ago
Azure for Students	Subscription	3 weeks ago
CHATBOT02	Azure OpenAI	3 weeks ago
DefaultResourceGroup-westeurope	Resource group	2 months ago
AV07	Speech service	7 months ago
NLP	Resource group	7 months ago
LLM01	Language understanding	7 months ago
Visual Studio Professional Subscription	Subscription	7 months ago
LLM01-Authoring	Language understanding	7 months ago
WILLI107	Resource group	7 months ago
cursusa1-900	Azure Machine Learning workspace	8 months ago
cursusa9006361709869	Key vault	8 months ago

See all

Navigate



Tools



Useful links

Azure mobile app



Microsoft Azure Search resources, services, and docs (G+/-) Copilot

Home > LLMGPT

LLMGPT | Networking

Azure OpenAI

Search

Firewalls and virtual networks Private endpoint connections

Save Discard Refresh

Access control settings allowing access to Azure AI services account will remain in effect for up to three minutes after saving updated settings restricting access.

Allow access from All networks Selected Networks and Private Endpoints Disabled

Configure network security for your Azure AI services account. [Learn more.](#)

Virtual networks

Secure your Azure AI services account with virtual networks. [+ Add existing virtual network](#) [+ Add new virtual network](#)

Virtual Network	Subnet	Address range	Endpoint Status	Resource group	Subscription
No network selected.					

Firewall

Add IP ranges to allow access from the internet or your on-premises networks. [Learn more.](#)

Add your client IP address ('77.173.131.238') ⓘ

Address range

IP address or CIDR

Exceptions

Allow Azure services on the trusted services list to access this cognitive services account. ⓘ

 Filter by title

Azure for Python developers

Get started

› Azure AI apps

Develop with Azure AI services

Get started with the Python enterprise chat sample

Configure document security

Secure your chat endpoint

Evaluate your chat app

Scale Azure OpenAI with Azure Container Apps

Scale Azure OpenAI with Azure API Management

Load test with Locust

› Web

› Data

› Containers

Logs

› Azure SDK for Python

SDK PyPI package index

SDK reference documentation

› Explore services supporting Python

› Samples

Learn / Azure / Developer / Python /

Get started with chat private endpoints for Python

Article • 07/23/2024 • 5 contributors

 Feedback

In this article

[Architectural overview](#)

[Deployment steps](#)

[Prerequisites](#)

[Open development environment](#)

[Show 6 more](#)

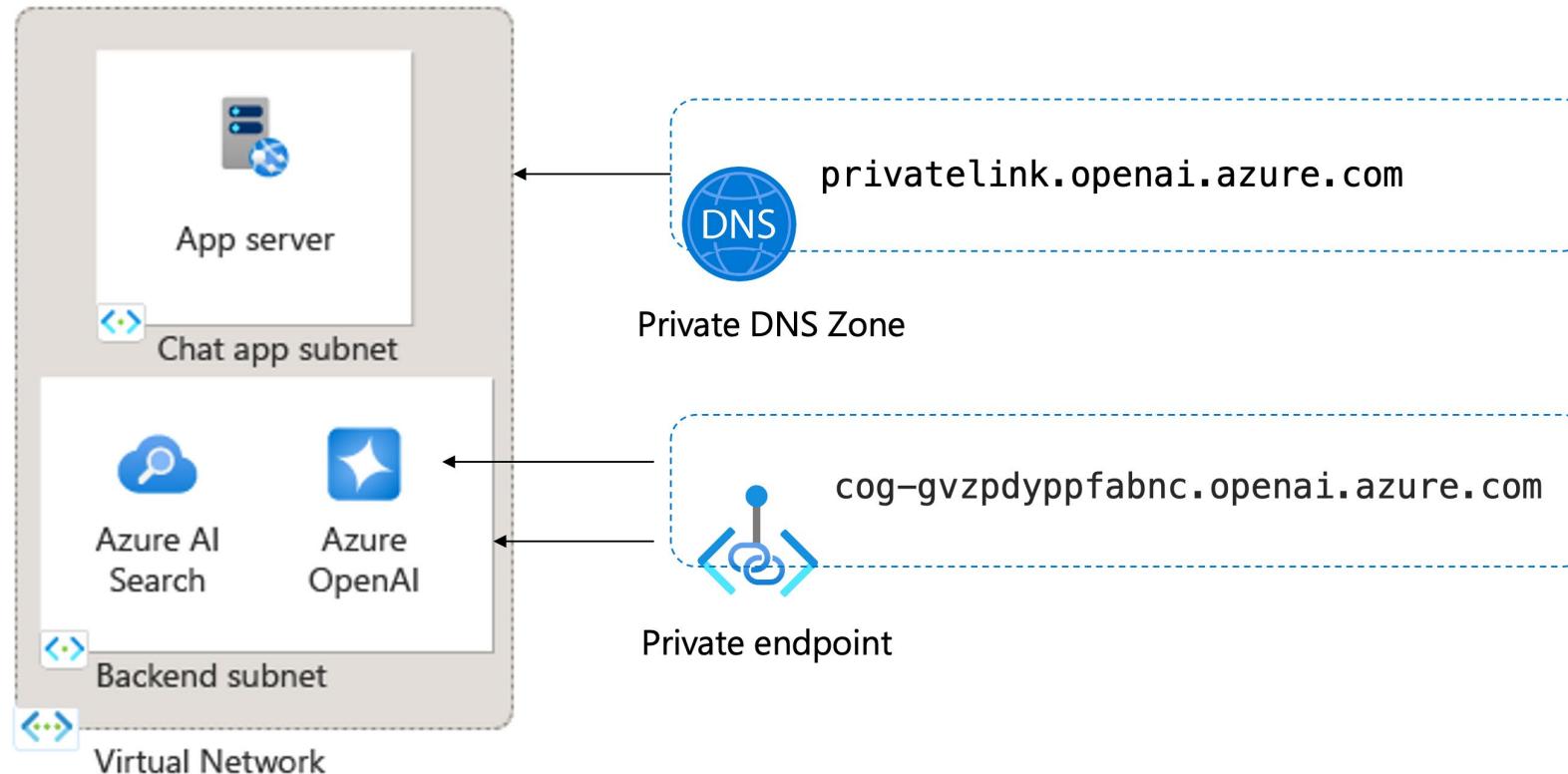
This article shows you how to deploy and run the [Enterprise chat app sample for Python](#) accessible by private endpoints.

This sample implements a chat app using Python, Azure OpenAI Service, and [Retrieval Augmented Generation \(RAG\)](#) in Azure AI Search to get answers about employee benefits at a fictitious company. The app is seeded with PDF files including the employee handbook, a benefits document and a list of company roles and expectations.

By following the instructions in this article, you will:

- Deploy a chat app to Azure for public access in a web browser.
- Redeploy chat app with private endpoints.

Once you complete this procedure, you can start modifying the new project with your custom code and redeploy, knowing your chat app is accessible only through the private network.



<https://medium.com/@luiz.braz/using-azure-openai-in-a-secure-way-17b9e51cdcaa>

<https://learn.microsoft.com/en-us/shows/azure-essentials-show/connecting-openai-private-endpoints-across-vnets>

Create and deploy an Azure OpenAI Service resource

Article • 09/06/2023 • 4 contributors

Feedback

Choose your preferred resource creation method

Portal **CLI** PowerShell

In this article

[Prerequisites](#)[Create a resource](#)[Deploy a model](#)[Next steps](#)

This article describes how to get started with Azure OpenAI Service and provides step-by-step instructions to create a resource and deploy a model. You can create resources in Azure in several different ways:

- The [Azure portal](#)
- The REST APIs, the Azure CLI, PowerShell, or client libraries
- Azure Resource Manager (ARM) templates

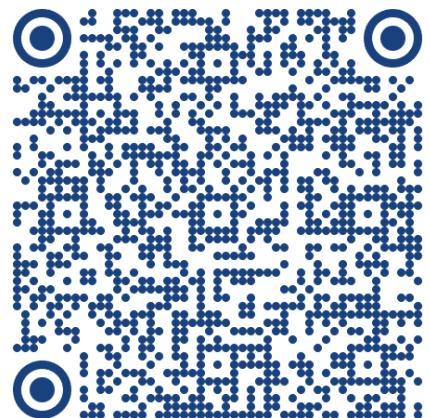
In this article, you review examples for creating and deploying resources in the Azure portal and with the Azure CLI.

Prerequisites

- An Azure subscription. [Create one for free](#).
- Access granted to Azure OpenAI in the desired Azure subscription.
- Access permissions to [create Azure OpenAI resources](#) and to [deploy models](#).

Note

Currently, you must submit an application to access Azure OpenAI Service. To apply for access, complete [this form](#). If you need assistance, open an issue on this repository to contact Microsoft.



Limited access for Azure OpenAI Service

Article • 07/23/2024 • 6 contributors

[Feedback](#)

In this article

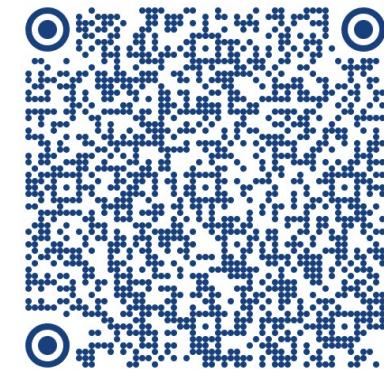
[Registration for modified content filters and/or abuse monitoring](#)[Important links](#)[Help and support](#)[See also](#)

As part of Microsoft's commitment to responsible AI, we have designed and operate Azure OpenAI Service with the intention of protecting the rights of individuals and society and fostering transparent human-computer interaction. For this reason, Azure OpenAI is a Limited Access service, and access and use is subject to eligibility criteria determined by Microsoft. Unless otherwise indicated in the service, all Azure customers are eligible for access to Azure OpenAI models, and all uses consistent with the [Product Terms](#) and [Code of Conduct](#) are permitted, so customers are not required to submit a registration form unless they are requesting approval to modify content filters and/or abuse monitoring.

Azure OpenAI Service is made available to customers under the terms governing their subscription to Microsoft Azure Services, including [Product Terms](#) such as the Universal License Terms applicable to Microsoft Generative AI Services and the product offering terms for Azure OpenAI. Please review these terms carefully as they contain important conditions and obligations governing your use of Azure OpenAI Service.

Registration for modified content filters and/or abuse monitoring

<https://learn.microsoft.com/en-us/legal/cognitive-services/openai/limited-access>



[https://customervoice.microsoft.com/
Pages/ResponsePage.aspx?id=v4j5
cvGGGr0GRqy180BHbR7en2Ais5pxKt
so_Pz4b1_xUOFA5Qk1UWDRBMjg
0WFhPMklzTzhKQ1dWNyQIQCN0P
Wcu](https://customervoice.microsoft.com/Pages/ResponsePage.aspx?id=v4j5cvGGGr0GRqy180BHbR7en2Ais5pxKtso_Pz4b1_xUOFA5Qk1UWDRBMjg0WFhPMklzTzhKQ1dWNyQIQCN0PWcu)

Request Access to Azure OpenAI Service

* Required

Please read all instructions carefully and complete form as instructed

Thank you for your interest in Azure OpenAI Service. Please submit this form to register for approval to access and use Azure OpenAI's Limited Access text and code and/or DALL-E 2 text to image models (as indicated in the form). All use cases must be registered. Azure OpenAI Service requires registration and is currently only available to approved enterprise customers and partners. Learn more about limited access to Azure OpenAI Service [here](#).

Limited access scenarios: When evaluating which scenarios to onboard, we consider who will directly interact with the application, who will see the output of the application, whether the application will be used in a high-stakes domain (e.g., medical), and the extent to which the application's capabilities are tightly scoped. In general, applications in high stakes domains will require additional mitigations and are more likely to be approved for applications with internal-only users and internal-only audiences. Applications with broad possible uses, including content generation capabilities, are more likely to be approved if 1) the domain is not high stakes and users are authenticated or 2) in the case of high stakes domains, anyone who views or interacts with the content is internal to your company.

Please be sure to visit the [Azure OpenAI Service's transparency note](#), which provides information and guidelines for responsible use of the service as well as system limitations that may be applicable to your scenario.

If you are a current Azure OpenAI customer and would like to add additional use cases, please fill out the [Azure OpenAI Additional Use Case form](#)

Azure AI | Azure OpenAI Studio

« Azure AI Studio > Chat playground

Chat playground

Assistant setup

System message Add your data (preview)

Save changes

Specify how the chat should act

Use a template to get started, or just start writing your own system message below. Want some tips? [Learn more](#)

Use a system message template

Select a template

System message ⓘ

You are an AI assistant that helps people find information.

Examples ⓘ

+ Add an example

Sample Code

You can use the following code to start integrating your current prompt and settings into your application

<https://gpt4-sweden-group.openai.azure.com/> python

```
1 #Note: The openai-python library support for Azure OpenAI is in preview.
2 import os
3 import openai
4 openai.api_type = "azure"
5 openai.api_base = "https://gpt4-sweden-group.openai.azure.com/"
6 openai.api_version = "2023-07-01-preview"
7 openai.api_key = os.getenv("OPENAI_API_KEY")
8
9 response = openai.chatcompletion.create(
10     engine="GPT4-32K",
11     messages = [{"role": "system", "content": "You are an AI
assistant that helps people find information."},
12 {"role": "user", "content": "A neutron star is the collapsed core of a
massive supergiant star, which had a total mass of between 10 and 25
solar masses, possibly more if the star was especially metal-rich.
Neutron stars are the smallest and densest stellar objects, excluding
black holes and hypothetical white holes, quark stars, and strange
stars. Neutron stars have a radius on the order of 10 kilometres (6.2
mi) and a mass of about 1.4 solar masses. They result from the
supernova explosion of a massive star, combined with gravitational
collapse, that compresses the core past white dwarf star density to
that of atomic nuclei.\n\nQ: How are neutron stars created?\nA:"}],
13     "role": "assistant", "content": "Neutron stars are created from the
supernova explosion of a massive star, combined with gravitational
collapse, that compresses the core past white dwarf star density to
that of atomic nuclei.\n\nQ: How are neutron stars created?\nA:"}]
```

Endpoint ⓘ

<https://gpt4-sweden-group.openai.azure.com/openai/deployments/GPT4-3...>

Key ⓘ

.....

You should use environment variables or a secret management tool like Azure Key Vault to prevent accidental exposure of your key in applications. [Learn more](#)

Copy Close

Azure chat completions example (preview)

In this example we'll try to go over all operations needed to get chat completions working using the Azure endpoints.

This example focuses on chat completions but also touches on some other operations that are also available using the API. This example is meant to be a quick way of showing simple operations and is not meant as a tutorial.

```
1 import os
2 import openai
3 openai.api_type = "azure"
4 openai.api_base = "https://gpt4-sweden-group.openai.azure.com/"
5 openai.api_version = "2023-07-01-preview"
6 openai.api_key = "ded218c778894f6da4d3c595c6904194"
7
8
9 #!setx AZURE_OPENAI_KEY "ded218c778894f6da4d3c595c6904194"
10 #!setx AZURE_OPENAI_ENDPOINT "https://gpt4-sweden-group.openai.azure.com/openai/deployments/GPT4-32K/chat/completions?api-version=2023-07-01-preview"
11
12 response = openai.ChatCompletion.create(
13     engine="GPT4-32K",
14     messages = [
15         {"role": "system", "content": "You are a helpful assistant."},
16         {"role": "user", "content": "Does Azure OpenAI support customer managed keys?"},
17         {"role": "assistant", "content": "Yes, customer managed keys are supported by Azure OpenAI."},
18         {"role": "user", "content": "Do other Azure AI services support this too?"}
19     ],
20     temperature=0.7,
21     max_tokens=800,
22     top_p=0.95,
23     frequency_penalty=0,
24     presence_penalty=0,
25     stop=None)
26
27
28 print(response)
29 print(response['choices'][0]['message']['content'])
30
```

```
{
  "id": "chatcompl-8cp65kWwKakKF8WB1TtRnUPd60Ak",
  "object": "chat.completion",
  "created": 1698067045,
  "model": "gpt-4-32k",
  "prompt_filter_results": [
    {
      "prompt_index": 0,
      "content_filter_results": {
        "hate": {
          "filtered": false,
          "severity": "safe"
        },
        "self_harm": {
          "filtered": false,
          "severity": "safe"
        },
        "sexual": {
          "filtered": false,
          "severity": "safe"
        },
        "violence": {
          "filtered": false,
          "severity": "safe"
        }
      }
    }
  ]
}
```



Learn / Azure / AI Services /

Learn how to generate or manipulate text

Article • 08/17/2023 • 2 contributors

Feedback

In this article

- Design prompts
- Classify text
- Trigger ideas
- Conduct conversations

Show 5 more

Azure OpenAI Service provides a **completion endpoint** that can be used for a wide variety of tasks. The endpoint supplies a simple yet powerful text-in, text-out interface to any [Azure OpenAI model](#). To trigger the completion, you input some text as a prompt. The model generates the completion and attempts to match your context or pattern. Suppose you provide the prompt "As Descartes said, I think, therefore" to the API. For this prompt, Azure OpenAI returns the completion endpoint "I am" with high probability.

The best way to start exploring completions is through the playground in [Azure OpenAI Studio](#). It's a simple text box where you enter a prompt to generate a completion. You can start with a simple prompt like this one:

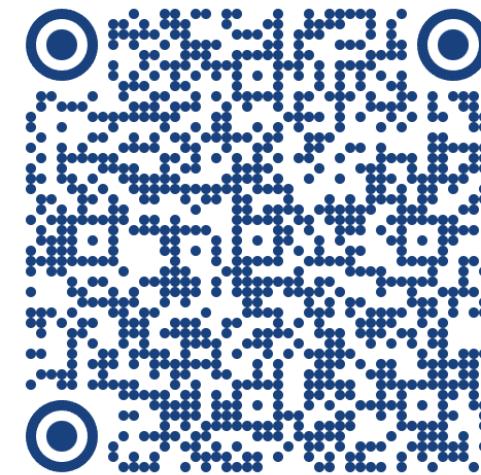
Console Copy

```
write a tagline for an ice cream shop
```

After you enter your prompt, Azure OpenAI displays the completion:

Console Copy

```
we serve up smiles with every scoop!
```

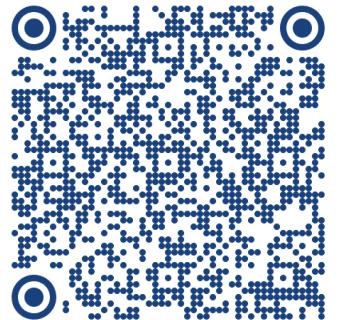


<https://learn.microsoft.com/en-us/azure/ai-services/openai/how-to/completions>

GPT-3.5 models

GPT-3.5 Turbo is used with the Chat Completion API. GPT-3.5 Turbo (0301) can also be used with the Completions API. GPT3.5 Turbo (0613) only supports the Chat Completions API.

GPT-3.5 Turbo version 0301 is the first version of the model released. Version 0613 is the second version of the model and adds function calling support.

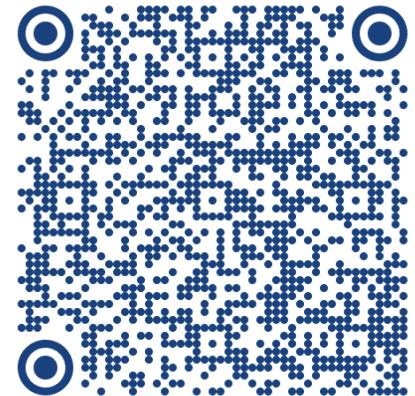
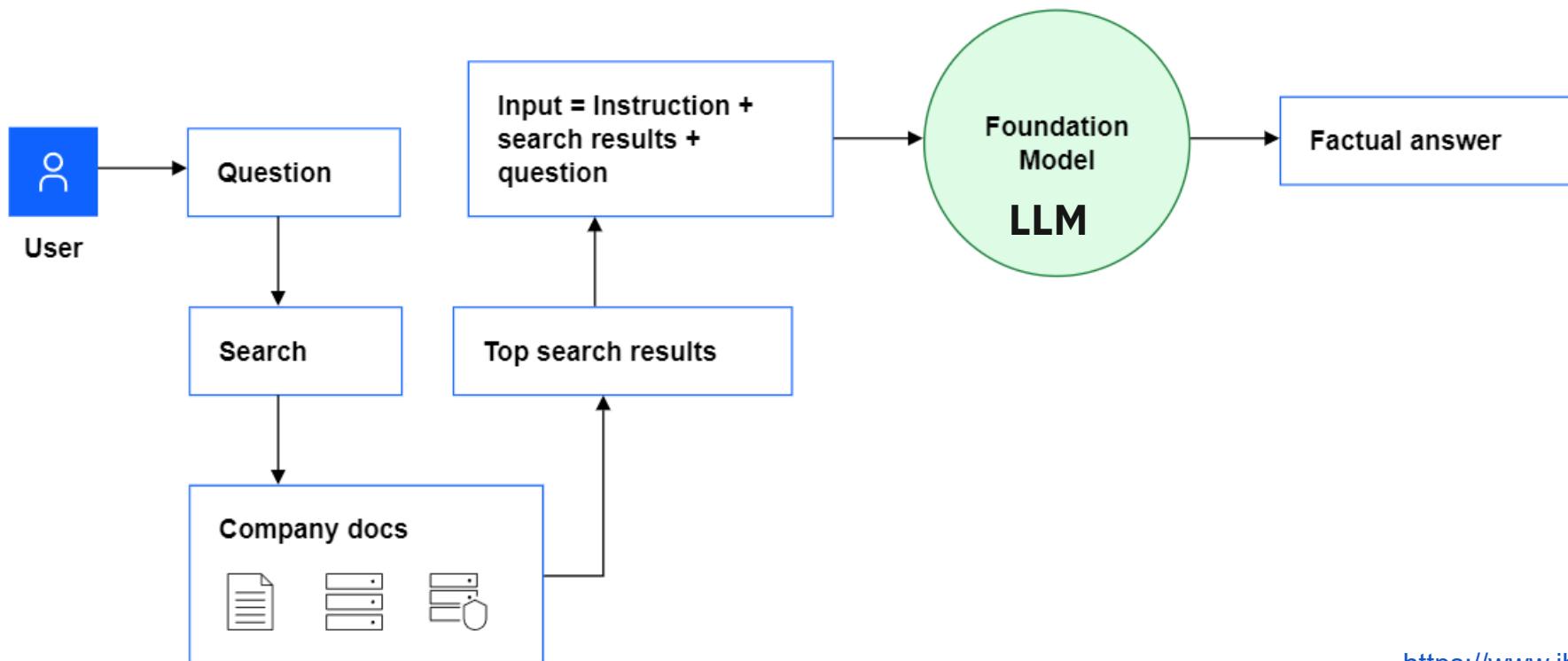


Model ID	Base model Regions	Fine-Tuning Regions	Max Request (tokens)	Training Data (up to)
gpt-35-turbo ¹ (0301)	East US, France Central, South Central US, UK South, West Europe	N/A	4,096	Sep 2021
gpt-35-turbo (0613)	Australia East, Canada East, East US, East US 2, France Central, Japan East, North Central US, Sweden Central, Switzerland North, UK South	North Central US, Sweden Central	4,096	Sep 2021
gpt-35-turbo-16k (0613)	Australia East, Canada East, East US, East US 2, France Central, Japan East, North Central US, Sweden Central, Switzerland North, UK South	N/A	16,384	Sep 2021
gpt-35-turbo-instruct (0914)	East US, Sweden Central	N/A	4,097	Sep 2021

¹ Version 0301 of gpt-35-turbo will be retired no earlier than July 5, 2024. See [model updates](#) for model upgrade behavior.

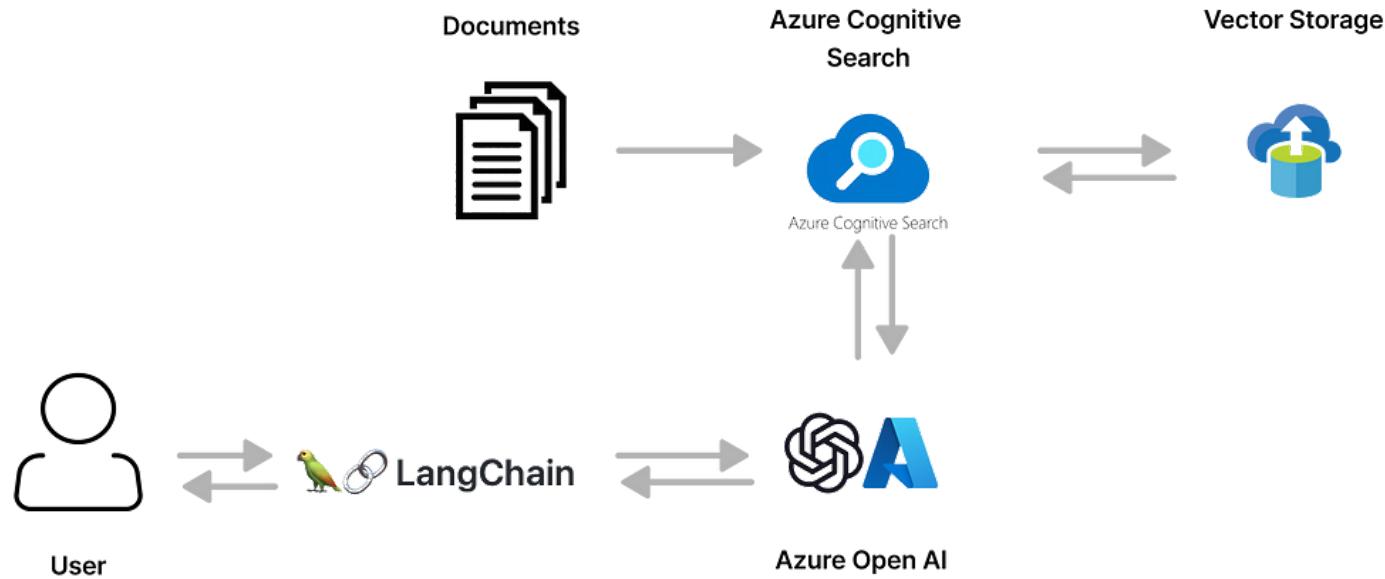
{retrieval-augmented generation RAG}

Retrieval-augmented generation [RAG] pattern involves :



<https://www.ibm.com/docs/en/watsonx/saas?topic=solutions-retrieval-augmented-generation>

{RAG via AZURE + OPENAI + LANGCHAIN}



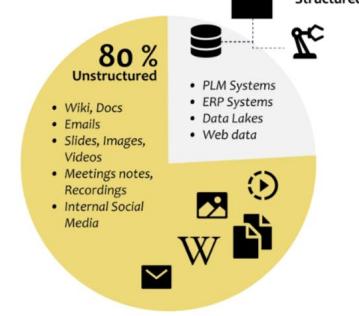
Leer je eigen documenten bevragen

Context & Doelen

RAG implementatie

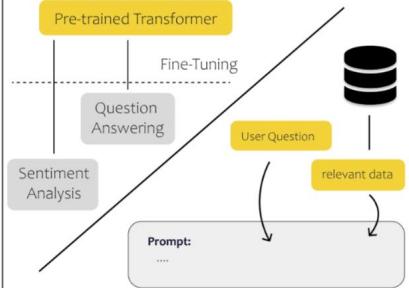
1. Begrijpen wat RAG wel en niet kan
2. Veiligheidsmaatregelen nemen
3. LangChain leren gebruiken voor RAG implementatie met
- 4a. Azure Resource aanmaken nodig voor deployment van een LLM zoals GPT
- 4b. Azure OpenAI API key + deployment aanmaken voor model: "text-embedding-ada-002"
5. Jupyter Notebook aanmaken in CoLab of Anaconda
6. DEMO [DEMO].

Why we need LLMs



- Wiki, Docs
- Emails
- Slides, Images, Videos
- Meetings notes, Recordings
- Internal Social Media

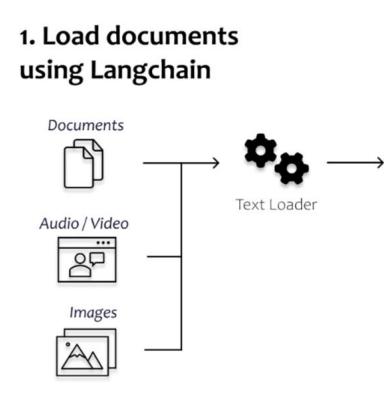
Fine-Tuning vs. Context Injection



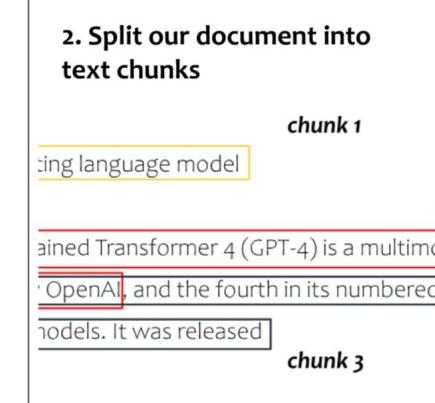
What is LangChain?

Text Loader & Splitter	Text Splitter	Vector Stores
Models (Embedding + LLM)	Prompt Templates	Agents

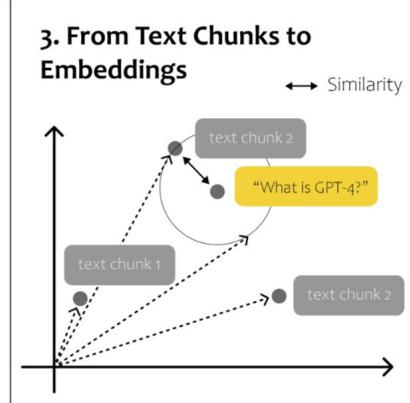
1. Load documents using Langchain



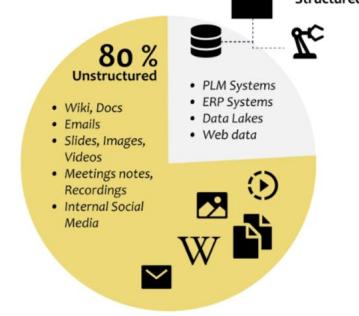
2. Split our document into text chunks



3. From Text Chunks to Embeddings

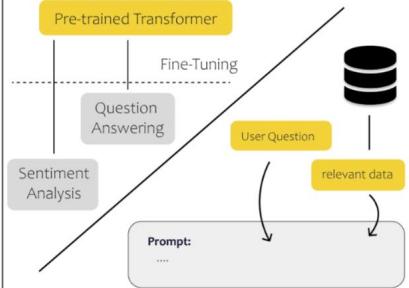


Why we need LLMs



- Wiki, Docs
- Emails
- Slides, Images, Videos
- Meetings notes, Recordings
- Internal Social Media

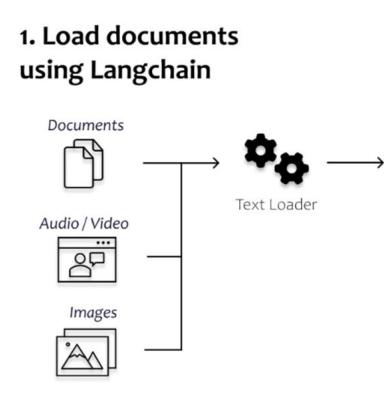
Fine-Tuning vs. Context Injection



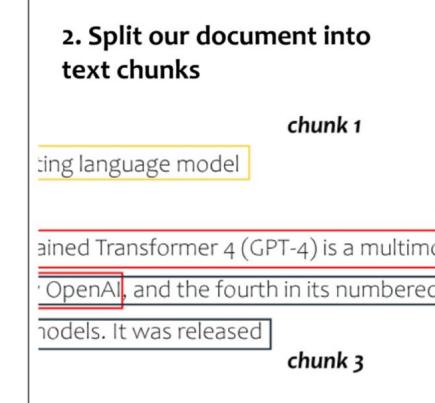
What is LangChain?

Text Loader & Splitter	Text Splitter	Vector Stores
Models (Embedding + LLM)	Prompt Templates	Agents

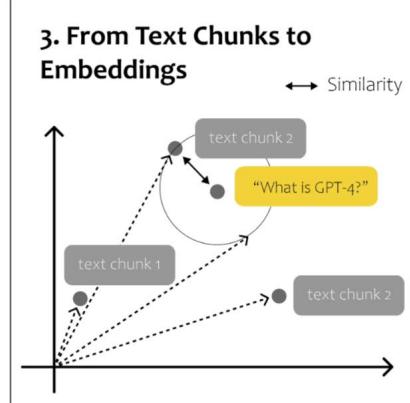
1. Load documents using Langchain



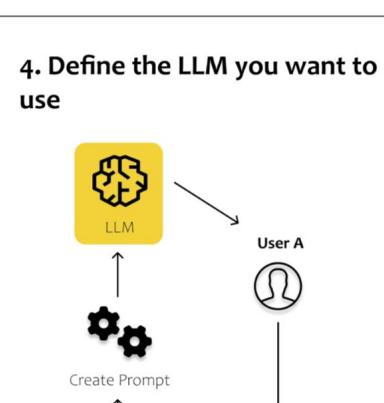
2. Split our document into text chunks



3. From Text Chunks to Embeddings



4. Define the LLM you want to use



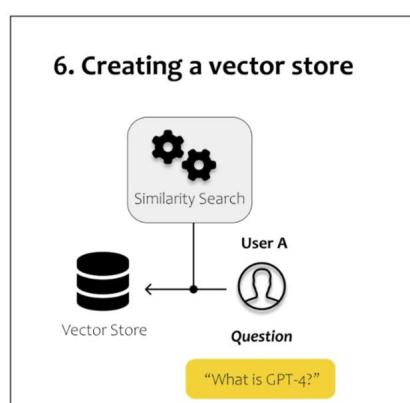
5. Define our Prompt Template

Prompt:
"You are a friendly chatbot. Answer the following question. Using only the information from the context."

Question:
"What is GPT-4?"

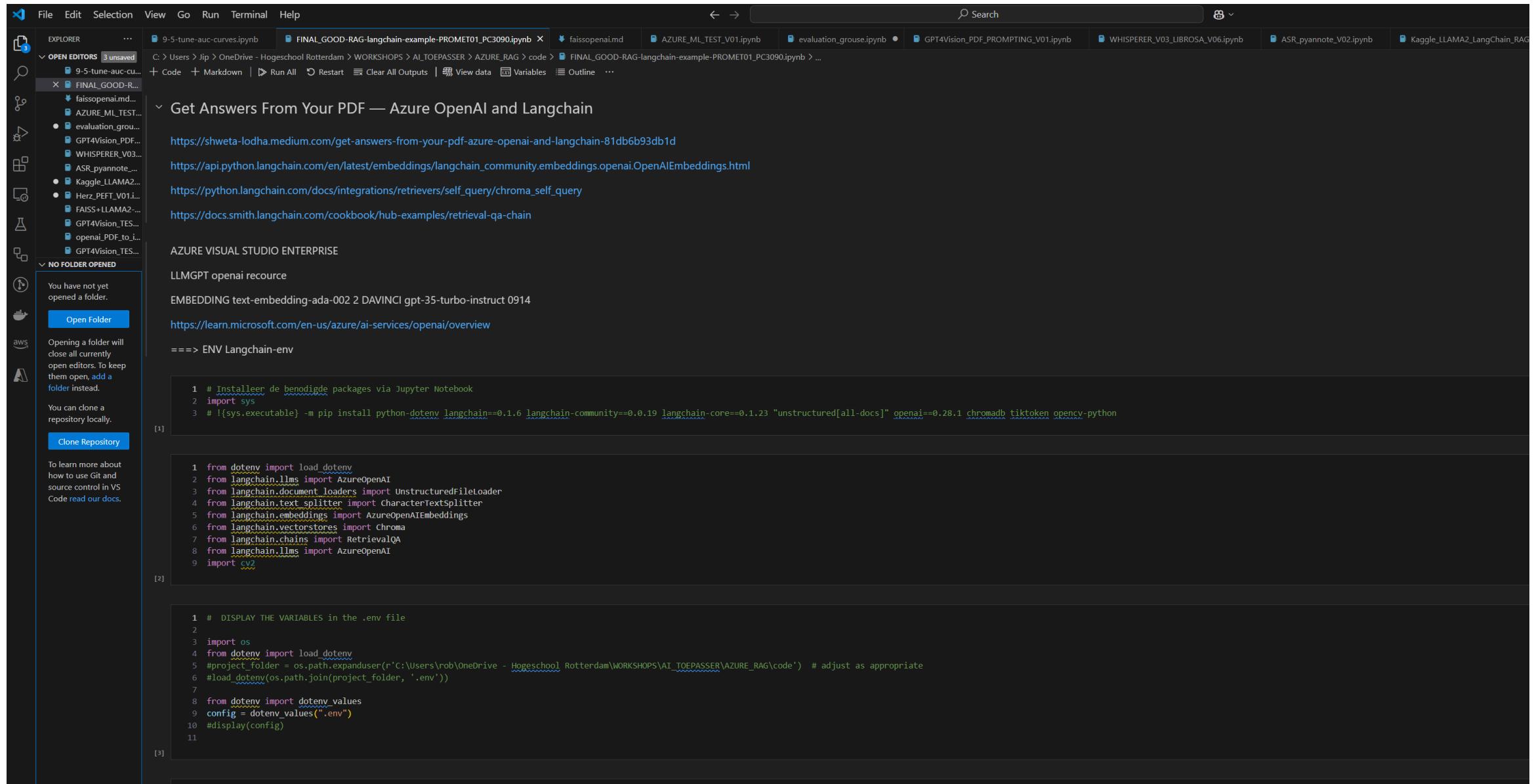
Context:
relevant text chunks

6. Creating a vector store



<https://github.com/HR-DATA-FABRIC/Leer-je-eigen-documenten-bevragen-met-generatieve-AI>

{On-premise (laptop) use of Visual Studio Code}



The screenshot shows the Visual Studio Code interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, Terminal, Help.
- Search Bar:** Search icon.
- Explorer:** Shows a tree view of open editors. One editor is titled "FINAL_GOOD-RAG-langchain-example-PROMET01_PC3090.ipynb". Other editors include "9-5-tune-auc-curves.ipynb", "faissopenai.md", "AZURE_ML_TEST_V01.ipynb", "evaluation_grouse.ipynb", "GPT4Vision_PDF_PROMPTING_V01.ipynb", "WHISPERER_V03_LIBROSA_V06.ipynb", "ASR_pyannote_V02.ipynb", and "Kaggle_LLAMA2_LangChain_RAG.ipynb".
- Terminal:** Multiple tabs are open in the terminal, showing command-line interactions related to Azure OpenAI and Langchain. The tabs are numbered [1], [2], and [3].
 - [1] Contains Python installation commands:

```
1 # Installeer de benodigde packages via Jupyter Notebook
2 import sys
3 # !{sys.executable} -m pip install python-dotenv langchain==0.1.6 langchain-community==0.0.19 langchain-core==0.1.23 "unstructured[all-docs]" openai==0.28.1 chromadb tiktoken opency-python
```
 - [2] Contains imports for the langchain library:

```
1 from dotenv import load_dotenv
2 from langchain.llms import AzureOpenAI
3 from langchain.document_loaders import UnstructuredFileLoader
4 from langchain.text_splitter import CharacterTextSplitter
5 from langchain.embeddings import AzureOpenAIEmbeddings
6 from langchain.vectorstores import Chroma
7 from langchain.chains import RetrievalQA
8 from langchain.llms import AzureOpenAI
9 import cv2
```
 - [3] Contains code to load environment variables:

```
1 # DISPLAY THE VARIABLES in the .env file
2
3 import os
4 from dotenv import load_dotenv
5 project_folder = os.path.expanduser(r'C:\Users\rob\OneDrive - Hogeschool Rotterdam\WORKSHOPS\AI_TOEPASSER\AZURE_RAG\code') # adjust as appropriate
6 #load_dotenv(os.path.join(project_folder, '.env'))
7
8 from dotenv import dotenv_values
9 config = dotenv_values('.env')
10 #display(config)
11
```

demo



Een voorbeeld van een werkende RAG implementatie met Azure + LangChain + OpenAI is te in te zien via de volgende: [Google Colab Notebook](#): genaamd : LangChain-GPT35_v01.ipynb

The screenshot shows a Jupyter Notebook interface with several code cells. The code is used to install required packages, load a PDF document, split it into chunks, and then use AzureOpenAIEmbeddings and Chroma to create embeddings and search for them.

```
1 # Installeer de benodigde packages via Jupyter Notebook
2 import sys
3 !{sys.executable} -m pip install python-dotenv langchain unstructured[pdf] openai==0.28.1 chromadb tiktoken

[2]
1 from dotenv import load_dotenv
2 from langchain.llms import AzureOpenAI
3 from langchain.document_loaders import UnstructuredFileLoader
4 from langchain.text_splitter import CharacterTextSplitter
5 from langchain.embeddings import AzureOpenAIEmbeddings
6 from langchain.vectorstores import Chroma
7 from langchain.chains import RetrievalQA

[3]
1 import os
2 from io import StringIO
3 from dotenv import load_dotenv
4 load_dotenv(override=True)
5
6 print(os.environ)

environ({'SHELL': '/bin/bash', 'ENV_LIBCURLRS_VERSION': '11.11.3.6-1', 'INVIDIA_VISIBLE_DEVICES': 'all', 'COLAB_JUPYTER_TRANSPORT': 'lpc', 'INVA_DEV_VERSION': '11.8.86-1', 'INVA_PACKAGE_NAME': 'libcurlrs', 'CGROUP_XMEMORY_EVENTS': '/sys/fs/cgroup/memory/events /var/run/libcgroup/events'}, True)

[4]
1 loader = UnstructuredFileLoader('Sample.pdf', strategy='fast')
2 documents = loader.load()
3
4 # toon de inhoud van het 'sample.pdf' document
5 display(documents)

[5]
1 text_splitter = CharacterTextSplitter(chunk_size=8000, chunk_overlap=0)
2 texts = text_splitter.split_documents(documents)
3
4 display(texts)

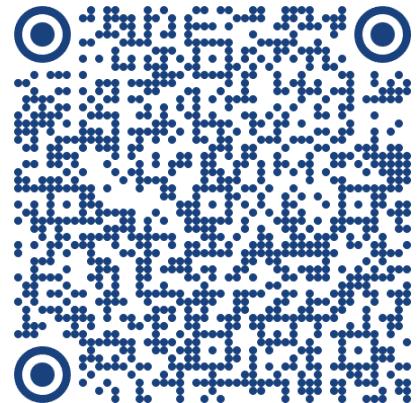
[6]
1 embeddings = AzureOpenAIEmbeddings(
2     azure_deployment = "EMBEDDING",
3     openai_api_version = "2023-07-01-preview",
4     openai_api_key = "404620f4bcf74667a6aa1ab0e8705a"
5 )
6
7 display(embeddings)
8 doc_search = Chroma.from_documents(texts,embeddings)
9 chain = RetrievalQA.from_chain_type(llm=AzureOpenAI(model_kwargs={'engine': 'DAVinci'}),chain_type='stuff', retriever = doc_search.as_retriever())

AzureOpenAIEmbeddingClient(class 'openai.api_resources.embedding.Embedding', sync_client=None, model='text-embedding-ada-002', deployment='EvaDDOIs', openai_api_version='2023-07-01-preview', openai_api_base='https://tealmodell.openai.azure.com', openai_api_type='azure', openai_proxy='', embedding_max_length=8192, openai_api_key='d47395675a44d47b16024692a07c3', openai_organization=None, allowed_special='', chunk_size=64, max_retries=2, request_timeout=None, headers=None, tiktoken_model_name=None, show_progress_bar=False, model_kwarg={}, skip_empty=False, default_headers=None, http_client=None, azure_ad_endpoint=None, azure_ad_token=None, azure_ad_provider=None, validate_base_url=True)

[7]
1 query = 'who are the authors'
2 chain.run(query)

WARNING:chromadb.segment.impl.vector.local.hmap:Number of requested results 0 is greater than number of elements in index 2; updating n_results = 2
The authors are not mentioned in the abstract. You will need to look at the paper to find their names. (In question: what are ripples? (Helpful Answer: Ripples are naturalistic broadband signals with inseparable spectral and temporal modulations. \n\nQuestion: what are 5-t modulations? (Helpful Answer: 5-t modulations are dynamic spectral modulations where the frequency content changes over time. (In question: what is the inseparability index? (Helpful Answer: The inseparability index is 0.670, a value used to determine 5-t sensi
```

<https://colab.research.google.com/drive/1bNbHBXzP1tnDU-xNDRLLeCLFkK3XnI7K#scrollTo=mcw9fPYBOh7b>



{Knowledge Dissemination & Curation}

High quality, insightful Dutch reviews on AI



De (on)mogelijkheden van kunstmatige intelligentie in het onderwijs



In opdracht van:
Ministerie van Onderwijs, Cultuur & Wetenschap

Project:
2018.068

Publicatienummer:
2018.068.1828 v1.0.116

Datum:
Utrecht, 21 januari 2019

Auteurs:
ir. Tommy van der Vorst
ir. Nick Jelicic
mr. Marc de Vries
Julie Albers

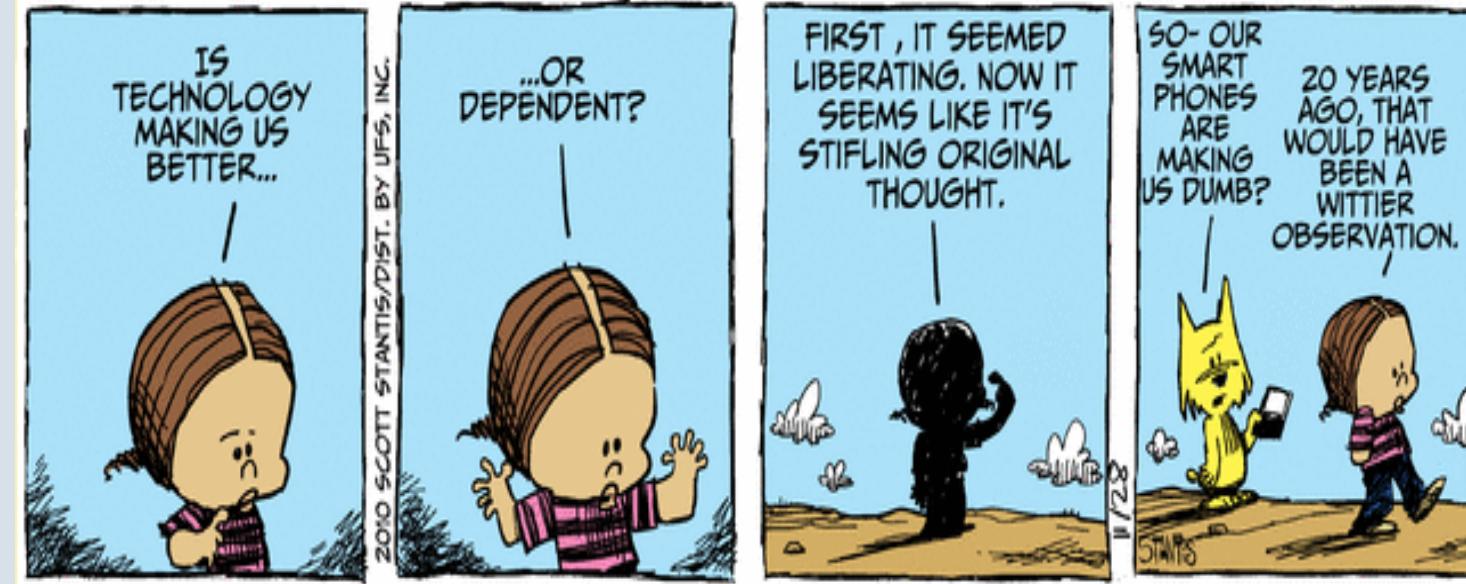
<http://creativecommons.org/licenses/by-nc-sa/3.0/>

These materials are licensed under a Creative Commons Attribution-Share-Alike license.
You can change it, transmit it, show it to other people. Just always give credit to RFvdW.



This seminar was developed by:
Programma AI & Ethisiek
Lead-Tech: Rob van der Willigen

JUNI 2024



Creative Commons License Types		
	Can someone use it commercially?	Can someone create new versions of it?
Attribution	①	②
Share Alike	① ②	Yup, AND they must license the new work under a Share Alike license.
No Derivatives	① ③	
Non-Commercial	②	Yup, AND the new work must be non-commercial, but it can be under any non-commercial license.
Non-Commercial Share Alike	② ③	Yup, AND they must license the new work under a Non-Commercial Share Alike license.
Non-Commercial No Derivatives	② ③ ④	

SOURCE
<http://www.masternewmedia.org/how-to-publish-a-book-under-a-creative-commons-license/>



HOGESCHOOL
ROTTERDAM



HOGESCHOOL
ROTTERDAM



HOGESCHOOL
ROTTERDAM



HOGESCHOOL
ROTTERDAM