



# ENGINEERING MEANINGFUL HUMAN CONTROL

## HOE REALISEREN WE DE ROL VAN DE MENS IN AUTOMATISCHE (\*AI) SYSTEMEN

Prof.dr.ir. Inald Lagendijk / TU Delft, the Netherlands

# Rekenkamer: nauwelijks aandacht voor ethiek bij algoritmes overheid

26 januari 2021 17:43

Aangepast: 26 januari 2021 18:07

Er gaat van alles mis met het gebruik van algoritmes door de overheid, concludeert de Algemene Rekenkamer in een kritisch rapport. Op tal van ministeries wordt hiervoorbeeld geen aandacht besteed aan de mogelijkheid van discriminatie.

Op slechts drie prioriteitenlijstjes gesprekken de ministeries. Waaronafhankelijk was niet het eerste doel.

Het doel was dat de overheid is gevoerd door ministeries van

## Wat is een algoritme?

Een algoritme is een reeks op grote schaal uitgewerkte stappen die kunnen worden volgt.

Zo kan een computer sporen, door te wijzen.

## Geen overzicht

Pas na het bezoeken veel ministeries blijkt dat er weinig bekend is over welke algoritmes er eigenlijk al gebruikt worden.

rtlnieuws

## Rekenkamer: ministers weten niet welke algoritmes ze gebruiken

Jan Fred van Wijnen 26 Jan

fd.

EUROPA

# EU legt kunstmatige intelligentie gedragsregels op

Mathijs Schippers  
Brussel

Een mijlpaal, noemt Eurocommissaris Margrethe Vestager het. De Deense presenteerde gisteren het langverwachte Brusselse voorstel voor het reguleren van kunstmatige intelligentie.

Het is 's werelds eerste poging om de ongewenste gevolgen van de zich in volle vaart ontluikende technologie aan banden te leggen. En Brussel hoopt dat goed voorbeeld zal doen volgen.

'Vertrouwen is noodzakelijk bij kunstmatige intelligentie en niet iets dat leuk is om erbij te hebben', zei Vestager. Ze bestrijdt dat Europa het risico loopt terrein te verliezen bij de ontwikkeling van toepassingen gebaseerd op kunstmatige

intelligentie door als enige beperkingen op te werpen. Ze liet juist weten bang te zijn dat mensen de beloftevolle technologie links laten liggen als er geen paal en perk gesteld wordt aan de negatieve aspecten ervan.

Kunstmatige intelligentie verwijst naar het leer vermogen van computers. Die kunnen snel grote hoeveelheden data verwerken, wat handig kan zijn bij het bepalen van iemands kredietwaardigheid of geschiktheid voor een baan. Maar kwaadwillenden kunnen de systemen zo inrichten dat bij de besluiten bijvoorbeeld sprake is van etnisch profileren en andere methodes die op gespannen voet staan met de fundamentele Europese waarden, zoals het recht op privacy.

Het gros van de toepassingen is volgens de Europese Commissie risicoloos. Daar zijn geen aanvul-

lende eisen voor nodig of slechts beperkte, zoals bij het gebruik van chatbots. In dat geval wil de Commissie dat de gebruiker weet dat hij met een computer praat.

Maar voor sommige toepassingen is meer vereist. Die kennen een 'hoog risico' van misbruik, zoals computersystemen die worden ingezet door overheden die willen vaststellen of iemand voor een uitspraak in aanmerking komt. Voordat zo'n systeem mag worden gebruikt, moet een test onder meer uitwijzen of het transparant en

**Systemen met een hoog risico op misbruik moeten op transparantie worden getest**

nauwkeurig is en of menselijk toezicht mogelijk is.

Daarnaast is er nog een aantal toepassingen dat volgens de Commissie ronduit verboden moet worden, zoals systemen die het gedachtegoed van iemand proberen te manipuleren of de kwetsbaarheid van kinderen exploiteren. Bedrijven die zich niet aan de regels houden, riskeren boetes tot €30 mln of 6% van de mondiale omzet, al naar gelang wat hoger is.

Er is bijval en kritiek voor het plan. 'Een evenwichtig startpunt', zegt Lawrence Kercknawi van public-affairsbureau Political Intelligence in Brussel, dat opkomt voor de belangen van de techsector. 'Het is goed dat Brussel een streep in het zand zet en zegt: dit vinden we onaanvaardbaar', zegt Bart Schermer van het Nederlandse public-affairsbureau Considerati.

fd.

Het Financiële Dagblad  
Donderdag 22 april 2021

# Uber's self-driving operator charged over fatal crash

16 September 2020

The back-up driver of an Uber self-driving car that killed a pedestrian has been charged with negligent homicide.

BBC  
NEWS

# Het AI Debat

- AI technologie brengt veel voordelen.
- Beloftevolle technologie vergt ook inperken van negatieve aspecten.
- Er zijn tientallen (high level) richtlijnen op ethisch, juridisch en sociaal vlak.
- En nu?
- Wat is eigenlijk het probleem? Verantwoordelijkheid? Voor wat?



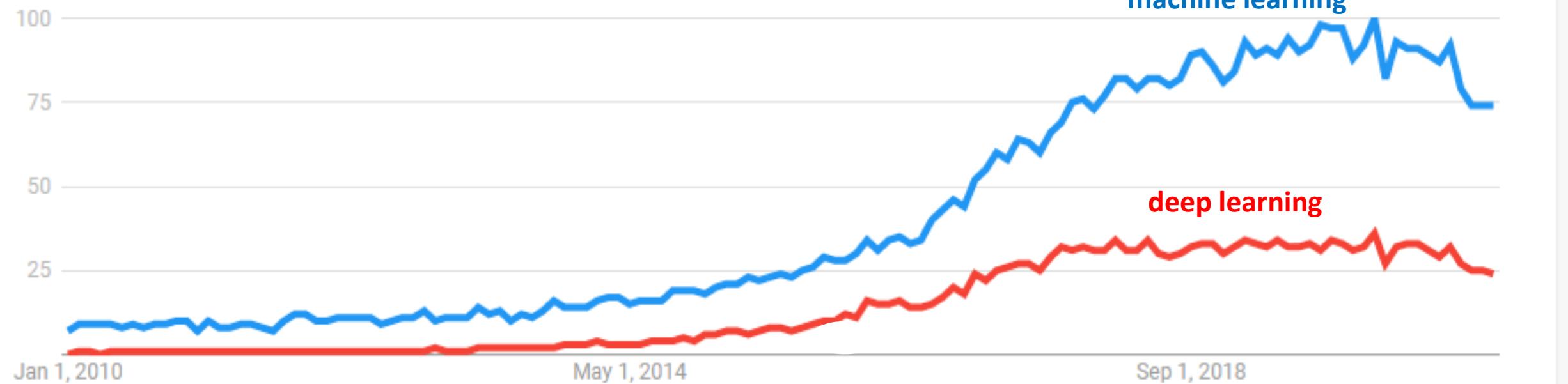
Google AI About Responsibilities Research Economic  
Overview Principles Responsible AI practices Review Process Publications  
**Objectives for AI applications**  
We will assess AI applications in view of the following objective.  
We believe that AI should:  

1. Be socially beneficial.
2. Avoid creating or reinforcing unfair bias.
3. Be built and tested for safety.
4. Be accountable to people.
5. Incorporate privacy design principles.
6. Uphold high standards of scientific excellence.
7. Be made available for uses that accord with these principles.

# Machine Leren: de AI Succesmotor

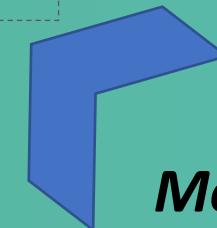
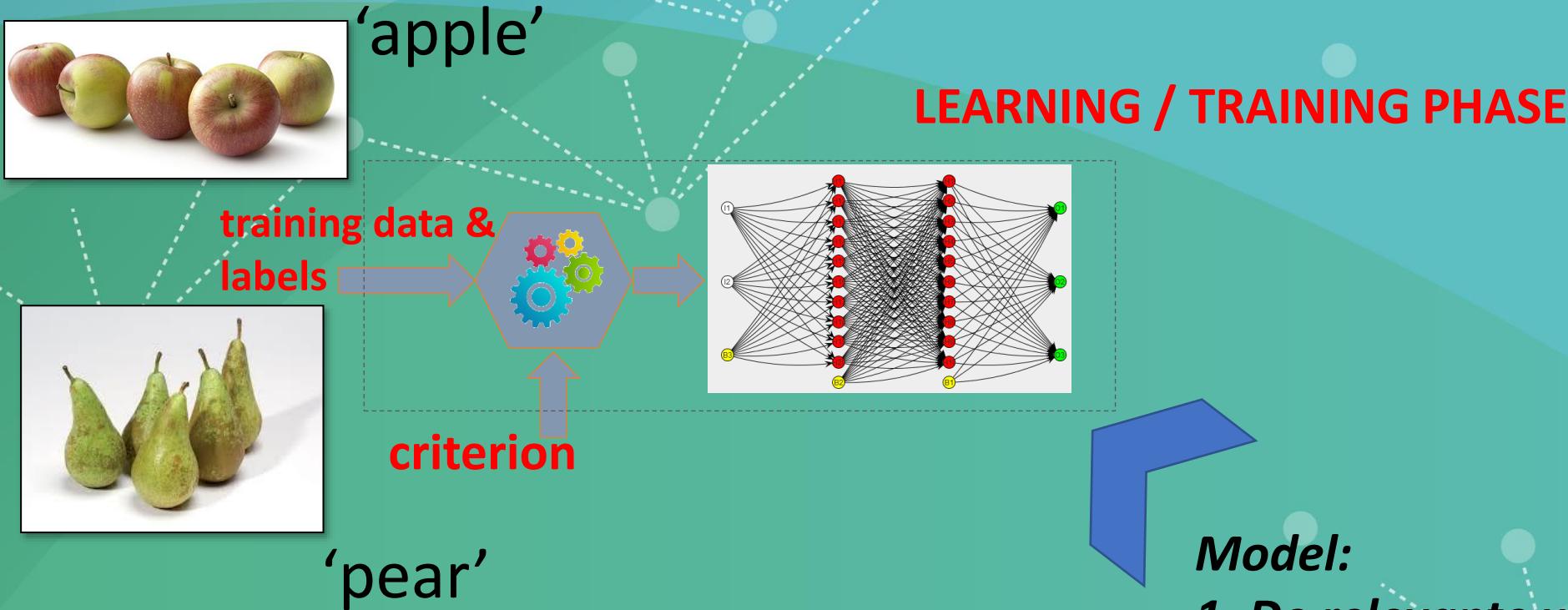


## Google Trends



*'leervermogen op basis van veel data'*

# Machine Leren: de AI Succesmotor



**Model:**

- 1. De relevante wereld kennis*
- 2. Relevante criteria voor prestatie*

# Machine L



training data &  
labels



training data &  
labels

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool

## Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

Perceptron (P) Feed Forward (FF) Radial Basis Network (RBF)

Recurrent Neural Network (RNN)

Long / Short Term Memory (LSTM)

Gated Recurrent Unit (GRU)

Auto Encoder (AE)

Variational AE (VAE)

Denoising AE (DAE)

Sparse AE (SAE)

Markov Chain (MC)

Hopfield Network (HN)

Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)

Deep Convolutional Network (DCN)

Deconvolutional Network (DN)

Deep Convolutional Inverse Graphics Network (DCIGN)

Generative Adversarial Network (GAN)

Liquid State Machine (LSM)

Extreme Learning Machine (ELM)

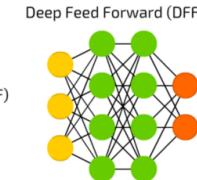
Echo State Network (ESN)

Deep Residual Network (DRN)

Kohonen Network (KN)

Support Vector Machine (SVM)

Neural Turing Machine (NTM)

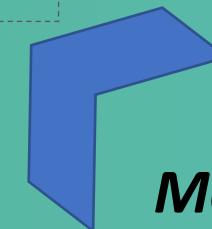
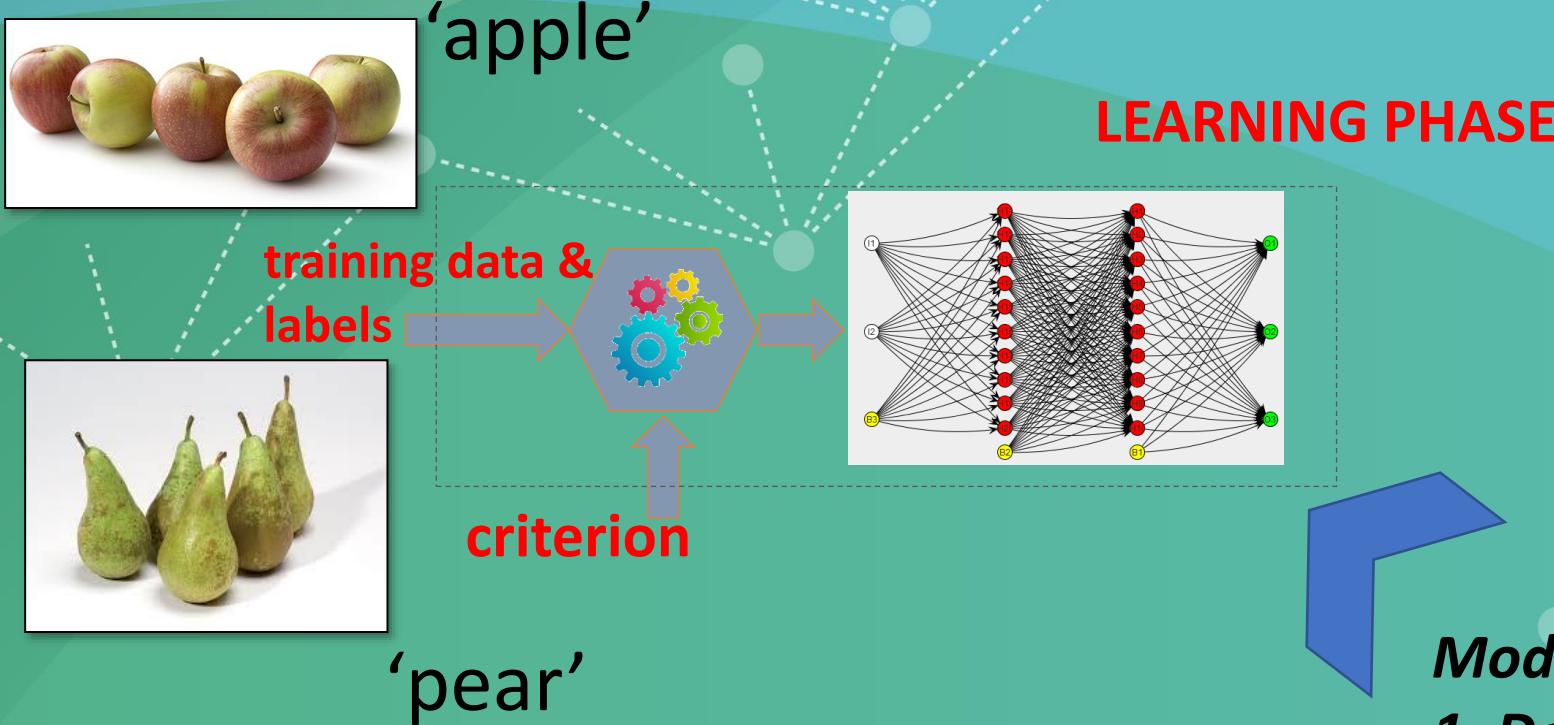


smotor



:  
*elevante wereld kennis  
vante criteria voor prestatie*

# Machine Leren: de AI Succesmotor

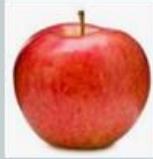


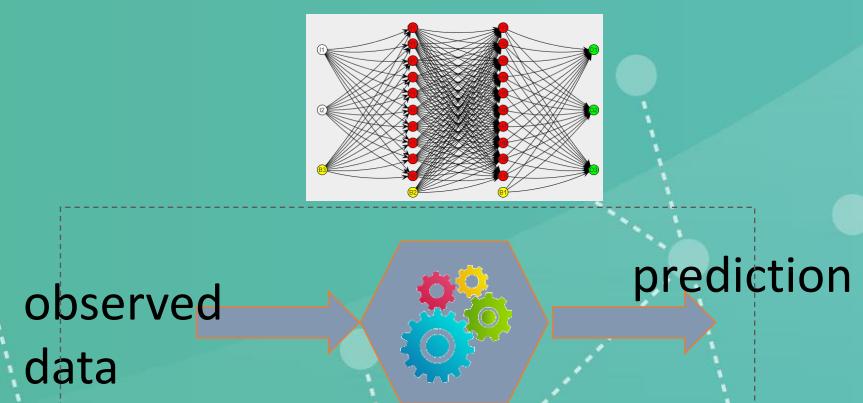
**Model:**

- 1. De relevante wereld kennis*
- 2. Relevante criteria voor prestatie*

# Machine Leren: de AI Succesmotor

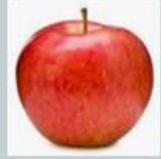
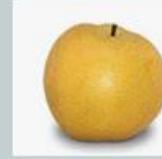


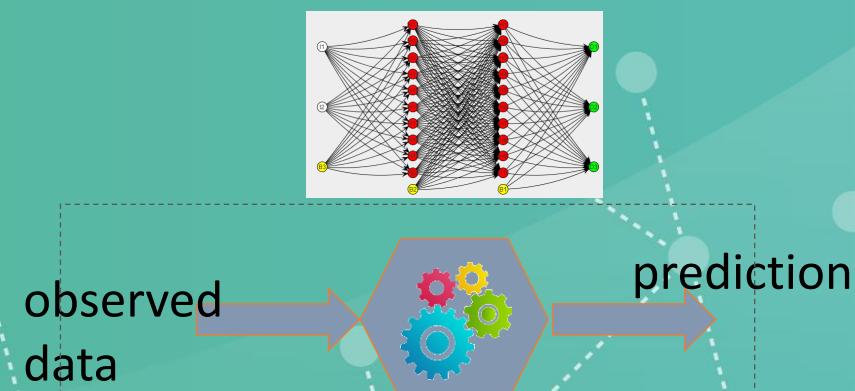
		ACTUAL	
		Apple	Pear
PREDICTED	Apple		
	Pear		



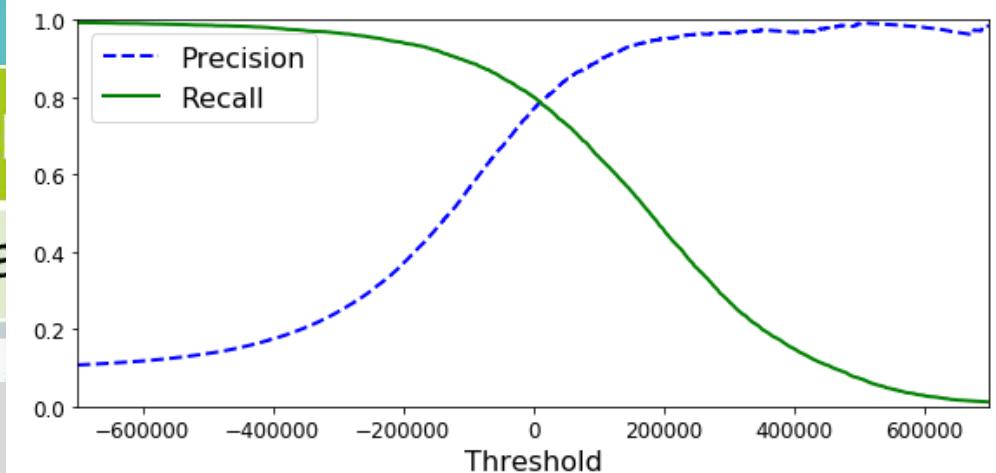
**Criterion:**

- False positives
- False negatives
- Precision
- Recall
- F1 measure

		ACTUAL	
		Apple	Pear
PREDICTED	Apple		
	Pear		



		ACTUA	
		Apple	Pea
	Apple		
	Pea		



## Antwoord 1: AI / Algoritmen maken ‘altijd’ fouten

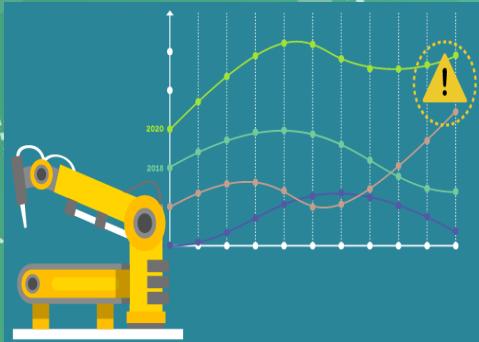
- Zijn alle fouten even erg (gebruikt criterium)?
- Kunnen wij (of AI) weten dat er een fout gemaakt wordt?
- Kunnen we begrijpen waarom (explainable, transparent)?
- Wie is verantwoordelijk (in de keten van betrokkenen)?
- ODD: operational design domain.



# Data? Criterium? Errors? Impact?



## Predictive maintenance



## Life expectancy prediction



## Job suitability prediction



## Sexual orientation prediction



**Criterion** Estimate of repair costs

**Data**  
\* vibration spectrum  
\* resource usage  
\* current signature

**Advantages** Maintenance planning  
**Risks** Reduced value of expertise

Risk score

\* blood pressure  
\* creatine  
\* white blood cell count

Informed treatment decisions  
Root cause unclarity

Avoid hiring unsuitable

\* facial emotion  
\* voice timbre  
\* vocabulary

High volume selection  
Reduced self presentation

Reproduce human ability

\* facial features  
\* morphology  
\* grooming

Security, marketing  
Stigmatization

# Data? Criterium? Errors? Impact?

Predictive maintenance

Life expectancy prediction

Job suitability prediction

WRONG WAY

## Framing and Formalism trap

- Formulate goal of the system as an **optimization criterion** (Aification).
- What are you **optimizing** for? How much of the world does the data/model capture? (reward hacking and tunnel vision).

Advantages

Maintenance planning

Risks

Reduced value of expertise

Informed treatment decisions

Root cause unclarity

High volume selection

Reduced self presentation

Security, marketing

Stigmatization

	Qualified person	Not-qualified person
Invited for interview	100	0
Not-invited for interview	0	100

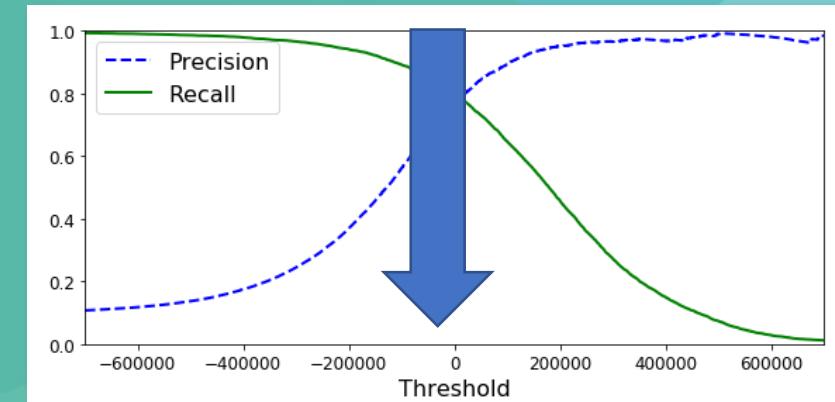


**'Unlikely ...'**

	Qualified person	Not-qualified person
Invited for interview	90	15
Not-invited for interview	10	85



**REWARD = 0.87**



# UNFAIR!

ALL

	Qualified person	Not-qualified person
Invited for interview	90	15
Not-invited for interview	10	85



TECH

SUBGROUP A

	Qualified person	Not-qualified person
Invited for interview	95	20
Not-invited for interview	5	80

SUBGROUP B

	Qualified person	Not-qualified person
Invited for interview	75	10
Not-invited for interview	25	90

# WHAT IS FAIR?

**ALL**

	Qualified person	Not-qualified person
Invited for interview	90	15
Not-invited for interview	10	85



TECH

**SUBGROUP A**

	Qualified person	Not-qualified person
Invited for interview	95	20
Not-invited for interview	5	80

**SUBGROUP B**

	Qualified person	Not-qualified person
Invited for interview	75	10
Not-invited for interview	25	90

**ALL**

	Qualified person	Not-qualified person
Invited for interview	90	15
Not-invited for interview	10	85



TECH

**SUBGROUP A**

	Qualified person	Not-qualified person
Invited for interview	95	20
Not-invited for interview	5	80

**SUBGROUP B**

	Qualified person	Not-qualified person
Invited for interview	75	10
Not-invited for interview	25	90

ALL

	Qualified person	Not-qualified person
Invited for interview	90	15
Total	105	30



## Antwoord 2: Vertrouwen op alleen numerieke criteria is problematisch

- Slechts dát wordt geoptimaliseerd.
- Onverwachte neveneffecten (reward hacking).
- Normen, waarden, wetten versus numeriek criterium.
- Startpunt is dus: het ideale criterium bestaat niet.
- Maar je hebt er wel één nodig. Wie is verantwoordelijk?

# Autonomous System for Good

Law and values: autonomy, beneficence, justice

Instrumental values: privacy, explainability, ...

**the new tower of Babel**



**many decisions are human**

Deep learning from data. Reasoning. Optimization, ...

# Autonomous System for Optimal Performance

Lower costs, better diagnosis, fewer malfunctioning, more mobility, ...

# Meaningful Human Control



In order for humans to **have**, and be **able** to take control over the decisions of an autonomous system.

# Meaningful Human Control



In order for humans to have, and be able to take control over the decisions of an autonomous system.

Tracking



human-AI system should be responsive to the **human (moral) reasons** relevant in the circumstances.

*Captures 'criterium issue'*

# Meaningful Human Control



In order for humans to have, and be able to take control over the decisions of an autonomous system.

## Tracking



human-AI system should be responsive to the **human (moral) reasons** relevant in the circumstances.

*Captures 'criterium issue'*

## Tracing



human-AI system behavior, capabilities, and possible effects in the world should be traceable to a proper moral understanding on the part of at least **one relevant human agent** who designs or interacts with the system.

*Captures 'errors issue'*

# Autonomous System for Good

Law and values: autonomy, beneficence, justice

Instrumental values: privacy, explainability, ...



Deep learning from data. Reasoning. Optimization, ...

# Autonomous System for Optimal Performance

Lower costs, better diagnosis, fewer malfunctioning, more mobility, ...

# Autonomous System for Good

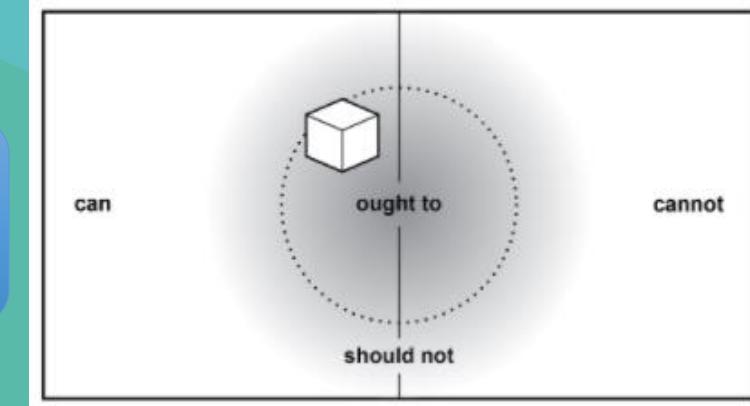
Law and values: autonomy, beneficence, justice

Instrumental values: privacy, explainability, ...

Responsibilities

*Methodology → Moral Operational Design Domain*

Deep learning from data. Reasoning. Optimization, ...



# Autonomous System for Optimal Performance

Lower costs, better diagnosis, fewer malfunctioning, more mobility, ...

# Autonomous System for Good

Law and values: autonomy, beneficence, justice

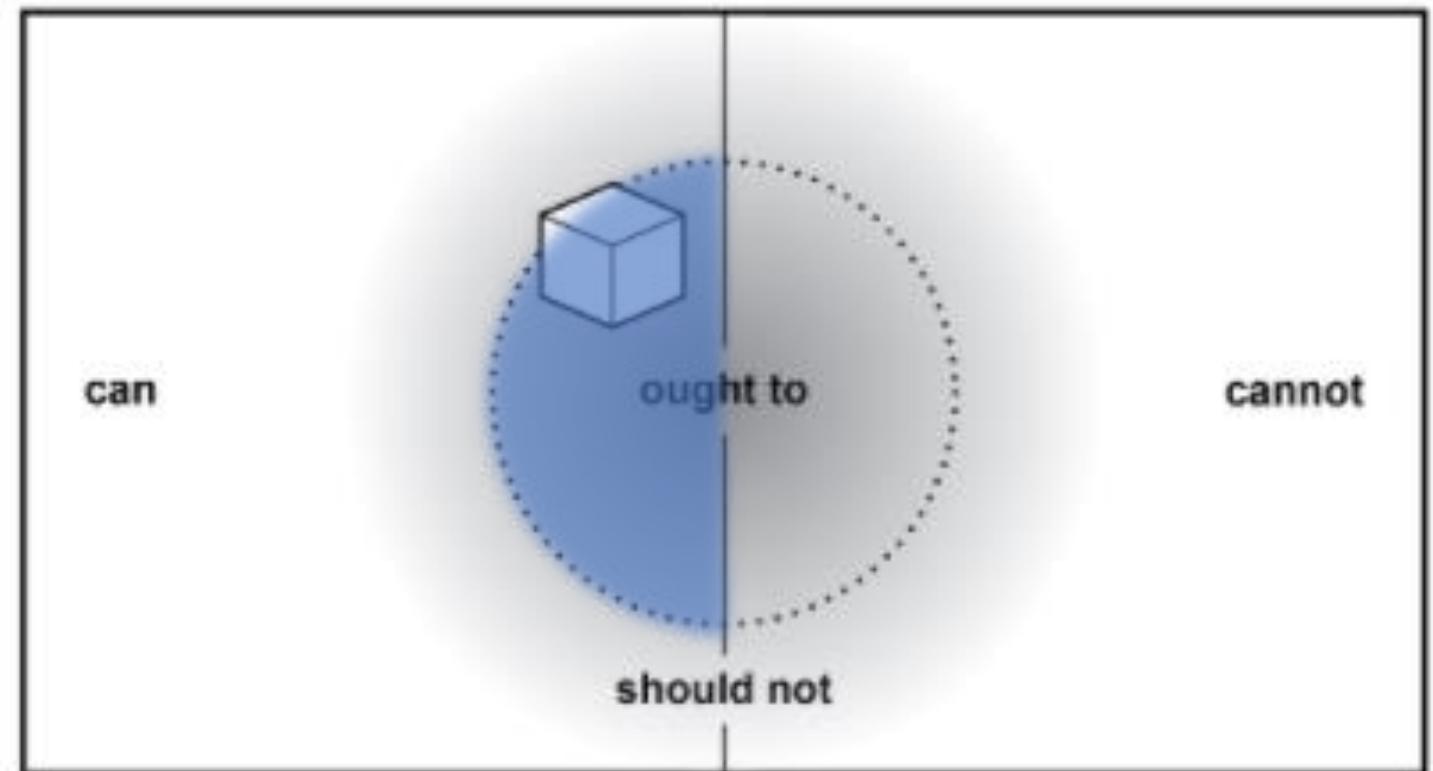
Responsibilities

*Methodology*

Instr

D

Deep le



# Autonomous System for Optimal Performance

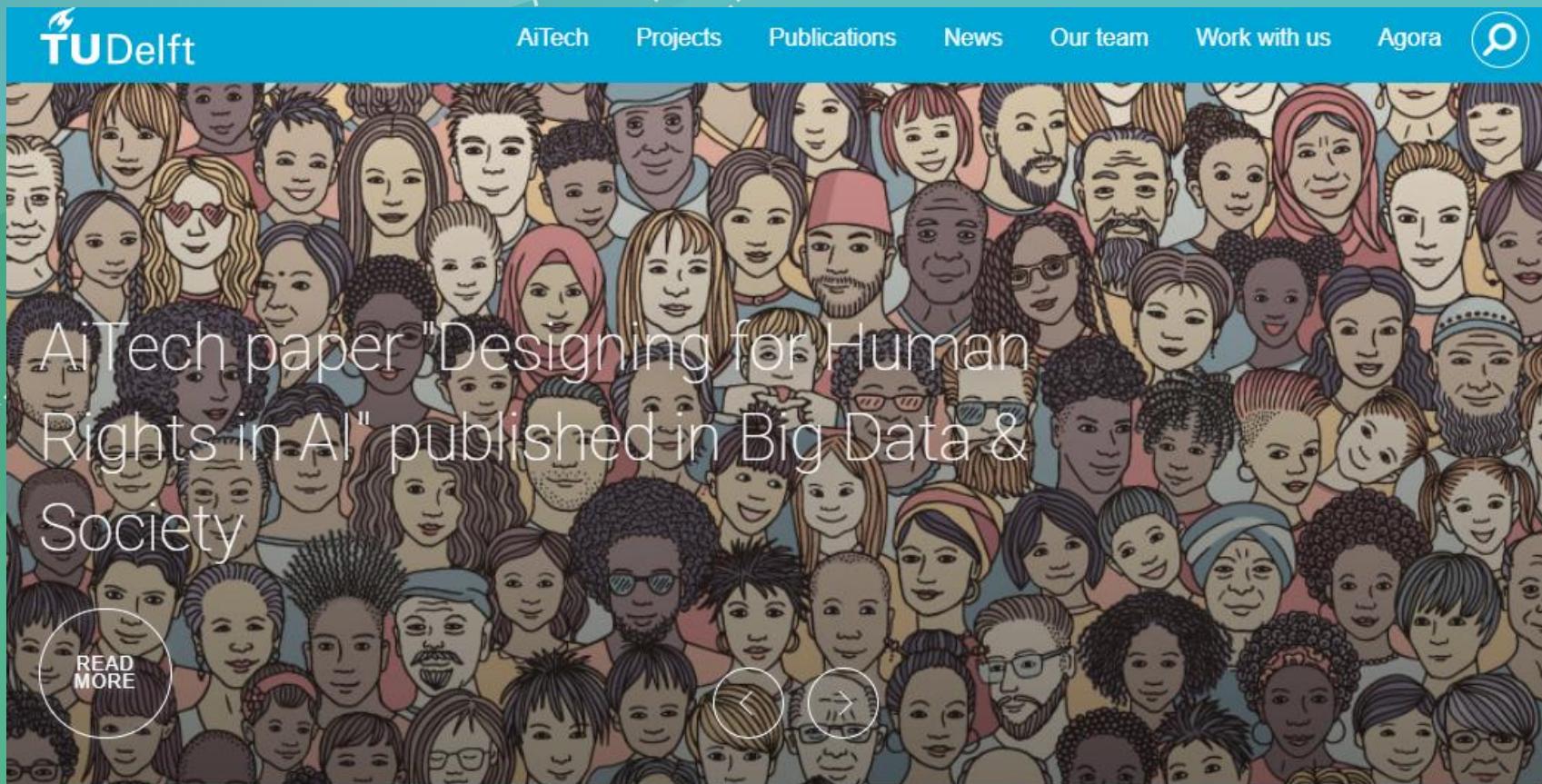
Lower costs, better diagnosis, fewer malfunctioning, more mobility, ...

		<b>Qualified person</b>	<b>Not-qualified person</b>
<b>Inside M-ODD</b>	Invited for interview	75	1
<b>Outside M-ODD</b>	Human action	23	29
<b>Inside M-ODD</b>	Not-invited for interview	2	70

# Properties to Realize

- Alignment of ability, authority and responsibility.
- Adequate and compatible mental models.
- Actions of AI agents are explicitly linked to human decisions.
- And so:
  - *Computational optimization is not always the solution.*
  - *Design AI systems for minimal invasive effect (use only when and where needed).*
  - *AI systems should eventually be aware of their own limitations.*

Visit us at <https://www.tudelft.nl/aitech/>



AiTech paper "Designing for Human Rights in AI" published in Big Data & Society

READ MORE

AiTech is TU Delft's multidisciplinary research program on awareness, concepts, and design & engineering of autonomous technology under meaningful human control

The website header features the TU Delft logo and navigation links for AiTech, Projects, Publications, News, Our team, Work with us, and Agora, along with a search icon.

### Why meaningful human control?

Today's engineers create systems that are ever more equipped with artificial intelligent technologies. Autonomous behavior of cars, robots, and decision support algorithms is becoming a reality. Our vision is that scientists should not only research the technology that makes

### Our 'how to' approach

Meaningful human control is particularly important in cases of failures or conflicts with the normative foundations of society, social conventions, and human acceptability. We believe these challenges demand a multidisciplinary effort, bringing together researchers across



Home &gt; News &gt; Netherlands AI Coalition pleased that backing from the National Growth Fund will let it achieve its aims

## Share

12/04/2021



## Netherlands AI Coalition pleased that backing from the National Growth Fund will let it achieve its aims

News



The AiNed investment programme for artificial intelligence run by the Netherlands AI Coalition (NL AIC) is getting 276 million euros as one of the innovation proposals that the National Growth Fund will be backing. The decision that the advisory committee chaired by Jeroen Dijsselbloem has made about the first round of the National Growth Fund was adopted by the council of ministers today. It includes honouring the 5 proposals for research and innovation that Mona Keijzer, the State Secretary for Economic Affairs and Climate Change, submitted on behalf of the collaborating companies, centres of expertise and authority bodies, including the NL AIC.

Kees van der Klaauw, the coalition manager at NL AIC: "We're delighted to have been told that our AiNed programme has been put forward by the advisory committee and that Phase 1 of the AiNed

- **ELSA Labs** zijn participatieve innovatie omgevingen waarin AI en data technologie en toepassingen in kaart gebracht en gevalideerd worden op een manier dieborgt dat alle relevante groepen betrokkenen de toepassingen als zinvol, haalbaar, verantwoord en wenselijk ervaren