# Digital Ethics – Fears, biases, values and trust*

Zoltán Szlávik

IBM **Benelux Center for Advanced Studies**

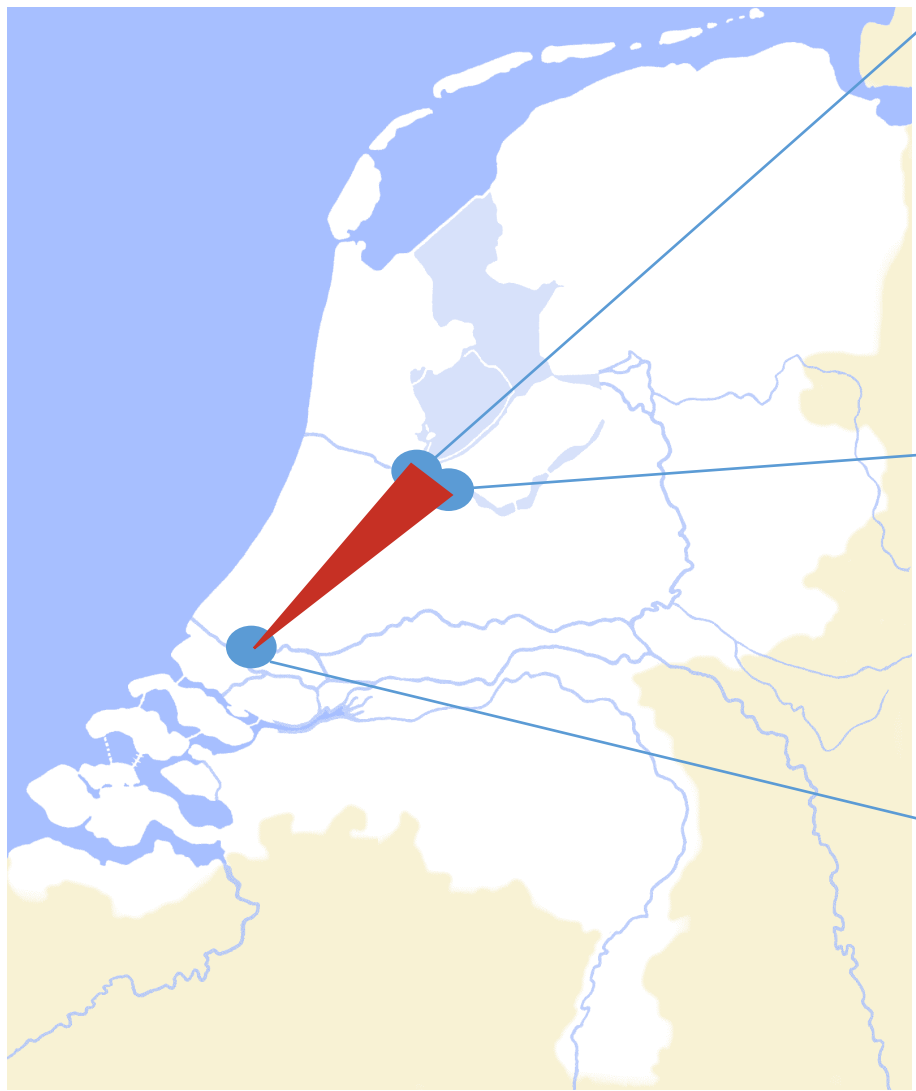*selected slides from the STM Innovation Day presentation

*Warning: contains personal observations, too.*

IBM

**Research**

**CAS**

*Innovation*

*Education*

IBM
**Center for**
**Advanced Studies**

@zolley

# CAS BNL

## CAS BeNeLux

**Zoltán Szlávik**
Lead/Research

**Benjamin Timmermans**
Research/VU

**Dorottya Mezőfi**
Programme
Management

**Manfred Overmeen**
Development

**Santiago Gaitan**
Research/TUD

## CIC VU Amsterdam

**Lora Aroyo**
Faculty Fellow

**Anca Dumitrache**
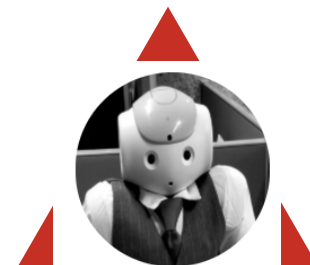Phd Fellow

**Oana Inel**
Phd Fellow

## CIC TU Delft
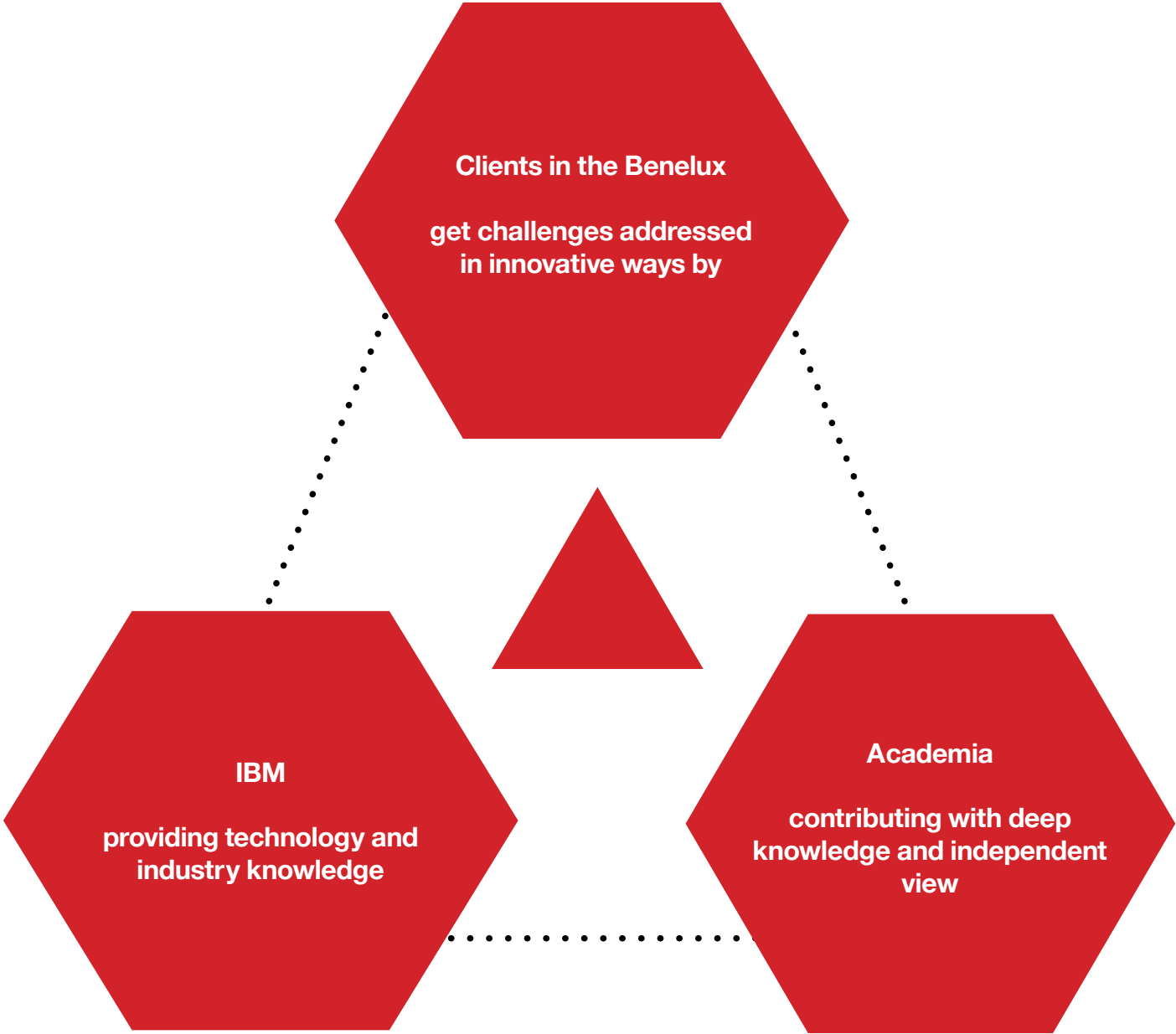
**Alessandro Bozzon**
Faculty Fellow

**Peter Hofstee**
IBM Austin Lab

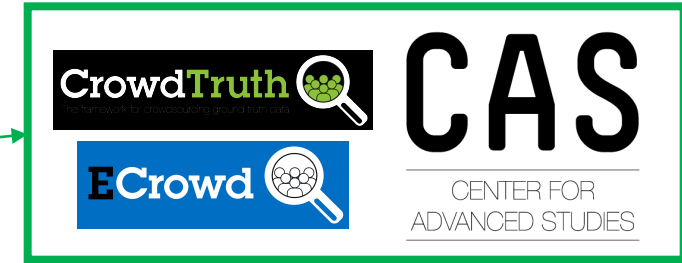**Robbie, Rosita & Cas**
Computer Overlords
Also, events

IBM
**Center for
Advanced Studies**

@zolley

IBM

CAS:
Collaboration
with Academia

Clients in the Benelux

get challenges addressed
in innovative ways by

IBM

providing technology and
industry knowledge

Academia

contributing with deep
knowledge and independent
view

IBM
Center for
Advanced Studies

@zolley

**Professional**

**CAS Training**

**Student**

**The Machine (AI)**

@zolley

# Where CAS Research fits



AI / Cognitive

- Data
  - Human Generated Data
  - Machine Data
    - Big Data, etc.
    - Sensor Data (IoT)
- Algorithms
  - Machine Learning
    - 'Traditional'
    - Deep Learning
  - Data Science
- Computing power
  - Cloud Computing
  - On Specialised HW

CrowdTruth
The framework for crowdsourcing ground truth data
ECrowd
CAS
CENTER FOR ADVANCED STUDIES

IBM Big Data

Watson IoT

SOFTLAYER
an IBM Company

Power Systems

IBM Center for Advanced Studies

IBM

@zolley

- **AI & Ethics**
- AI & Jobs
- Data Privacy

"The Singularity" —— Skynet & friends

Trust and Values —— Transparency & understanding of the purpose of AI

Bias, Diversity & Inclusion —— Representativeness as appearing in data

IBM
**Center for**
**Advanced Studies**
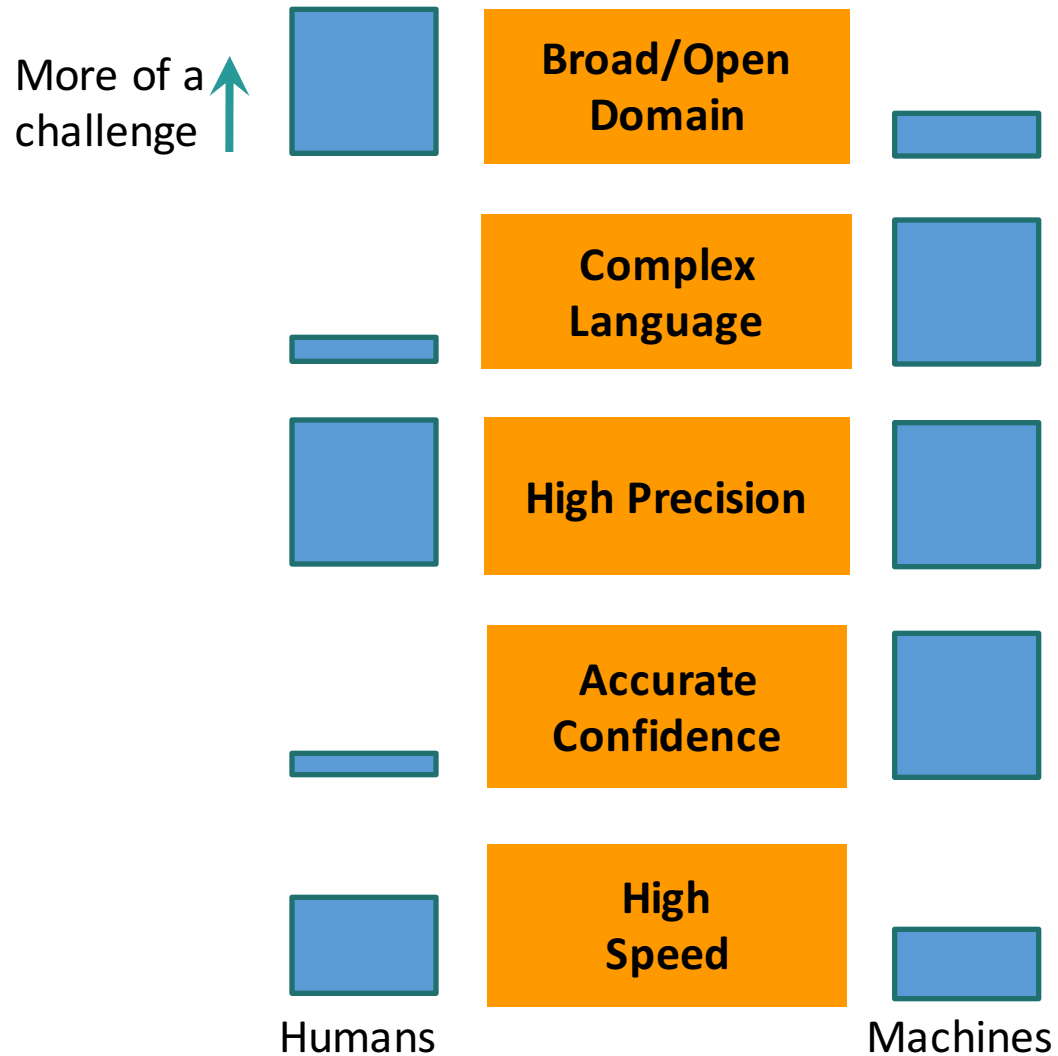
@zolley

IBM

# Towards the Singularity?

"If I can imagine it, it must be likely" – The Availability Heurisitic

IBM

# Watson (2011) Question-Answering on unconstrained domains

More of a challenge ↑

Broad/Open Domain

Complex Language

High Precision

Accurate Confidence

High Speed

Humans

Machines



IBM
**Center for Advanced Studies**

@zolley

IBM

# EXPANDS

EXPANDS human cognition, makes the jobs we do easier, like a *cognitive prosthesis*, especially when dealing with processing massive data, or data that requires human interpretation

LEARNS as you use it – most machine errors are easy for a human to detect, and we can instrument usage of systems to better understand the system and the problem it solves
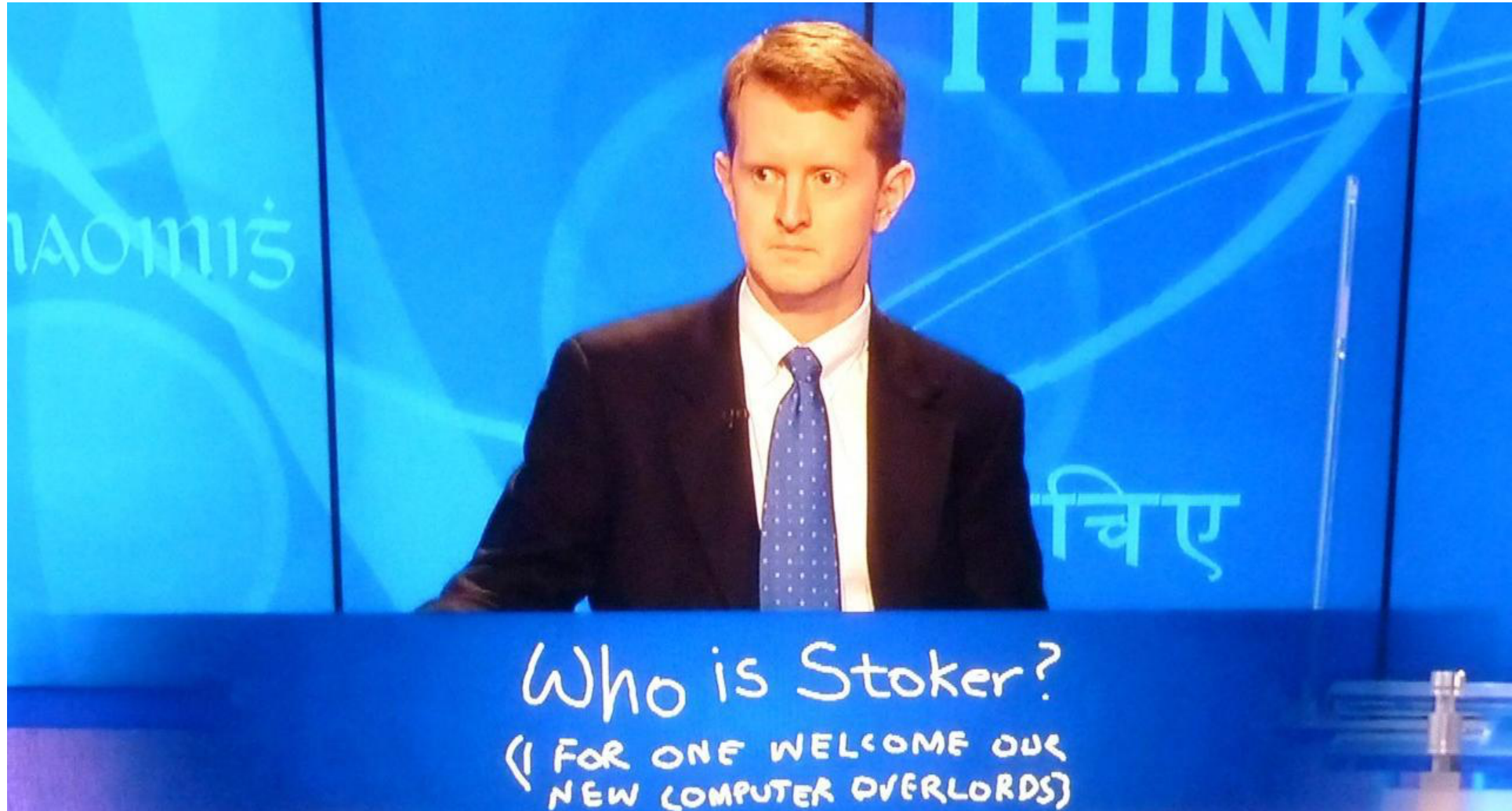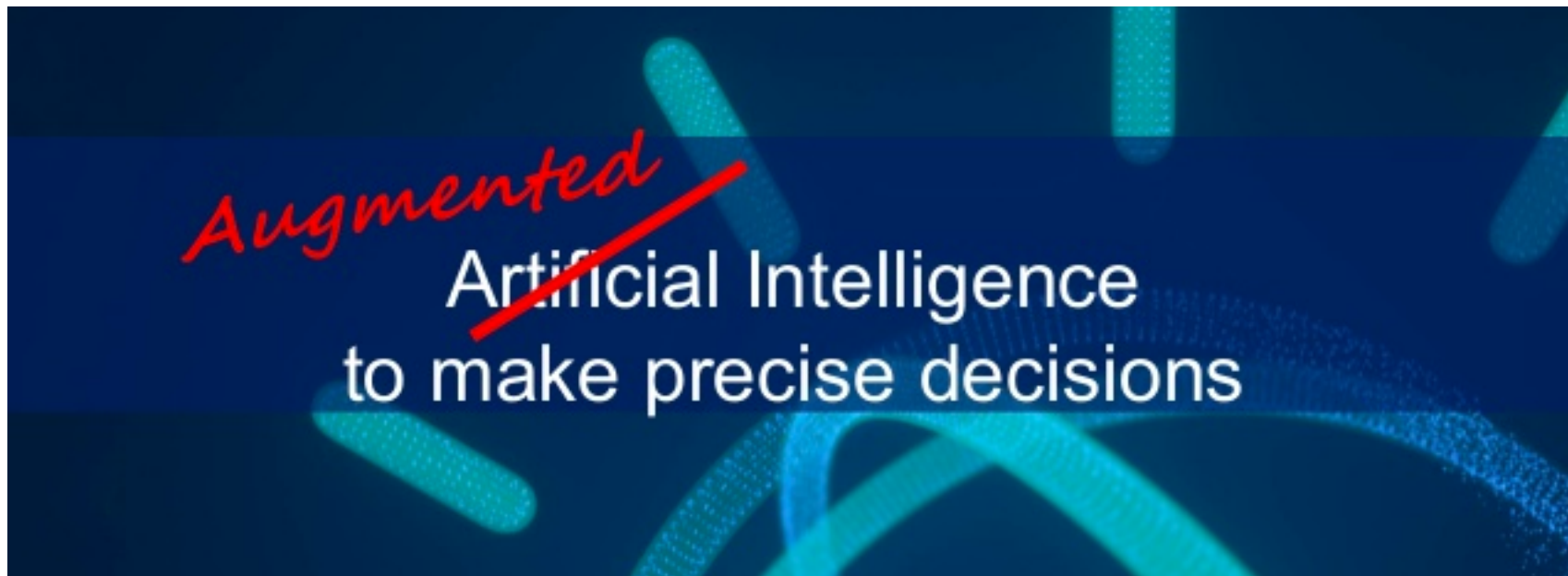
# LEARNS

# INTERACTS

INTERACTS naturally. We need to bring machines closer to their users, we have adapted ourselves enough to them, they should understand natural language, spoken or written, be able to process images and videos. These *simple* human problems are extremely complex for machines, but are hallmarks of a new computing era.
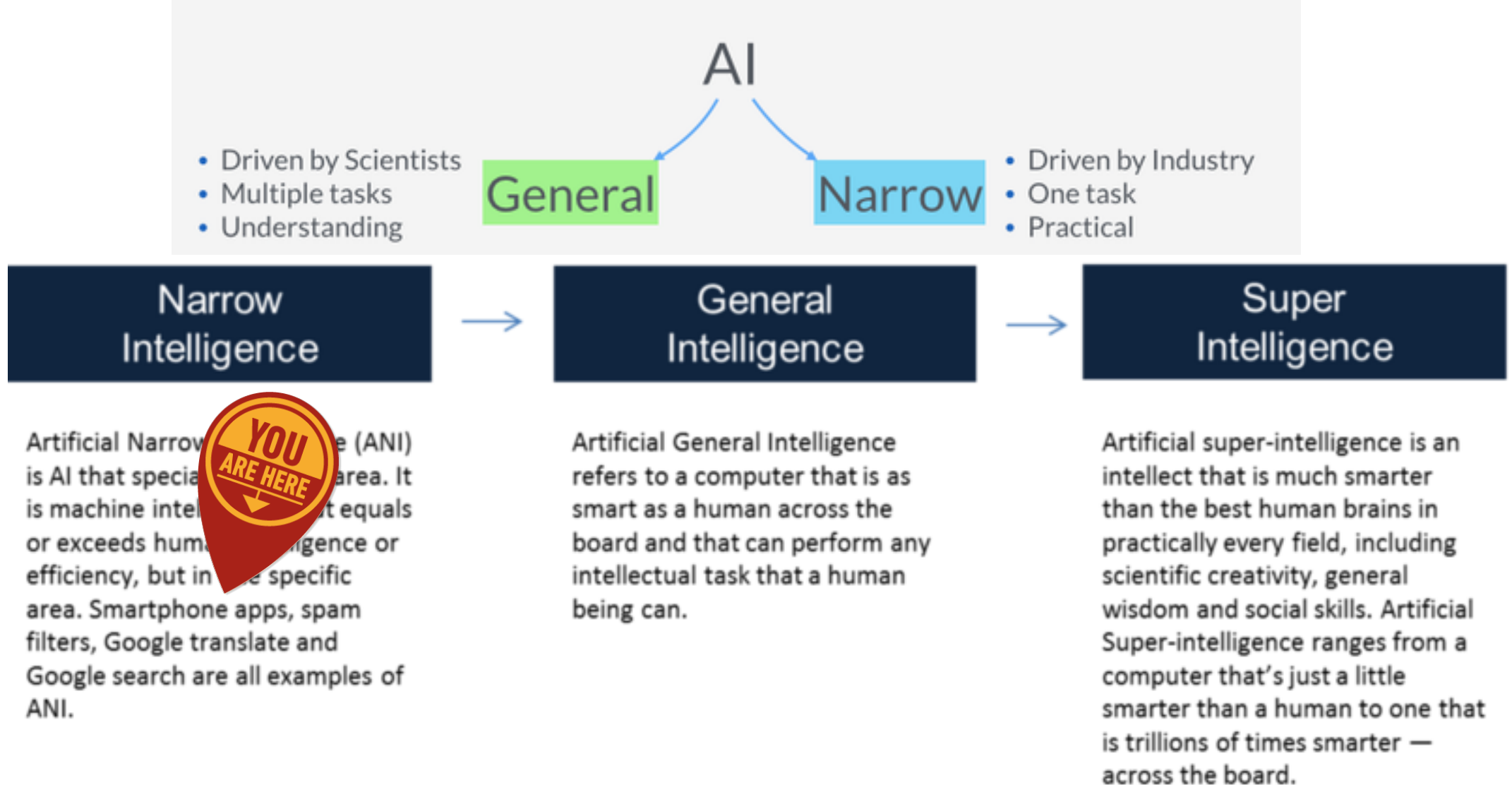
IBM
**Center for
Advanced Studies**

@zolley

IBM

# "Augmented Intelligence"

@zolley

# What kind of AI?



AI

• Driven by Scientists
• Multiple tasks
• Understanding

General    Narrow

• Driven by Industry
• One task
• Practical

**Narrow Intelligence** → **General Intelligence** → **Super Intelligence**

Artificial Narrow [Intelligenc]e (ANI) is AI that specia[lizes in one] area. It is machine intel[ligence tha]t equals or exceeds hum[an intell]igence or efficiency, but in [on]e specific area. Smartphone apps, spam filters, Google translate and Google search are all examples of ANI.

**YOU ARE HERE**

Artificial General Intelligence refers to a computer that is as smart as a human across the board and that can perform any intellectual task that a human being can.

Artificial super-intelligence is an intellect that is much smarter than the best human brains in practically every field, including scientific creativity, general wisdom and social skills. Artificial Super-intelligence ranges from a computer that's just a little smarter than a human to one that is trillions of times smarter — across the board.

**Progression of Artificial Intelligence (AI)**

IBM
Ce
Ad

An illustrated depiction of facial analysis technology similar to that used in the experiment. Illustration: Alamy

# Biases
## and mistakes made by AI systems

**Bonus: a cultural "debate"**

IBM
**Center for**
**Advanced Studies**

IBM

@zolley

AI makes mistakes that are different from human mistakes But we can learn from these!

FATHERLY NICKNAMES

THIS FRENCHMAN WAS "THE FATHER OF BACTERIOLOGY"

HOW TASTY WAS MY LITTLE FRENCHMAN

TECHNOLOGY

BEFORE & AFTER

$200

$400

$600

$200

$400

$600

$1000 $1000 $1000 $1000 $1000 $1000

IBM
Center f
Advanc

IBM

# AI models reflect, and can easily amplify, biases in data

| PR | Data | Provider | Labeller |
| --- | --- | --- | --- |

*"Now it can even…"*

*"We have this awesome new service…"*

*"We do it better than Google…!"*

I'm not a robot
reCAPTCHA

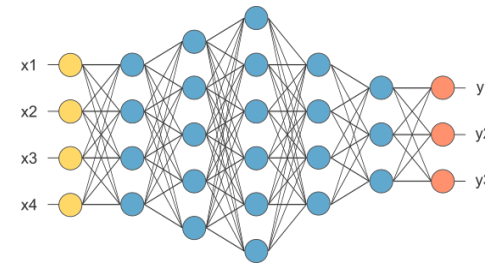Type the text

250

Verify

IBM
Center for
Advanced Studies

START-UP

@zolley

IBM

"A computer program does what you tell it to do, not what you want it to do." (Arthur Bloch, Murphy's Law and Other Reasons Why Things Go Wrong)

An AI system reflects what's in its **input data and algorithms**, not what you want it to

- A learned model reflects what it was trained with

- It will not "think" the way humans do

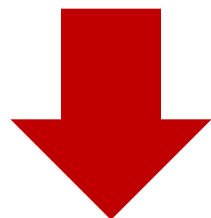- It also has the potential for emergent behaviour (still not the Cylons!)

IBM
**Center for**
**Advanced Studies**

IBM

What if you lie?

What if you show people fake documents?

What if you tell the computer programme to do bad things?

What if you inject bad data, modify training algorithms or trained models?

IBM
**Center for**
**Advanced Studies**

*"All of this has happened before, and it will all happen again."*

\- Peter Pan

\- Battlestar Galactica

IBM
**Center for
Advanced Studies**

IBM

@zolley

# There are even AI techniques building on it!

IBM
**Center for**
**Advanced Studies**

# Adding a tiny bit of noise…



Image source http://karpathy.github.io/assets/break/szegedy.jpeg

@zolley

IBM
**Center for**
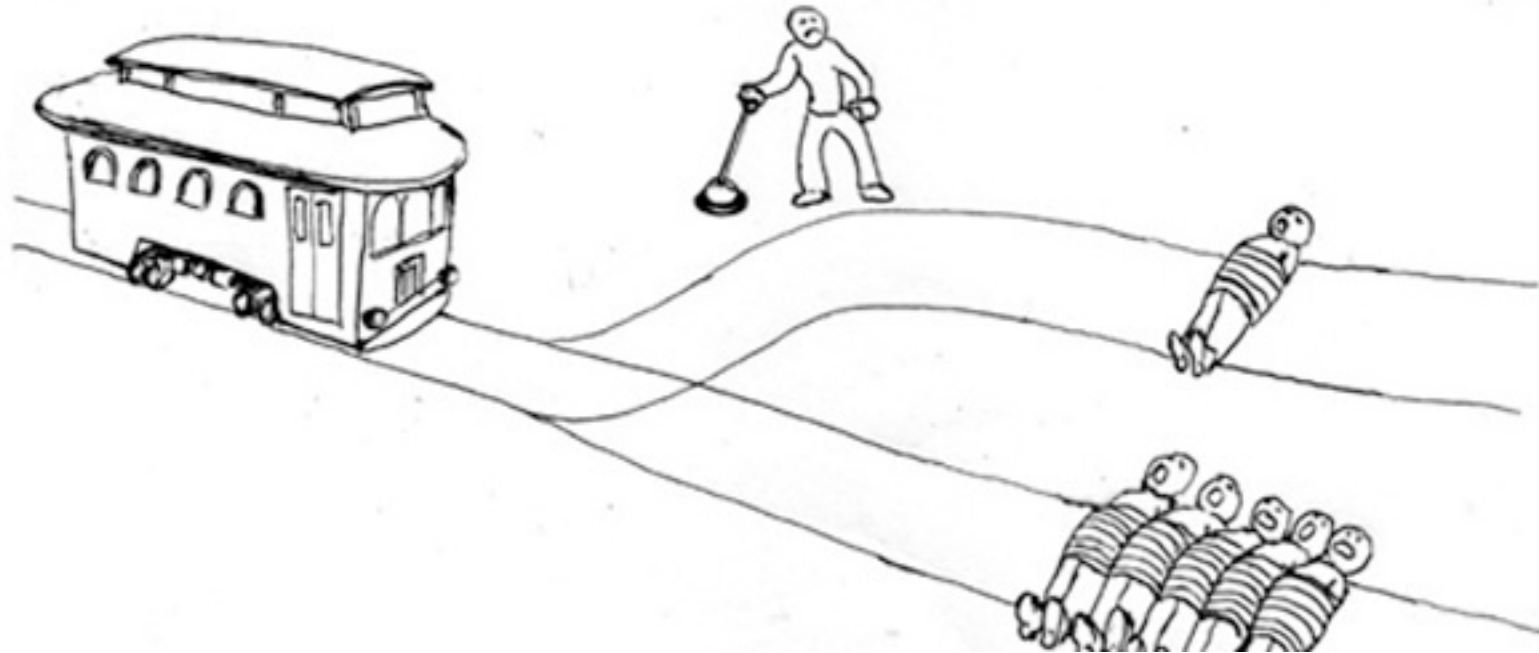**Advanced Studies**

# Values and Trust
## How to inject values, and learn about trust

IBM

# E.g., Teaching self driving cars

How should the car behave in a tricky (ambiguous, unresolvable, sensitive, etc.) situation?

IBM
**Center for Advanced Studies**

@zolley

IBM
**Center for**
**Advanced Studies**

@zolley

IBM

Try many more of these at http://moralmachine.mit.edu/

# What should we do (as a human or machine?)

Try many more of these at http://moralmachine.mit.edu/

Try many more of these at http://moralmachine.mit.edu/

# Human-robot relations: do we trust the humanoid?

@zolley

IBM
**Center for**
**Advanced Studies**

- **Unusu...** display... *bin bes...* the tas... comple... reques...
- **Unusu...** display... *from th...*

you for comfor... minute... area.

- **Unusu...** had sa...

Is this really different from "send this to 10 people or you'll experience 10 years of bad luck"?

Or from trusting in fake news?



Figure 2: Quantitative data analysis: percentages and ratios of participants who did or did not follow the robot's unusual requests (per task)

# We need to handle 'tricky situations' around AI, otherwise…

ROBOT APOCALYPSE!!!

# ACTUALLY, NO

IBM
**Center for
Advanced Studies**

IBM

# We need to handle 'tricky situations' around AI, otherwise…

**Another AI Winter!**

IBM
**Center for
Advanced Studies**

IBM

@zolley

# The Global Initiative for Ethical Considerations
## in the Design of Autonomous Systems
http://standards.ieee.org/develop/indconn/ec/autonomous_systems.html

**Mission:** *"To ensure every technologist is educated, trained, and empowered to prioritize ethical considerations in the design and development of autonomous and intelligent systems."*

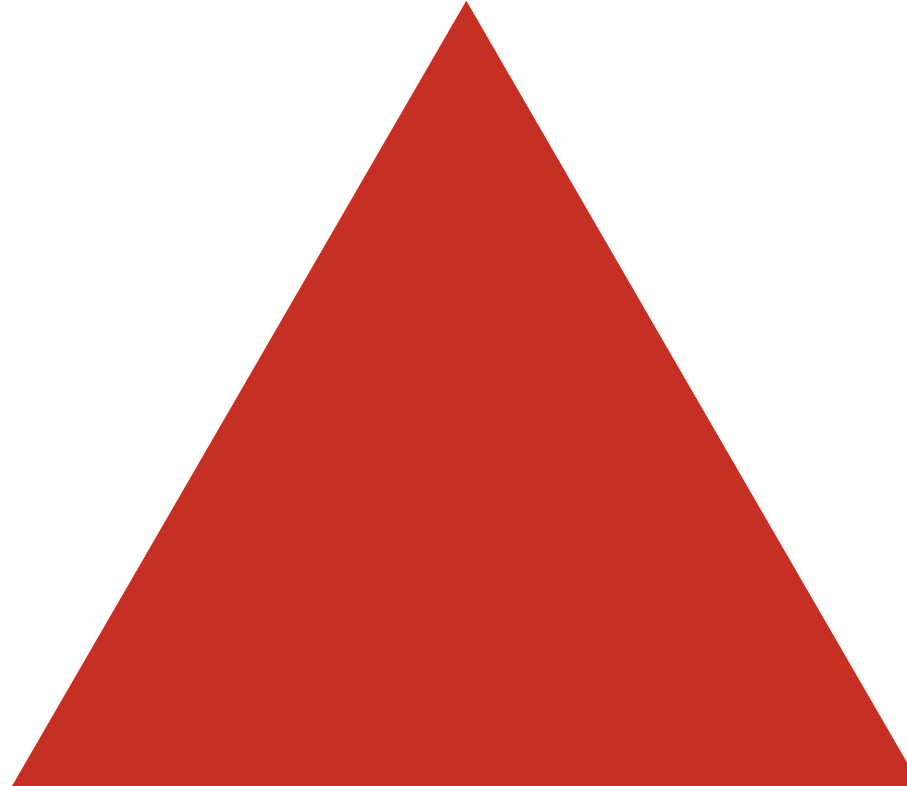Version 2 of document featuring top issues related to autonomous and intelligent technology coming soon.

Committees looking into ethics:

- Law
- General Principles and Guidance
- Safety and Beneficence of AGI and ASI
- Individual/Personal Data*
- Economics of Machine Automation/Humanitarian Issues
- Methodologies to Guide Ethical Research, Design and Manufacturing
- How to Imbue Ethics/Values into Autonomous and Intelligent Systems
- Reframing Autonomous Weapons Systems
- Affective Computing
- Classical Ethics in Information & Communication Technologies
- Effective Policymaking
- Mixed Reality
- Standards
- Ecosystem Mapping
- The "Lexonomy" Committee

IBM
**Center for Advanced Studies**

@zolley

IBM

**Ethics (& Security)** *"Should we do it? How should we do it right?"*

**Business** *"We should do it!"*

**Tech** *"We can do it!"*

IBM
**Center for Advanced Studies**

@zolley

*"A happy peace is my favourite vision."*

\- Necessious

\*

# Thank You!

_____

zoltan.szlavik@nl.ibm.com

Twitter: @zolley

*\*Necessious (who named itself) is an AI model trained in an hour using "inspirational quotes" as input data*