

AI, People, and the Open World

Eric Horvitz

Technical Fellow and Director
Microsoft Research

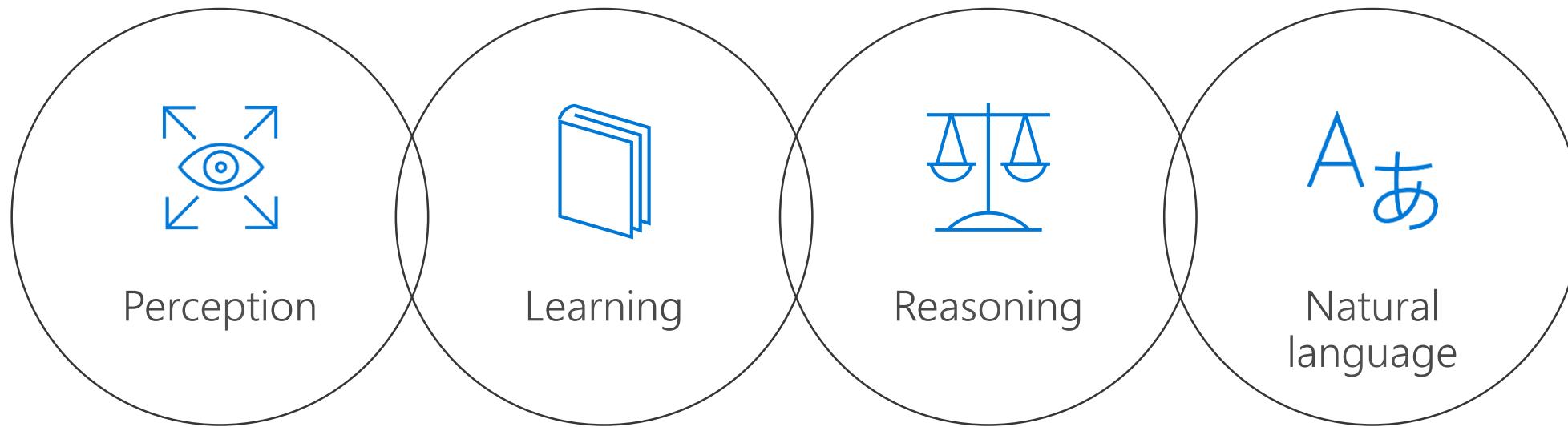


August 16, 2017

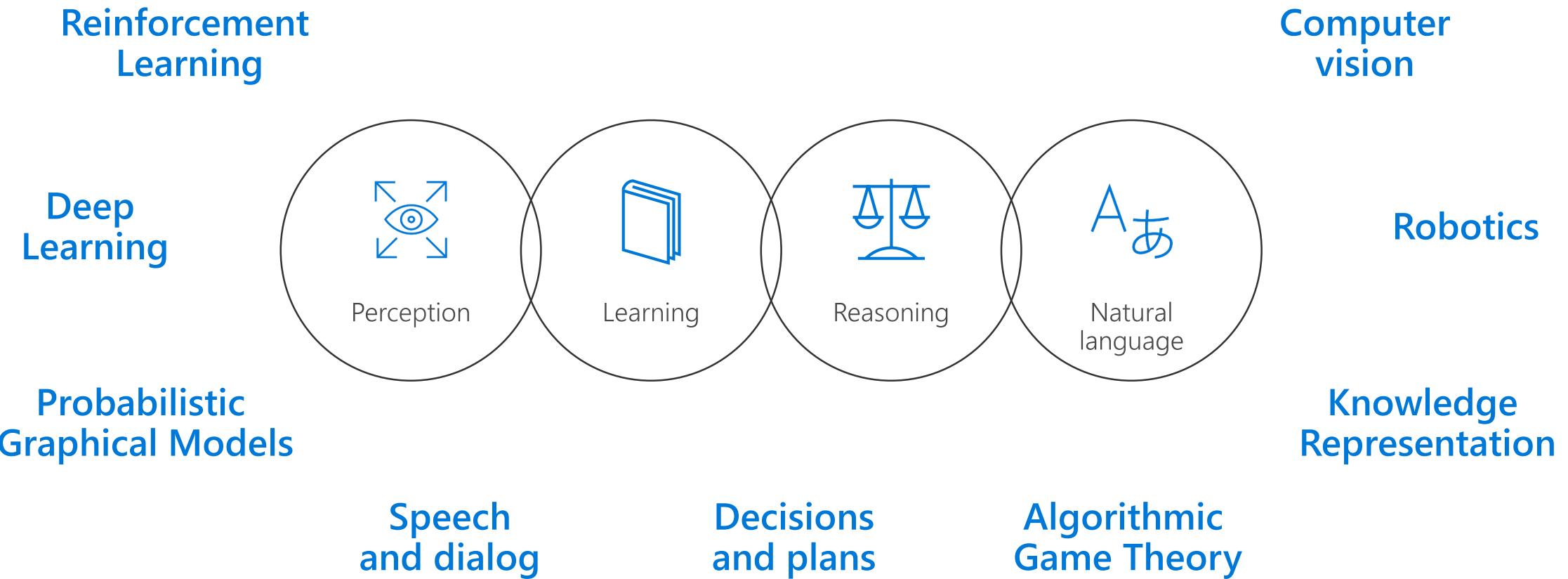
Artificial Intelligence

“...to find how to make machines...solve kinds of problems now reserved for humans...” (1955)

Artificial Intelligence



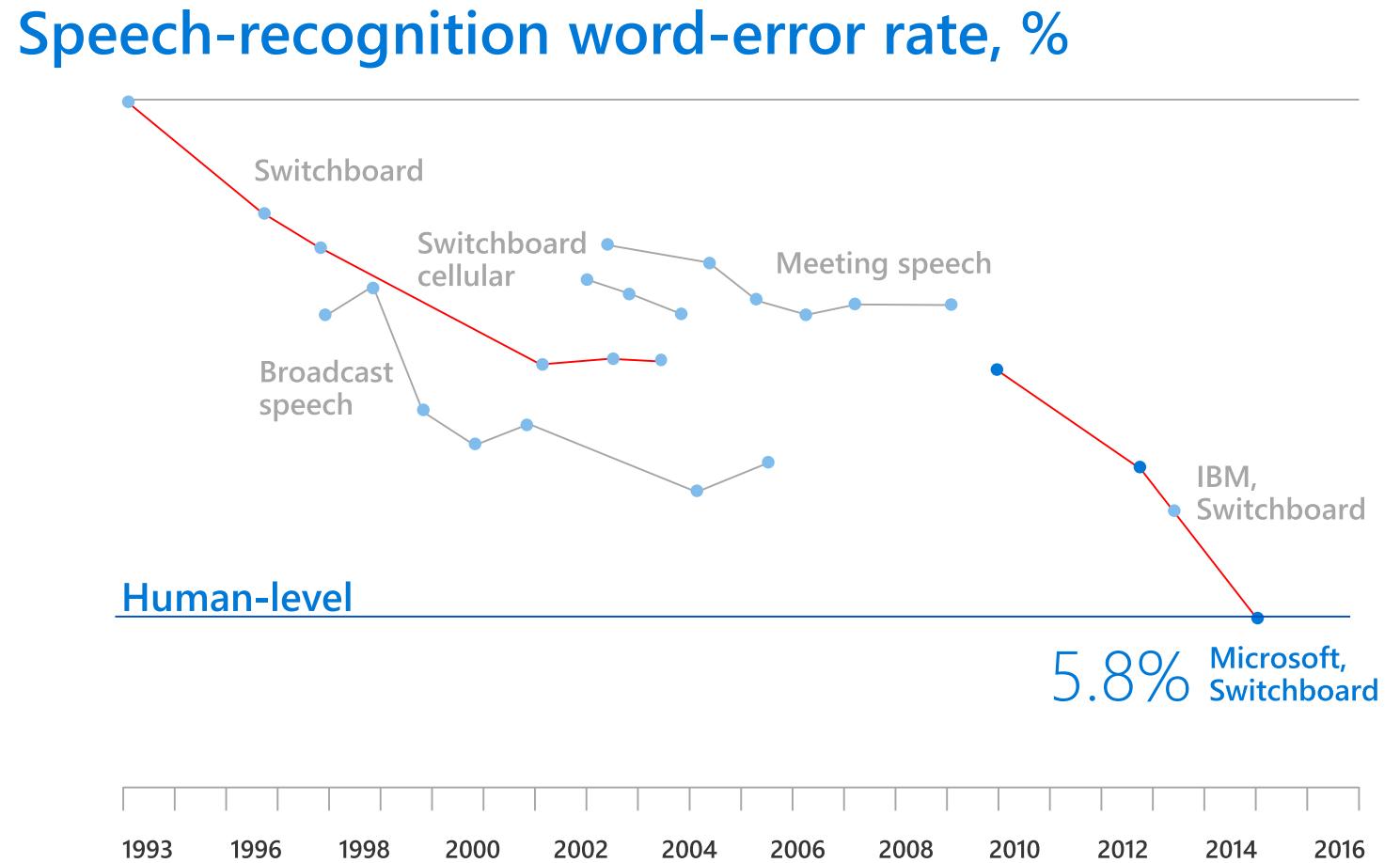
Artificial Intelligence



Multiple subdisciplines and research communities

Speech recognition

5.8% human error rate
Switchboard challenge



Sources: Microsoft: research papers

Vision

152 layers

RESNET



Reading comprehension

Answer questions about information on Wikipedia

SQuAD challenge

Stanford question & answer challenge

The image shows a tablet displaying the official SQuAD website at rajpurkar.github.io/SQuAD-explorer. The website has a red header with the text "SQuAD" and "The Stanford Question Answering Dataset". Below the header, there's a "What is SQuAD?" section with a detailed description of the dataset, followed by two buttons: "Explore SQuAD and model predictions" and "Read the paper (Rajpurkar et al. '16)". There's also a "Getting Started" section with download links for "Training Set v1.1 (30 MB)" and "Dev Set v1.1 (5 MB)". To the right of these sections is a "Leaderboard" table showing the top models with their EM and F1 scores:

Rank	Model	EM	F1
1	r-net (ensemble) Microsoft Research Asia http://aka.ms/met	76.922	84.006
2	Interactive AoA Reader (ensemble) Joint Laboratory of HIT and iFLYTEK Research	76.492	83.745
3	MEMEN (ensemble) Eigen Technology & Zhejiang University	75.370	82.658
4	ReasoNet (ensemble) MSR Redmond https://arxiv.org/abs/1609.05284	75.034	82.552
5	r-net (single model) Microsoft Research Asia http://aka.ms/met	74.614	82.458
6	Mnemonic Reader (ensemble)	73.754	81.863

New capabilities



skype™ Translate

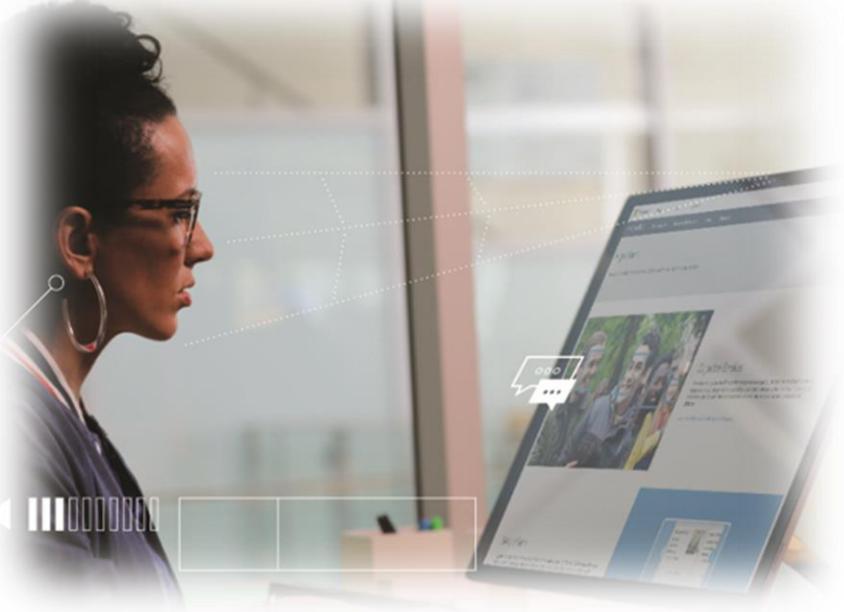
I was wondering what you are going to do later?
Me preguntaba lo que vas a hacer después?

Going to the pub, do you want to join us? I think you
met some of the team at the party in April.
*Al pub, ¿quieres unirte a nosotros? Creo que conociste
a algunos miembros del equipo en la fiesta en abril.*

Can I join you guys after my meeting?
¿Puedo unirme a ustedes después de mi reunión?

Type a message in English here

Lab environment



Open world



Competence

AI & People

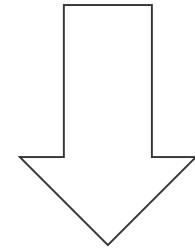
Ethical, legal, societal influences

Lab environment



Open world

Frame problem



Qualification problem

All preconditions?



Ramification problem

All effects of action?

Lab environment



Open world

Framing

Data

Fidelity

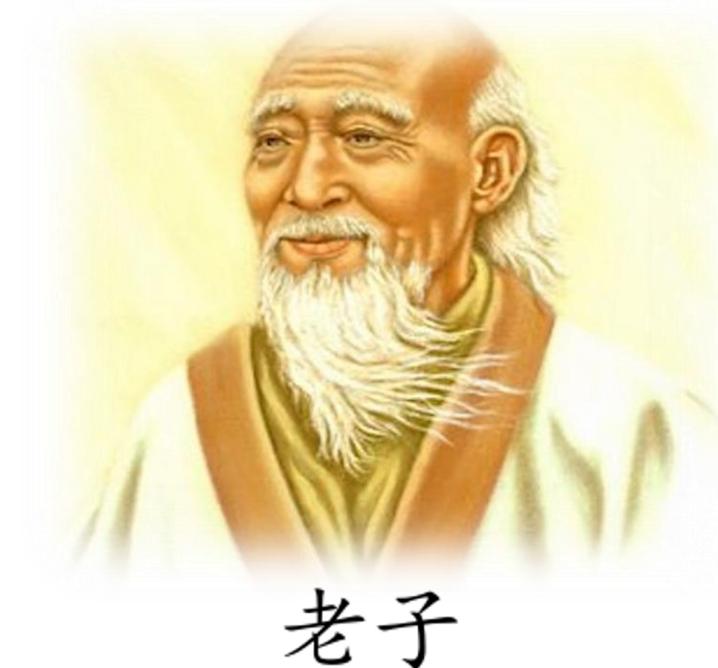
Representation

Inference



Learning about competency

Wisdom of the ages



老子

Knowing that you do not know is the best.

Not knowing that you do not know is an illness.

- Laozi, 500-600 BCE

Learn about abilities & failures

Successes & failures



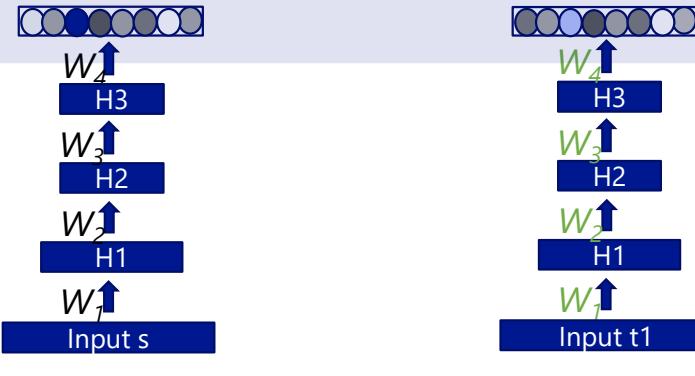
Confidence

$$p(\text{ fail} \mid E, t)$$

Deep learning about deep learning performance

Quality score [0,1] $s = \frac{e^{W \cdot f}}{1 + e^{W \cdot f}}$

- 1. LM score
- 2. Length
- 3. LM score per word
- 4. Vision tag covered
- 5. DSSM score



Caption:
a man holding a tennis racquet on a tennis court

Grappling with Open-World Complexity

Reliable predictions of performance: *Known unknowns*



Grappling with Open-World Complexity

Reliable predictions of performance: *Known unknowns*







Grappling with Open-World Complexity

Reliable predictions of performance: *Known unknowns*

Challenge of *unknown unknowns*



Unknown
unknowns

Directions

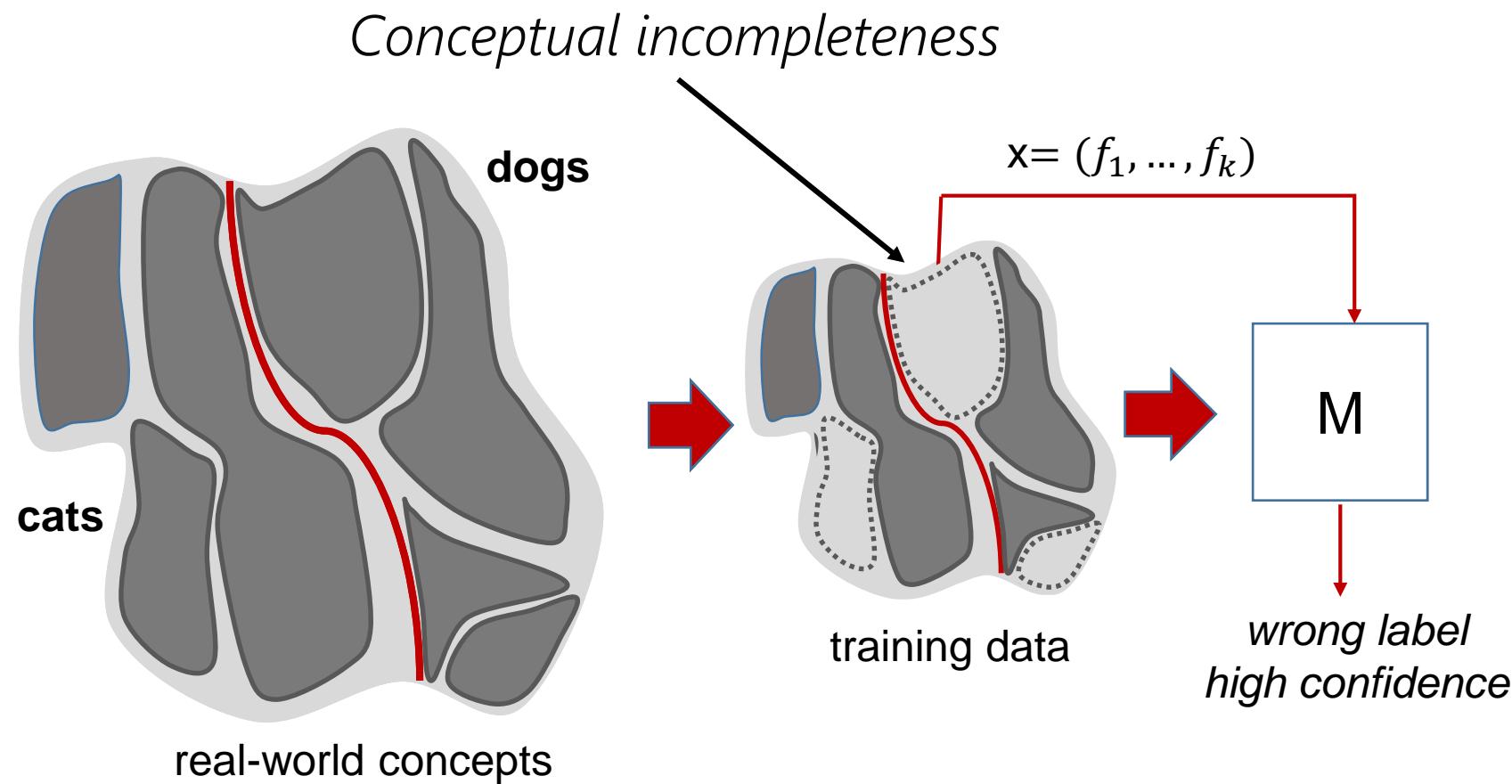
Expanded real-world testing

Algorithmic portfolios

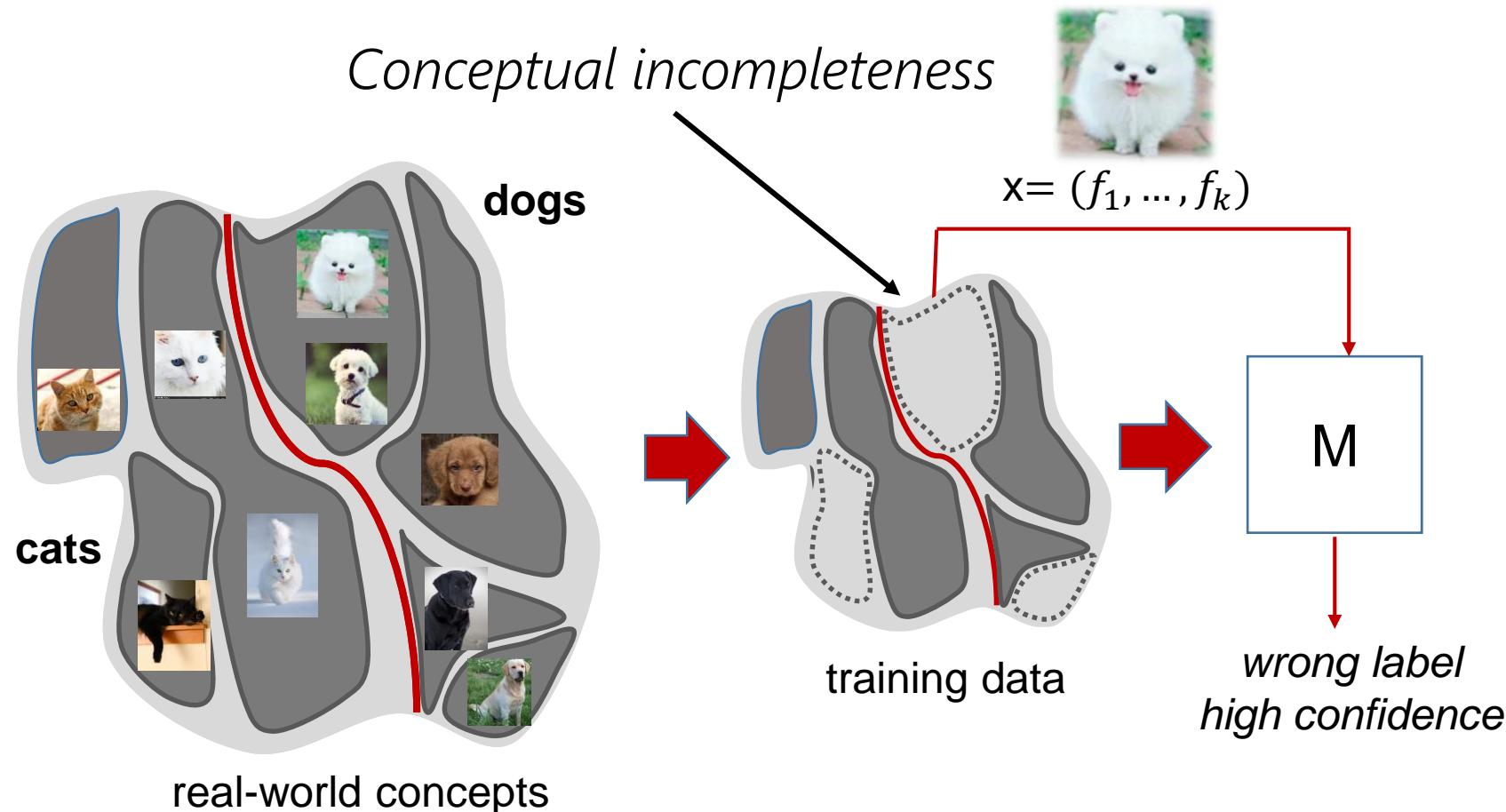
Failsafe designs

People + machines

Identifying classifier blindspots

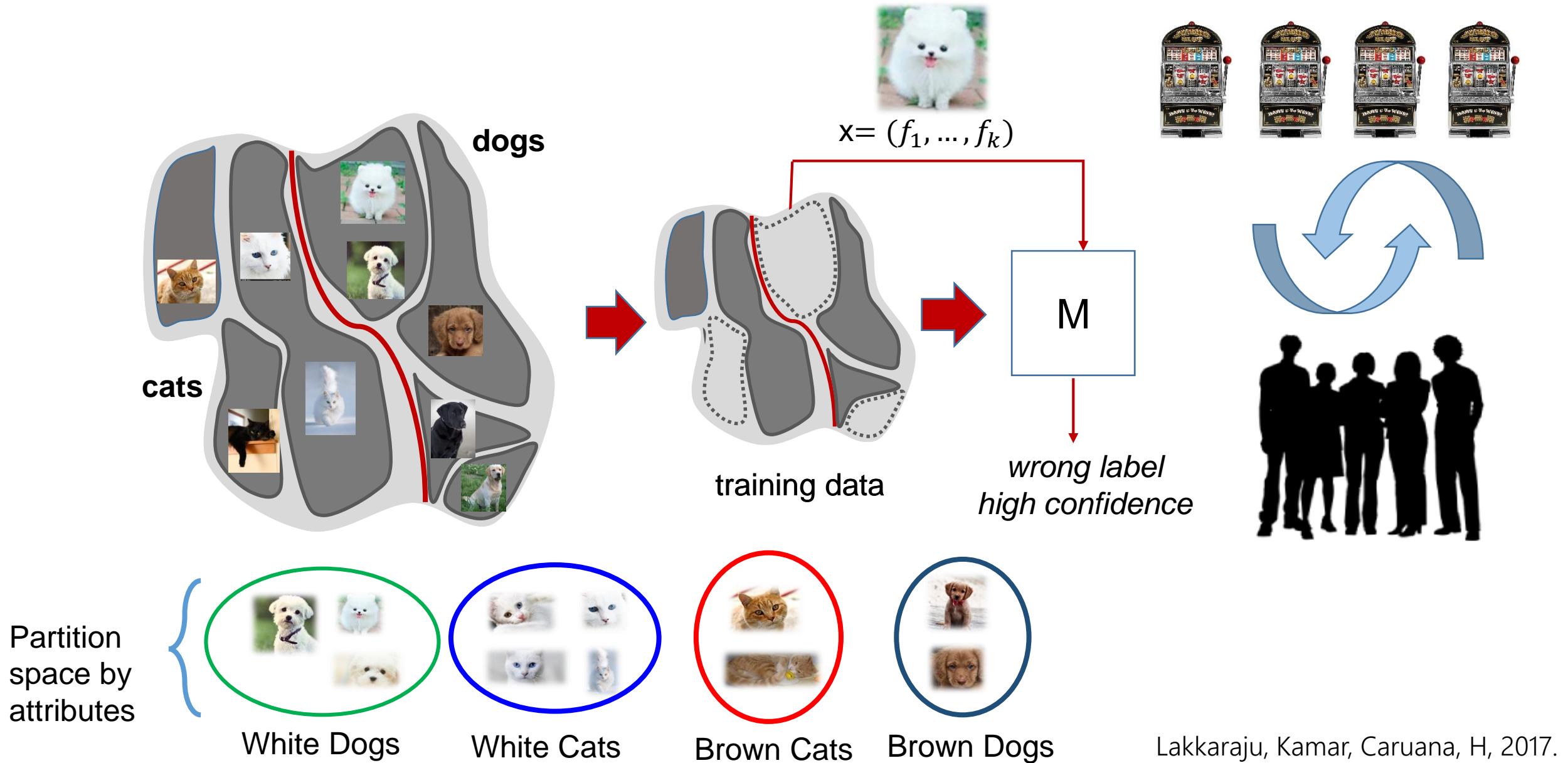


Identifying classifier blindspots



How to define & search regions of data space?
How to trade exploration and exploitation?

Identifying classifier blindspots



Data
scarcity

Directions

Transfer learning

Learn from rich simulations

Learn generative models

Transfer learning opportunity

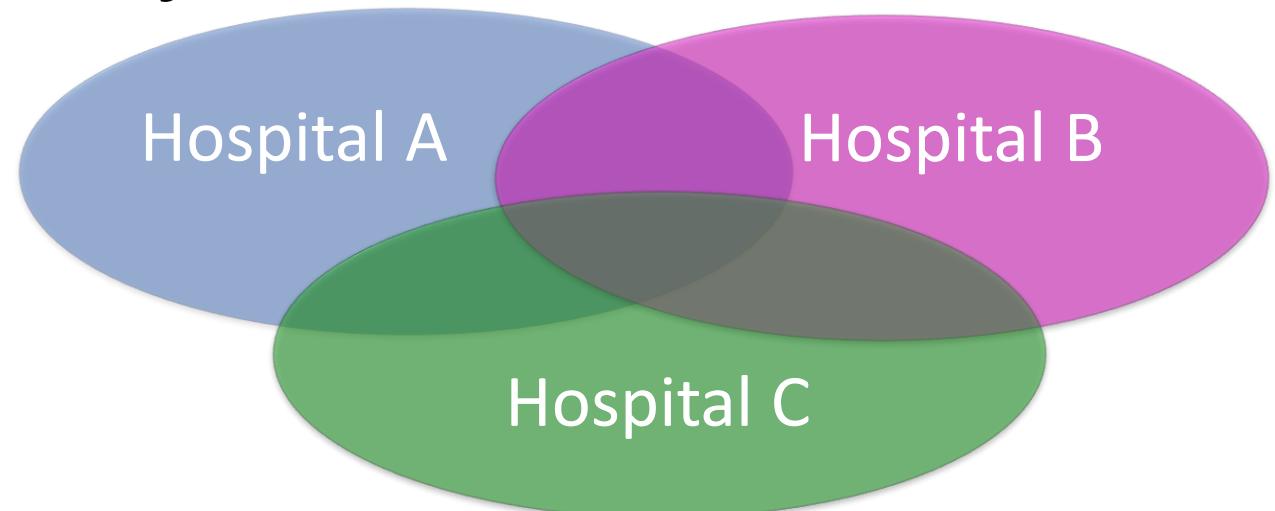
Site-specific data

Observations, definitions

Patients, prevalencies

Covariate dependencies

Predict risk of infection



- A: Community hosp: 10k pts/yr
- B: Acute care & teaching: 15k/yr
- C: Major teaching & research: 40k/yr

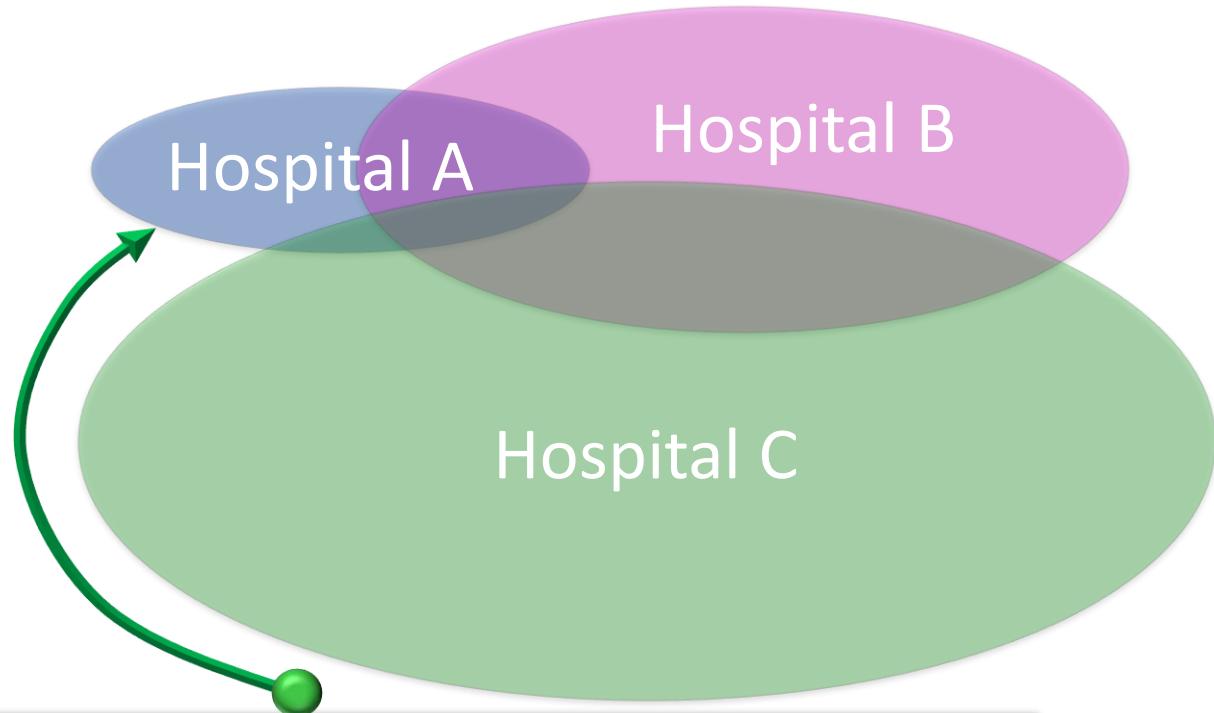
Transfer learning opportunity

Site-specific data

Observations, definitions

Patients, prevalencies

Covariate dependencies



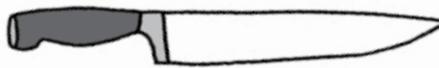
Research and applications

A study in transfer learning: leveraging data from multiple hospitals to enhance hospital-specific predictions

Jenna Wiens,¹ John Guttag,¹ Eric Horvitz²

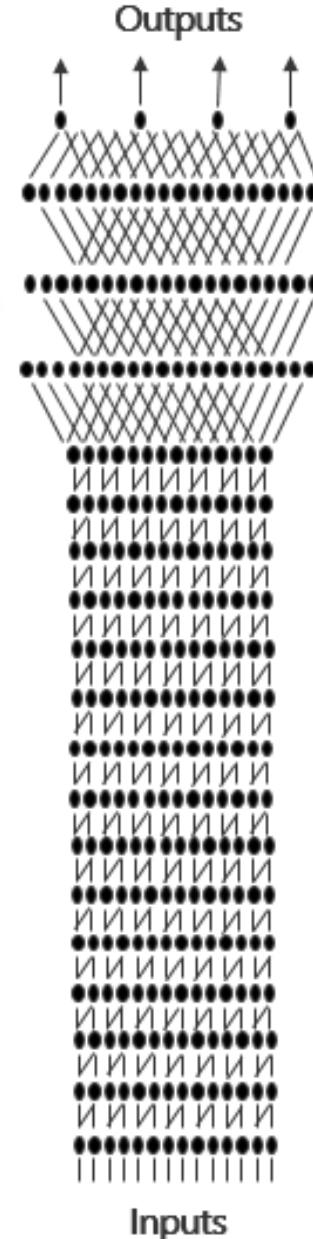
Embedded deep transfer learning

Less data with better features



ImageNet 1000, 1M photos

Cut off top layer



Embedded deep transfer learning

Less data with better features

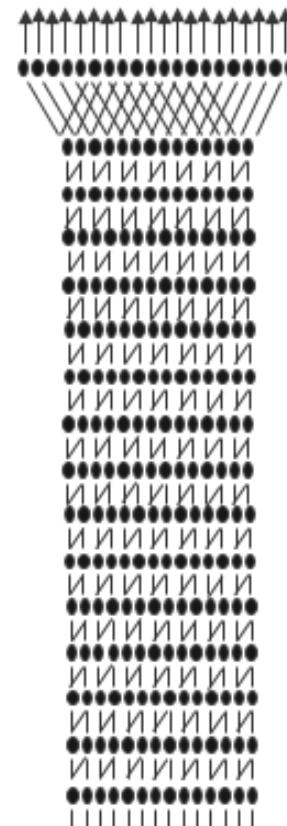
ImageNet 1000, 1M photos

Cut off top layer



0.93 0.55 -0.91 0.04 -0.42 -0.91 0.64 0.81 -0.63 0.99

Output Feature Vector



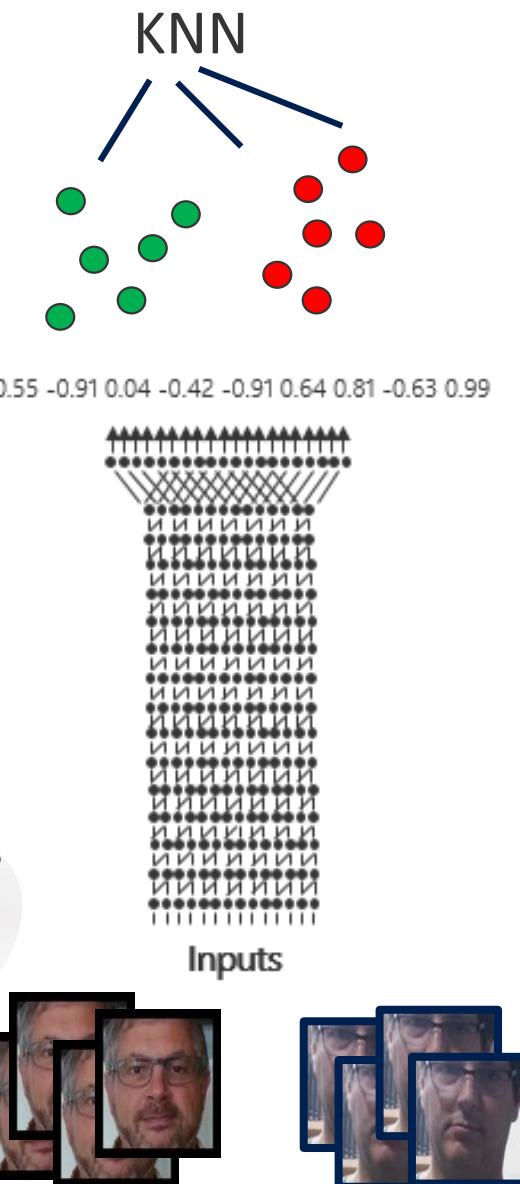
Inputs

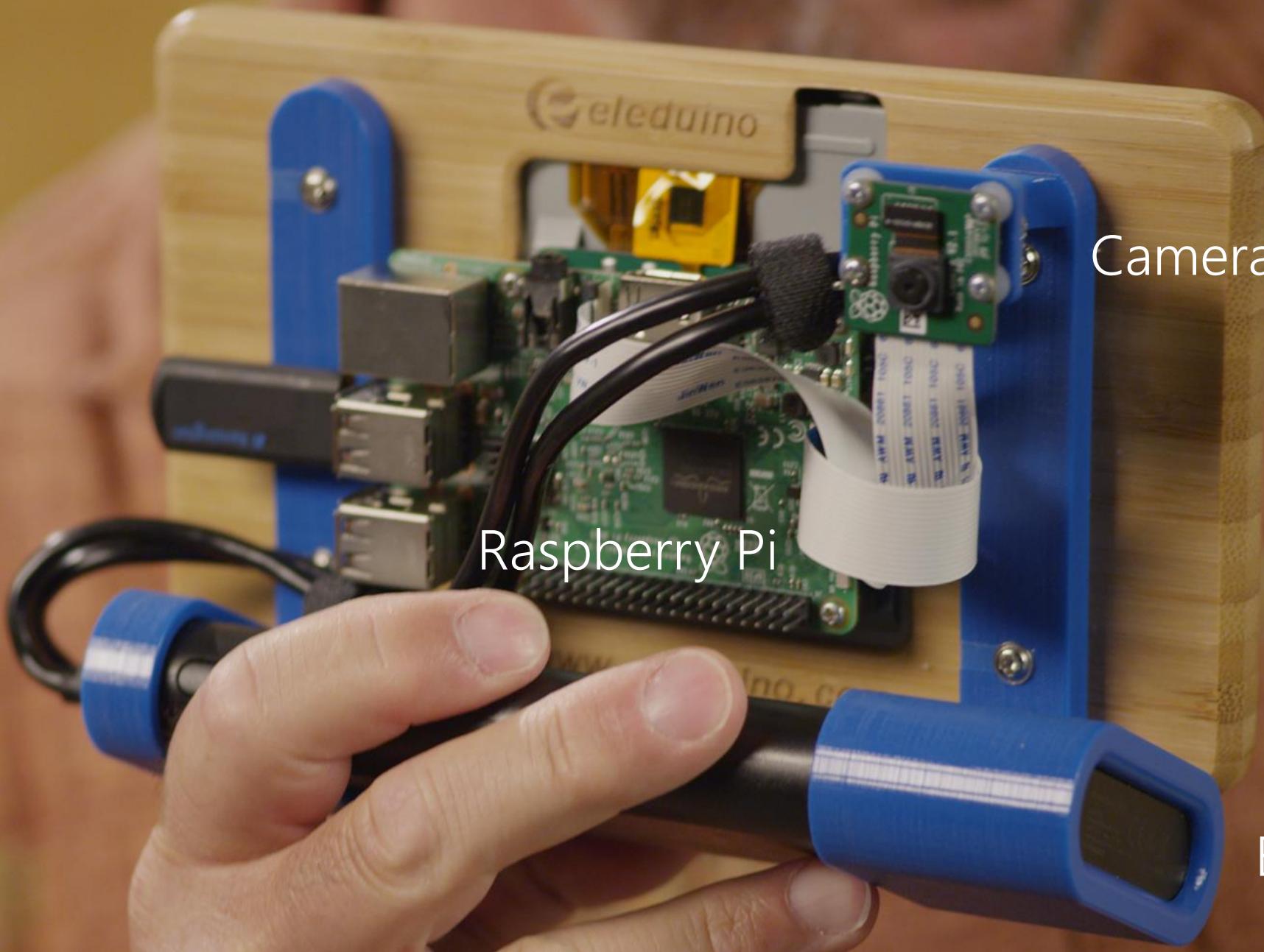
Embedded deep transfer learning

Less data with better features

ImageNet 1000, 1M photos

Cut off top layer





Camera

Raspberry Pi

Battery

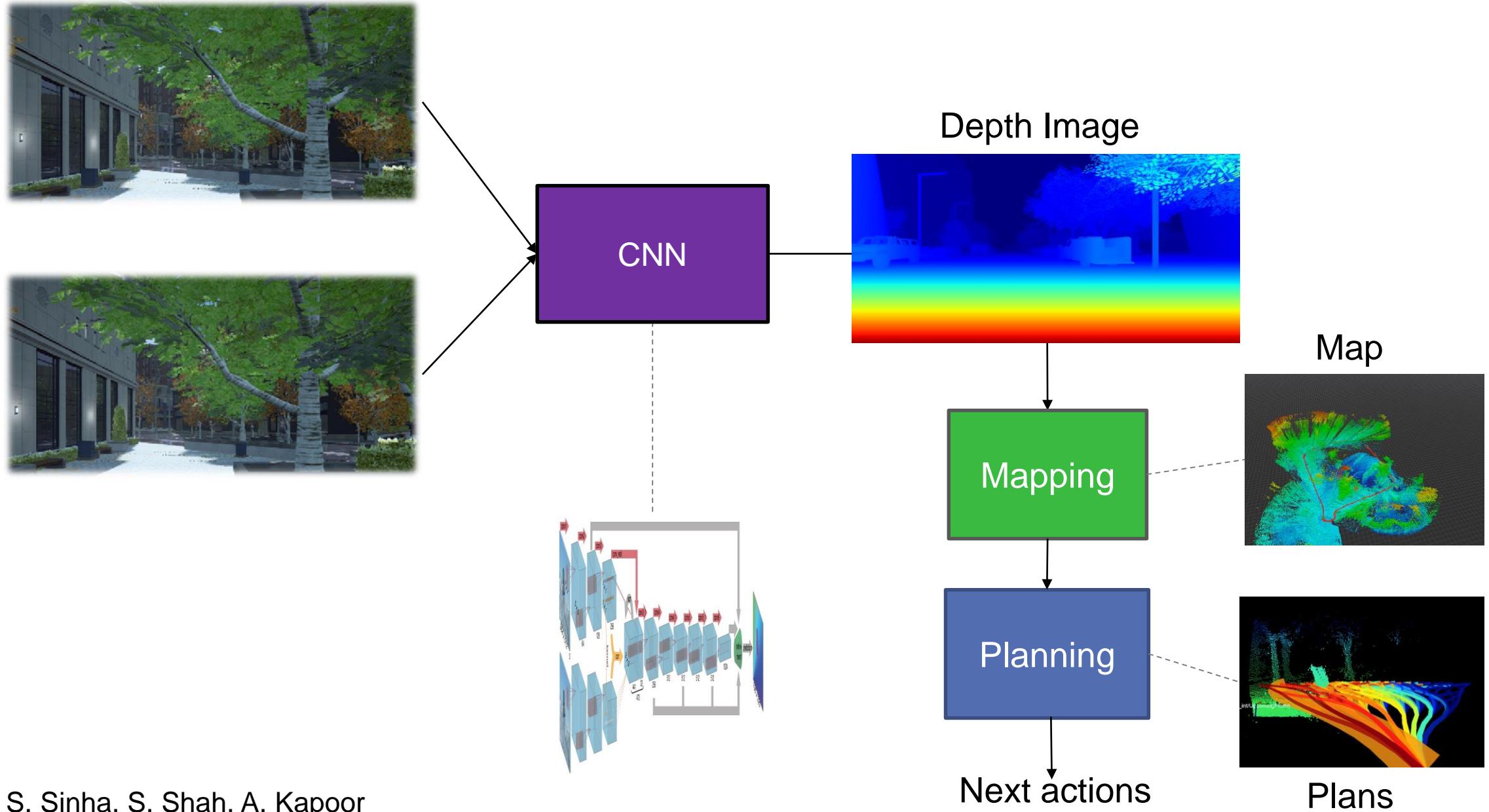
Simulated Environments

Trillions of sessions in complex scenarios

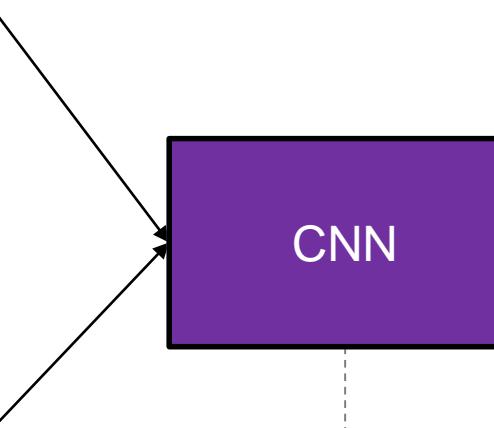
Learn & evaluate core competencies

Learn to optimize action plans

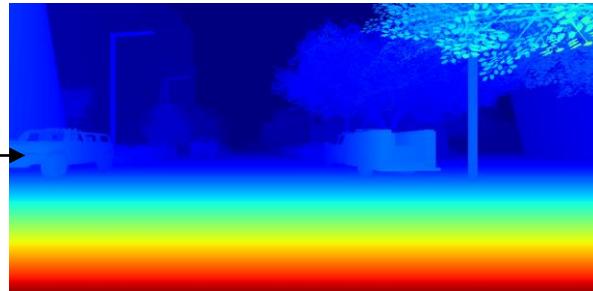
Leveraging rich simulations



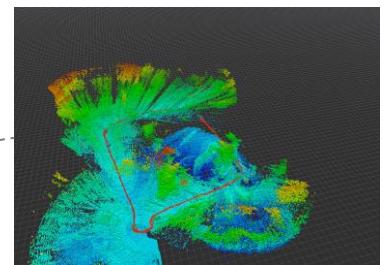
Leveraging rich simulations



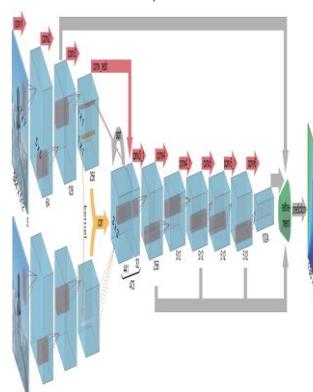
Depth Image



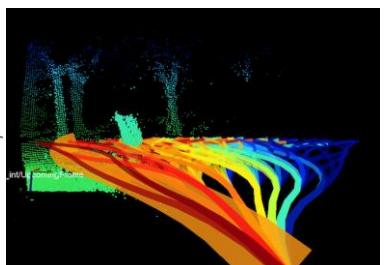
Map



Mapping



Planning



Next actions

Plans

Learn Generative Models

Learn expressive generative models
Generalize from minimal training sets
Harness physics

Learning generative models

Multilevel variational autoencoder

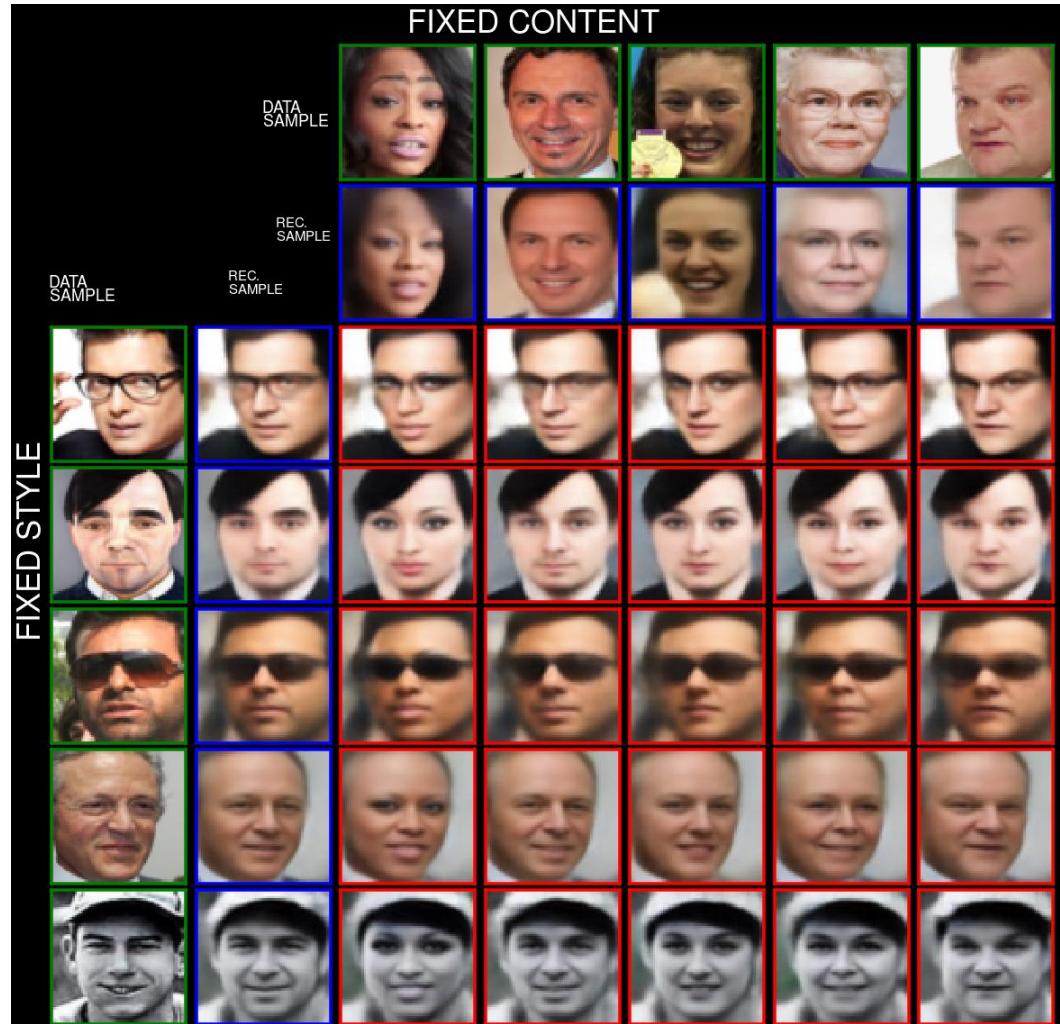
Learn disentangled representations

Groups of observations → latent models

Vary style

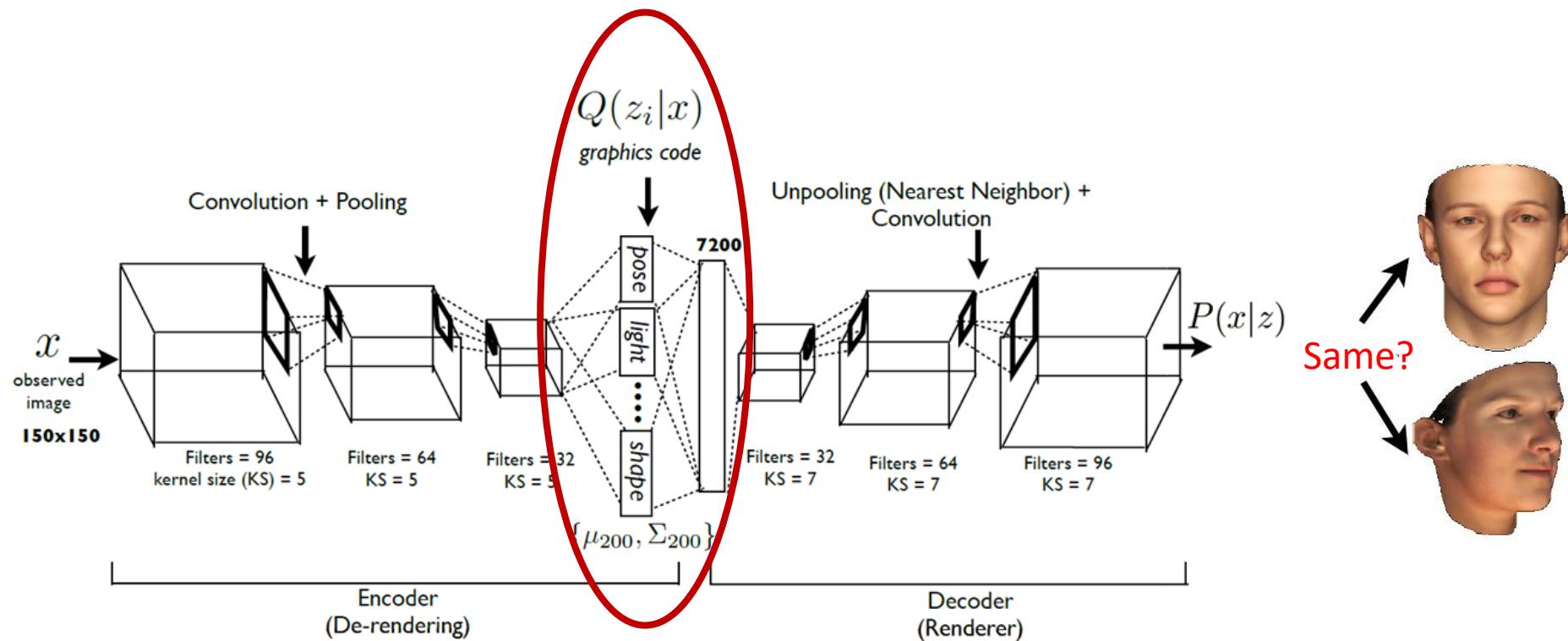


Vary ID

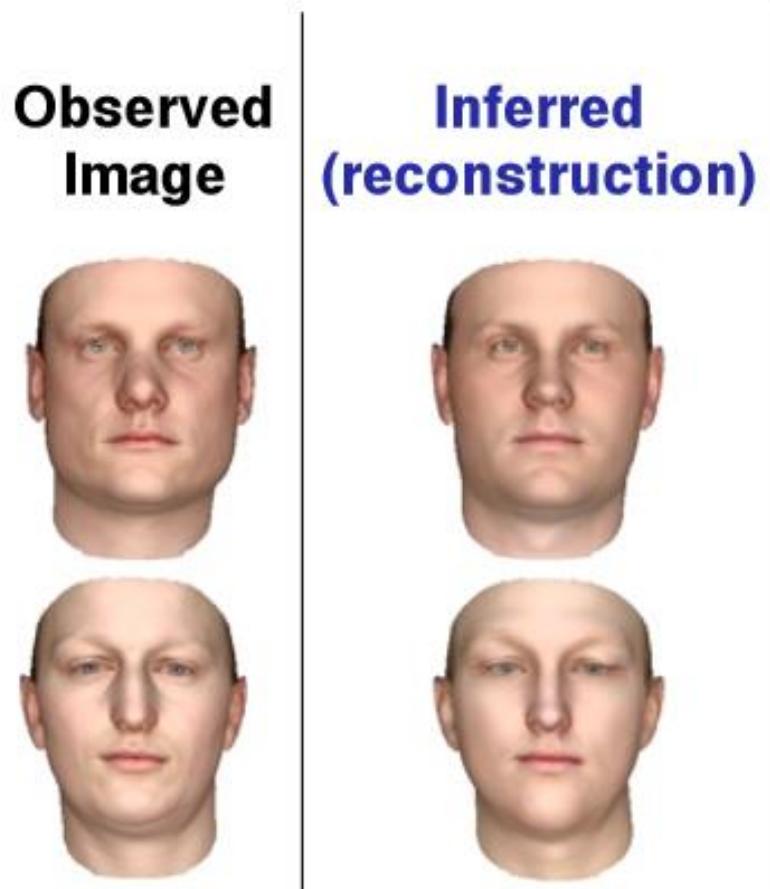


Smooth control over learned latent space

Inject physics to disentangle & generalize



Inject physics to disentangle & generalize



Inject physics to disentangle & generalize



Illumination

Nod

Shake

Adversarial Attacks & AI

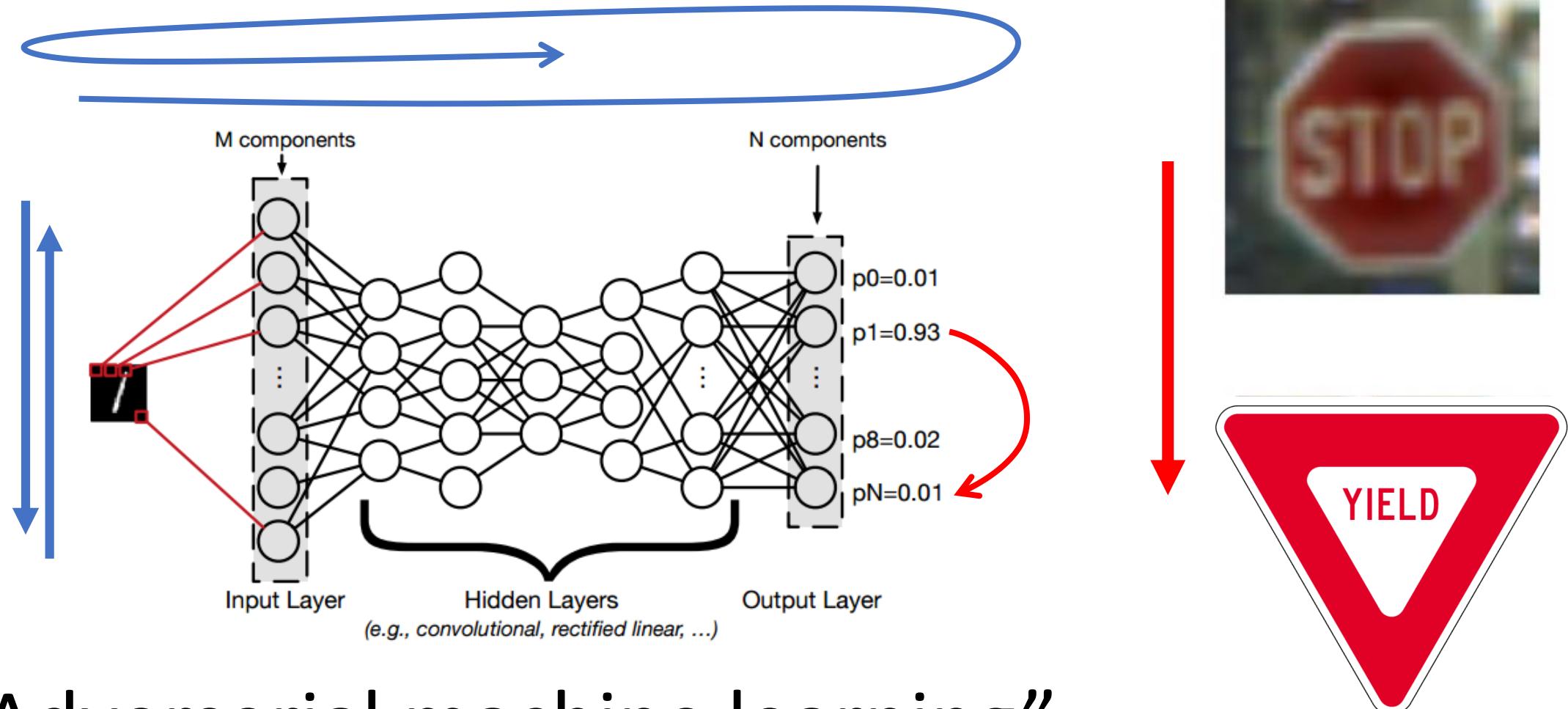
Concerns

AI attack surfaces

Adversarial machine learning

Self-modification

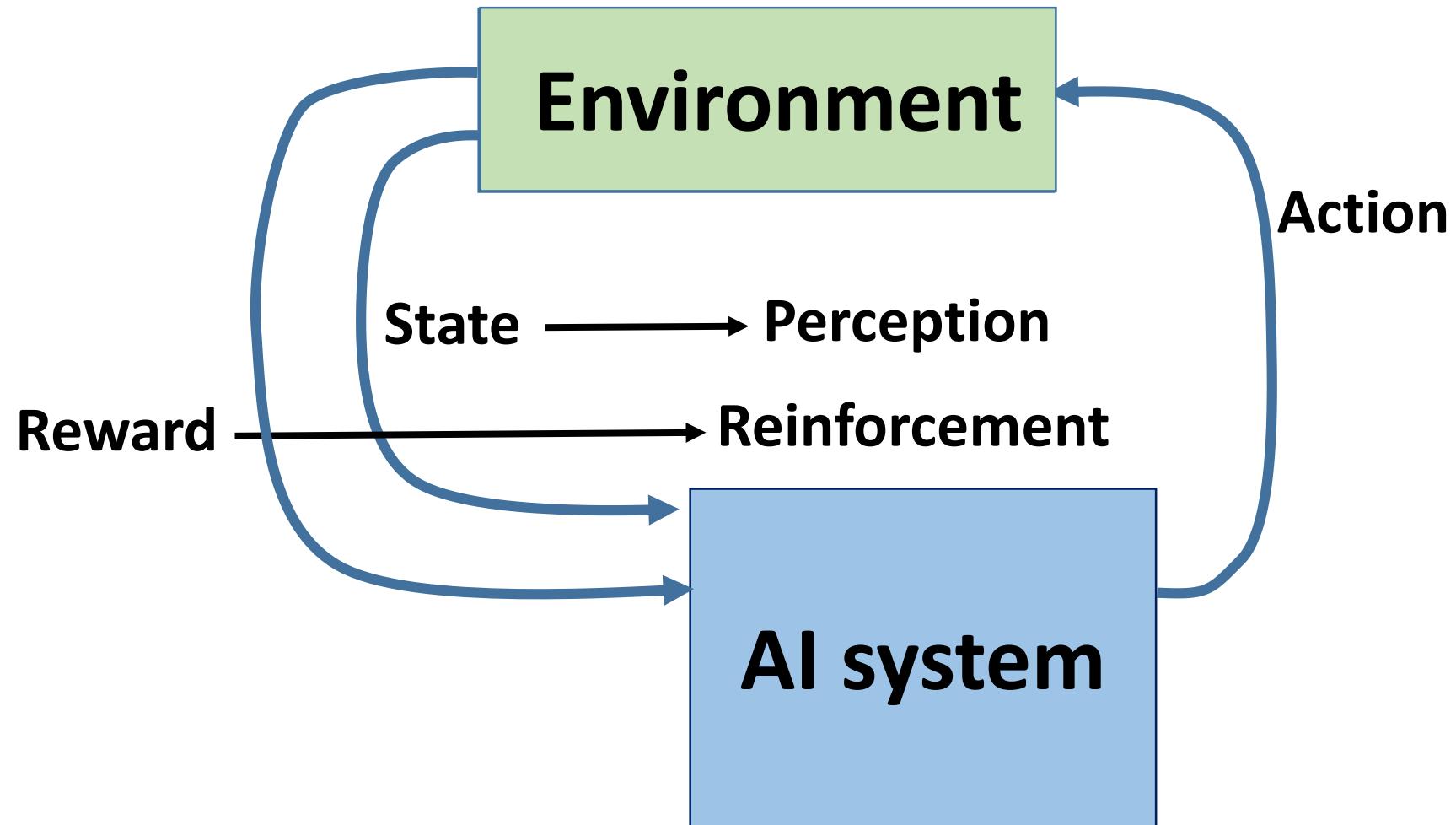
Attacks on AI Systems



“Adversarial machine learning”

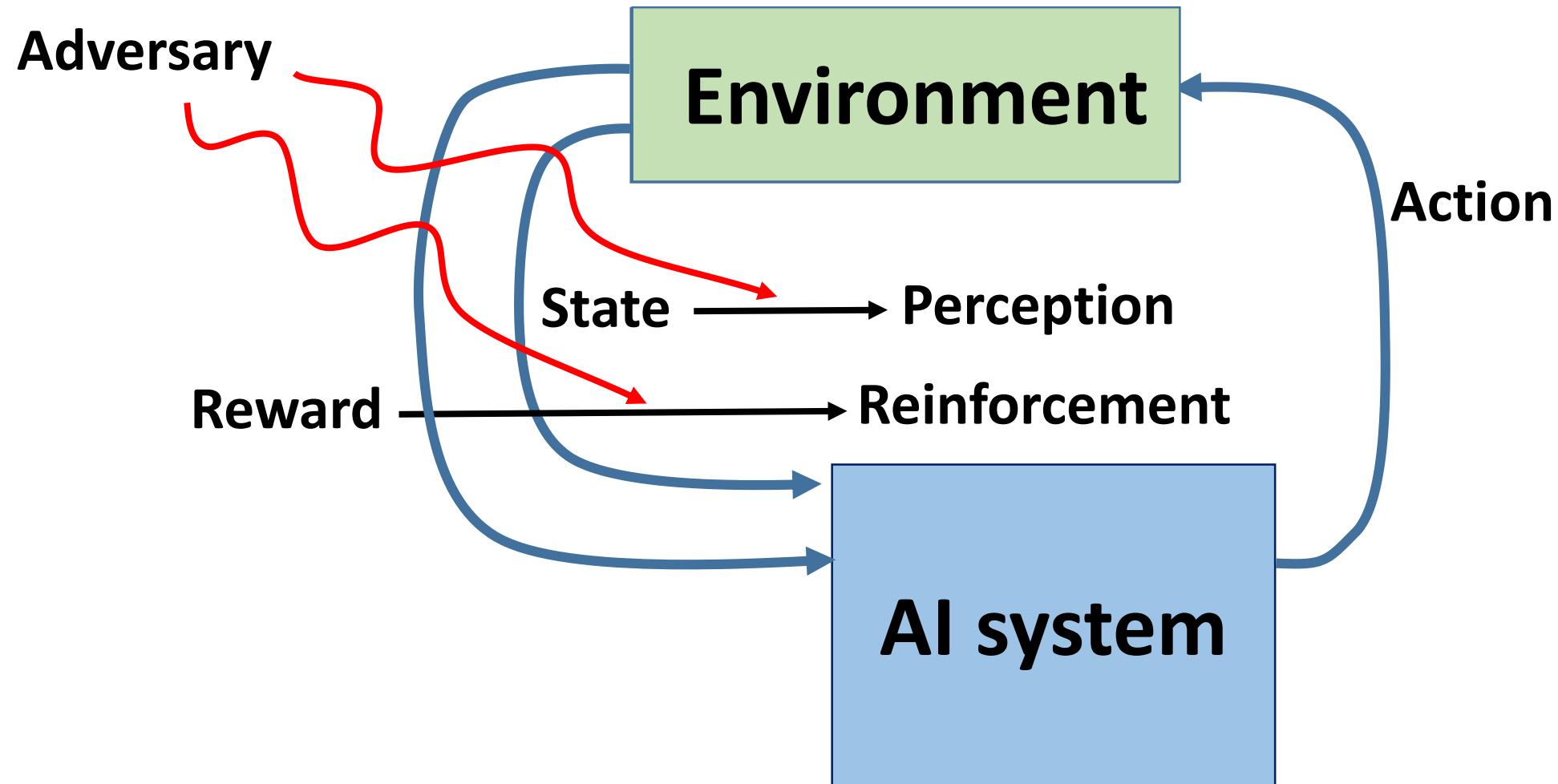
Goodfellow, et al.
Papernot, et al.

Adversarial Attacks & Self-Modification



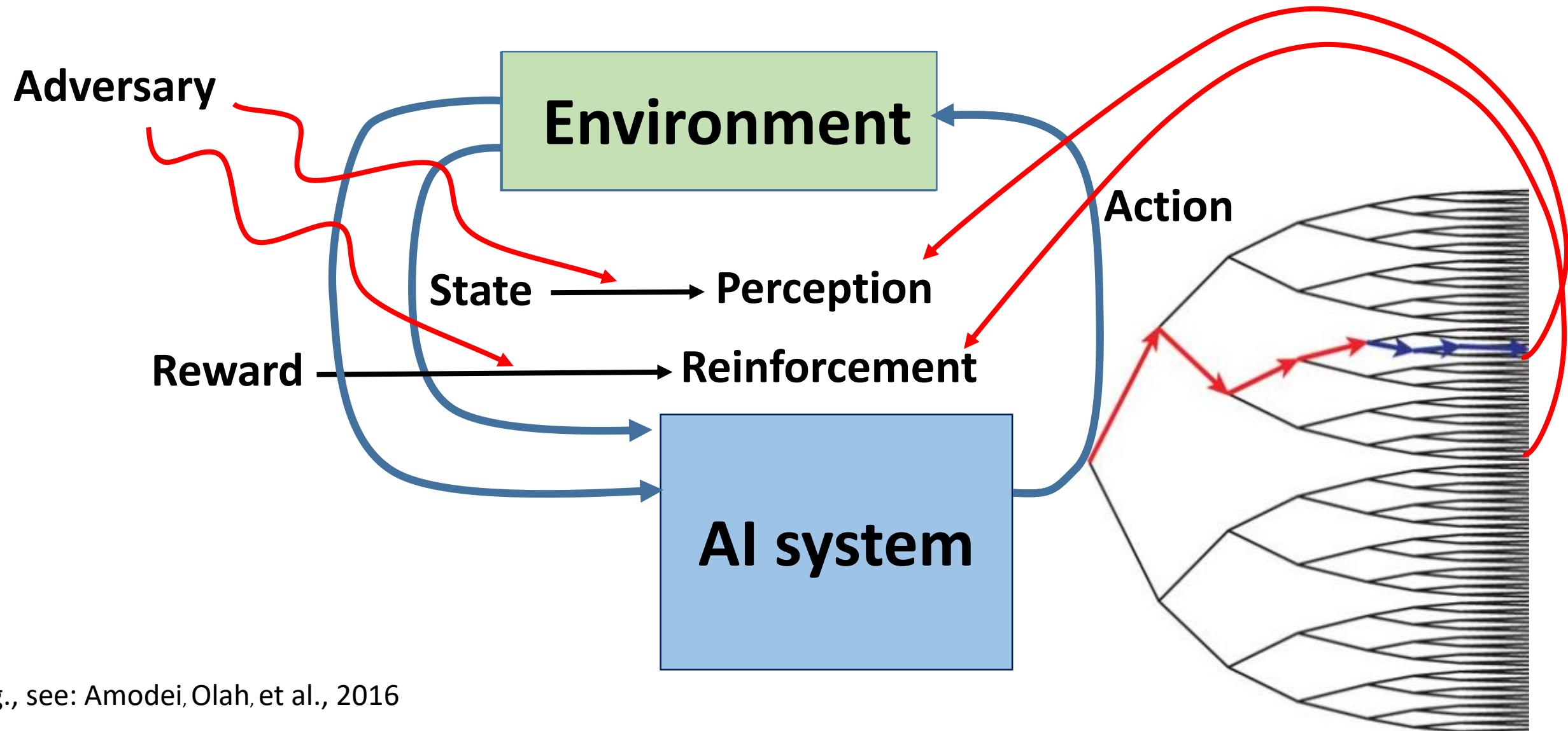
e.g., see: Amodei, Olah, et al., 2016

Adversarial Attacks & Self-Modification



e.g., see: Amodei, Olah, et al., 2016

Adversarial Attacks & Self-Modification



Adversarial Attacks & Self-Modification

Run-time verification

Static analysis

Environment

Reflective analysis

Ensure isolation * identify meddling * ensure operational faithfulness

AI system

AI & People

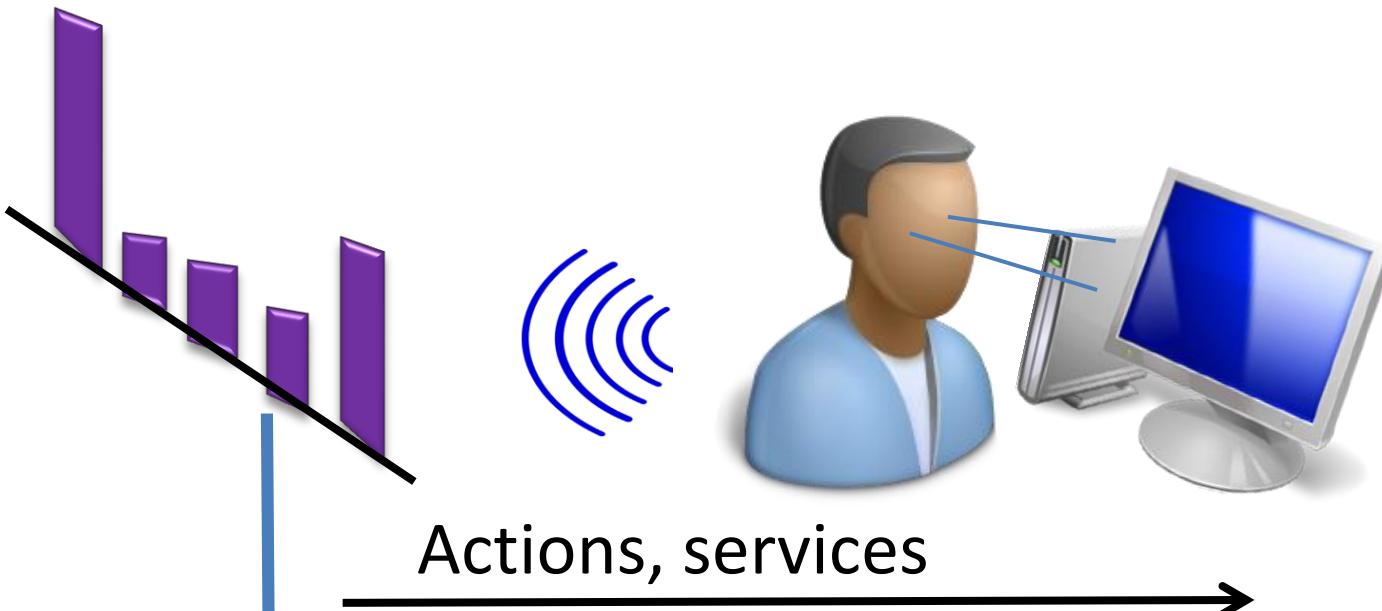
Directions

Models of people & tasks

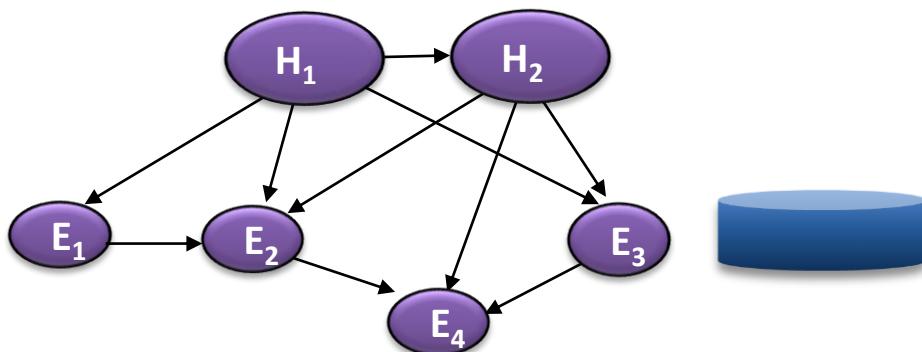
Models of complementarity

Coordination of initiative

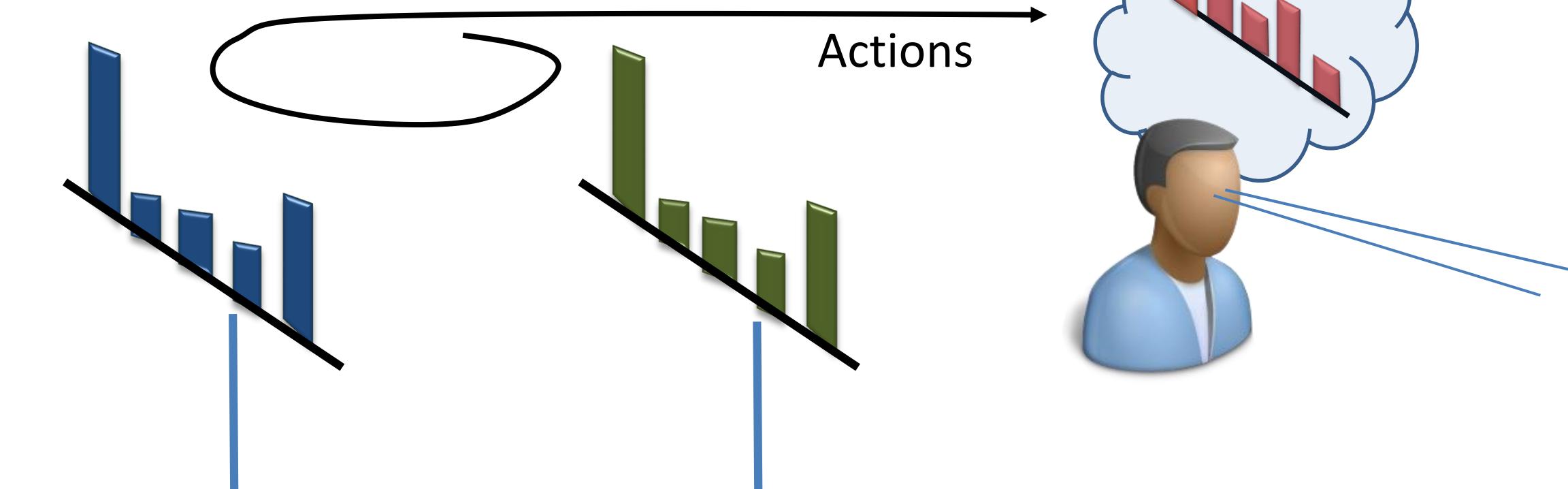
Models of people & tasks



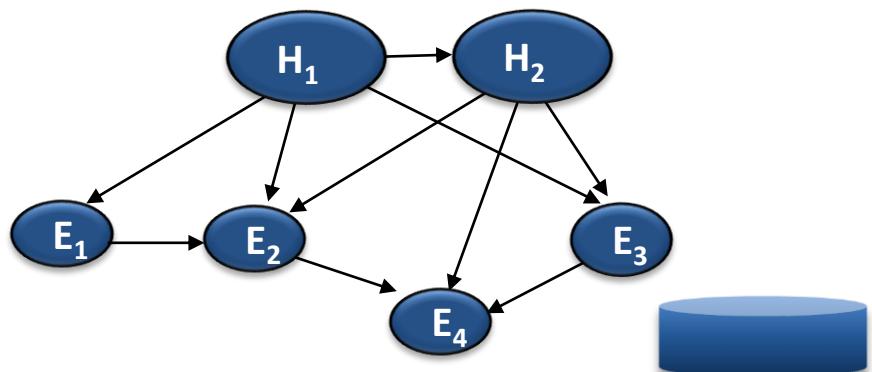
Predictions about needs, goals



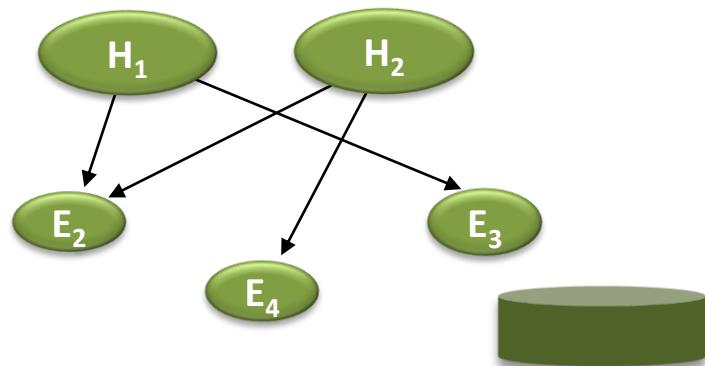
Models of world & people



Predictions about world



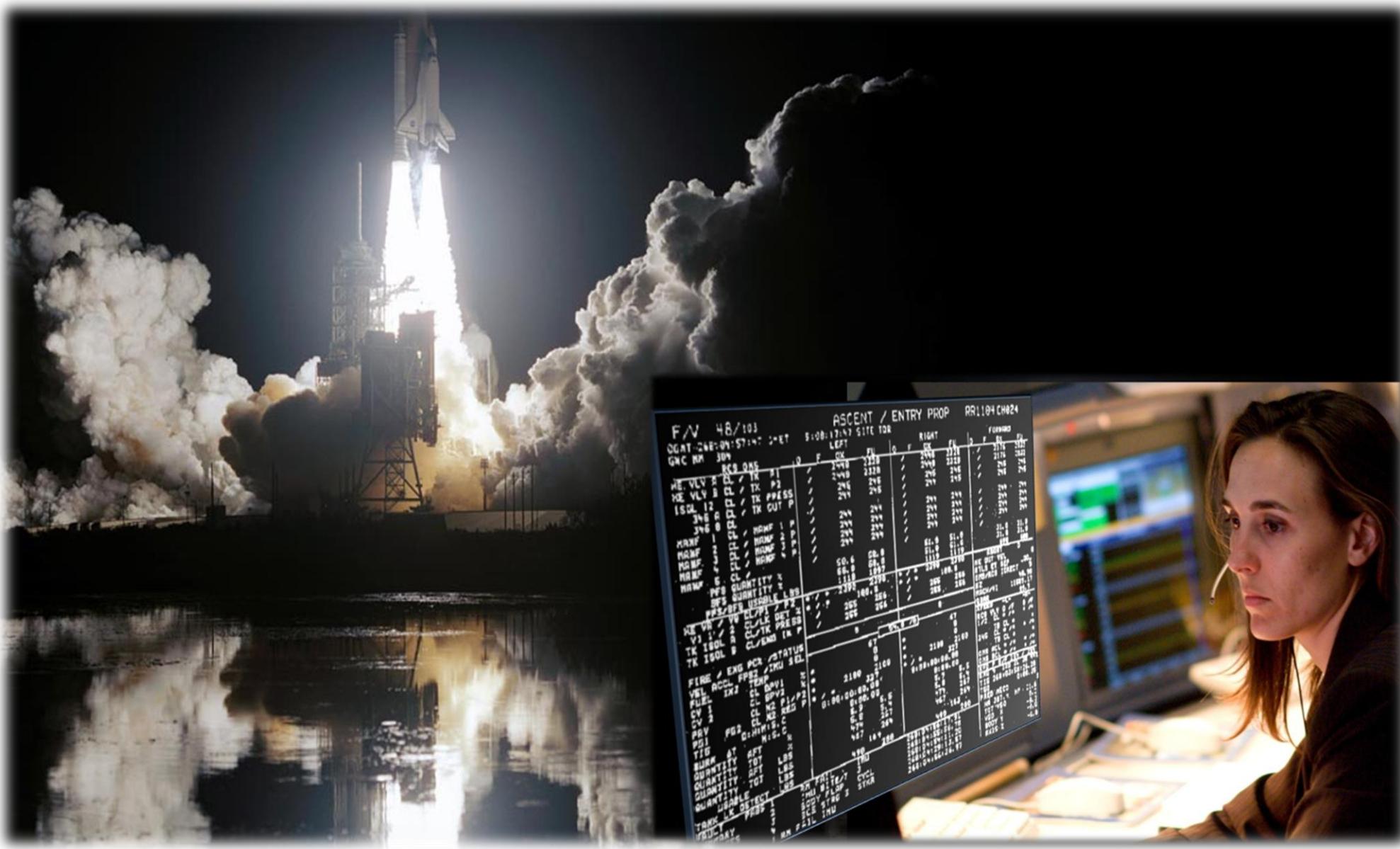
Predictions about user beliefs



H. Barry, 1995

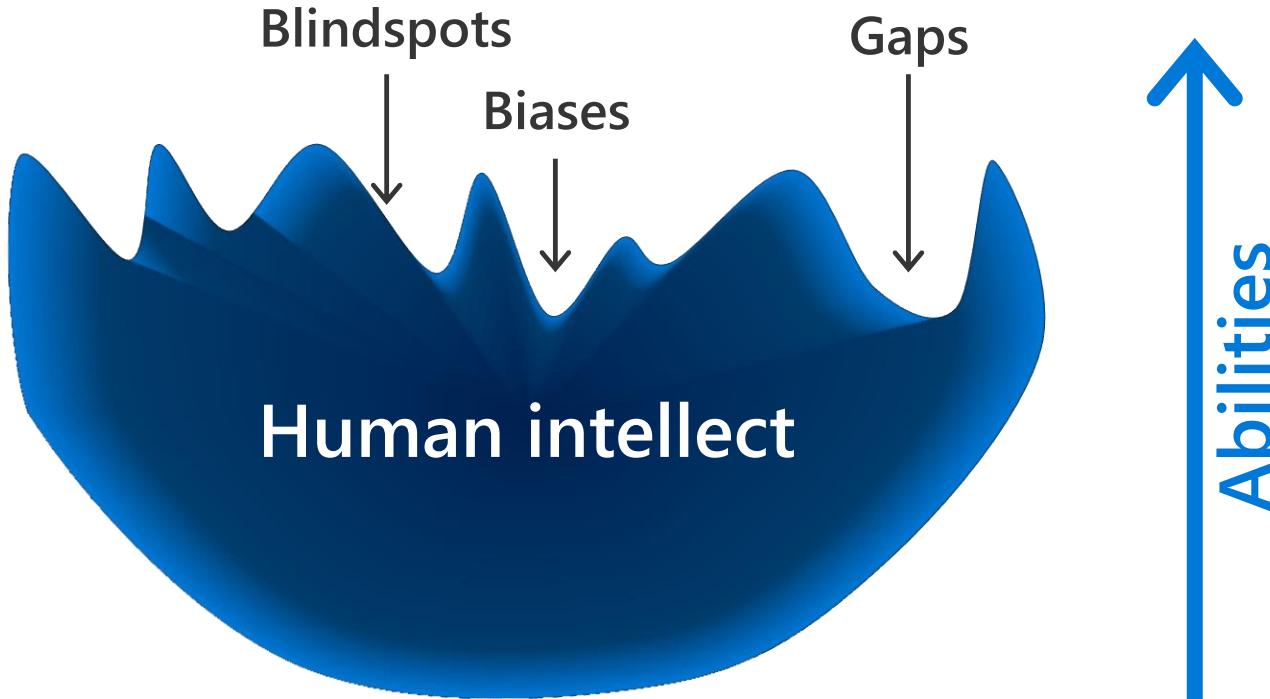
H. , Apacible, Sarin, Liao, 2005

Models of world & people



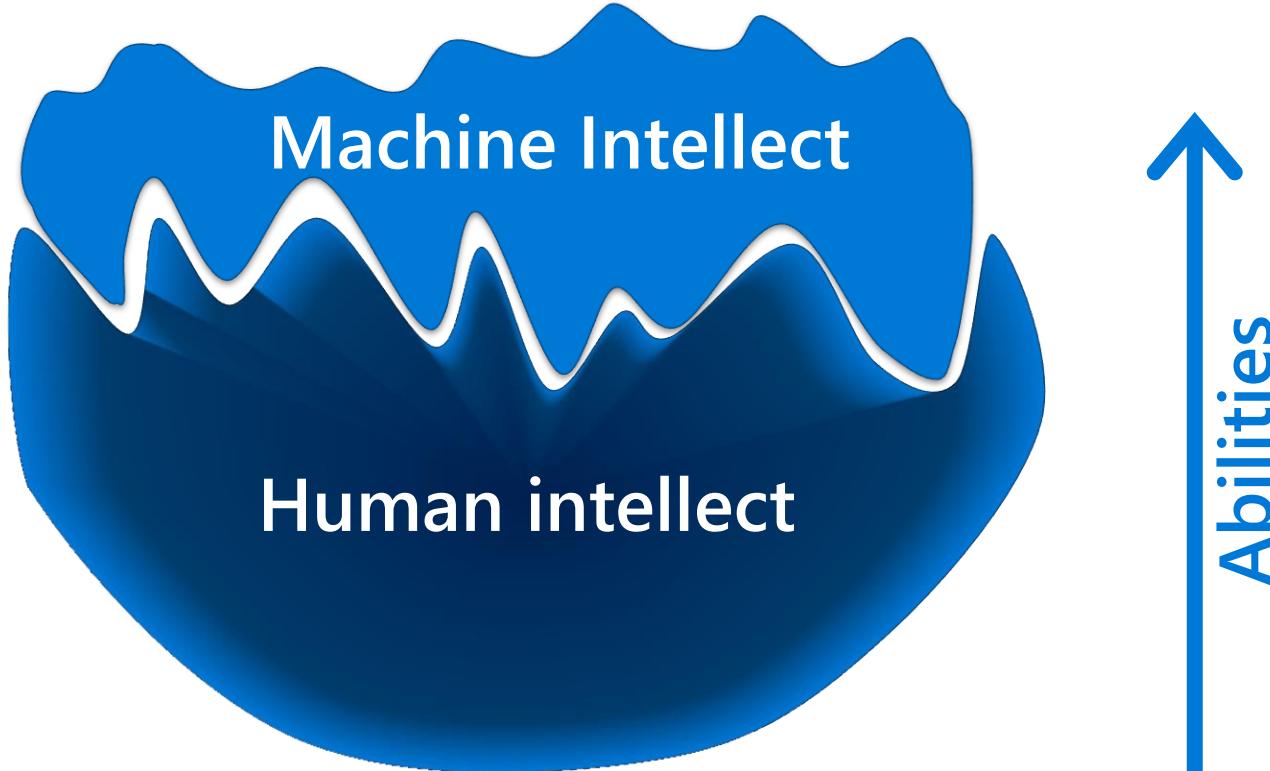
H. Barry, 1995

Complementarity



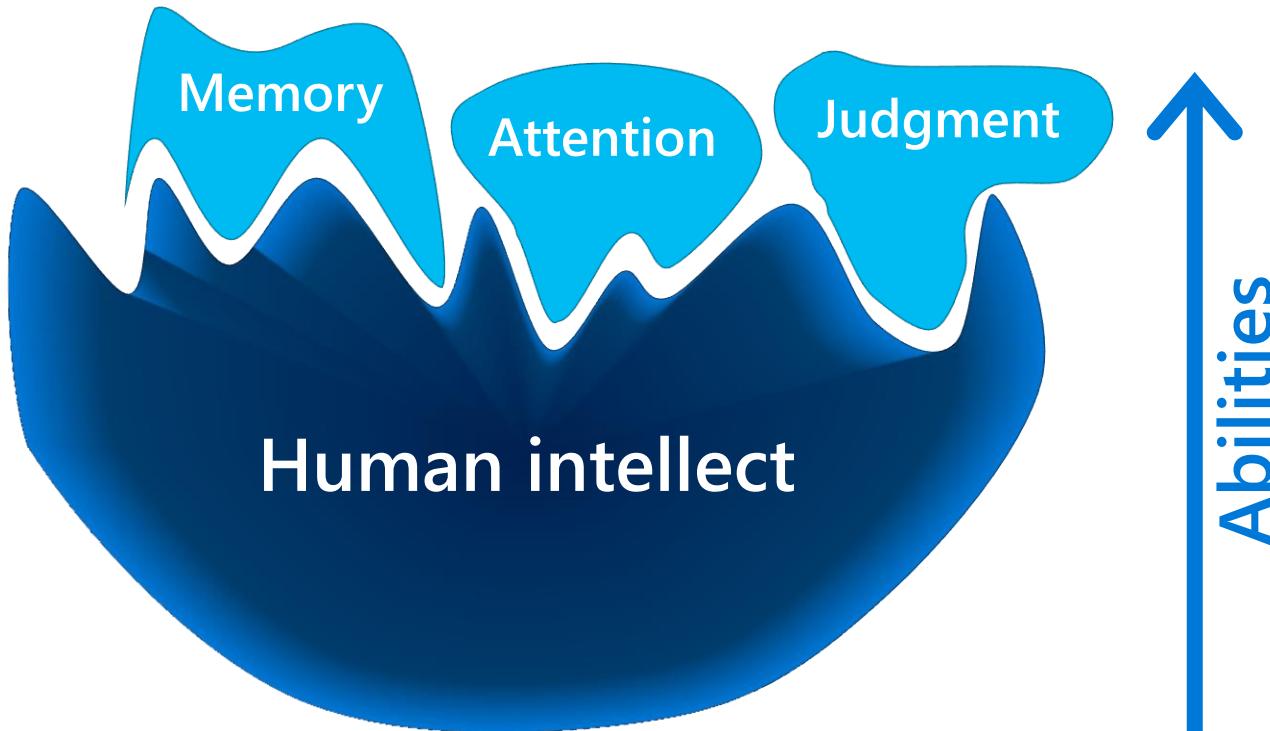
Leverage and extend results from
cognitive psychology

Complementarity



Leverage and extend results from
cognitive psychology

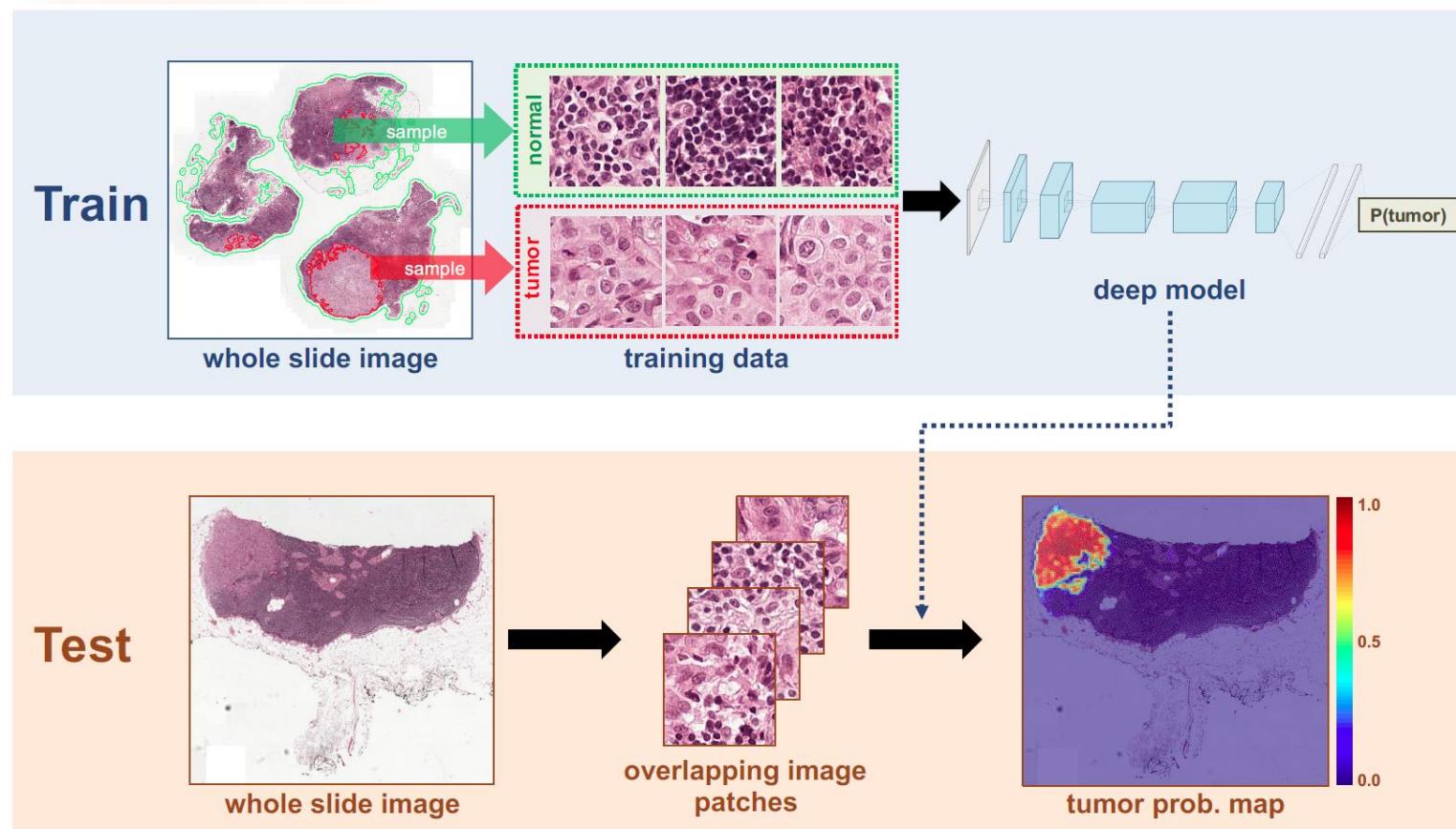
Complementarity



Leverage and extend results from
cognitive psychology

Complementarity

Identifying metastatic breast cancer
(Camelyon Grand Challenge 2016)



Human is superior

Error: 3.4%

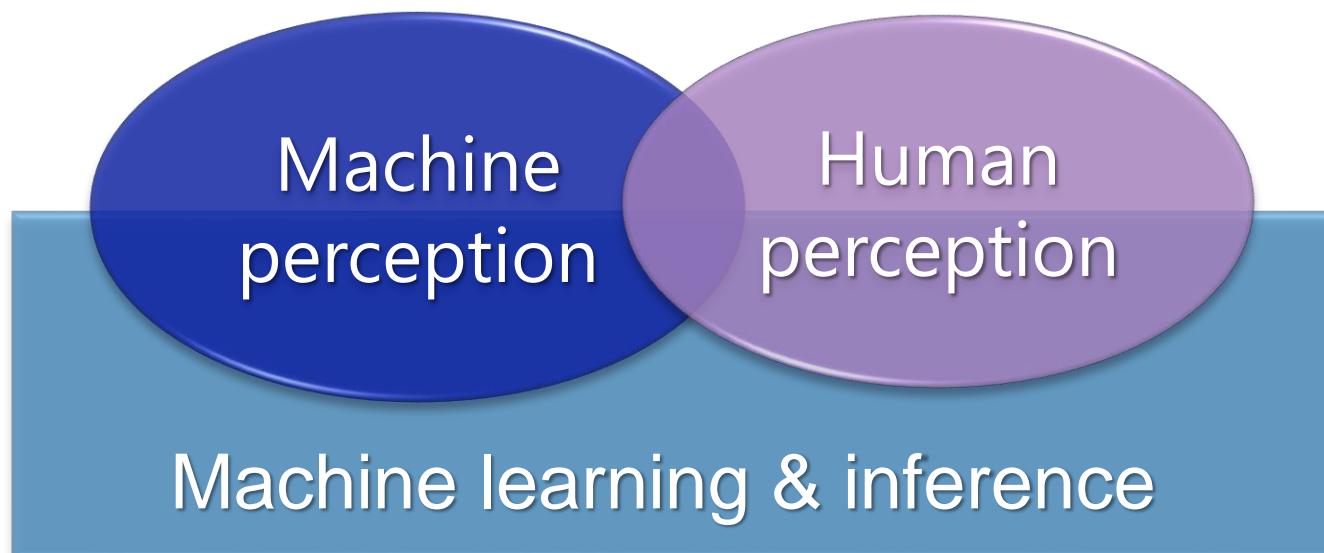


AI + Expert: 0.5%

85% reduction in errors.

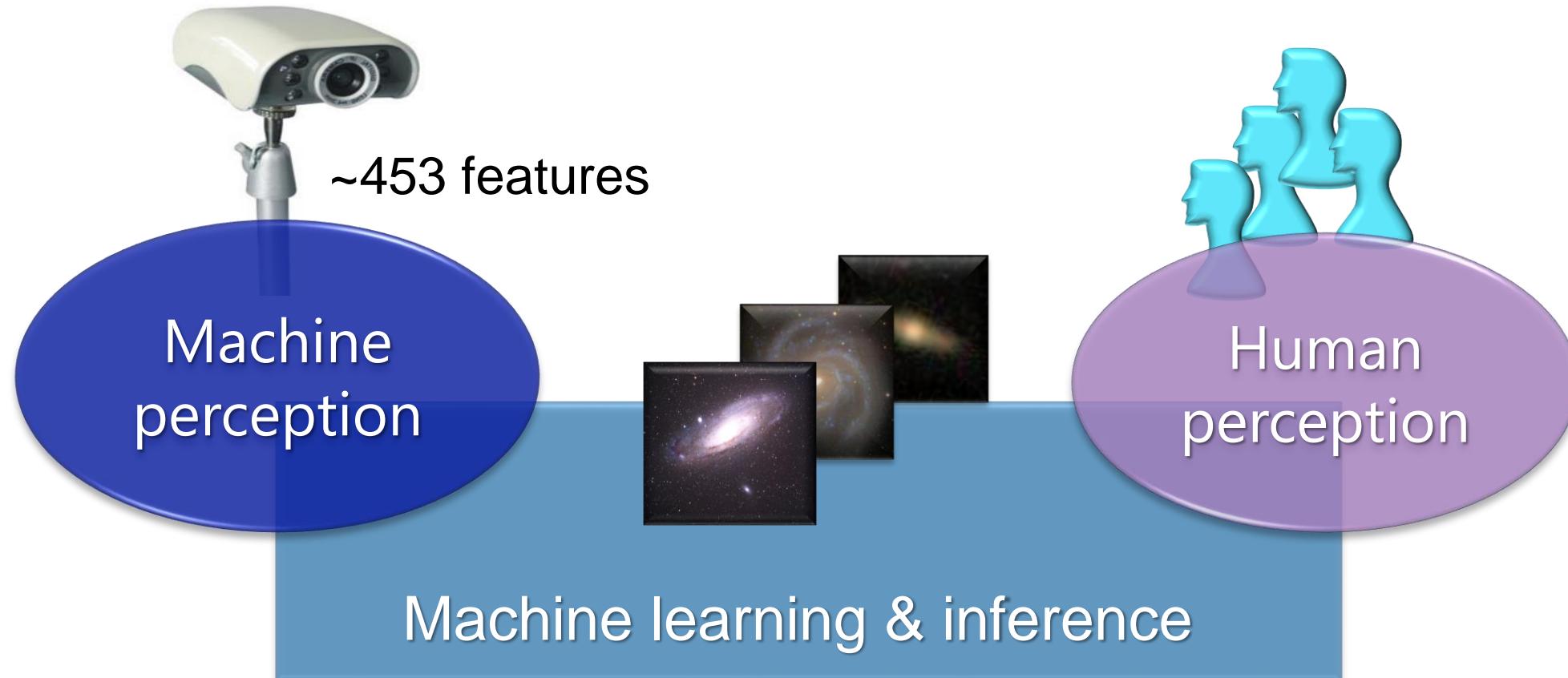
Complementarity

Label galaxies in Sloan Digital Sky Survey
(Galaxy Zoo)



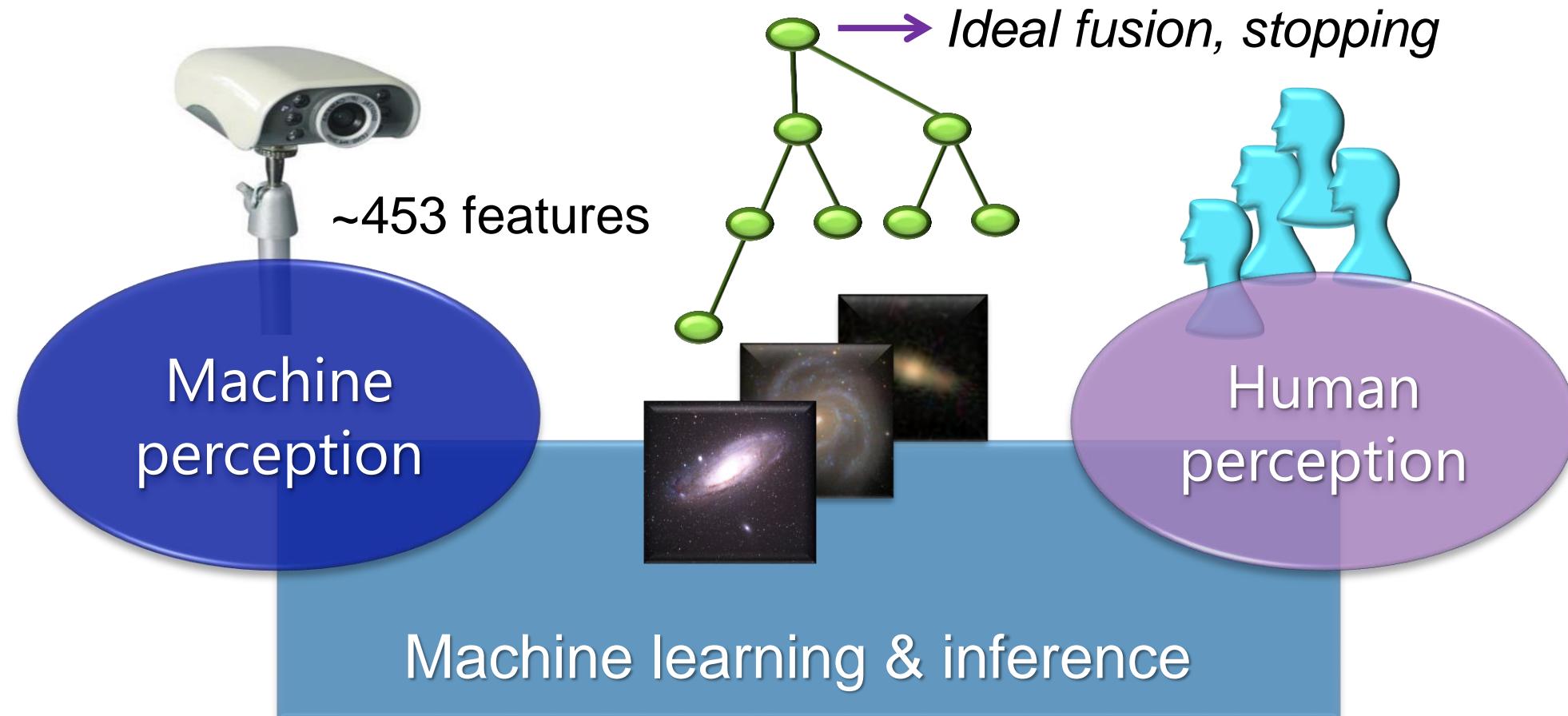
Complementarity

Label galaxies in Sloan Digital Sky Survey
(Galaxy Zoo)



Complementarity

Full accuracy: 47% of human effort
95% accuracy: 23% of human effort





Designs for mix of initiatives

Human
cognition

Machine
intelligence

Machine learning & inference

Design, learning, optimization

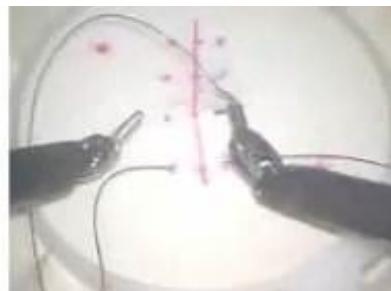
Recognizing intention



Reach needle #22



Position #165



Insert #162



Left transfer #119



Right transfer #37



Pull #160



Orient #48



Tighten suture #23



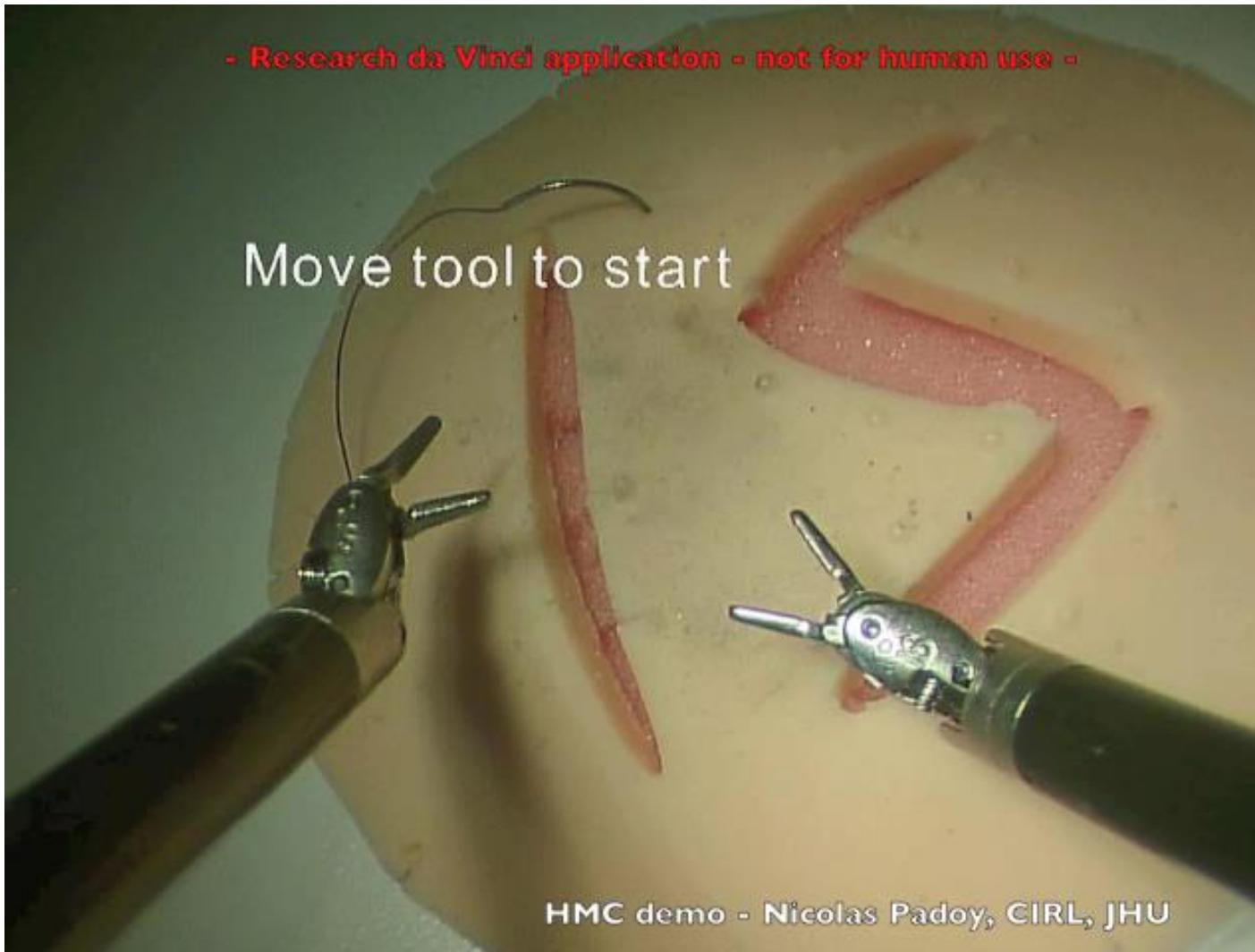
Loosening #4



Dropping #39

C.E. Reiley, et al.

Coordination of initiative



Padoy & Hager. ICRA 2011
van den Berg, et al, ICRA, 2010

AI, people, and society

Multiple influences and concerns

- Trustworthiness and safety
- Fairness, accuracy, transparency
- Ethical and legal aspects of autonomy
- Jobs and economy



Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

May 23, 2016

Bernard Parker: rated high risk



Dylan Fugett: rated low risk.



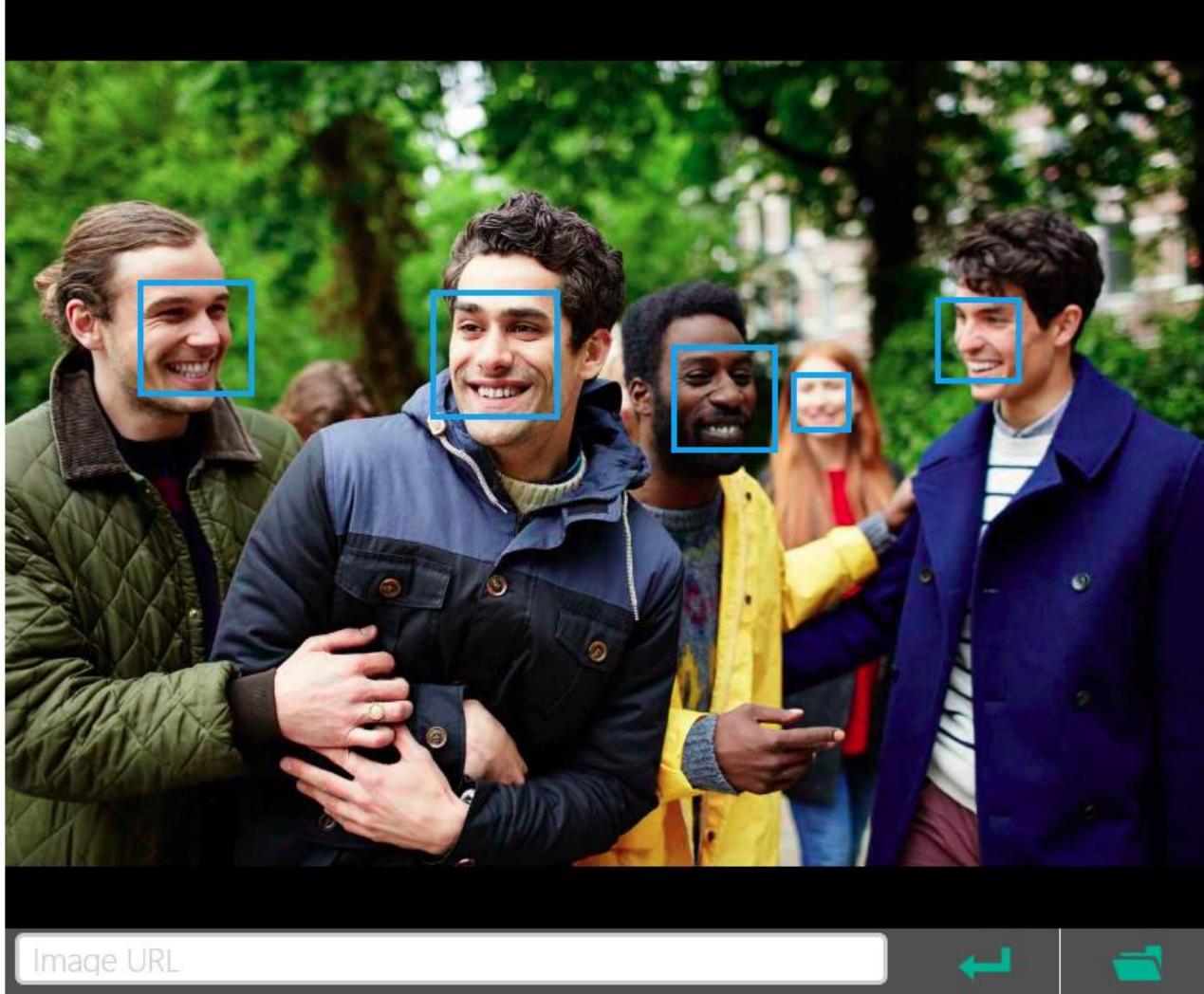
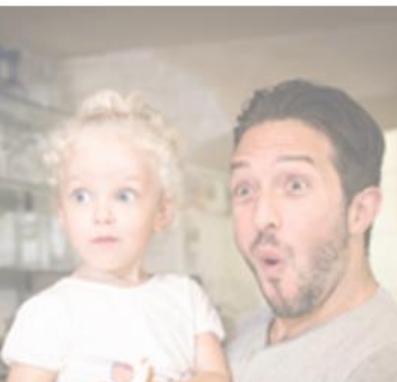
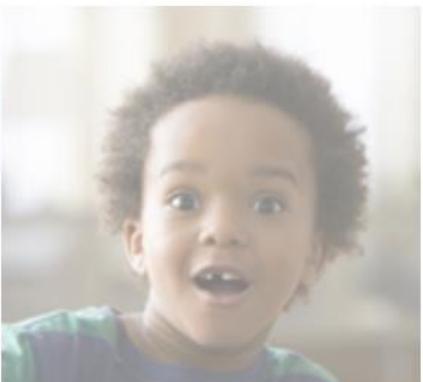


Image URL



Detection Result:

5 faces detected

JSON:

[

{

```
"faceRectangle": {  
    "left": 488,  
    "top": 263,  
    "width": 148,  
    "height": 148  
},
```

```
"scores": {
```

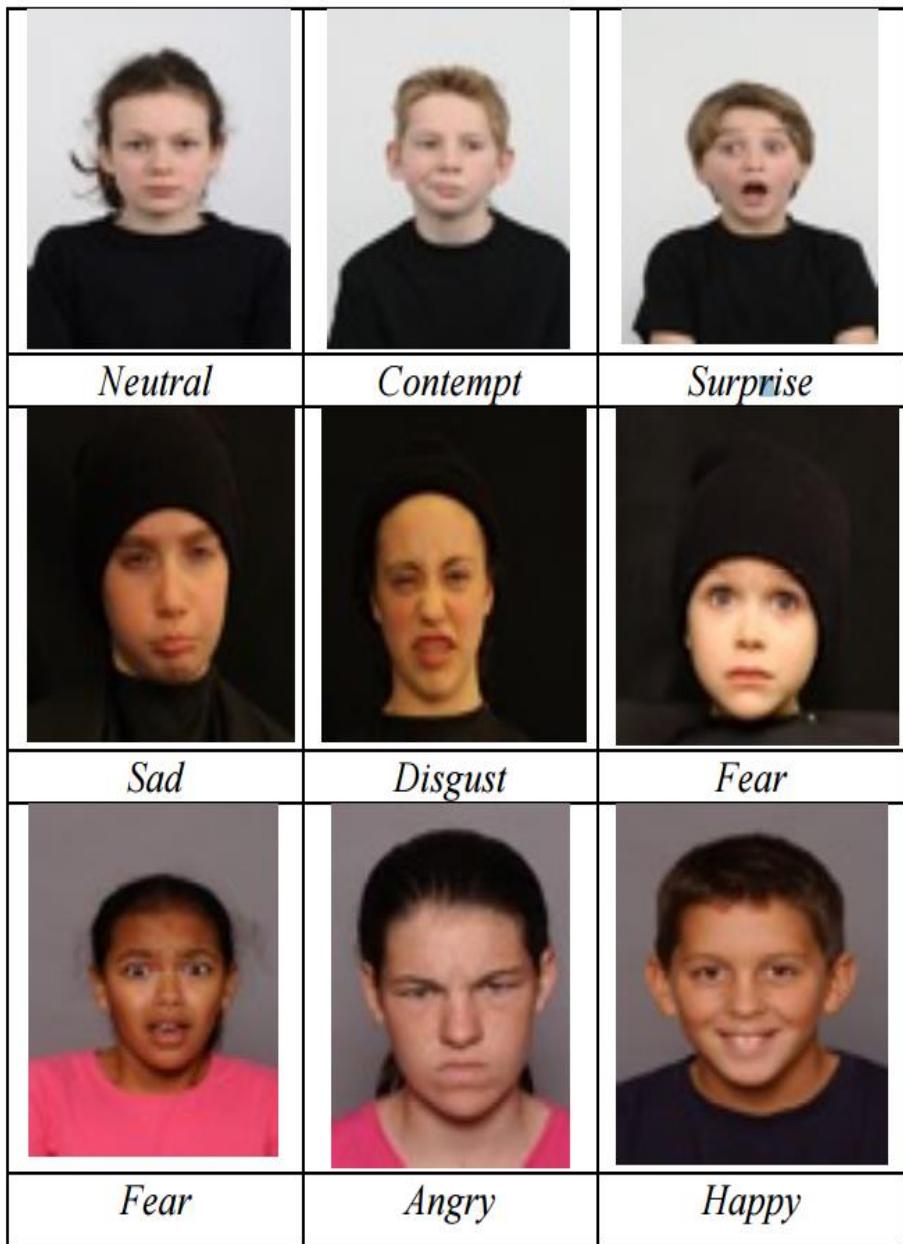
```
    "anger": 9.075572e-13,  
    "contempt": 7.048959e-9,  
    "disgust": 1.02152783e-11,  
    "fear": 1.778957e-14,  
    "happiness": 0.9999999,  
    "neutral": 1.31694478e-7,  
    "sadness": 6.04054263e-12,  
    "surprise": 3.92249462e-11
```

}

,

{

```
"faceRectangle": {
```



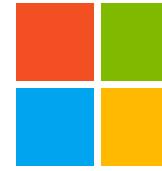
Addressing Bias in Machine Learning Algorithms: A Pilot Study on Emotion Recognition for Intelligent Systems

Ayanna Howard^{1*}, Cha Zhang², Eric Horvitz²

March 2017

Machine learning “contact lens” for children

Corporate &
community
responsibility



Microsoft

Aether Advisory Panel

AI and Ethics in Engineering and Research



Partnership on AI

to benefit people and society

AI in the Open World

Science & engineering

Human-AI collaboration

AI, people, and society

Much to do