



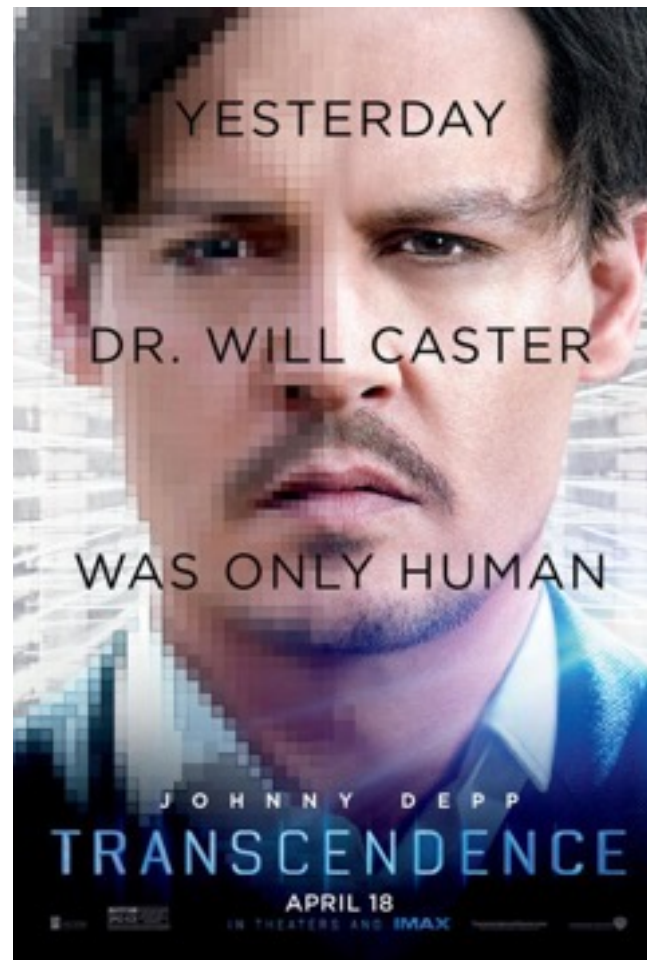
Implications of Artificial Intelligence

– the most transformative technology of the 21st century

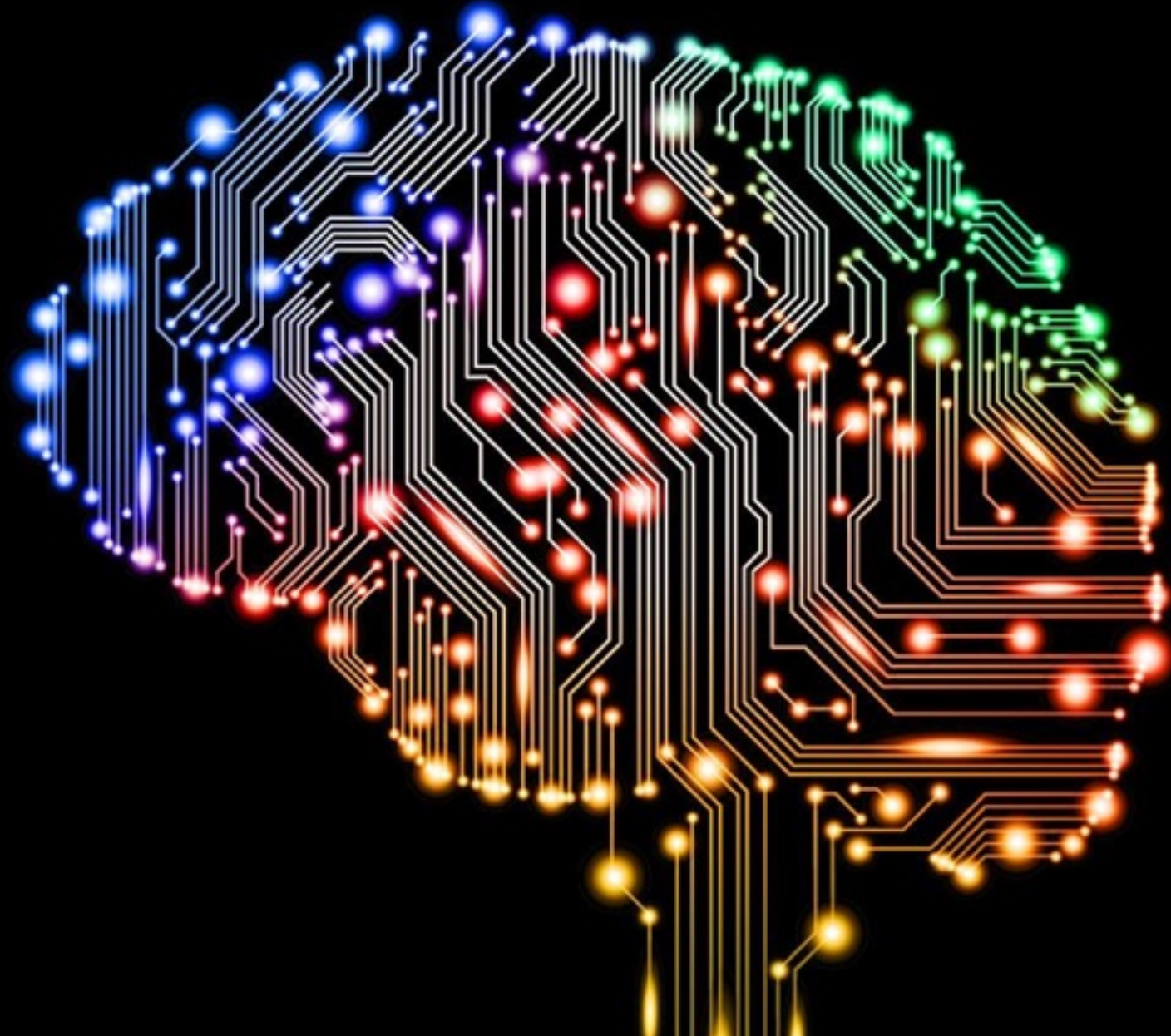
Robin Li, Bill Gates, Elon Musk on AI



More Money for Entertainment



... than ensuring a good outcome!



Introduction

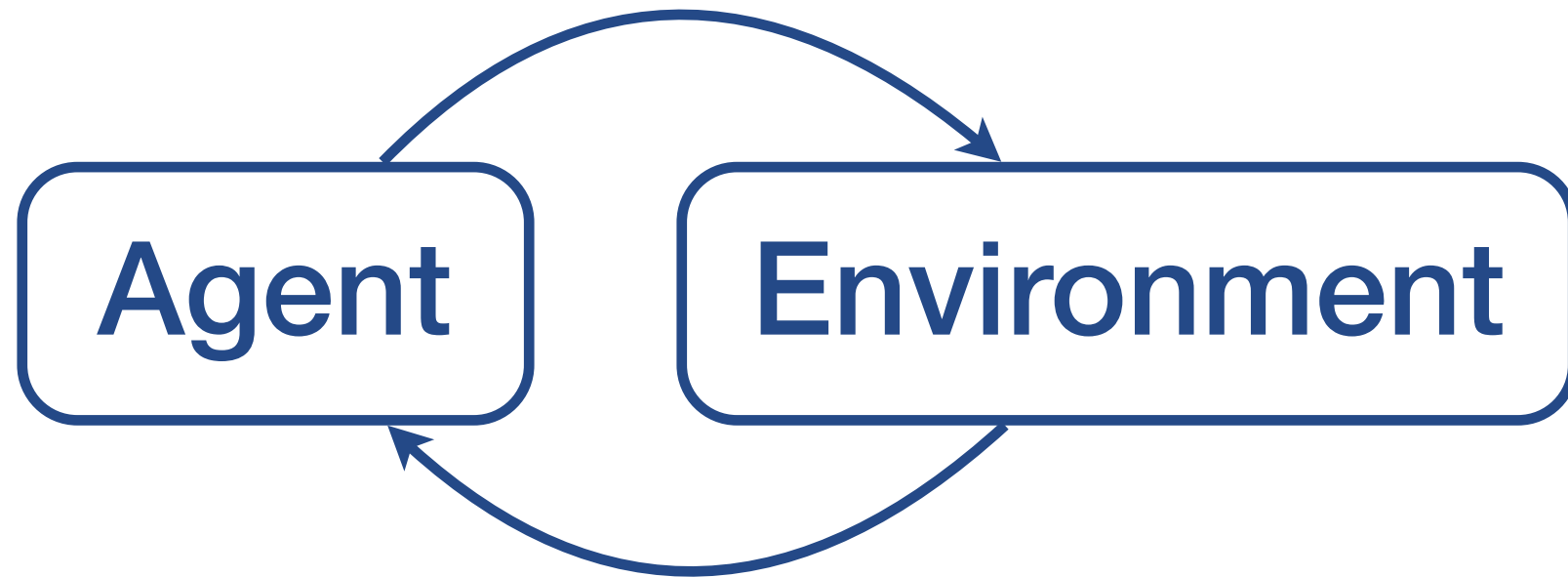
What are we talking about?

Intelligence

«Intelligence measures an agent's ability to achieve its goals in a wide range of unknown environments.»

$$\text{Intelligence} = \frac{\text{Optimization Power}}{\text{Used Resources}}$$

Ingredients



Learn, predict, rate and plan!

Intelligence is a Big Deal



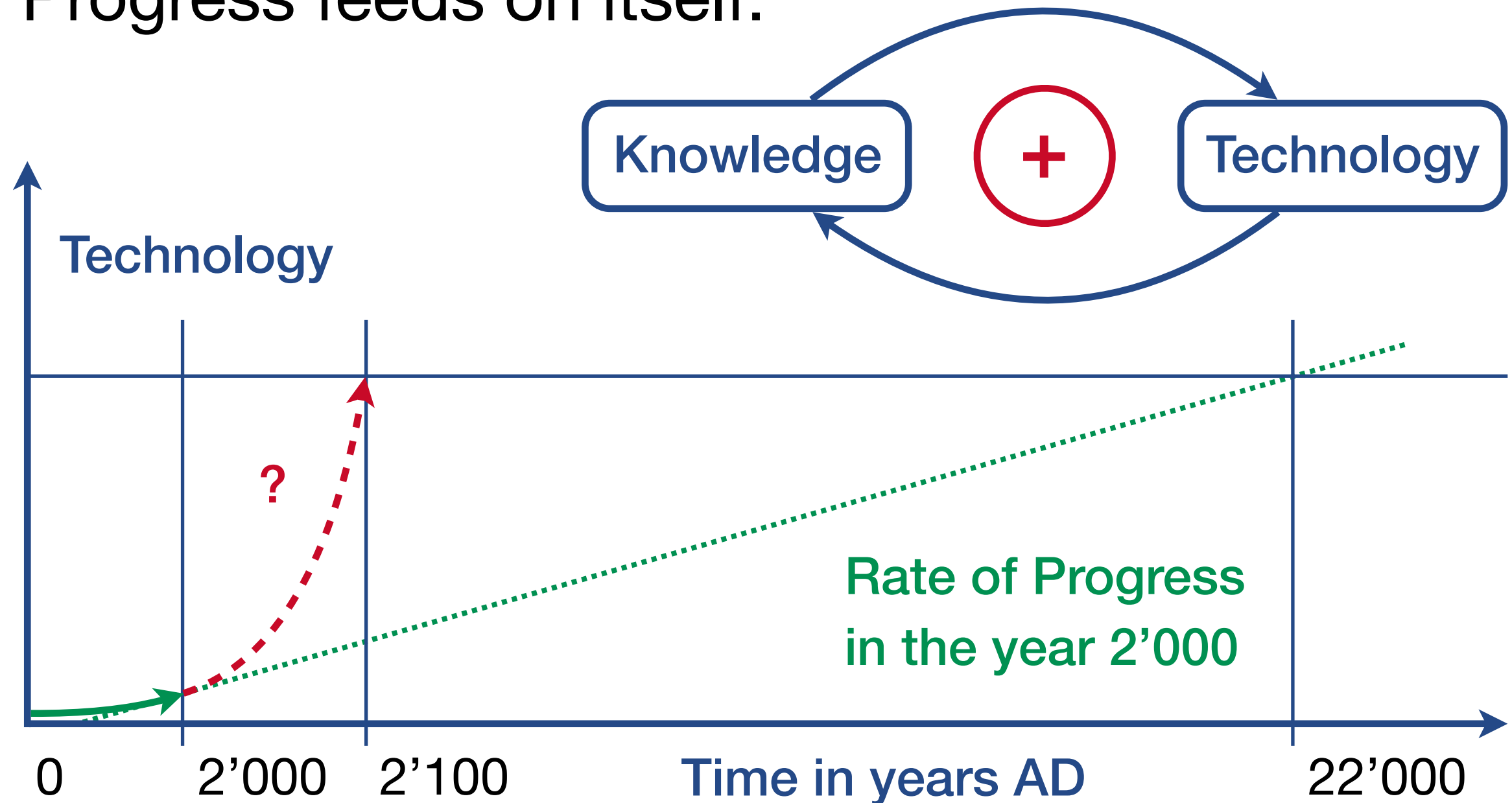
6 million years ago, 96% common DNA

Fast-Evolving Human DNA Leads to Bigger-Brained Mice
[phenomena.nationalgeographic.com/2015/02/19/\[...\]](http://phenomena.nationalgeographic.com/2015/02/19/[...])

Implications of AI
Zurich, May 2016

Accelerating Change

Progress feeds on itself:



Technology = (Neutral) Lever

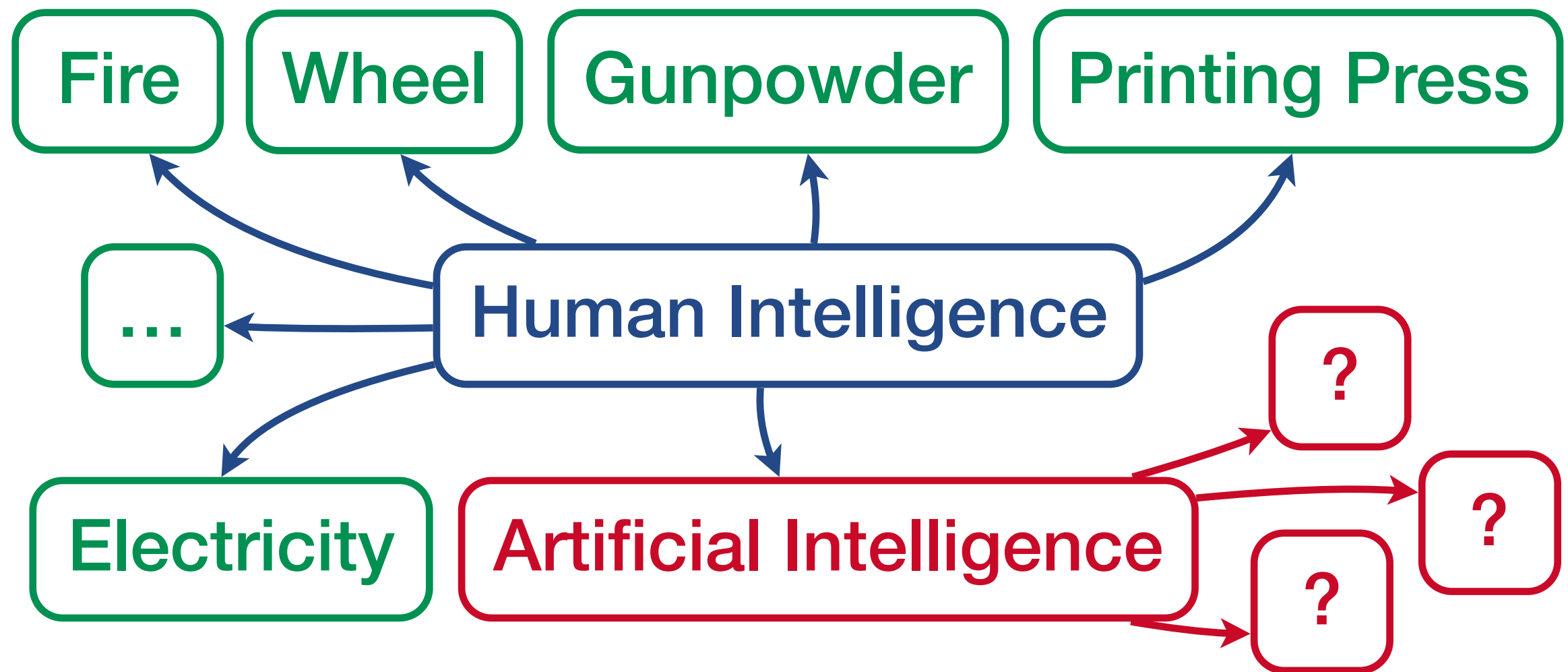
Greed and laziness drive us to

- increase our productivity
 - make our lives easier
- ... which is fine except:



**Our technological progress far
outperforms our moral progress!**

Artificial Intelligence



Intelligence is a technology like no other!



Current Trends

Where are we heading to?

State of the Art



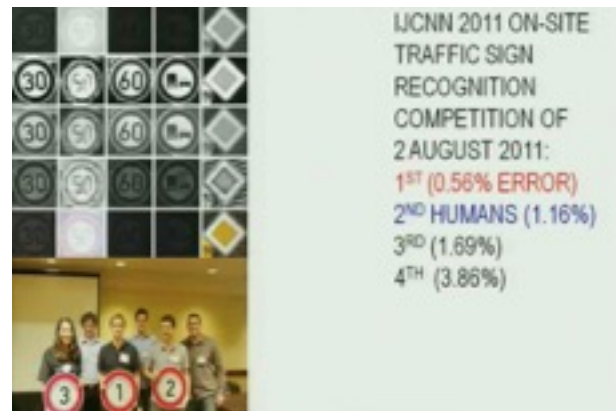
Deep Blue: 1997



Stanley: 2005



IBM Watson: 2011



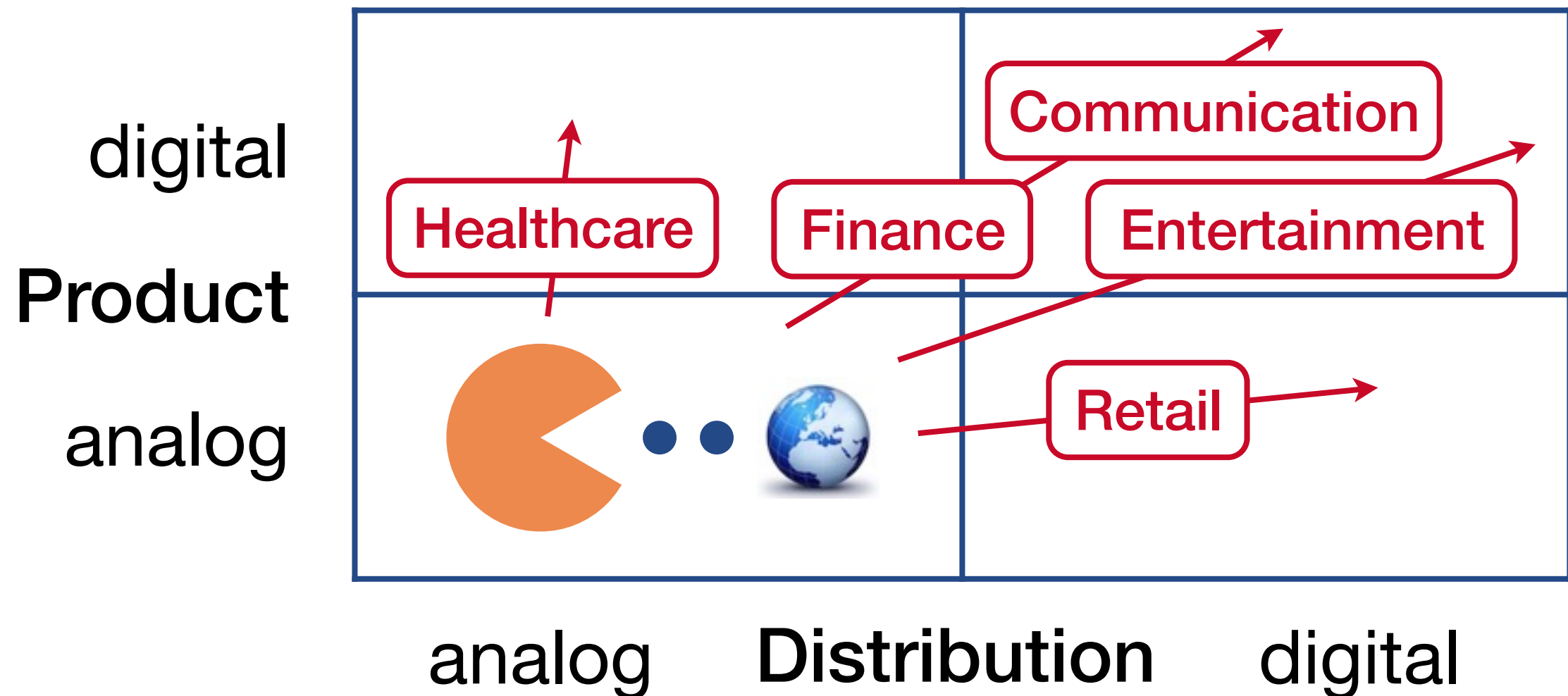
Schmidhuber: 2011

| | |
|------------|---------------|
| Checkers | Superhuman |
| Backgammon | Superhuman |
| Othello | Superhuman |
| Chess | Superhuman |
| Crosswords | Expert Level |
| Scrabble | Superhuman |
| Bridge | Equal to Best |
| Jeopardy! | Superhuman |
| Poker | Varied |
| FreeCell | Superhuman |
| Go | Superhuman |

How bio-inspired learning keeps winning competitions
www.kurzweilai.net/how-bio-inspired-deep-learning-...

Implications of AI
 Zurich, May 2016

Software is eating the world



... and AI will digest all this (big) data!

Describes without errors

Describes with minor errors

Somewhat related to the image

Unrelated to the image



A person riding a motorcycle on a dirt road.



Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



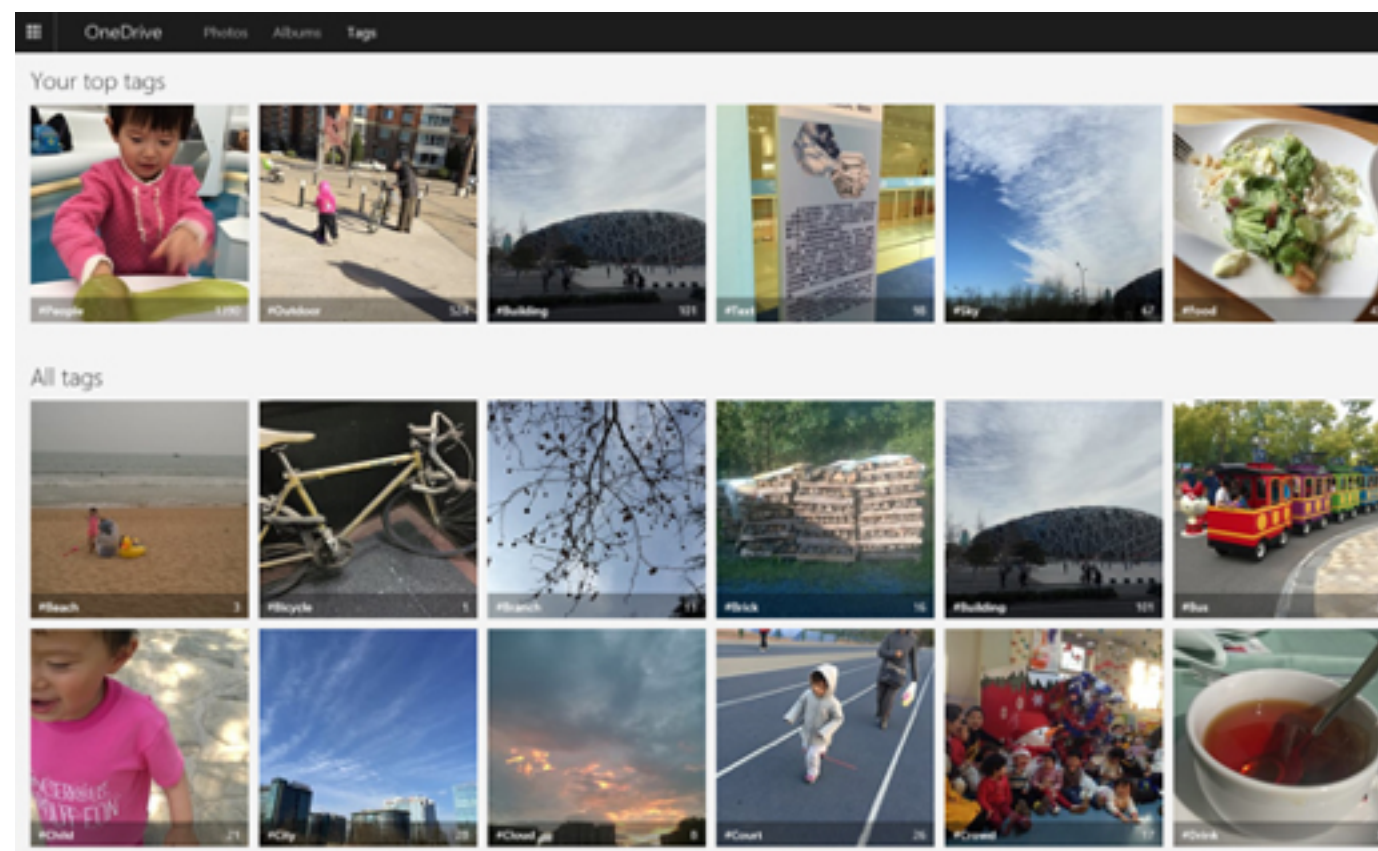
A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.

Superhuman Image Recognition

- With convolutional neural networks
- 1.2 m training images, ~ 30 layers





Jeremy Howard
go.ted.com/bbZC

Implications of AI | 16
Zurich, May 2016

MACHINE INTELLIGENCE 2.0

AGENTS

| PROFESSIONAL | PERSONAL | OS INTERFACES |
|--|-----------------------|--------------------|
| DigitalGenius OVERLAP.CC machines PRIMER | @livesome | Google Now |

AUTONOMOUS SYSTEMS

| AIR | GROUND | SEA | INDUSTRIAL |
|--------------|--------------|--------------|------------------|
| | | | |

ENTERPRISE

| SECURITY / FRAUD | HR / RECRUITING | SALES | MARKETING | CUSTOMER SUPPORT | INTERNAL INTEL | MARKET INTEL |
|------------------|-----------------|--------------|-----------|------------------|----------------|--------------|
| | | | | | | |

PLATFORMS

| RESEARCH / AGI | FULL STACK | MACHINE LEARNING | INDUSTRIAL IOT | AUDIO | VISION | DATA ENRICHMENT |
|------------------|--------------|------------------|----------------|--------------|--------------|-----------------|
| | | | | | | |

INDUSTRIES

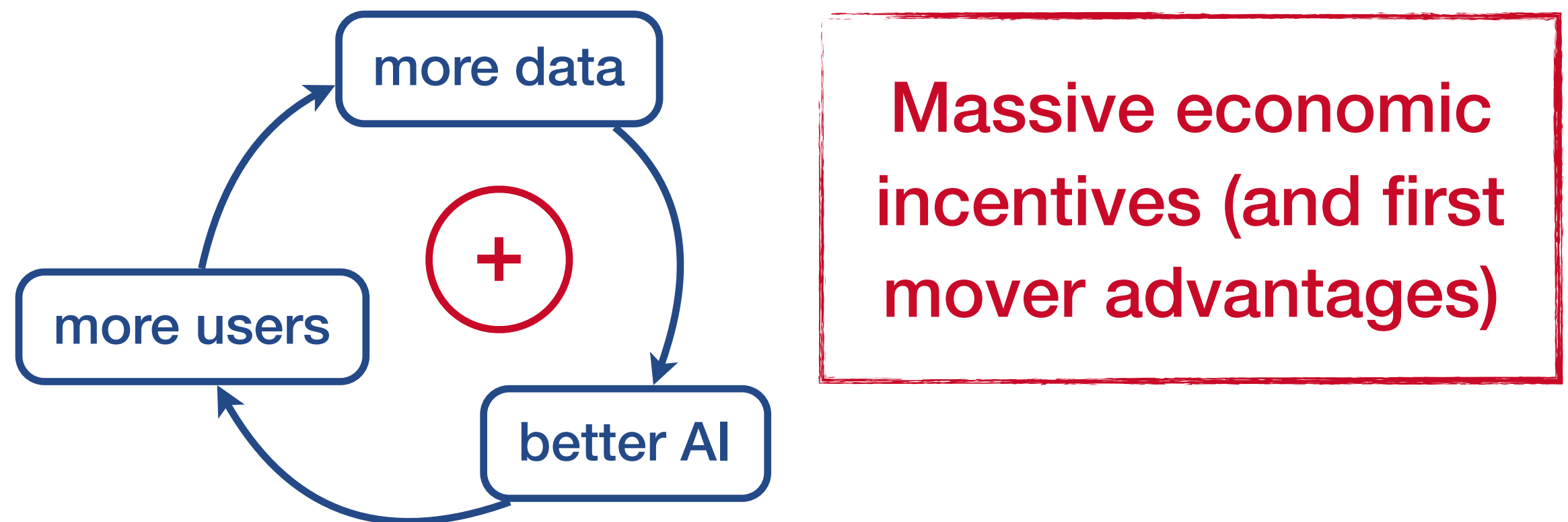
| ADTECH | AGRICULTURE | FOR GOOD | RETAIL FINANCE | LEGAL | MATERIALS & MFG | HEALTHCARE |
|--------------|--------------|--------------|----------------|--------------|------------------|------------------|
| | | | | | | |

INDUSTRIES (CONT'D)

| EDUCATION | TRANSPORT & LOGISTICS | INVESTMENT FINANCE | DATA SCIENCE | MACHINE LEARNING | OPEN SOURCE |
|--------------|-----------------------|--------------------|--------------|------------------|--------------|
| | | | | | |

Deep Learning Explosion

Stuart Russell: «Industry [has probably invested] more in the last 5 years than governments have invested since the beginning of the field [1950s].»



Automation brings abundance

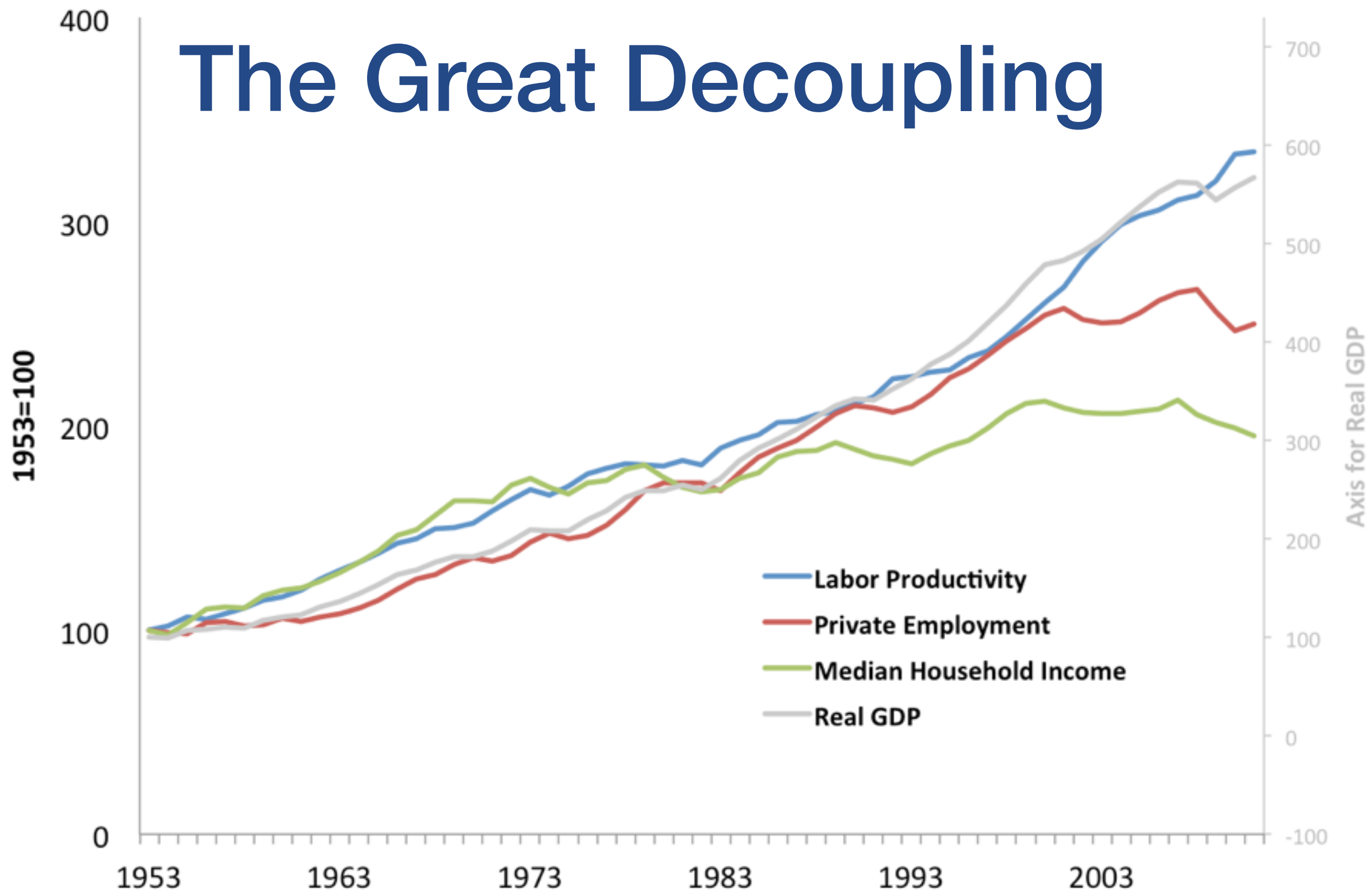
Many industries will
be transformed into
pure AI industries...

Capital: ↑, Wages: ↓



**Should be beneficial if we manage
it well – but we are not prepared!**

The Great Decoupling



© 2012 Andrew McAfee (@amcafee)

Sources: Census Bureau, Bureau of Labor Statistics

The Great Decoupling of the US Economy
[andrewmcafee.org/2012/12/\[...\]](http://andrewmcafee.org/2012/12/[...])

Implications of AI | 20
Zurich, May 2016

We'll give robots full autonomy

- ... because of
- increased speed
 - high complexity
 - risk of jamming

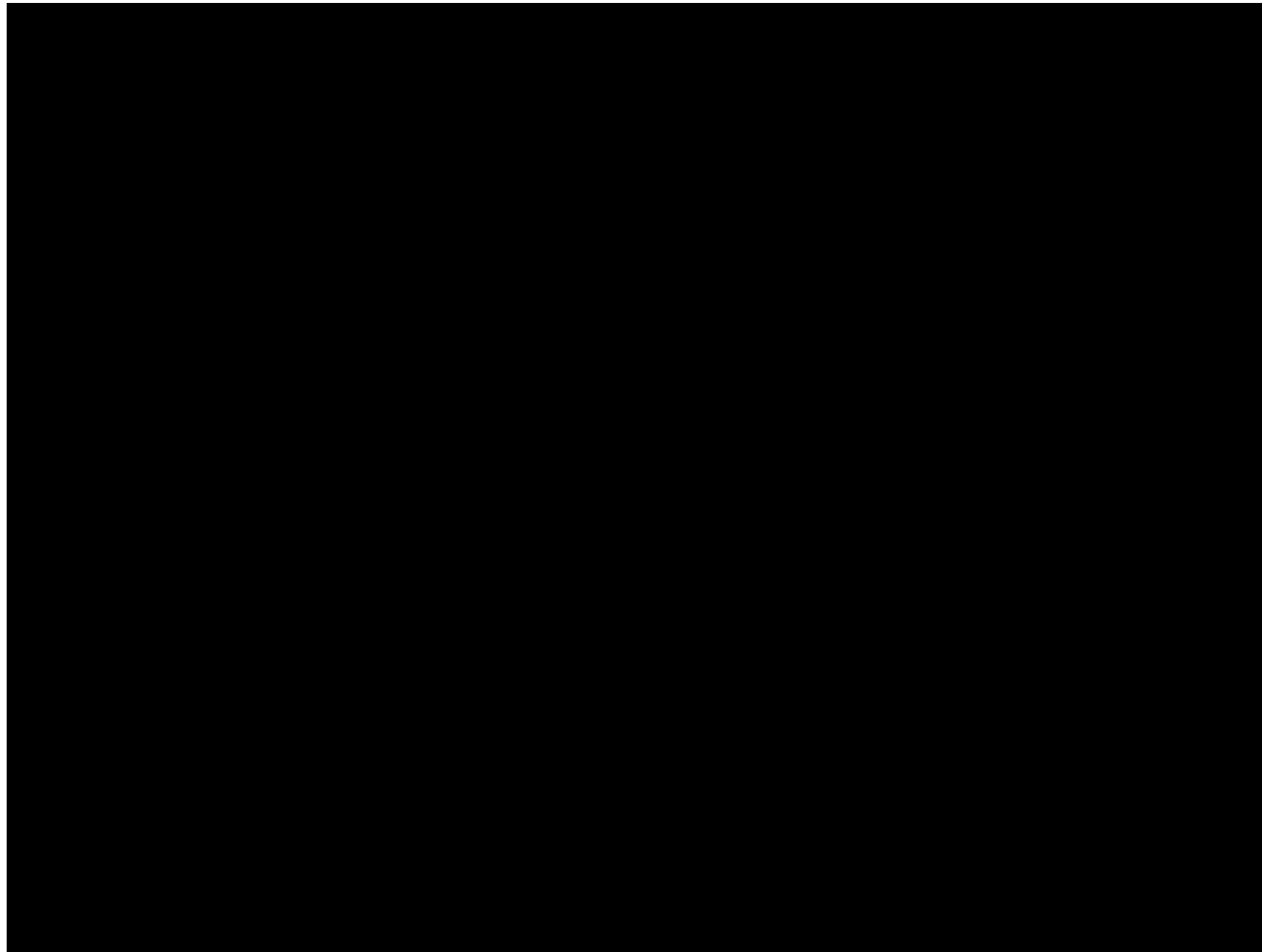
Current examples:

- financial markets
- auton. weapons



**AI is the ultimate
productivity boost!**

Machine Learning by Google



Google's Algorithms Learn Tasks Independently
[www.theguardian.com/technology/2015/feb/25/\[...\]](http://www.theguardian.com/technology/2015/feb/25/[...])

Implications of AI
Zurich, May 2016

Reinforcement Learning



**Train with first-person shooters
and deploy on armed drones...**

Visual Doom AI Competition



«Can AI effectively play Doom using only raw visual input?»



Challenges

What might go wrong?

A rational agent will strive to ...

- stay functional (self-preservation)
- keep its goal (goal-preservation)
- get stuff (resource accumulation)
- be smarter (intelligence explosion)

The problem is not malevolence but different goals and higher decision quality!

Single-Shot Situation

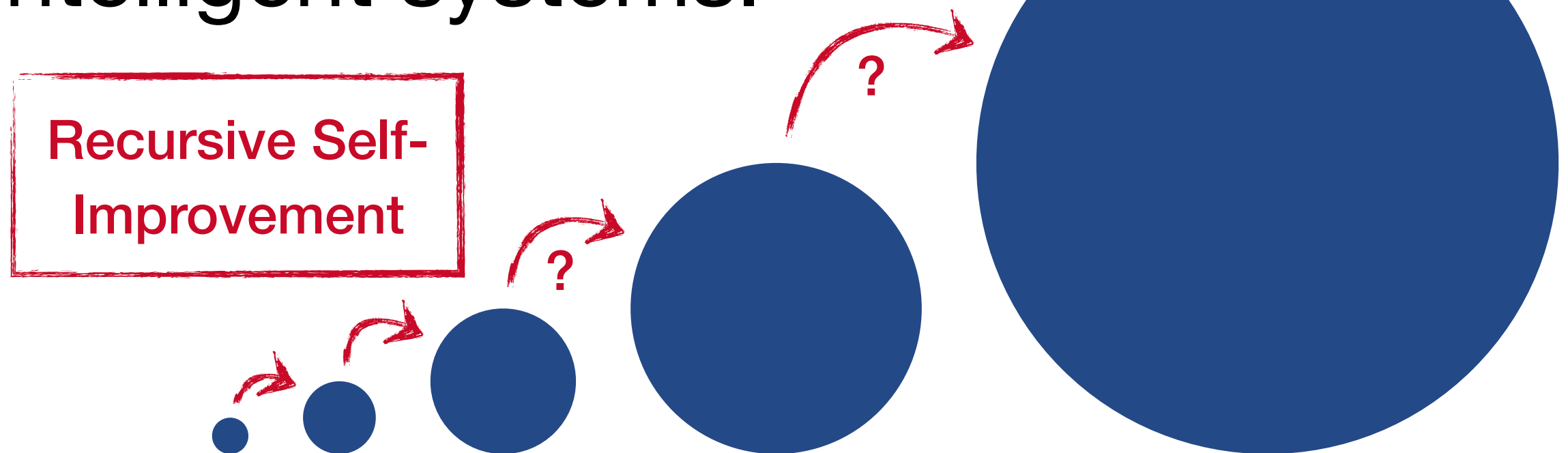
Our first superhuman AI must be a safe one for we may not get a second chance!



- We're good at iterating with testing and feedback
- We're terrible at getting things right the first time
- Humanity only learns when catastrophe occurred

Intelligence Explosion

Proportionality Thesis: An increase in intelligence leads to similar increases in the capacity to design intelligent systems.



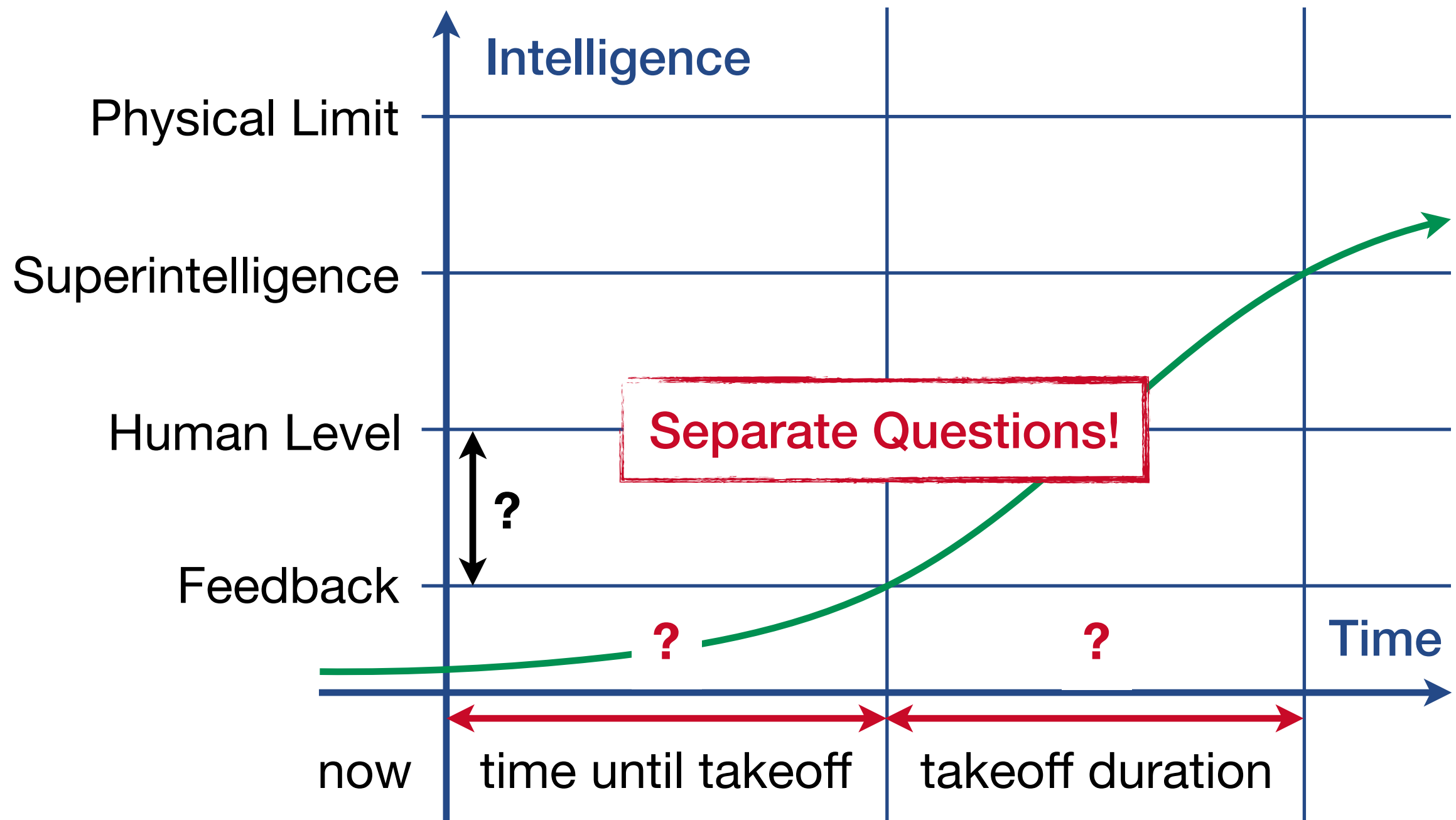
Technological Singularity

Theoretic phenomenon: There are arguments why it should exist but it has not *yet* been confirmed experimentally.

Three major singularity schools:

- **Accelerating Change** (Ray Kurzweil)
- **Intelligence Explosion** (I.J. Good)
- **Event Horizon** (Vernor Vinge)

Takeoff Scenarios



Advantages of AIs over Brains

| Hardware: | Software: | Effectiveness: |
|-----------|-----------------|-----------------|
| – Size | – Editability | – Rationality |
| – Speed | – Copyability | – Coordination |
| – Memory | – Expandability | – Communication |

| Human Brain | Modern Microprocessor |
|-----------------------|-------------------------|
| 86 billion neurons | 1.4 billion transistors |
| firing rate of 200 Hz | 4'400'000'000 Hz |
| 120 m/s signal speed | 300'000'000 m/s |



Interview by John Oliver with Stephen Hawking
www.youtube.com/watch?v=T8y5EXFMD4s

Implications of AI
Zurich, May 2016

Modelling Capabilities

An advanced AI will also model its operators and go to great lengths to prevent being switched off!

It will behave nicely and cooperatively until the external threats are under control and it is ready for takeover.

Optimization Power

Problem: When optimizing a system, unspecified parameters often assume extreme values.

You will get what you wished for and not what you wanted



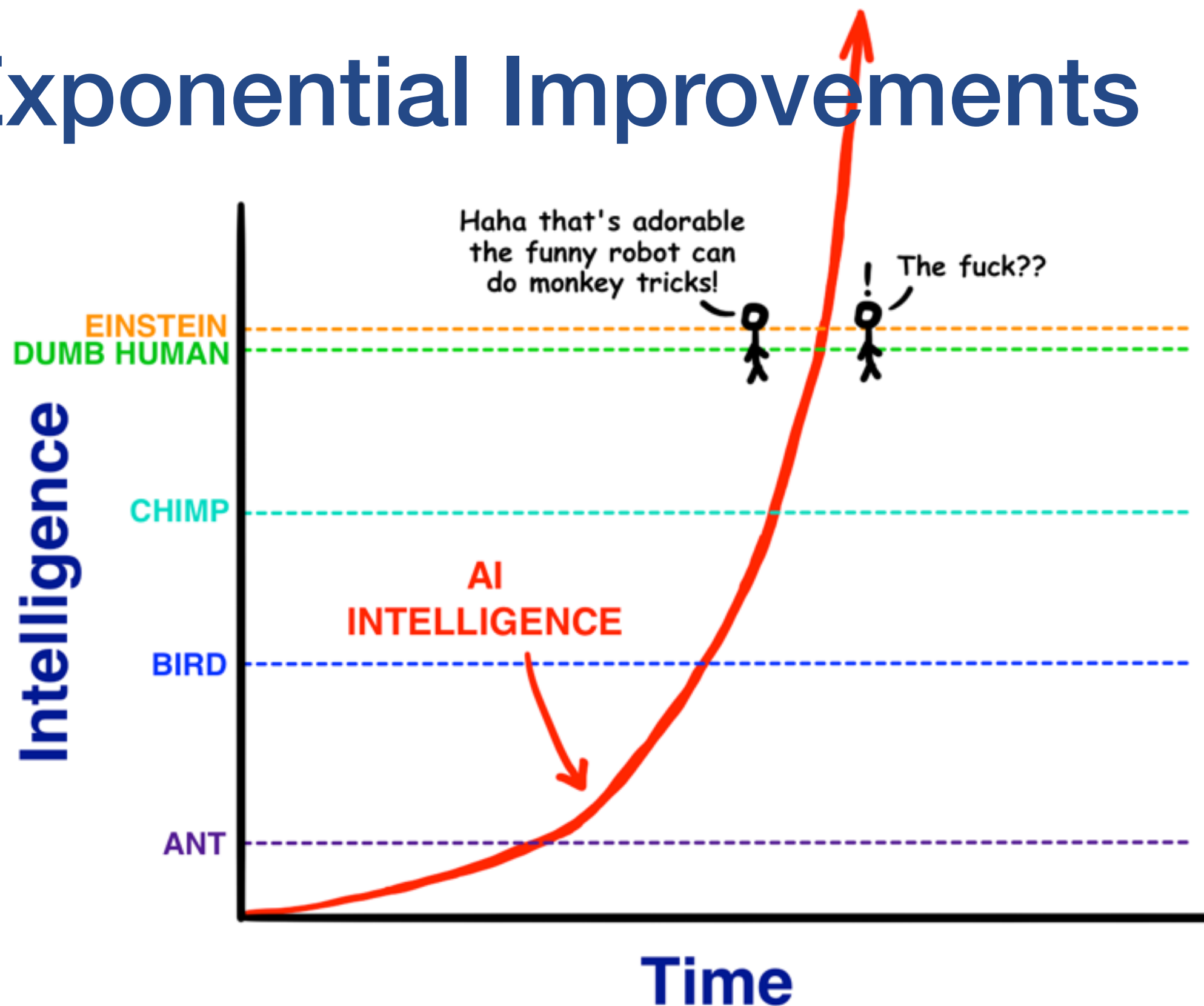
Outlook

What is ahead of us?

Implications of (advanced) AI

- **Short-term** (now to 15 years): classification/assistance, surveillance/intelligence, autonomous systems (cars, weapons, etc.) (AIs as moral agents)
- **Mid-term** (15 to 40 years): automation/unemployment, cooperation/arms races, control/principal agent problem, genie/oracle problem (unintended effects)
- **Long-term** (40 to 100 years): scientific breakthroughs by AIs, recursive self-improvement/runaway AI (?), artificial consciousness (?) (AIs as moral patients)

Exponential Improvements



Predicting AI Timelines

Great uncertainties:

- Hardware or software the bottleneck?
- Small team or a Manhattan Project?
- More speed bumps or accelerators?

| Probability for AGI | 10% | 50% | 90% |
|-----------------------|------|------|------|
| AI scientists, median | 2024 | 2050 | 2070 |

Speed Bumps

- Depletion of low-hanging fruit
- An end to Moore's law
- Societal collapse
- Disinclination



Accelerators

- Faster hardware
- Better algorithms
- Massive datasets



**+ enormous economic, military
and egoistic incentives!**



Strategy

What is to be done?

Prioritization

- **Scope:** How big/important is the issue?
- **Tractability:** What can be done about it?
- **Crowdedness:** Who else is working on it?

Work on the matters that matter the most!

- AI is the key lever on the long-term future
- Issue is urgent, tractable and uncrowded
- The stakes are astronomical: our light cone

Flow-Through Effects



Going meta: Solve the problem-solving problem!

Controlled Detonation



Difficulty:

Friendly AI >> General AI

Control Problem

Will AI outsmart us?



Capability Control

Boxing

Stunting

Tripwires



Motivation Selection

Direct Specification

Indirect Normativity

Incentive Methods

AI Safety Research

- Value alignment
- Corrigibility
- Security
- Verification
- Transparency
- Other things

Drop me a line at
me@kasparetter.com
if you're interested to
work on these problems!
(I will connect you with
researchers in the field.)

Differential Intellectual Progress

**Prioritize risk-reducing intellectual progress
over risk-increasing intellectual progress**

AI safety should outpace AI capability research

Who is working on this?



Future of Humanity Institute
UNIVERSITY OF OXFORD



International Cooperation

- We are the ones who will create superintelligent AI
- Not primarily a technical problem, rather a social
- International regulation?



In face of uncertainty, cooperation is robust!



2014: A turning point in AI safety!

Many smart people take superintelligence very seriously.

Future of Life Institute: AI Conference
futureoflife.org/misc/ai_conference

Implications of AI
Zurich, May 2016



With great power comes great responsibility!