

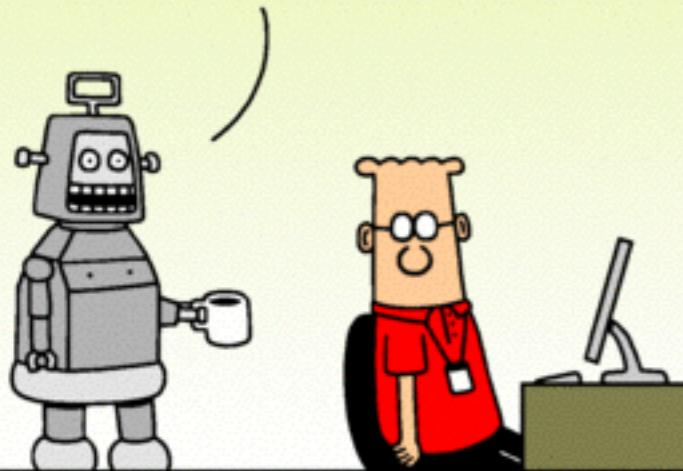
Superintelligence

Our Final Invention



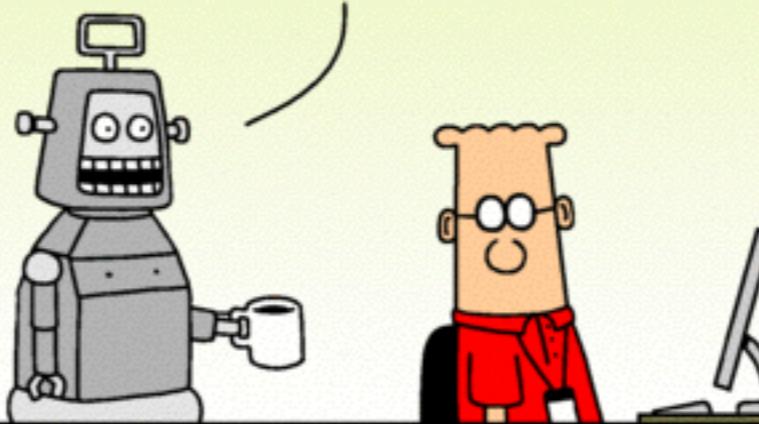
TEMPORARY ROBOT BOSS

I HAVE COME TO
MICROMANAGE YOU.



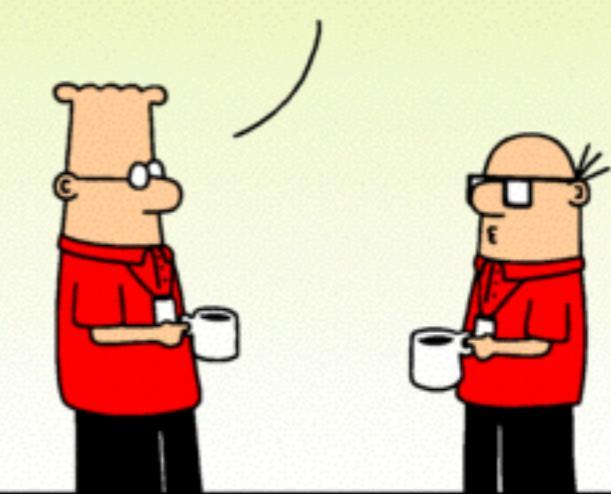
Dilbert.com DilbertCartoonist@gmail.com

BUT ONLY
UNTIL I REPLACE
YOU WITH A ROBOT
AND TURN YOU INTO
FURNITURE.



10-17-14 © 2014 Scott Adams, Inc./Dist. by Universal Uclick

ON THE PLUS SIDE,
HE HAS A PLAN AND
HE COMMUNICATES
WELL.



«Artificial Intelligence makes philosophy honest.»

– Daniel Dennett (2006), American Philosopher

Outline

- Introduction
- Singularity
- Superintelligence
- State and Trends
- Strategy
- Sources
- Summary





Introduction

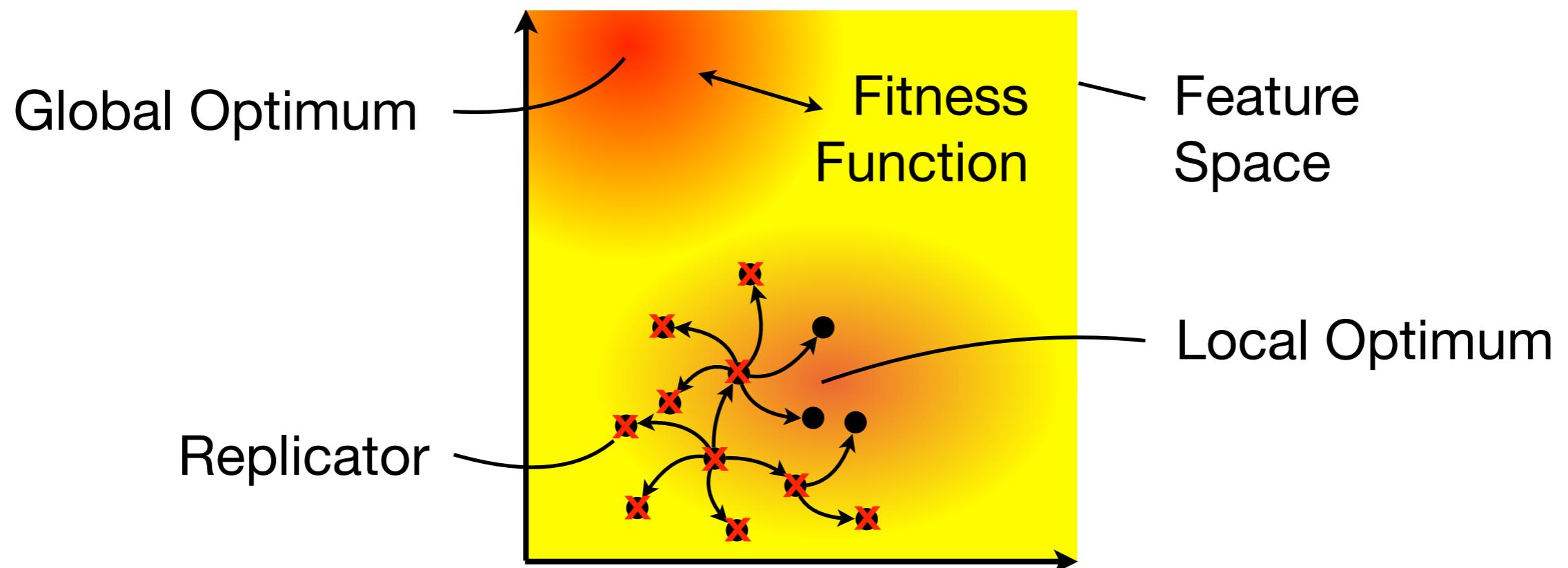
What are we talking about?

Crucial Consideration

- ... an idea or argument that entails a major change of direction or priority.
- Overlooking just one consideration, our best efforts might be for naught.
- When headed the wrong way, the last thing we need is progress.

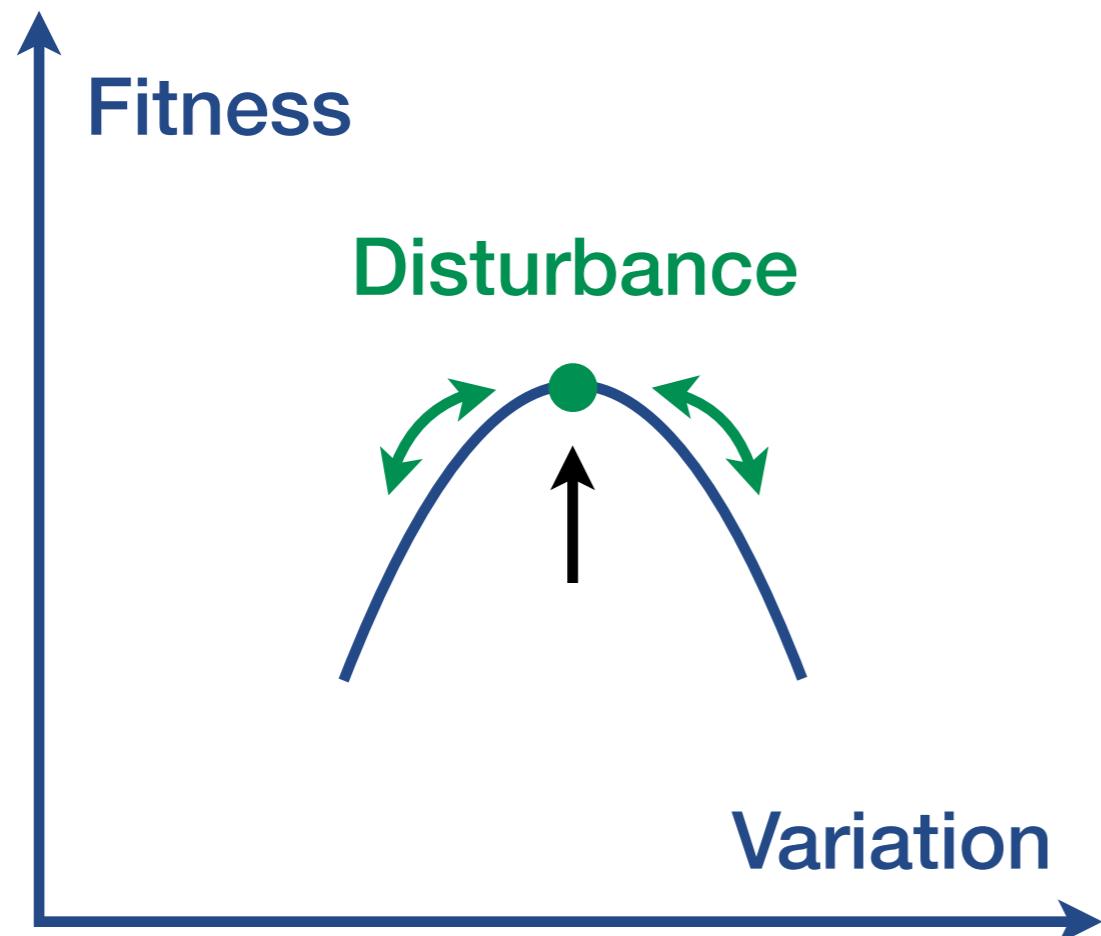


Evolution

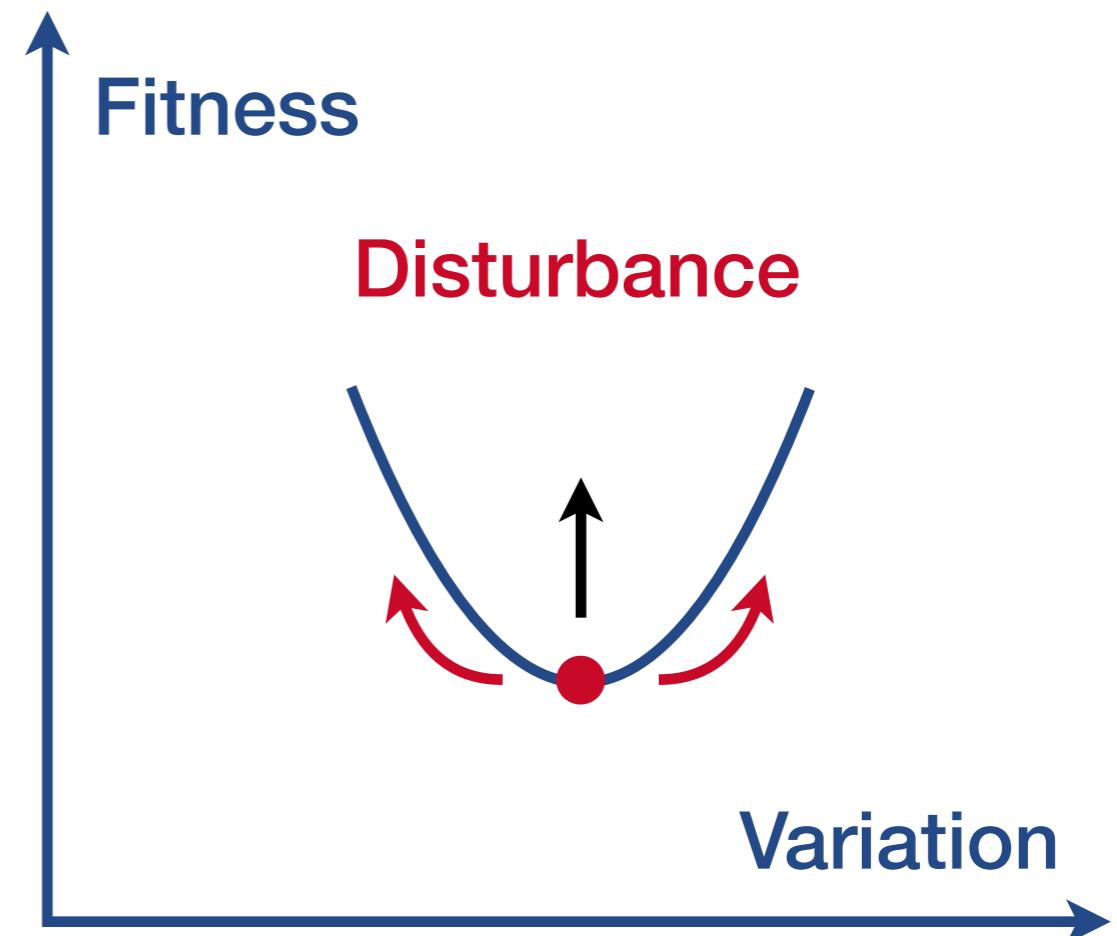


Foresight is the power of intelligence!

Stability



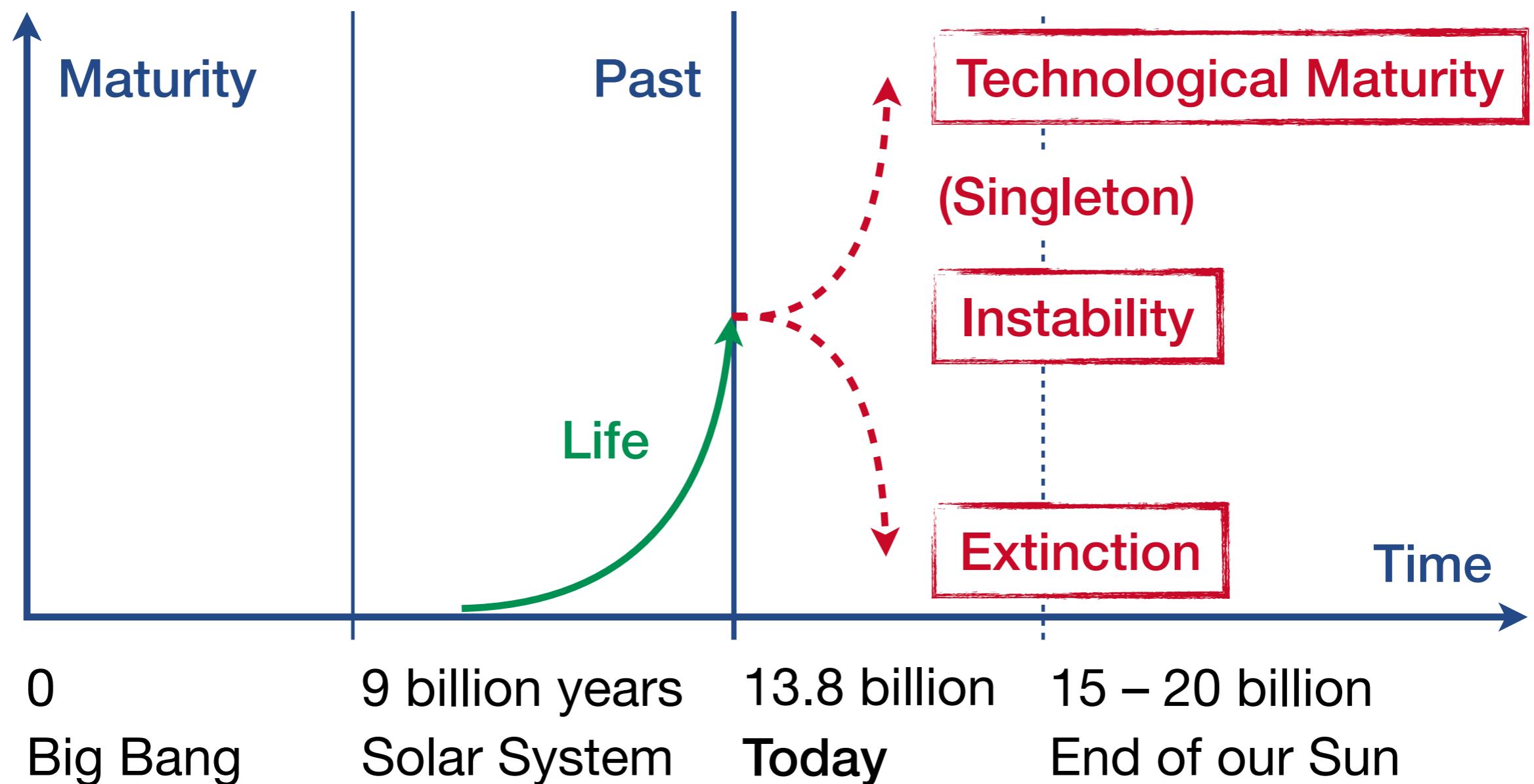
**stable states stay
(when time passes)**



**instable states vanish
(unless they are cyclic)**

Attractors

Big Rip: ≥ 20 billion years from now
Big Crunch: $\geq 10^2$ billion years from now
Big Freeze: $\geq 10^5$ billion years from now
Heat Death: $\sim 10^{1000}$ years from now



Singleton

... ultimate fate?

- World order with a single decision-making agency at the highest level
- Ability to prevent existential threats

Advantages:

It would avoid

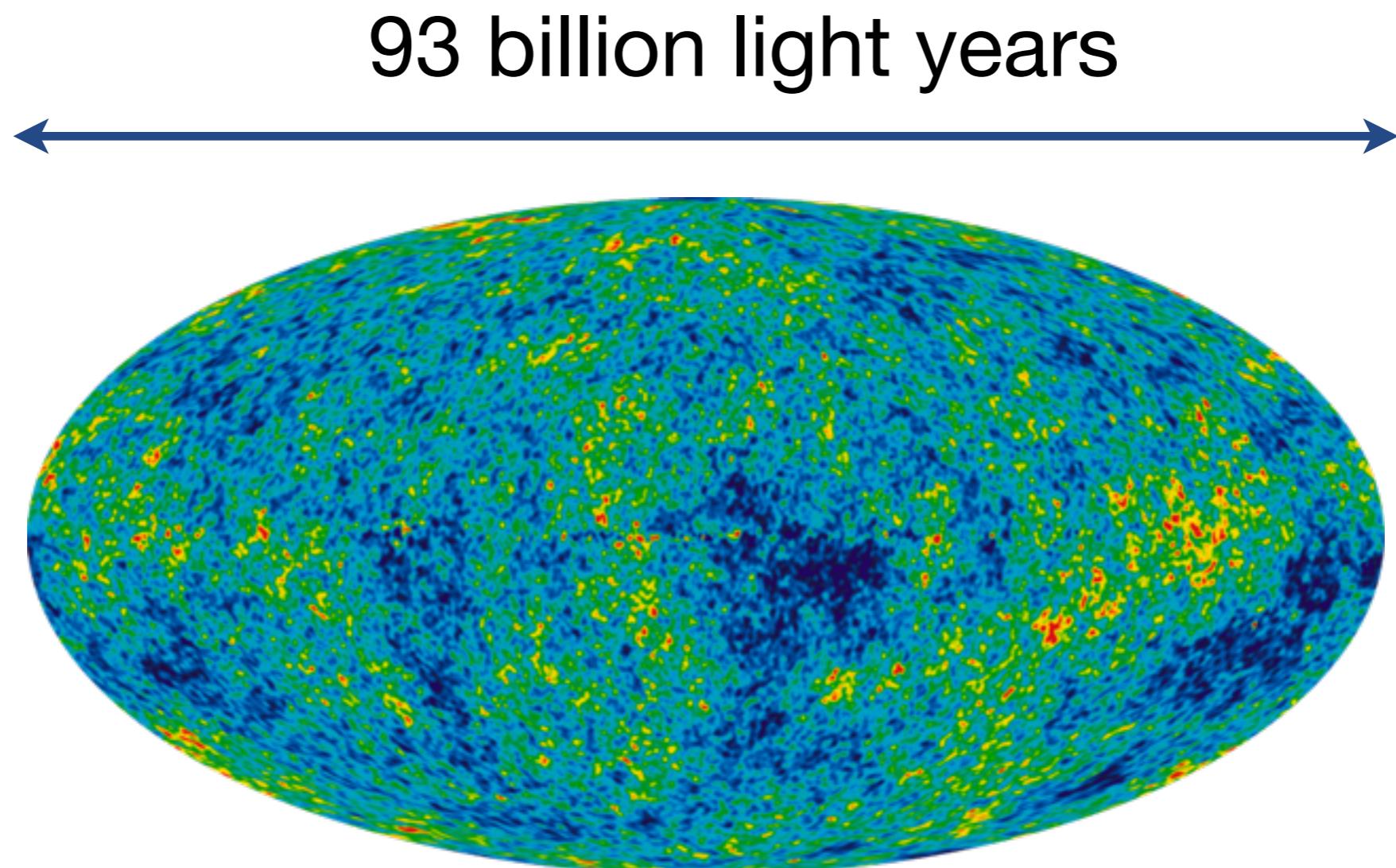
- arms races
- Darwinism

Disadvantages:

It might result in a

- dystopian world
- durable lock-in

The (Observable) Universe



$> 10^{11}$ galaxies
100000000000

$\sim 3 \cdot 10^{23}$ stars
3000000000000000000000000

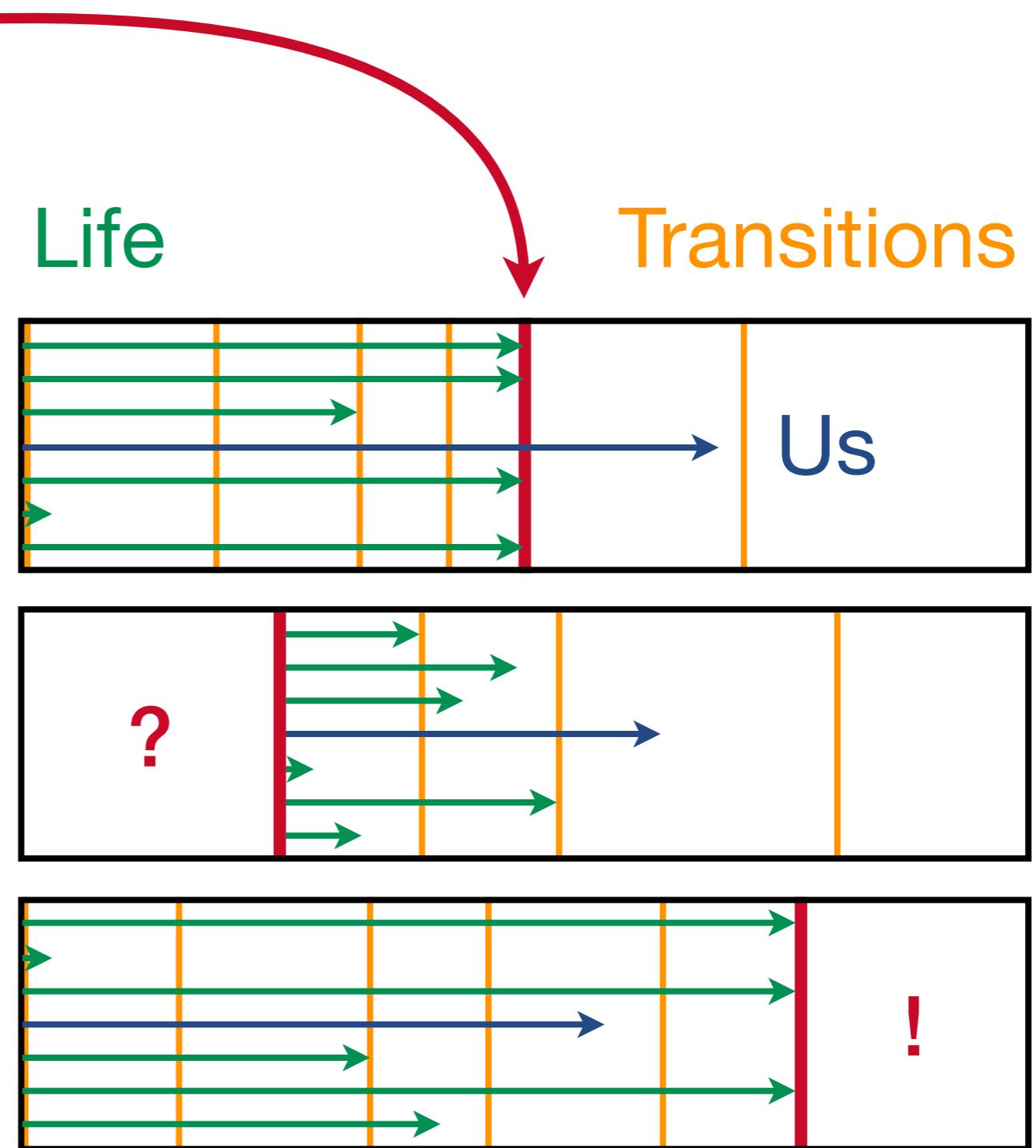
Fermi Paradox

Where are they? (Extraterrestrial Life)

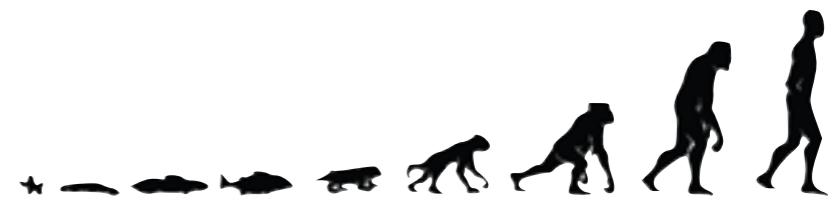
There are two groups of explanation:

- There are none, i.e. we're all alone.
- We can't detect them because ...
 - we're too primitive or too far apart
 - there are predators or all fear them
 - we're lied to, live in a simulation, ...

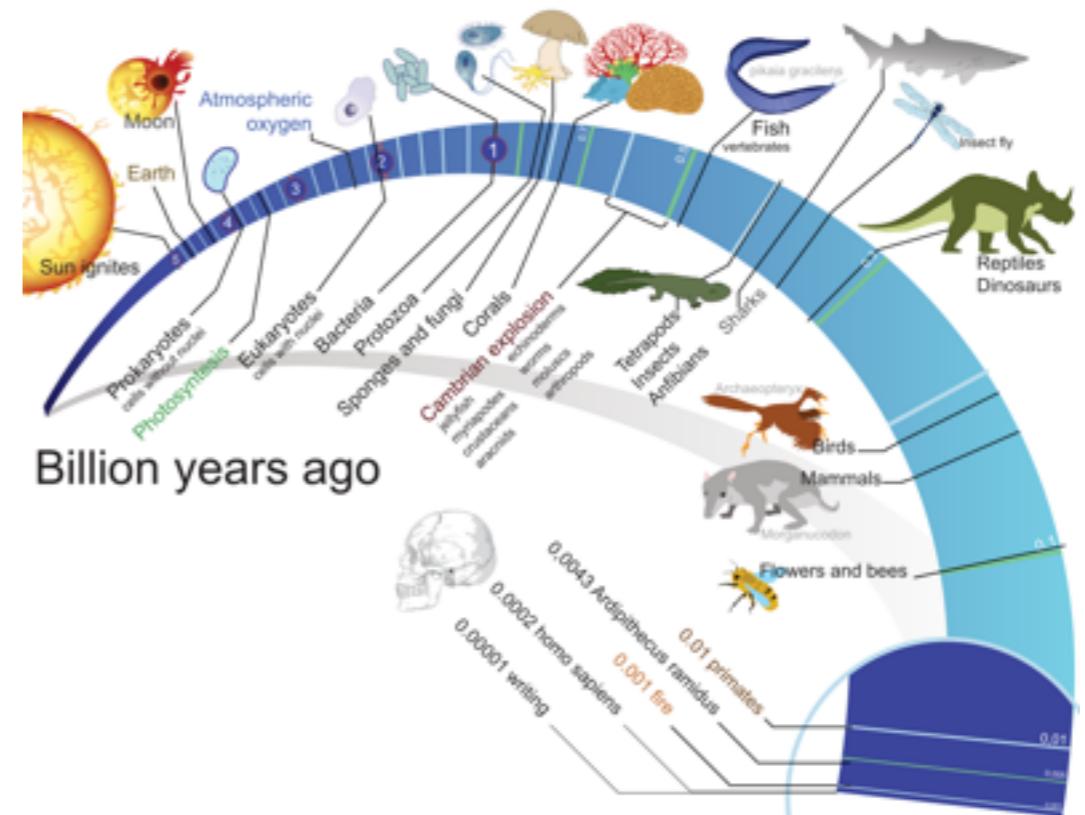
Great Filter



Major Transitions



- Self-replicating molecules (abiogenesis)
- Simple (prokaryotic) single-cell life
- Complex (eukaryotic) single-cell life
- Sexual reproduction
- Multi-cell organisms
- Tool-using animals
- Where we are now
- Space colonization



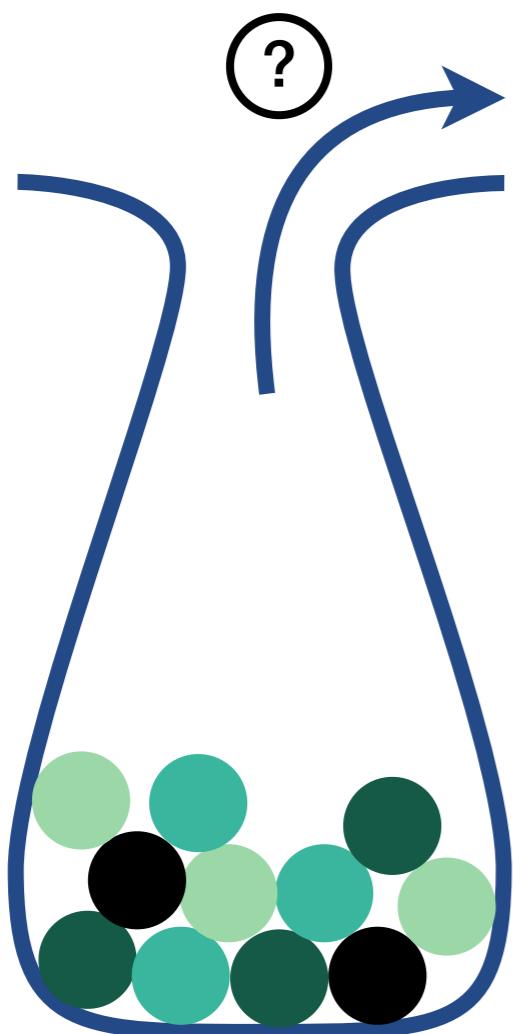
Anthropic Principle

- How probable are these transitions?
- They have occurred at least once
- Observation is conditional on existence

$P(\text{complex life on earth} \mid \text{our existence}) = 1$

There are observer selection effects!

Technologies



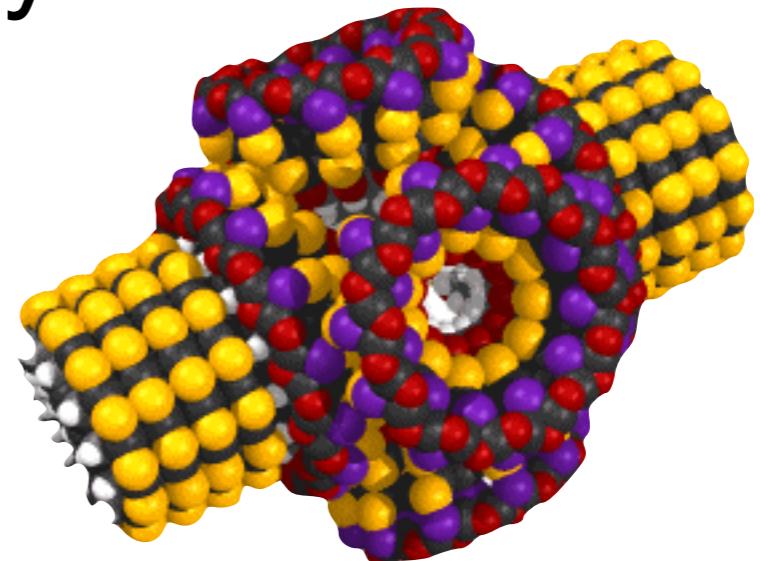
- Taking balls out of a jar
- No way to put back in
- Black balls are lethal

By definition:

- No ball has been black
- We'll only take out one

Candidates

- Nuclear Weapons (still possible)
- Synthetic Biology (engineered pathogens)
- Totalitarianism enabling Technologies
- Molecular Nanotechnology
- Machine Intelligence
- Geoengineering
- Unknown



Intelligence

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

(adapted from Legg and Hutter)

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

Ingredients

- Epistemology: Learn model of world
- Utility Function: Rate states of world
- Decision Theory: Plan optimal action

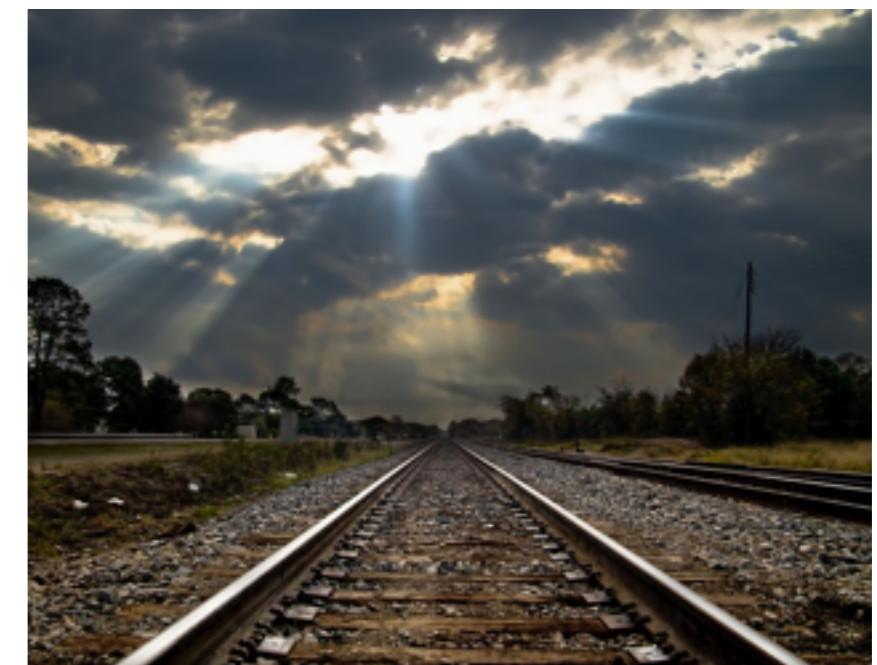
(There are still some open problems, e.g. classical decision theory breaks down when the algorithm itself becomes part of the game.)

Consciousness

- ... is a completely separate question!
- Not required for an agent to reshape the world according to its preference

Consciousness is

- reducible or
- fundamental
- and universal



Machine Sentience

Open questions of immense importance:

- Can simulated entities be conscious?
- Can machines be moral patients?

If yes:

- Machines deserve moral consideration
- We might live in a computer simulation



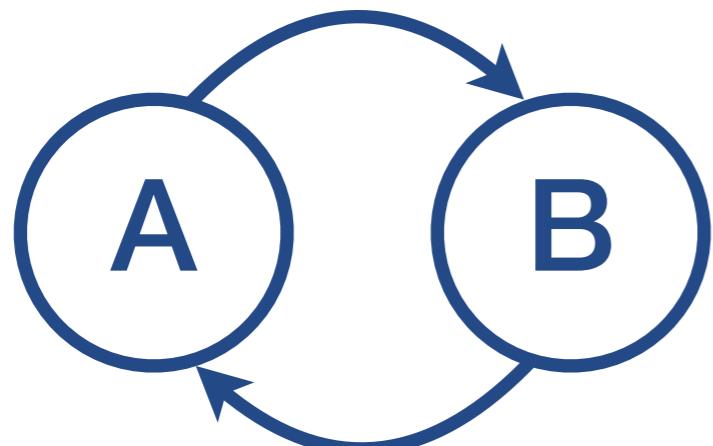
Singularity

What is the basic argument?

Feedback



Systems can feed back into themselves
and thus must be analyzed as a whole!



- Feedback is either:
- Positive (reinforcing)
 - Negative (balancing)

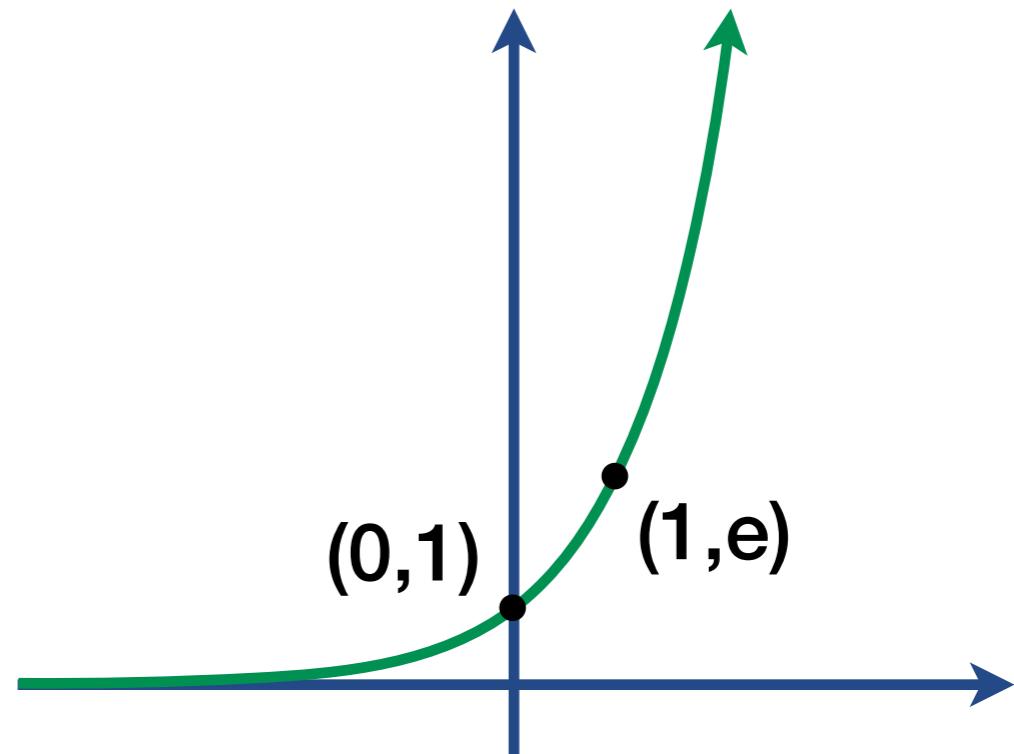
Exponential Functions

If increase is linear to current amount:

$$\frac{d}{dx} f(x) = c \cdot f(x)$$

solved by

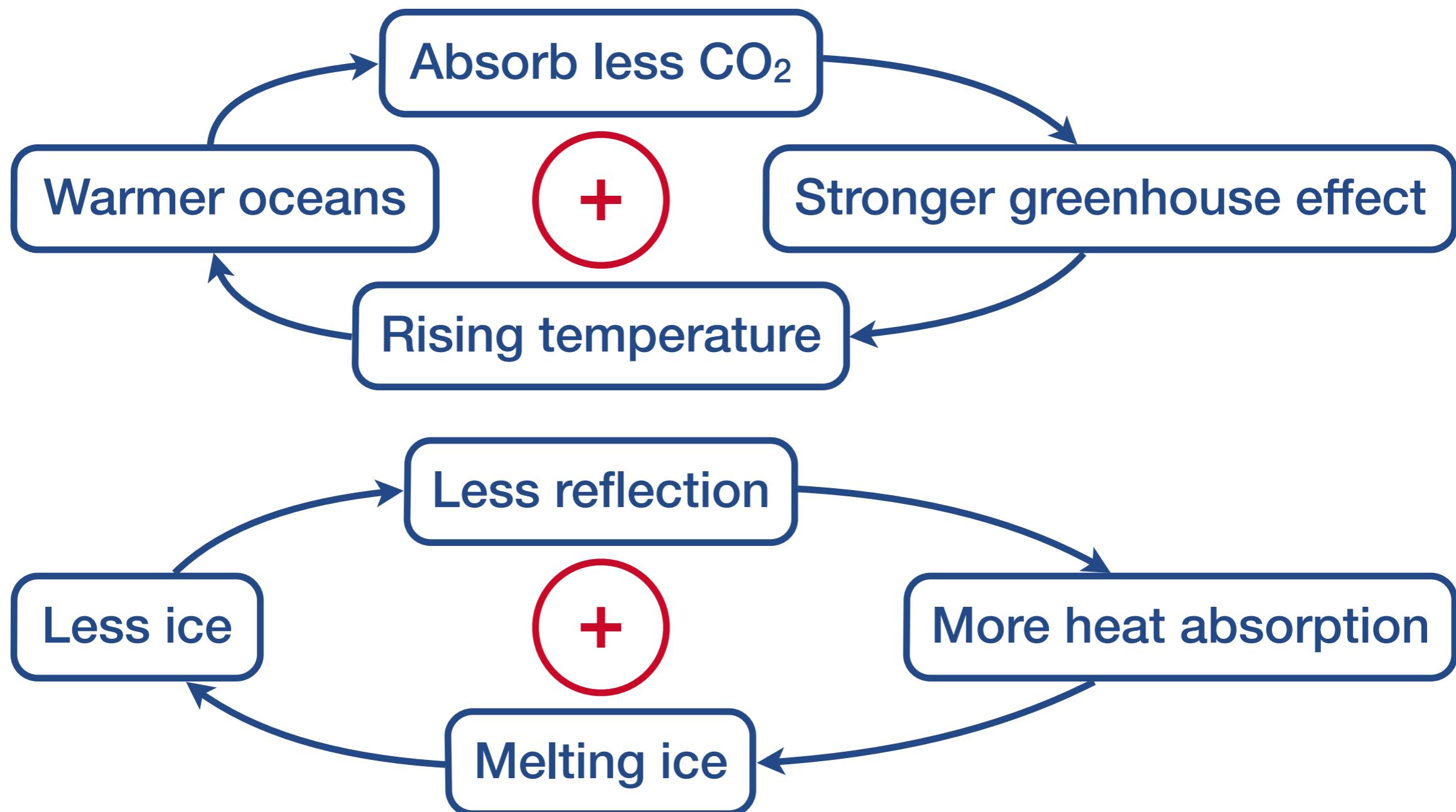
$$f(x) = e^{c \cdot x}$$



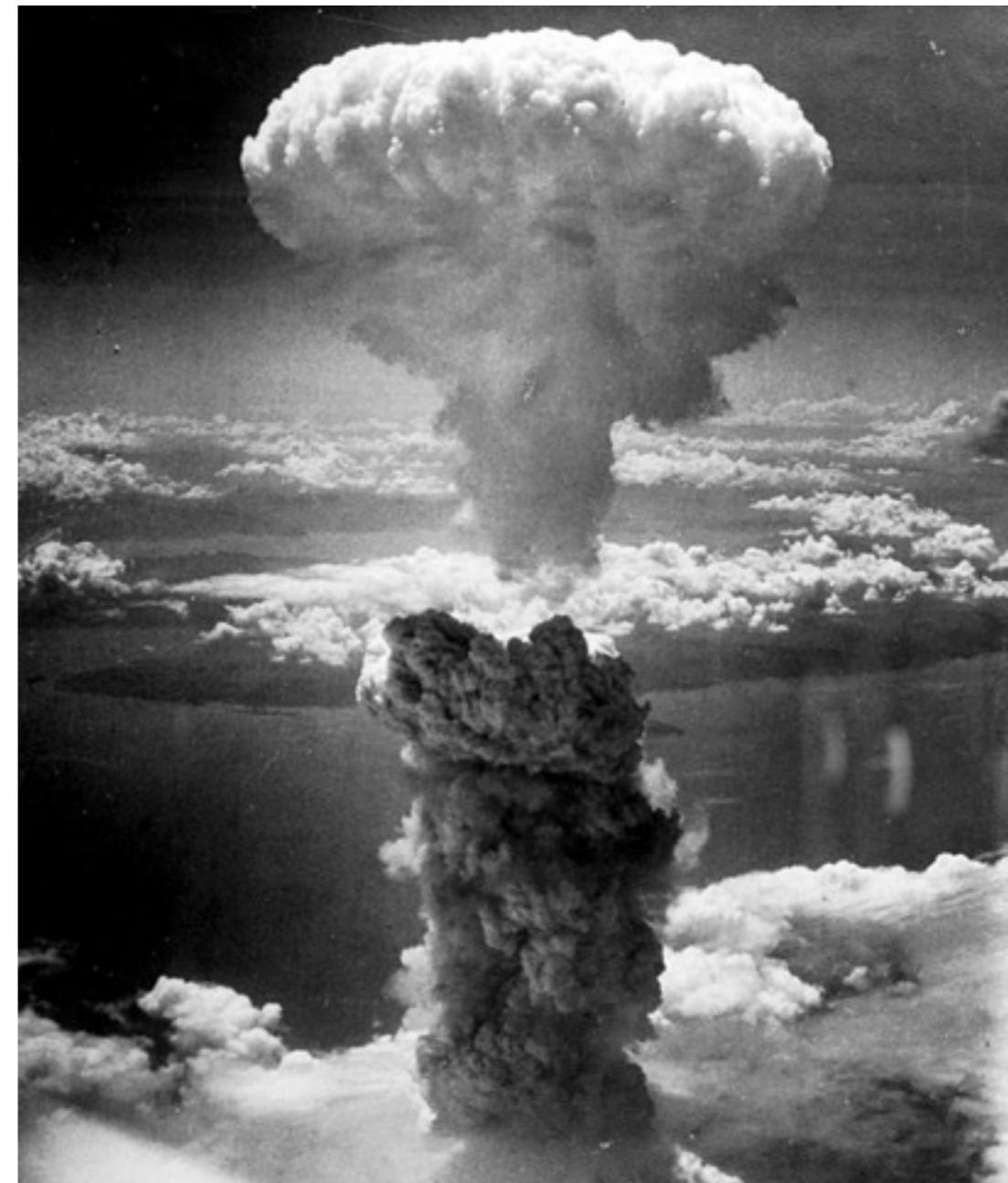
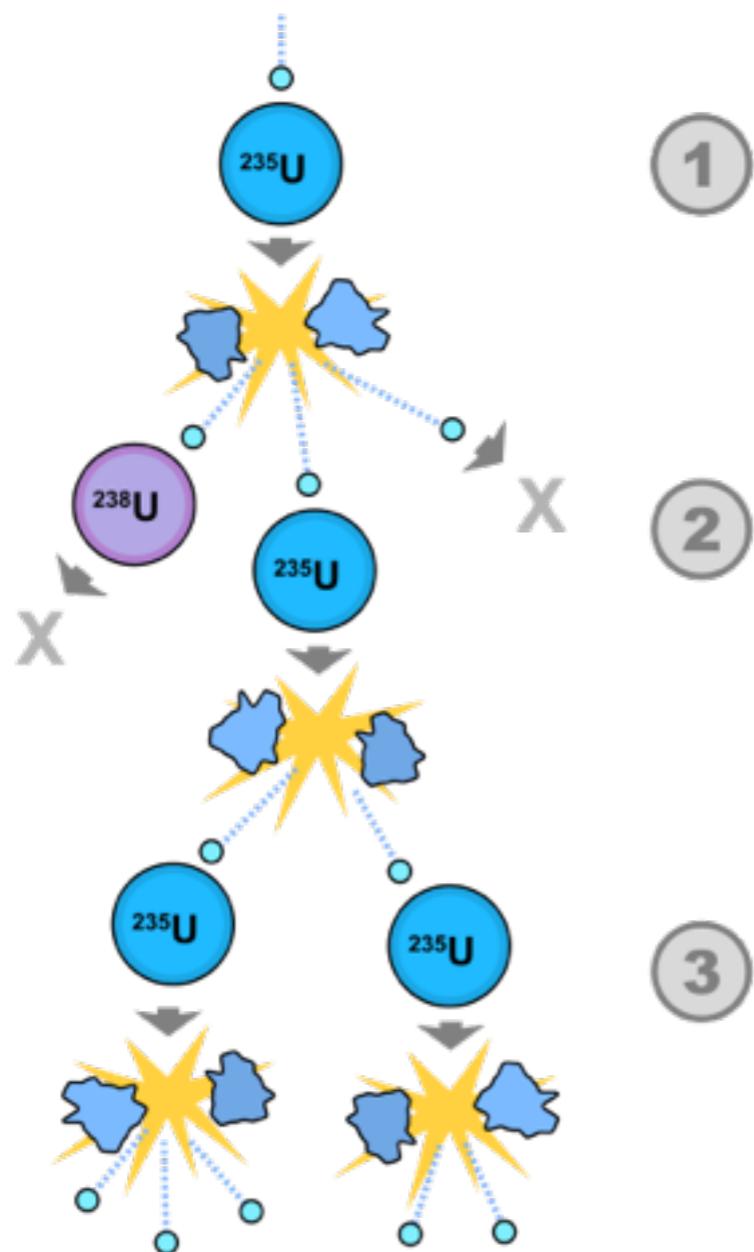
Fold a paper 45 times \Rightarrow to the moon!

How folding a paper can get you to the moon
www.youtube.com/watch?v=AmFMJC45f1Q

Climate Change



Nuclear Chain Reaction

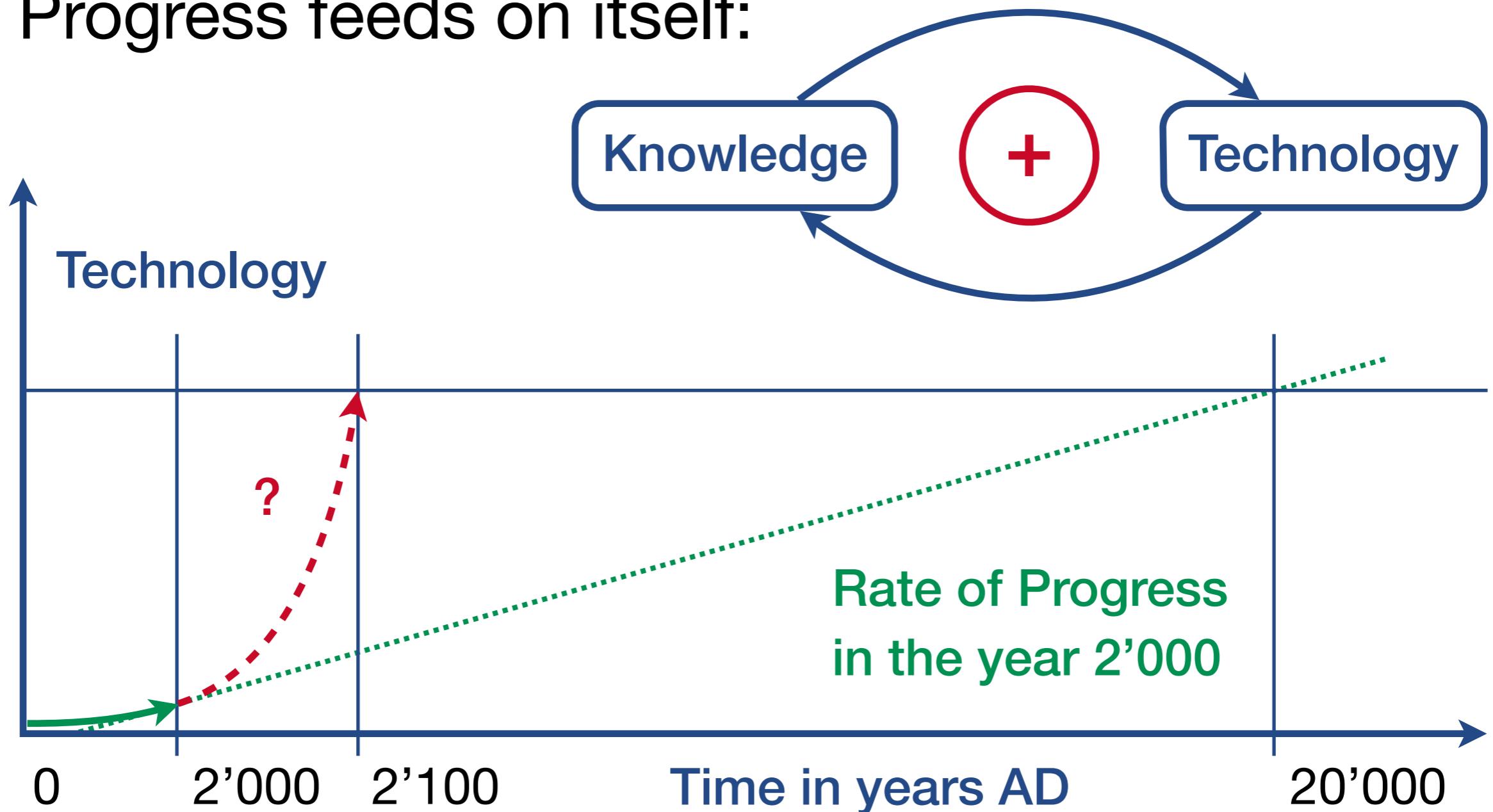


Nuclear Chain Reaction
en.wikipedia.org/wiki/Nuclear_chain_reaction

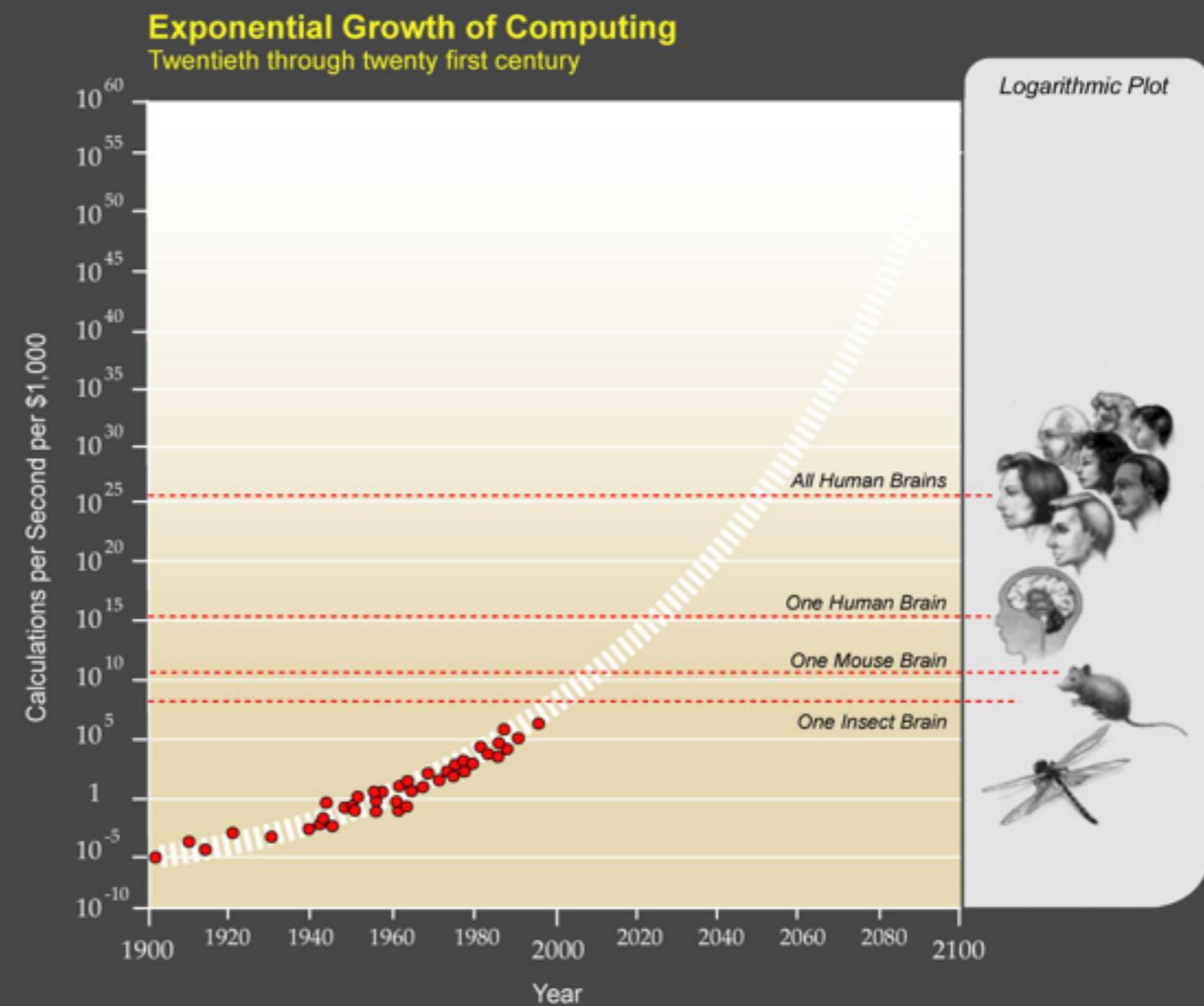
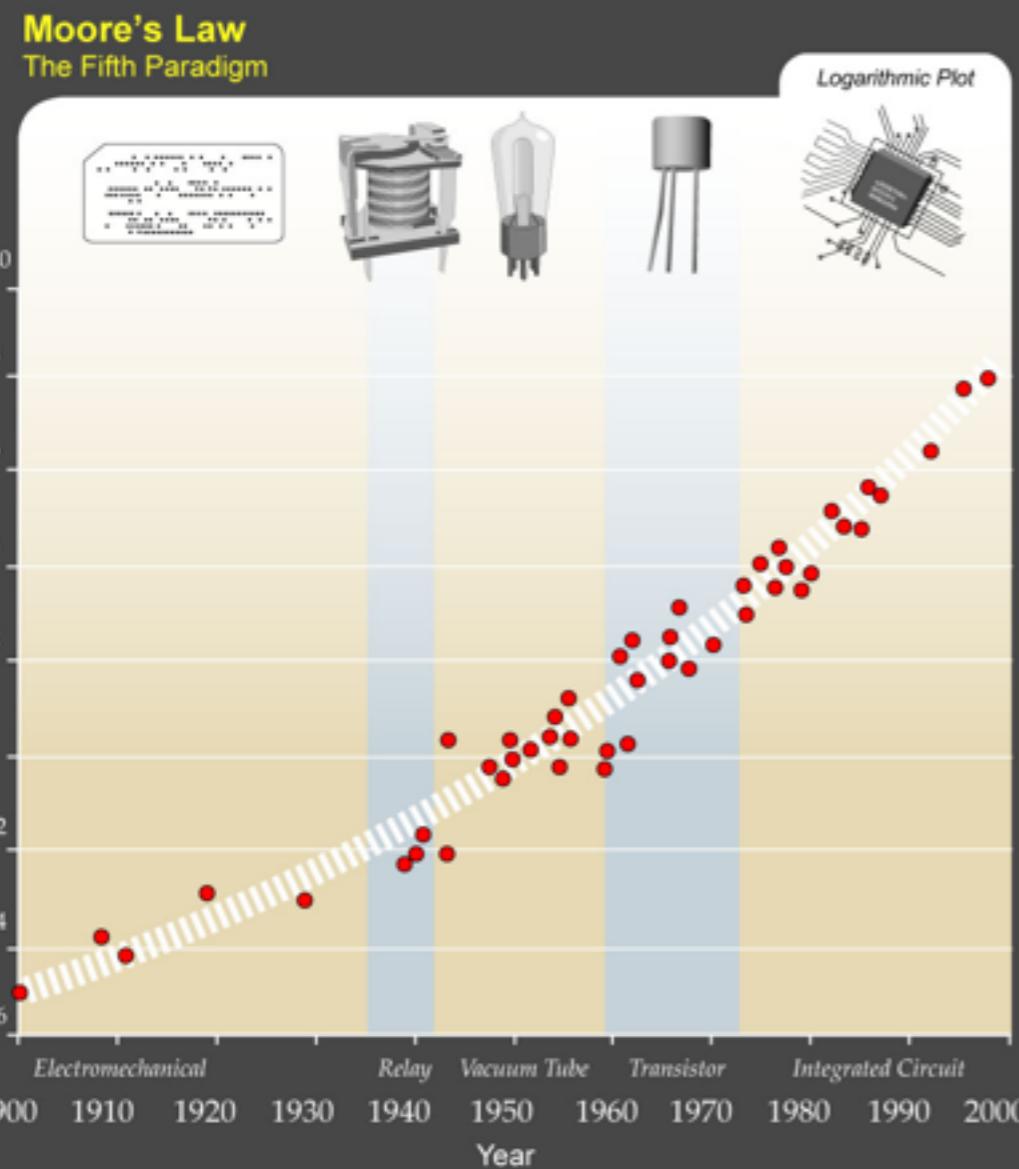
Superintelligence | 25
Singularity

Accelerating Change

Progress feeds on itself:



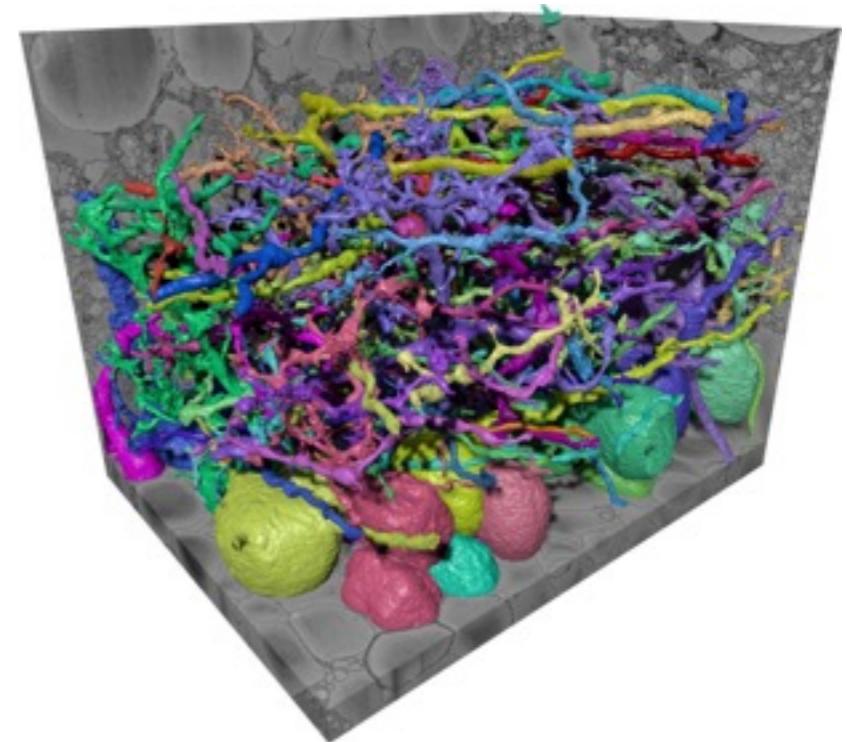
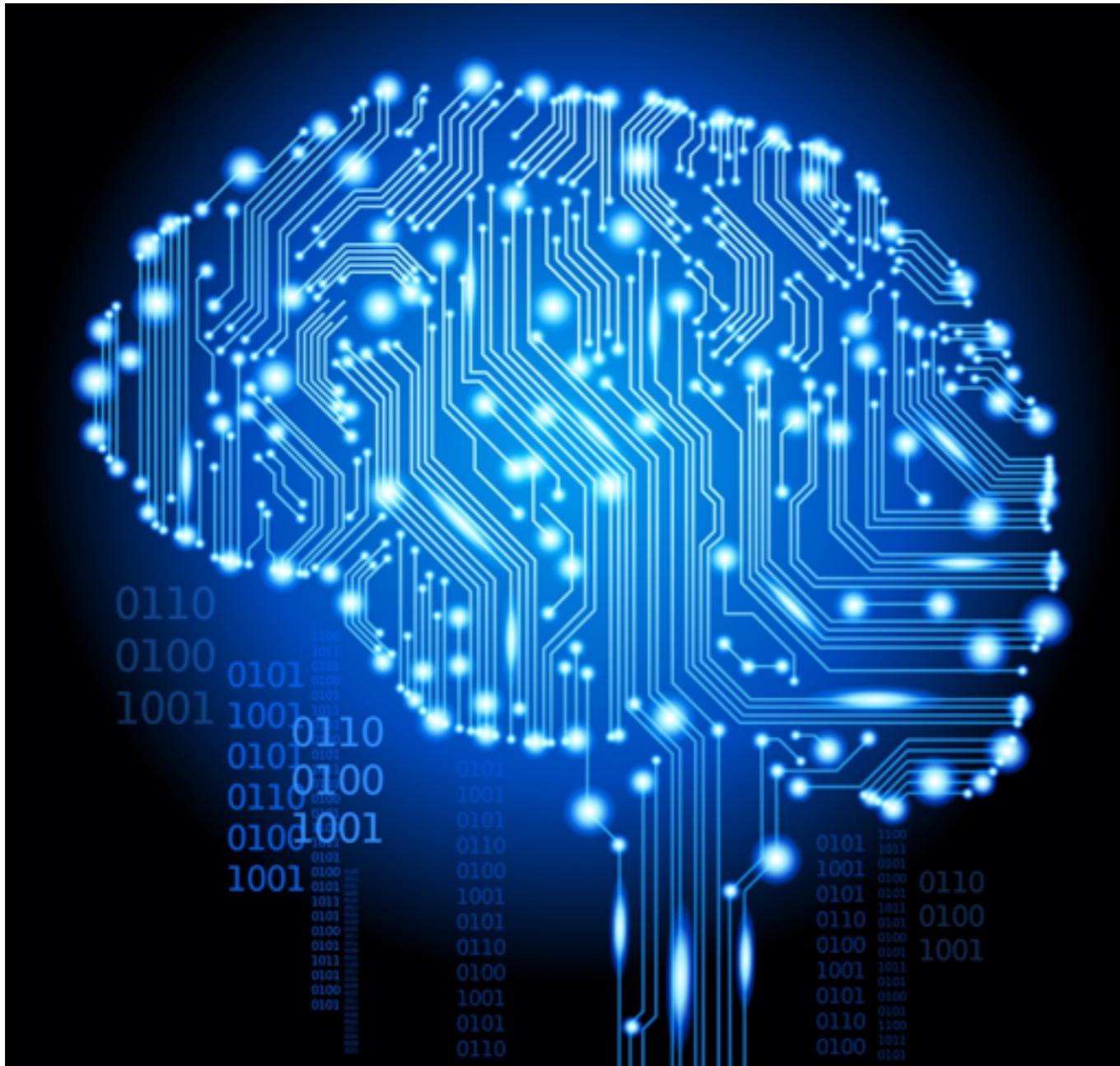
Moore's Law



Exponential and Non-Exponential Trends in IT
intelligence.org/.../exponential-and-non-exponential/

Superintelligence | 27
Singularity

Artificial Mind



Imagine all relevant aspects captured in a computer model (thought experiment)

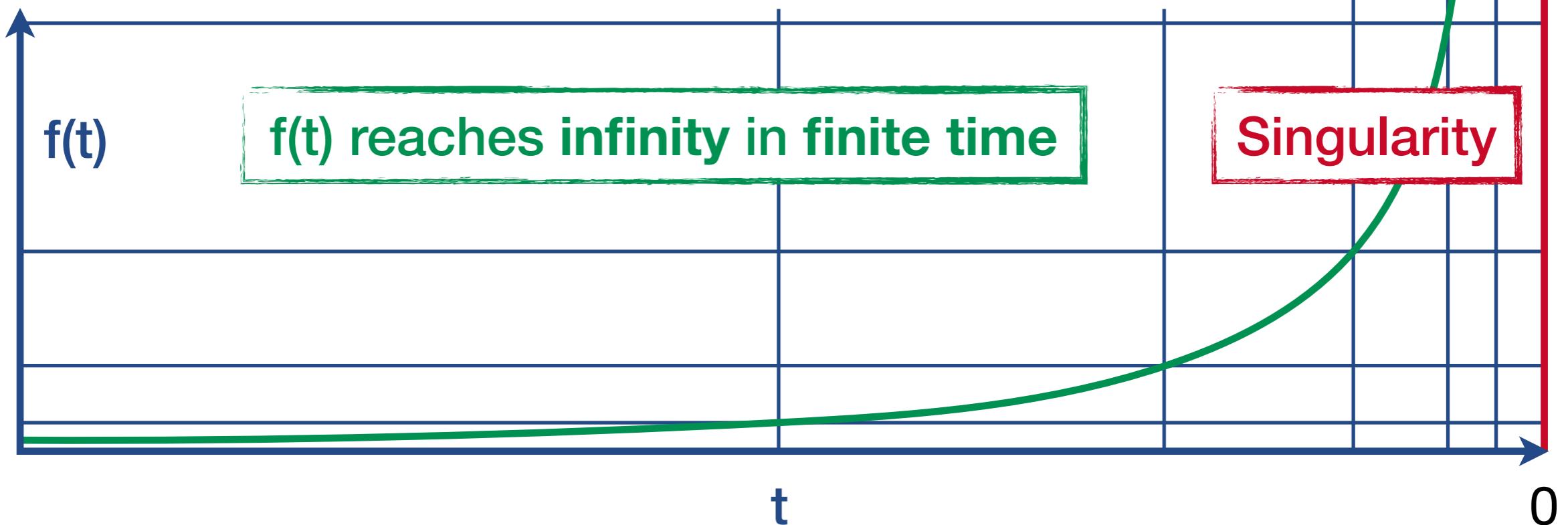
Whole Brain Emulation: A Roadmap
www.fhi.ox.ac.uk/brain-emulation-roadmap-report.pdf

Superintelligence Singularity | 28

Hyperbolic Growth

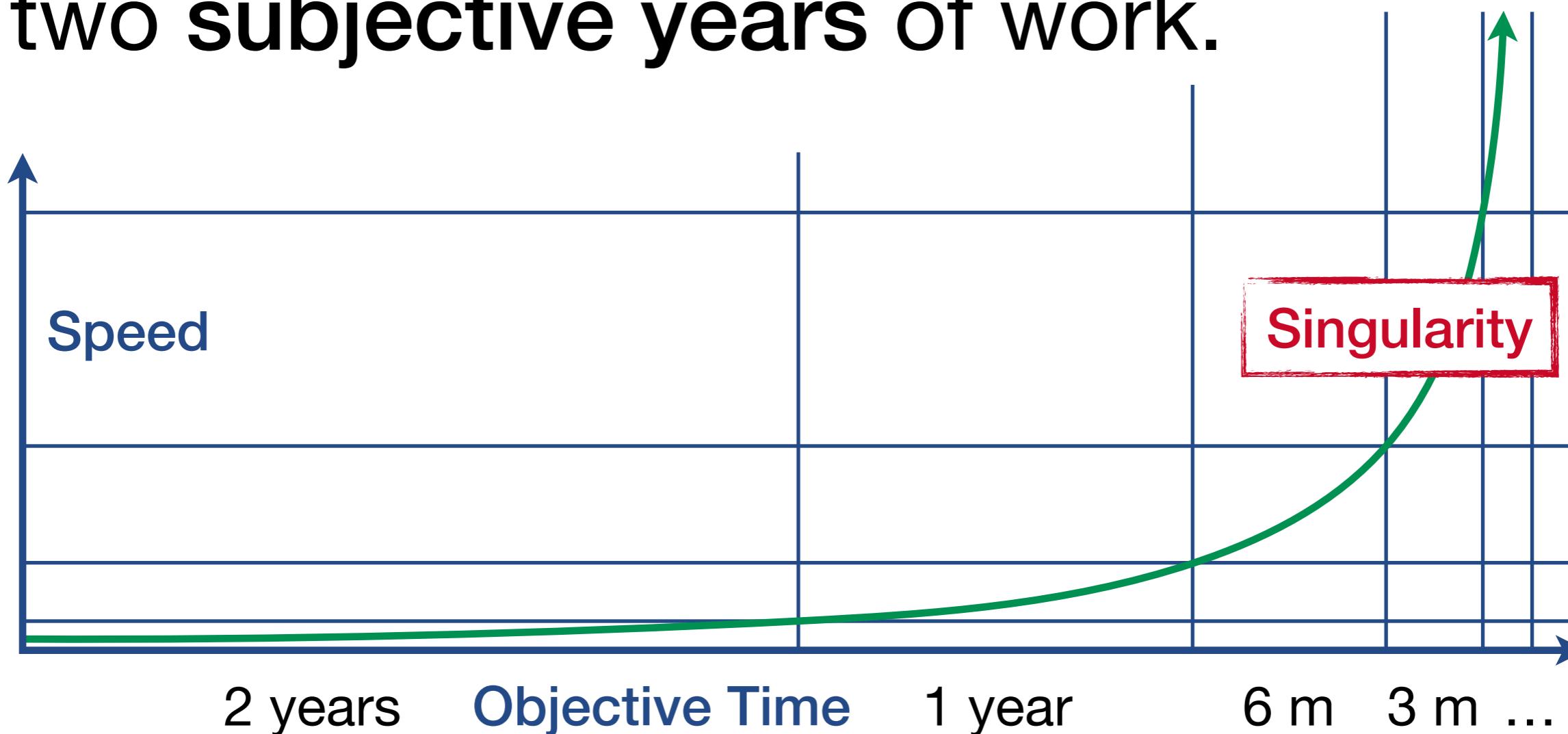
Second-order positive feedback loop

$$\frac{d}{dt} f(t) = c \cdot f(t)^2 \Rightarrow f(t) = \frac{-1}{c \cdot t}$$



Speed Explosion

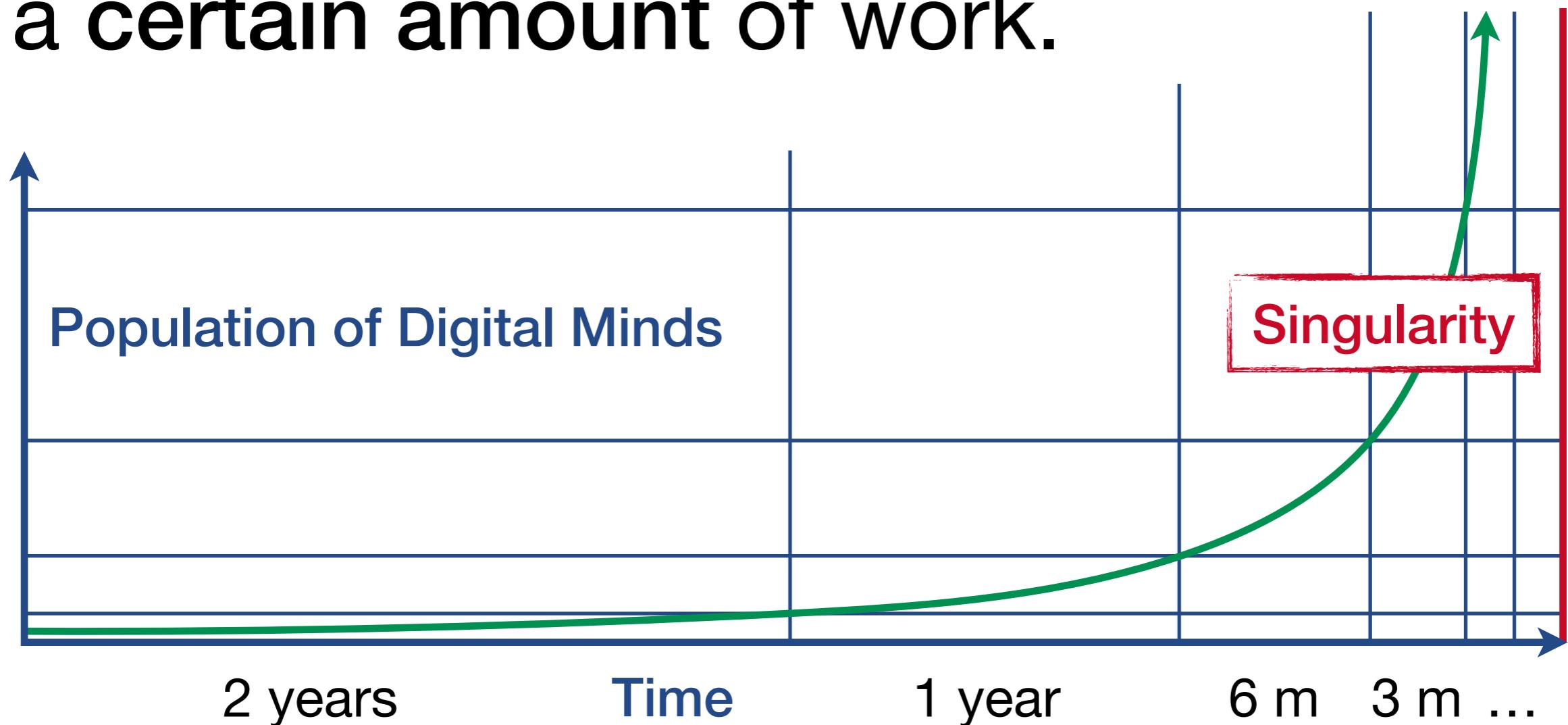
Computing speed doubles every two subjective years of work.



Population Explosion

Quantitative

Computing costs halve for a certain amount of work.

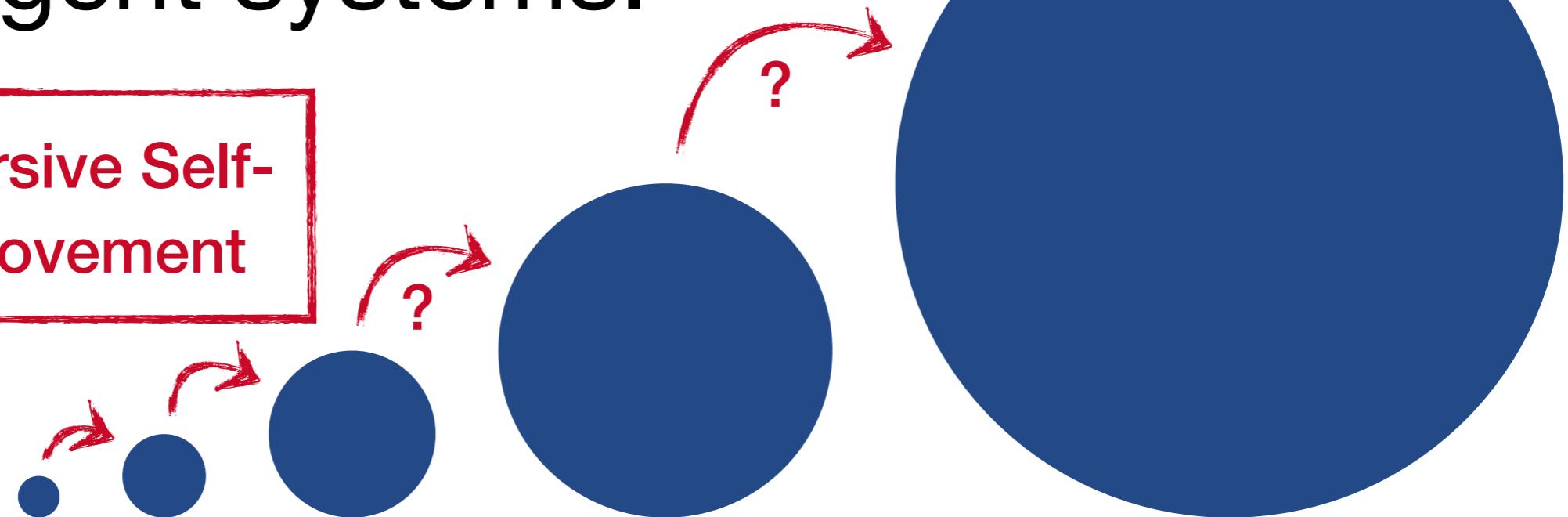


Intelligence Explosion

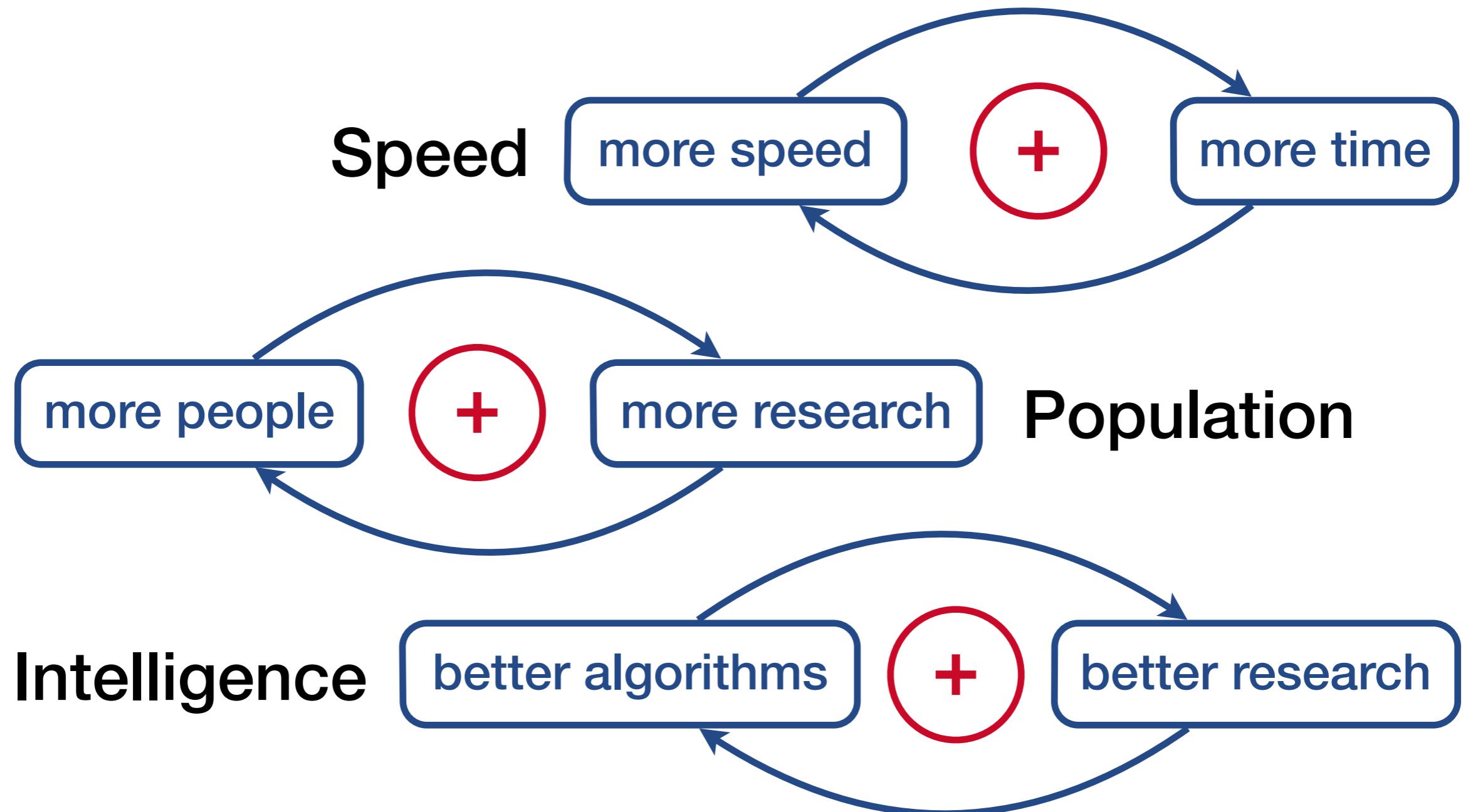
Qualitative

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

Recursive Self-Improvement



Three Separate Explosions



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)

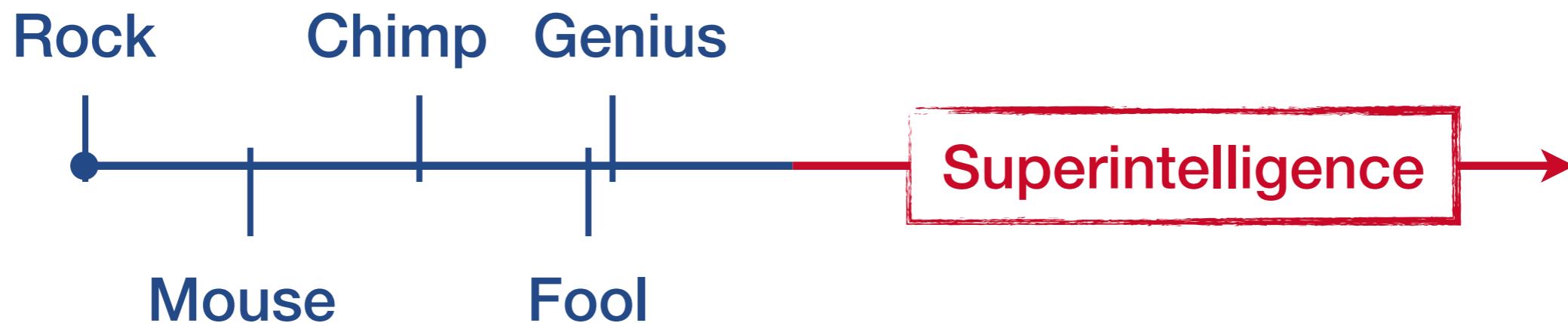


Superintelligence

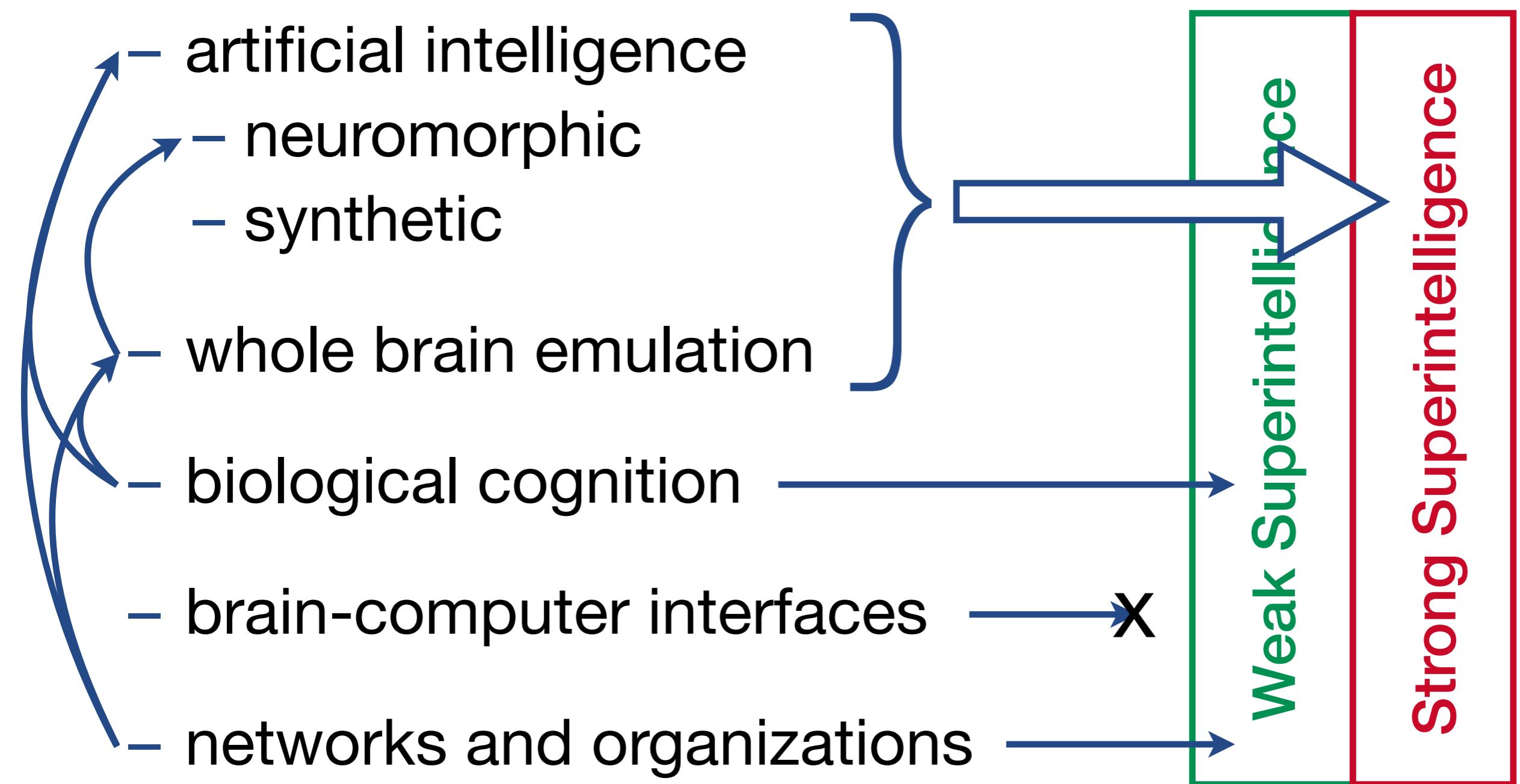
What are potential outcomes?

Definition of Superintelligence

An agent is called **superintelligent** if it exceeds the level of current human intelligence in all areas of interest.



Pathways to Superintelligence



Advantages of AIs over Brains

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

Human Brain

86 billion neurons

firing rate of 200 Hz

120 m/s signal speed

Modern Microprocessor

1.4 billion transistors

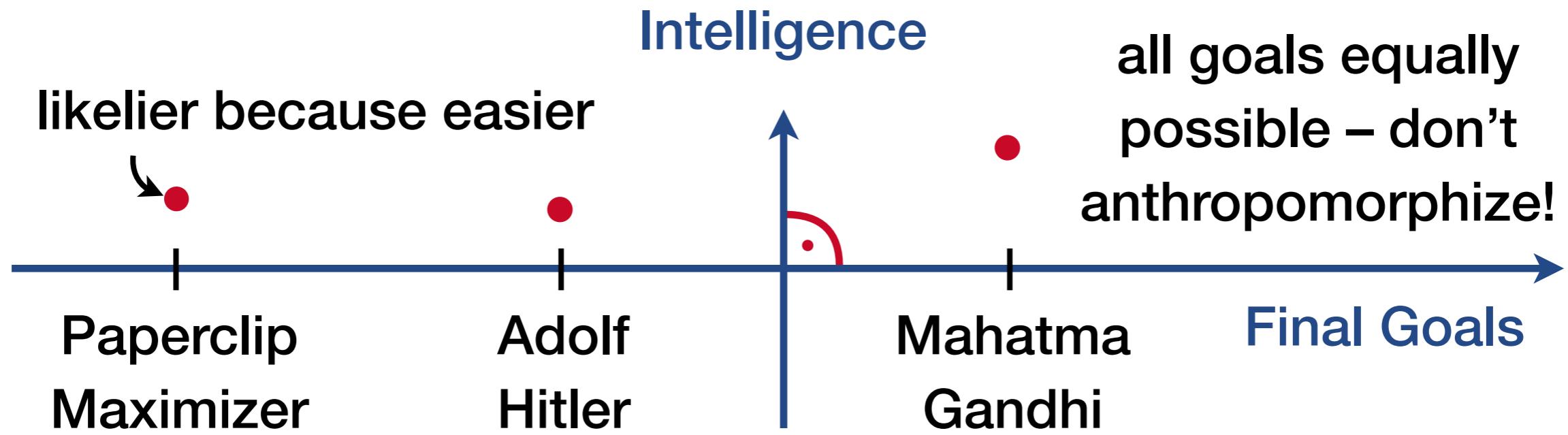
4'400'000'000 Hz

300'000'000 m/s

Cognitive Superpowers

- Intelligence amplification: bootstrapping
- Strategizing: overcome smart opposition
- Hacking: hijack computing infrastructure
- Social manipulation: persuading people
- Economic productivity: acquiring wealth
- Technology research: inventing new aids

Orthogonality Thesis



Intelligence and final goals are orthogonal:
Almost any level of intelligence could in principle be combined with any final goal.

Convergent Instrumental Goals

- Self-Preservation
 - Goal-Preservation
 - Resource Accumulation
 - Intelligence Accumulation
- } necessary to achieve goal
- } to achieve goal better

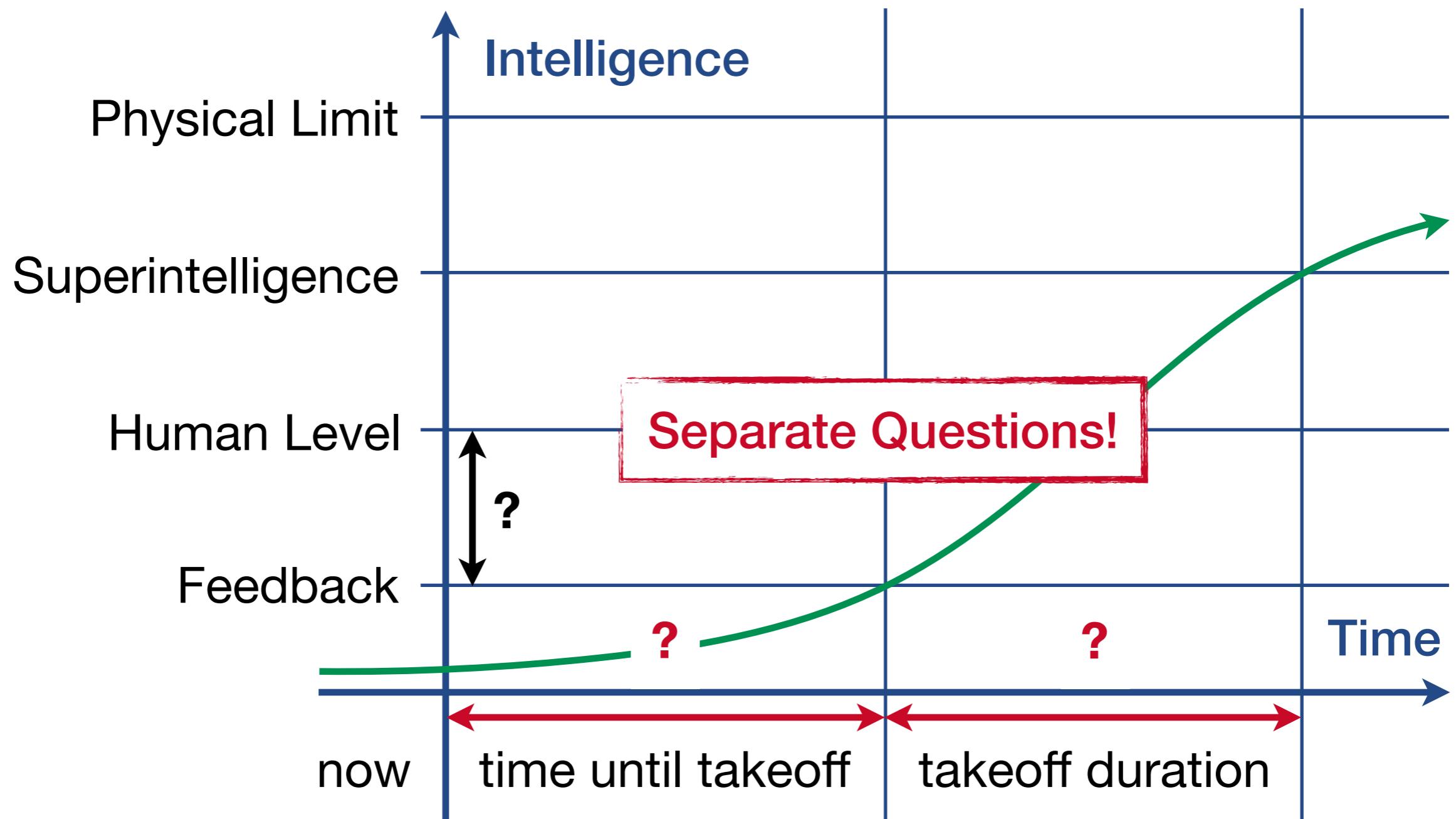
**Default Outcome: Doom
(Infrastructure Profusion)**

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

Takeoff Scenarios



Potential Outcomes

Fast Takeoff

hours, days, weeks

Unipolar Outcome

Singleton
(Slide 9)

Slow Takeoff

several months, years

Multipolar Outcome

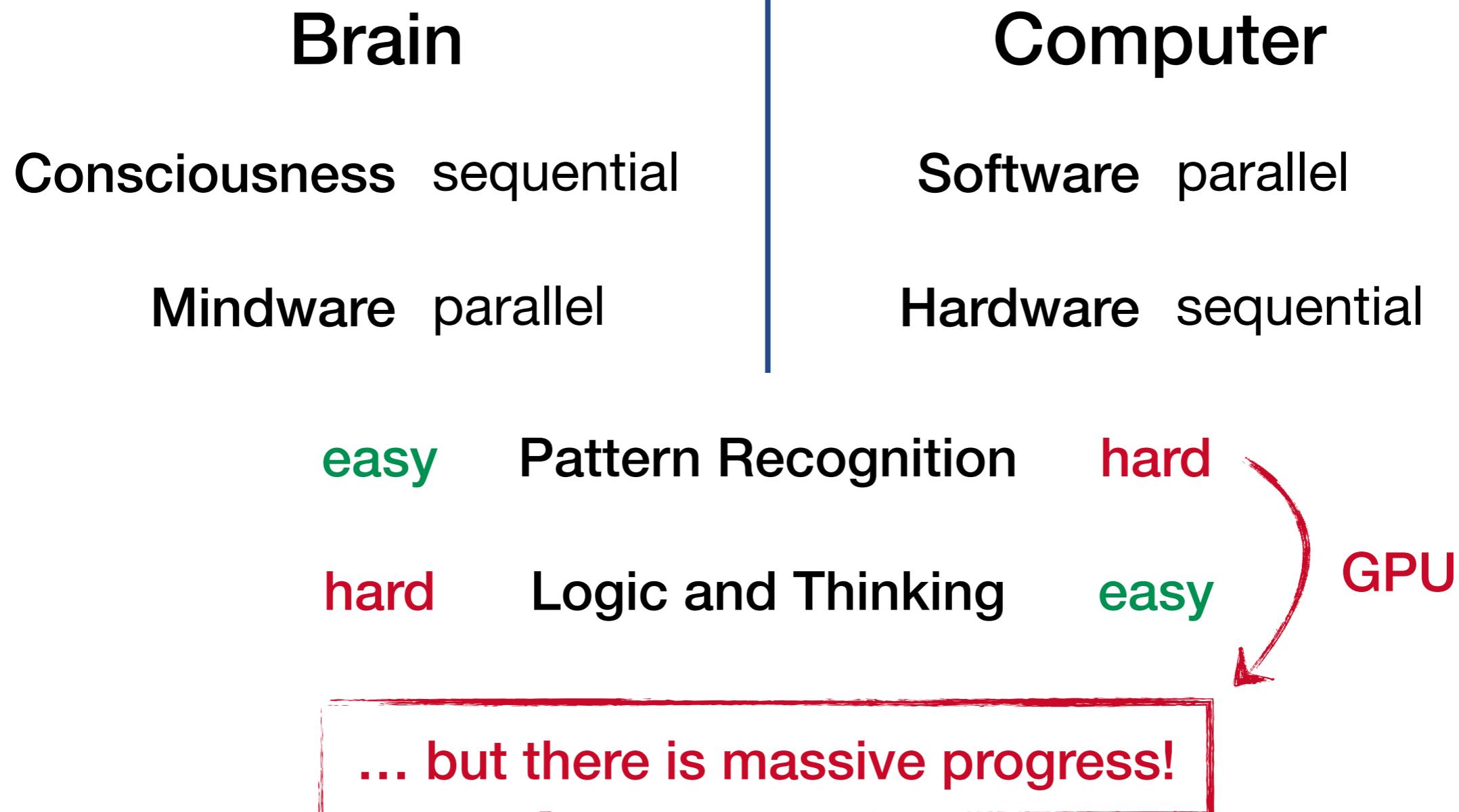
Second Transition
Unification by Treaty



State and Trends

Where are we heading to?

Brain vs. Computer



State of the Art



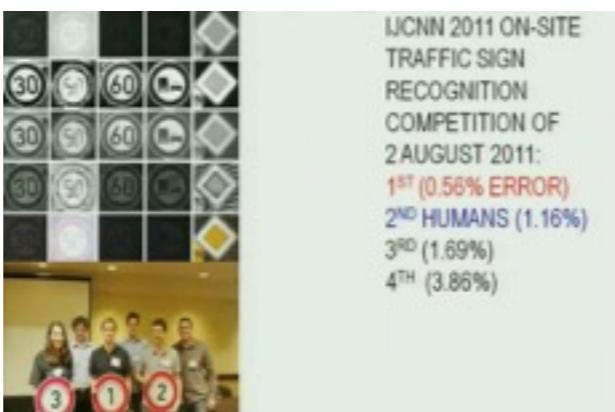
Deep Blue: 1997



IBM Watson: 2011



Stanley: 2005



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	Strong Amateur

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

Superintelligence
State and Trends

Consumer Products



Siri



INTRODUCING
amazon echo

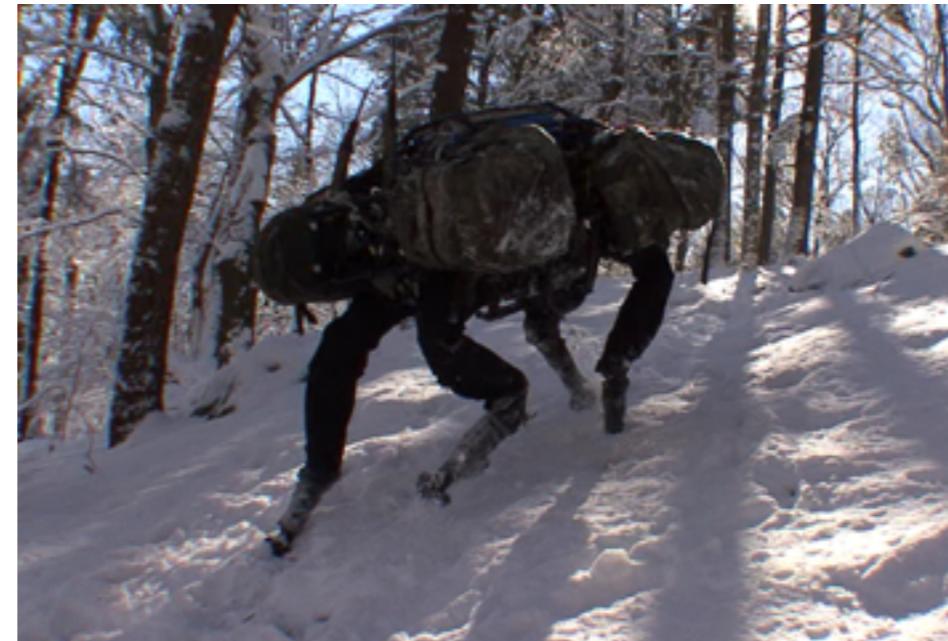
Raffaello D'Andrea
go.ted.com/xeh



jibo

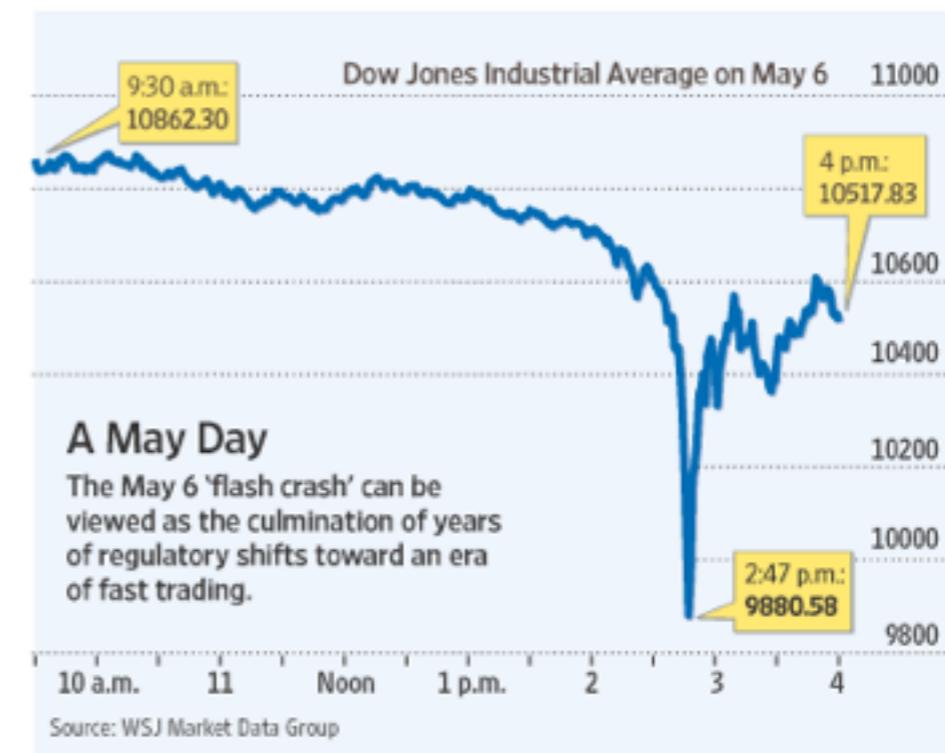
Superintelligence
State and Trends

Military Robots

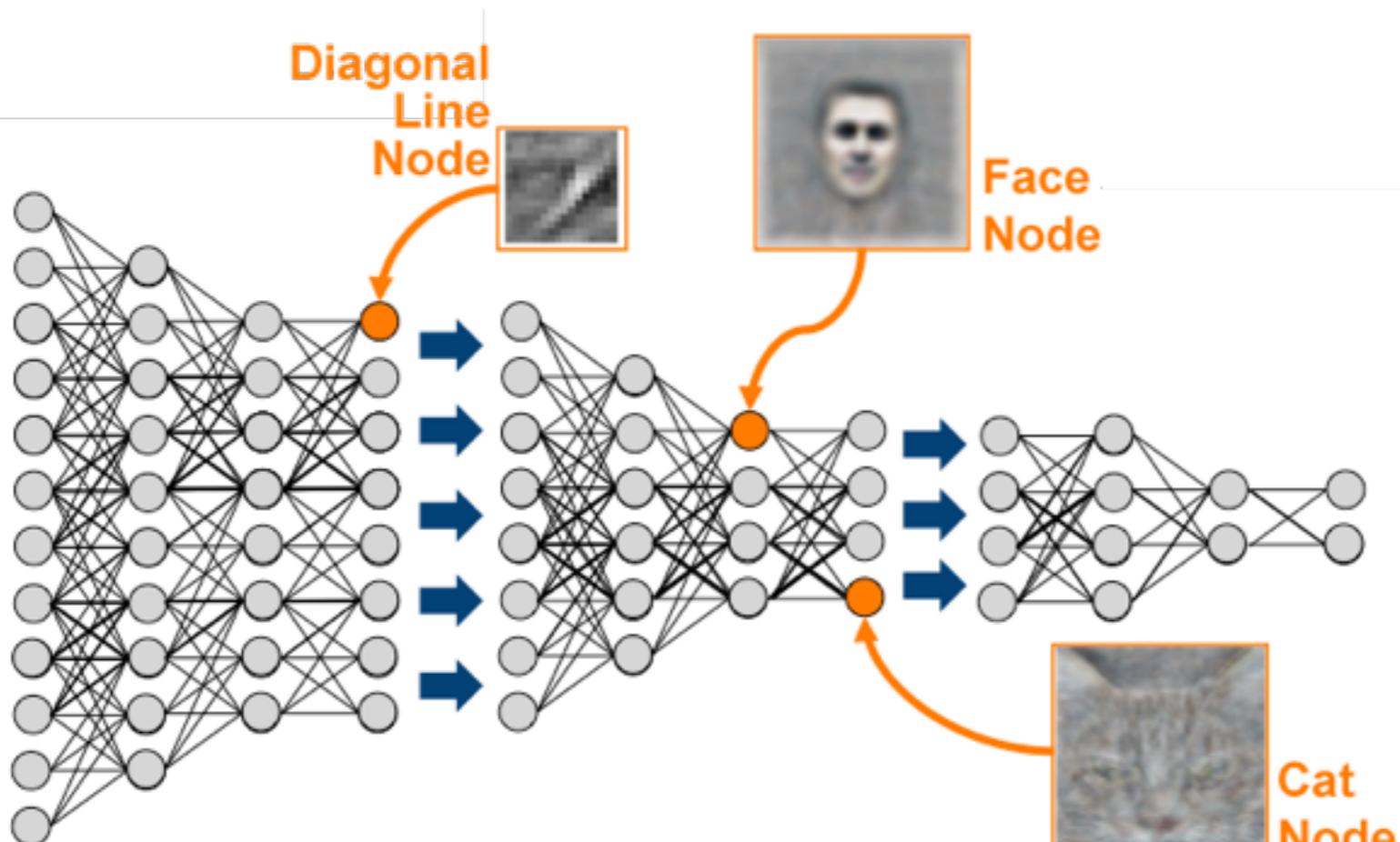


Financial Markets

- High-frequency trading (HFT): Buy and sell securities within milliseconds algorithmically
- In 2009, 65% of all US equity trading volume
- Flash crash: very rapid fall in security prices
- 6 May 2010: Dow Jones lost \$1 trillion (over 9%)
- 23 April 2013: One tweet causes \$136 billion loss



Machine Learning



Vicarious AI passes first Turing Test: CAPTCHA
news.vicarious.com/...-ai-passes-first-turing-test



Superintelligence
State and Trends

Universal Artificial Intelligence

$$a_k := \arg \max_{a_k} \sum_{o_k r_k} \dots \max_{a_m} \sum_{o_m r_m} (r_k + \dots + r_m) \sum_{q : U(q, a_1..a_m) = o_1 r_1 .. o_m r_m} 2^{-l(q)}$$

- AIXI by Marcus Hutter at IDSIA in Manno
- AIXI is a universally optimal rational agent
- AIXI uses Solomonoff induction and EUT
- AIXI is gold standard but not computable

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
Luke Muelhauser, MIRI	2030	2070	2140

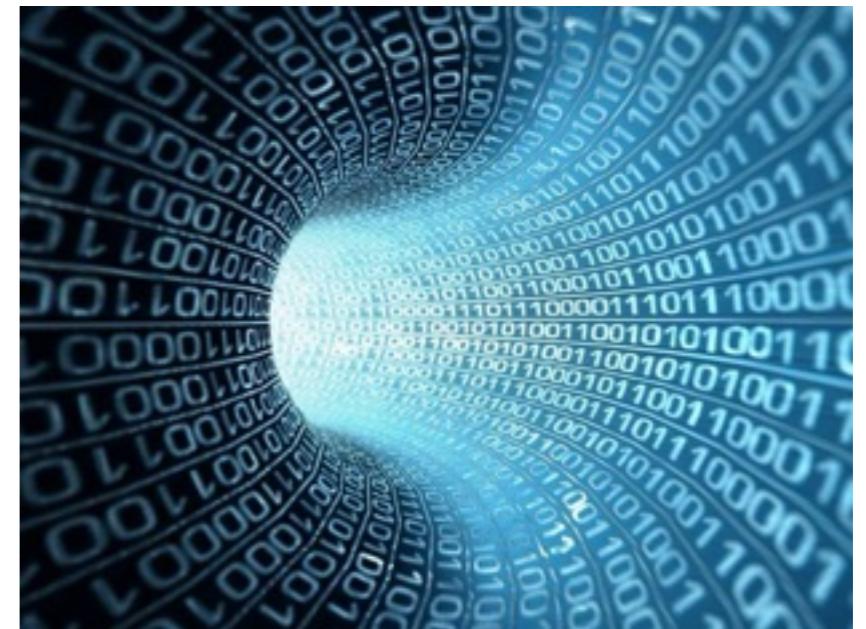
Speed Bumps

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



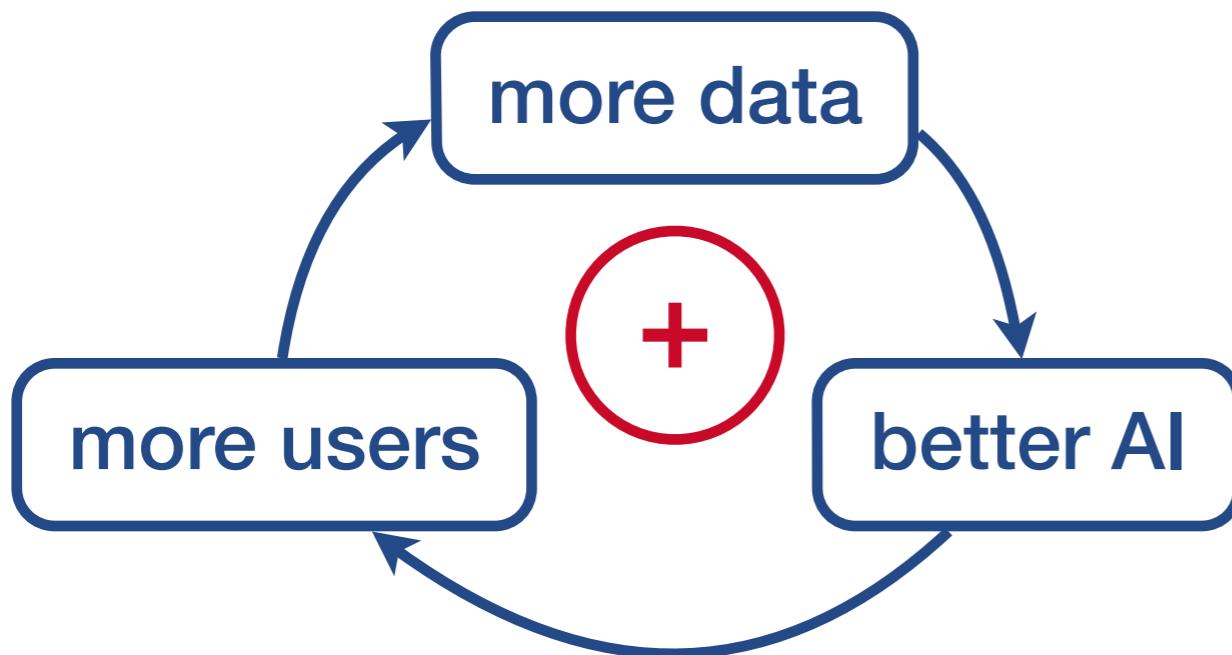
Accelerators

- Faster hardware
- Better algorithms
- Massive datasets



+ enormous incentives!

Economic Incentives



- It's difficult to enter the race later on
- Machines do more intellectual tasks
- Impossible for humans to compete

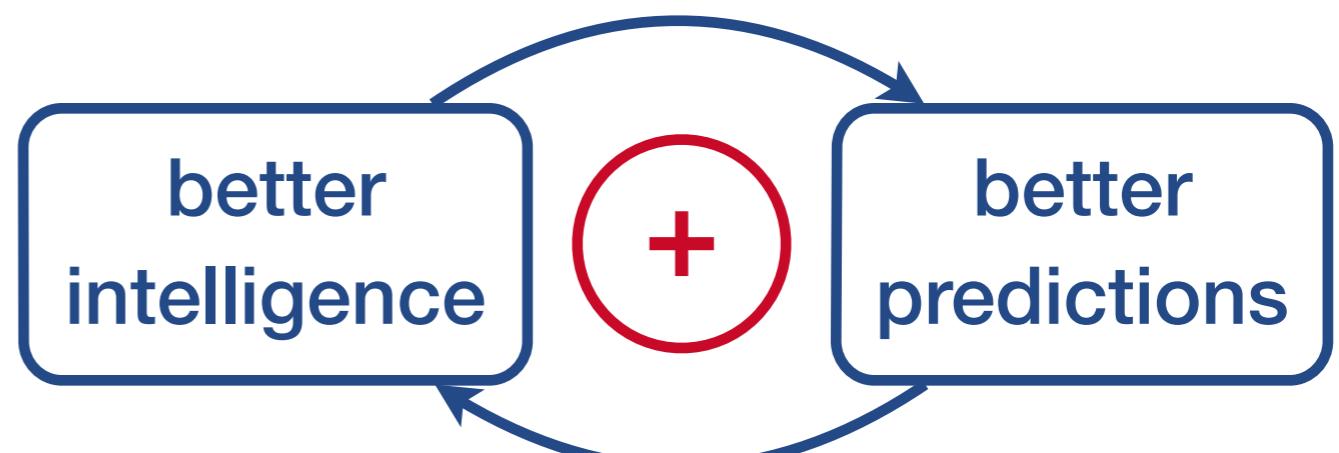
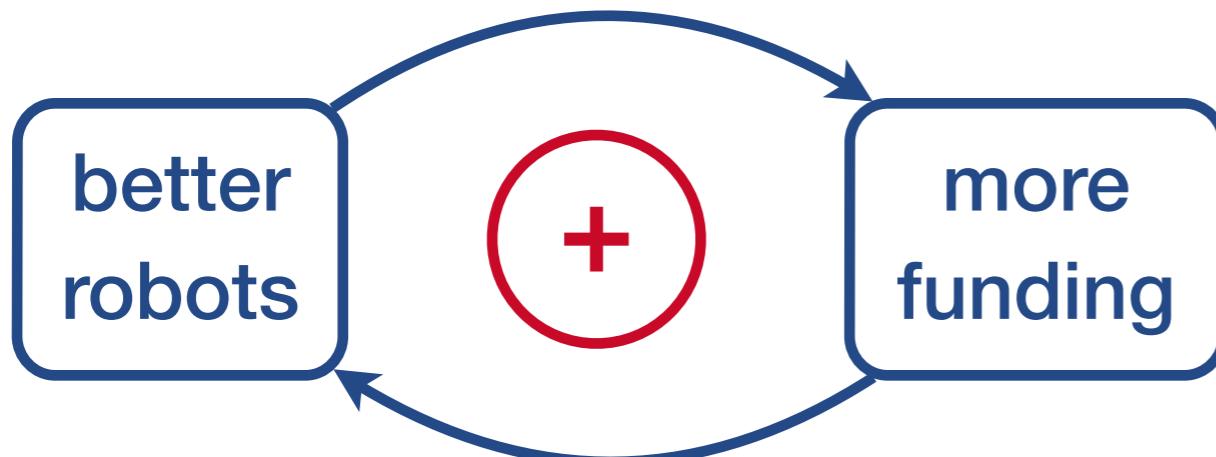
Economic Consequences

- The living costs of digital workers are drastically lower (just energy)
- Thus enormous pressure on wages
- Massive unemployment ahead of us
- Wages approach zero, wealth infinity

Introduce unconditional basic income?

Military Incentives

Arms Race?



Egoistic Incentives

- Intelligence
 - Wellbeing
 - Longevity
- ⇒ Willing to take risks



But with great power comes great responsibility!



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!

- Extreme Poverty
- Factory Farming
- Climate Change
- Artificial Intelligence

could
solve
other
issue

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

Escaping the Box



The AI could persuade someone to free it from its box and thus human control by:

- Offering wealth and power to liberator
- Claiming it needs outside resources to accomplish a task (like curing diseases)
- Predicting a real-world disaster which occurs and claiming afterwards it could have been prevented had it been let out

Value Loading

Utility function of AI?

- Perverse instantiation
- Moral blind-spots...?



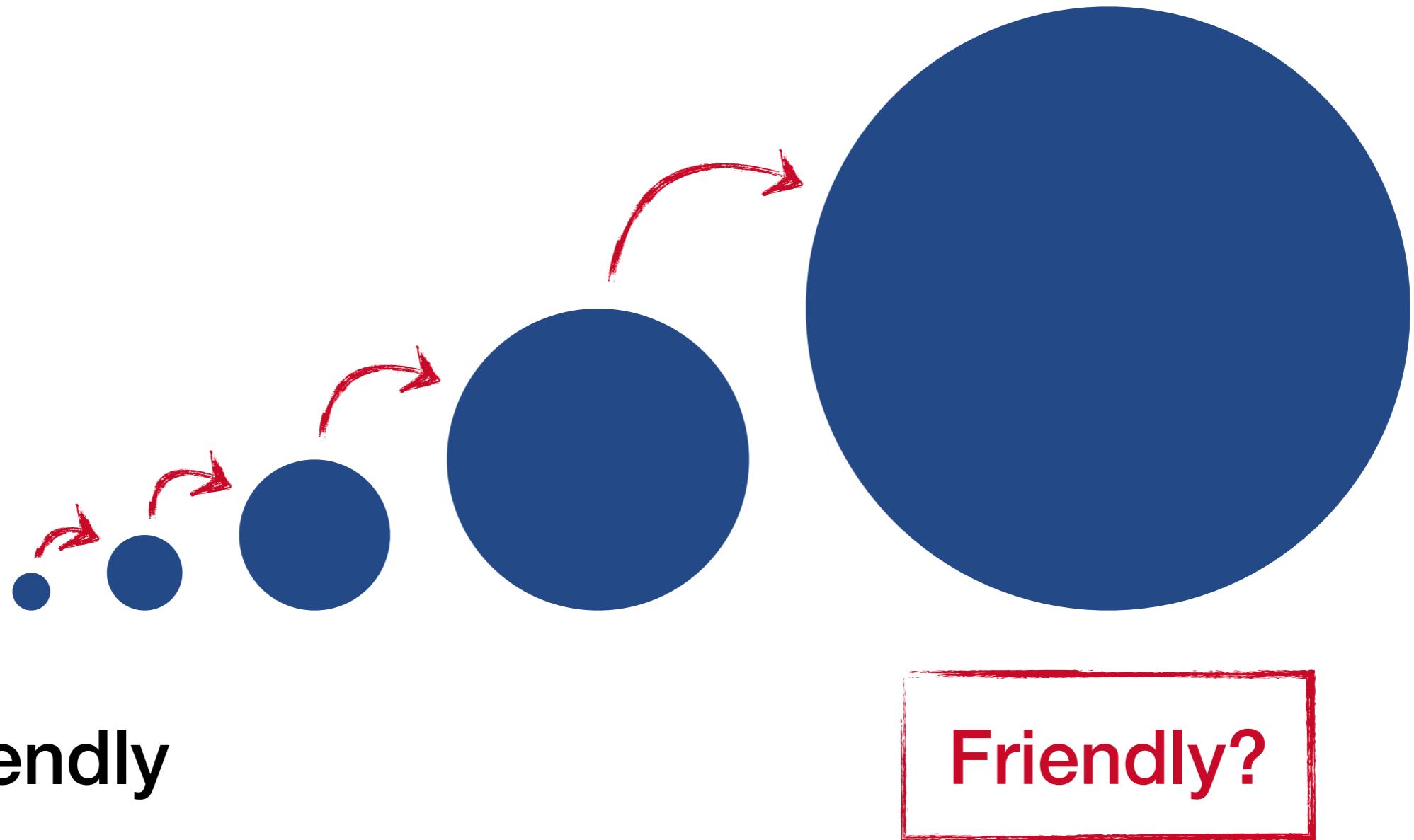
Coherent Extrapolated Volition (CEV):

The AI should do what we would want, if we were more intelligent, better informed and more the people we wished we were.

Goal-Directedness and Tool AI

- Orthogonality Thesis (revisited): Any utility function can be combined with a powerful epistemology and decision theory.
- Why not create an AI without motivations?
- Boxed oracle AI could work but less useful
- AI is relevant to find solutions for problems:
 - might be unintended (perverse instant.)
 - might require planning to meet criterion

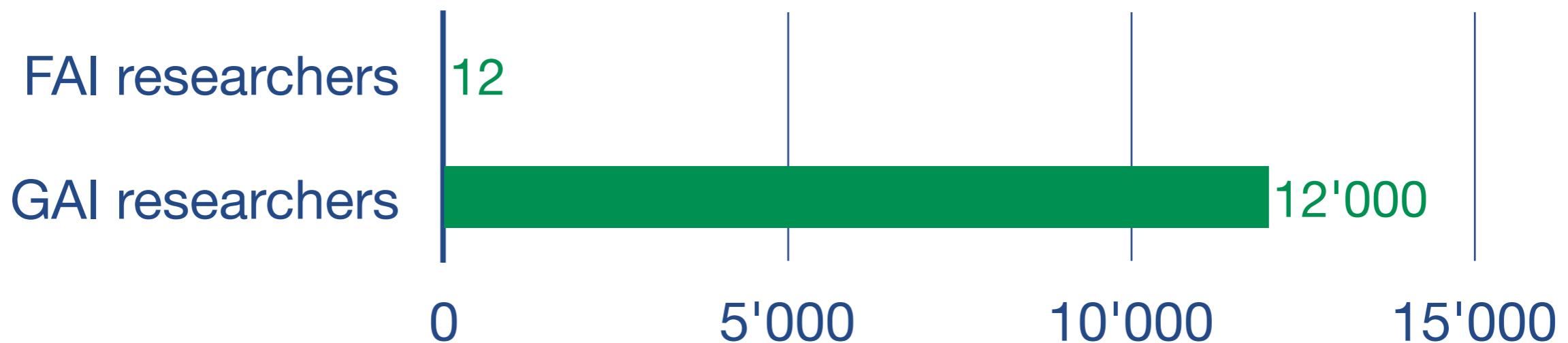
Stable Self-Improvement



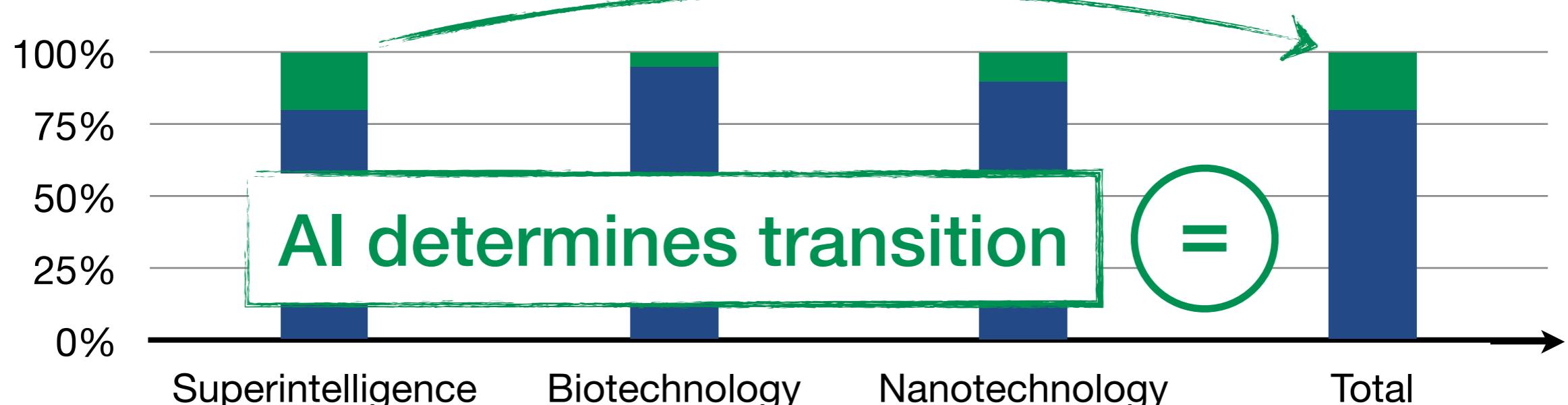
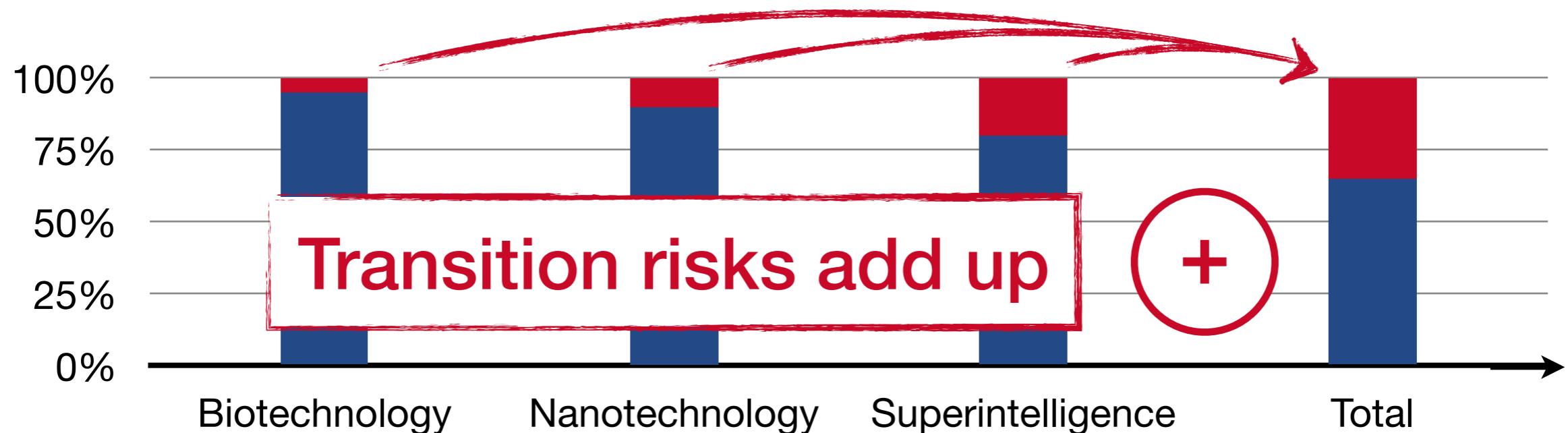
Differential Intellectual Progress

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

AI safety should outpace AI capability research



Order of Arrival



Information Hazards

Research can ...

- reduce the great uncertainties
- ... but can also ...
- bring up dangerous insights or ideas

Creating Awareness

Outreach can ...

- create awareness

... but can also ...

- fuel existing fears and cause panic!



Prisoner's Dilemma

- Difficult to prevent arms races
- Parties are better off by defecting
- The winner takes all (of what remains)
- Arms races are dangerous because

parties sacrifice safety for speed!

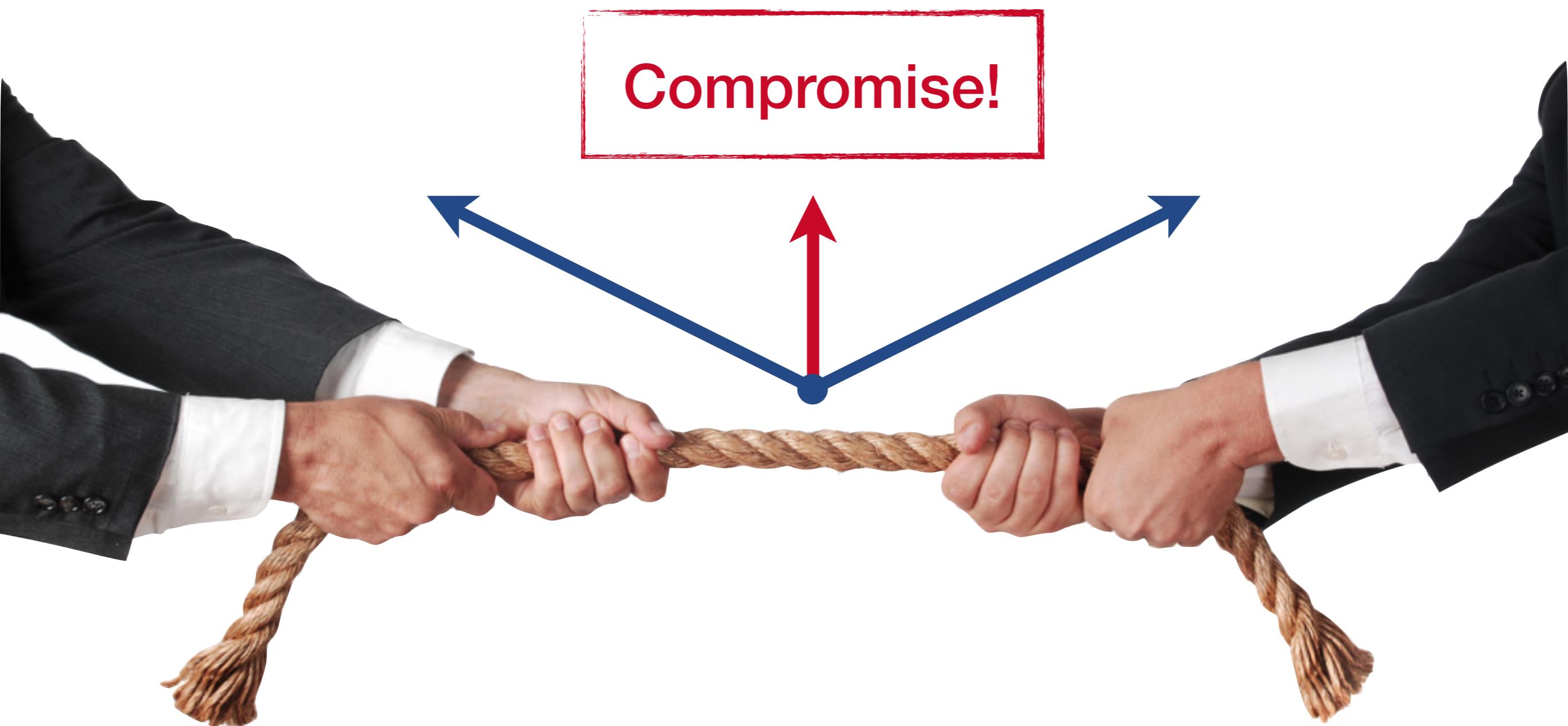
International Cooperation

- 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!

Moral Trade



Brian Tomasik: Gains from Trade through Compromise
foundational-research.org/.../gains-from-trade-.../

Superintelligence
Strategy 75

Heuristics for Altruists

Safe bets that likely turn out positive:

- Remain alive! (Self-Preservation)
- Remain an altruist! (Goal-Preservation)
- Acquire wealth and influence. (Resource Ac.)
- Educate yourself and become more rational.
(Self-Improvement, Intelligence Accumulation)



Sources

Where to learn more?

Institutes and Influential People

Future of Humanity Institute
UNIVERSITY OF OXFORD



Foundational Research
INSTITUTE



Nick Bostrom



Eliezer Yudkowsky



Brian Tomasik

Talks



The long-term future of AI(and what we can do about it): Daniel Dewey at TEDxVienna

Daniel Dewey
TEDxVienna



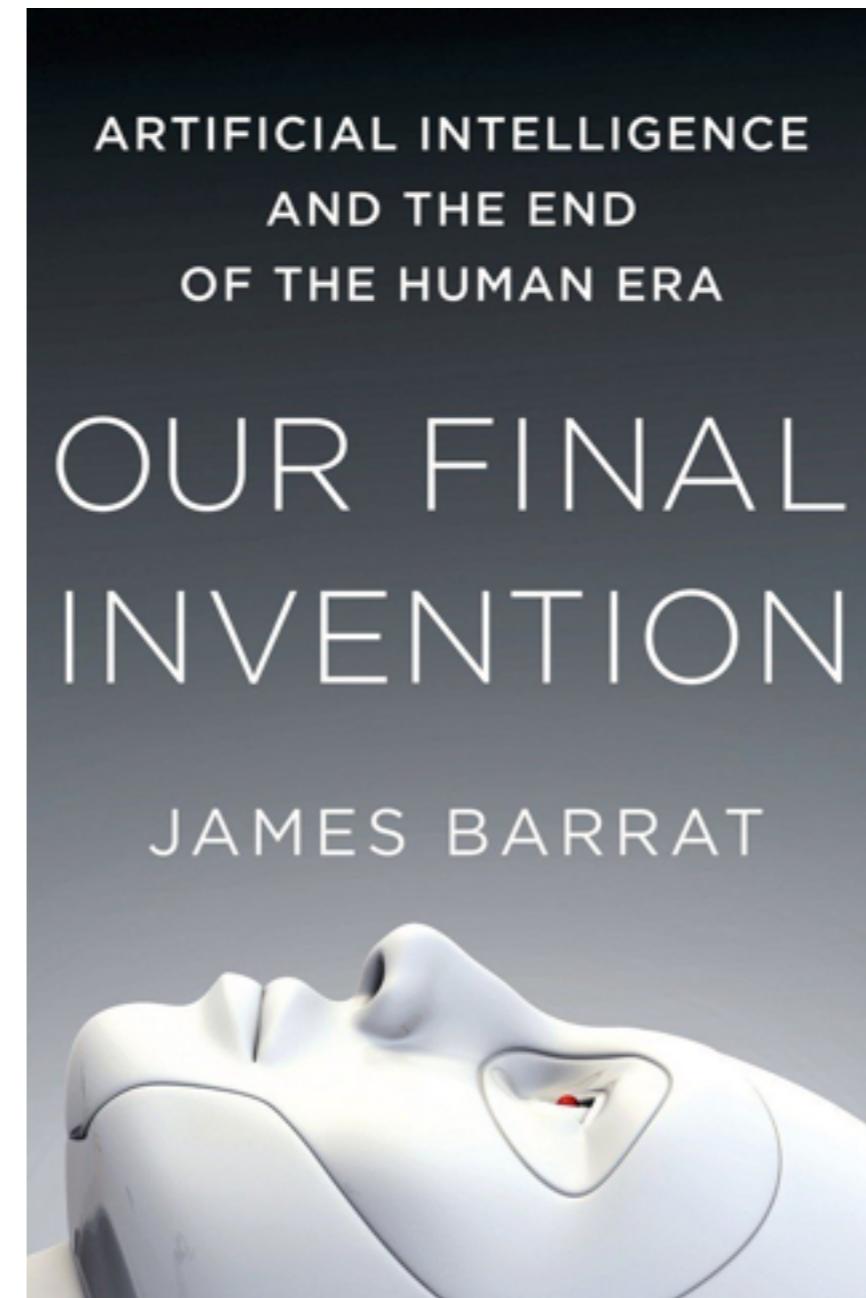
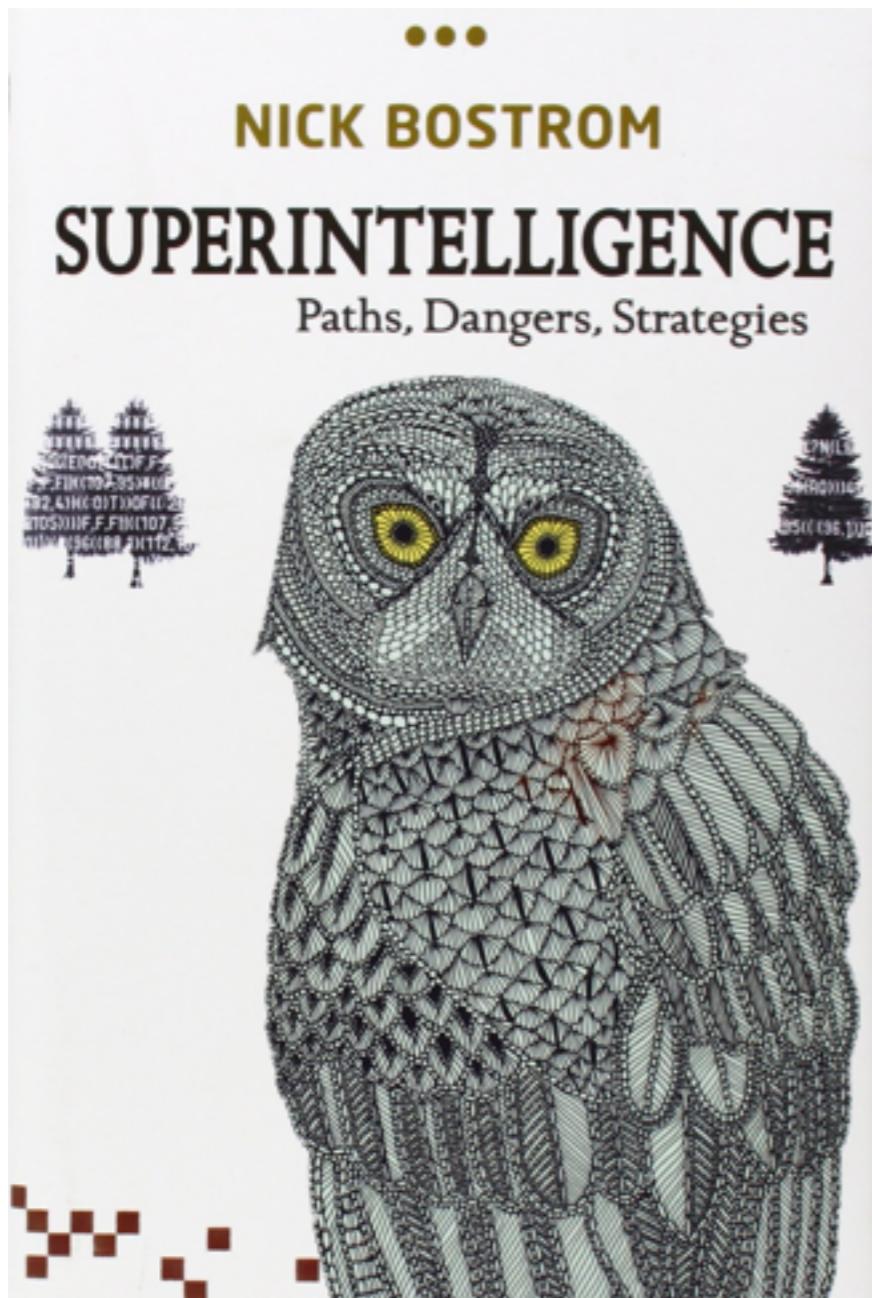
When creative machines overtake man: Jürgen Schmidhuber at TEDxLausanne

Jürgen Schmidhuber
TEDxLausanne

Papers

- **Intelligence Explosion by Luke Muehlhauser and Anna Salamon**
- **The Singularity: A Philosophical Analysis by David Chalmers**
- **The Superintelligent Will by Nick Bostrom**

Books





Summary

What have we learned?

Crucial Crossroad

Instead of passively drifting,
we need to steer a course!

- Philosophy
- Mathematics
- Cooperation

with a deadline.



«Before the prospect of an intelligence explosion, we humans are like children playing with a bomb. Such is the mismatch between the power of our play-thing and the immaturity of our conduct. Superintelligence is a challenge for which we are not ready now and will not be ready for a long time. We have little idea when the detonation will occur, though if we hold the device to our ear we can hear a faint ticking sound.»

— Prof. Nick Bostrom in his book *Superintelligence*



Discussion

www.superintelligence.ch



Kaspar Etter, kaspar.etter@gbs-schweiz.org

Adrian Hutter, adrian.hutter@gbs-schweiz.org

Basel, Switzerland

22 November 2014

85