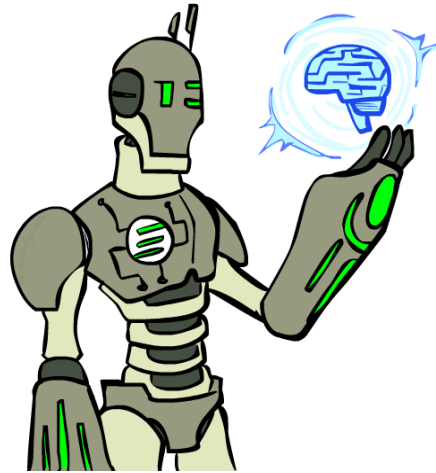


# Introduction to Artificial Intelligence

Lecture: Artificial General Intelligence

Prof. Gilles Louppe  
[g.louppe@uliege.be](mailto:g.louppe@uliege.be)

# Today\*

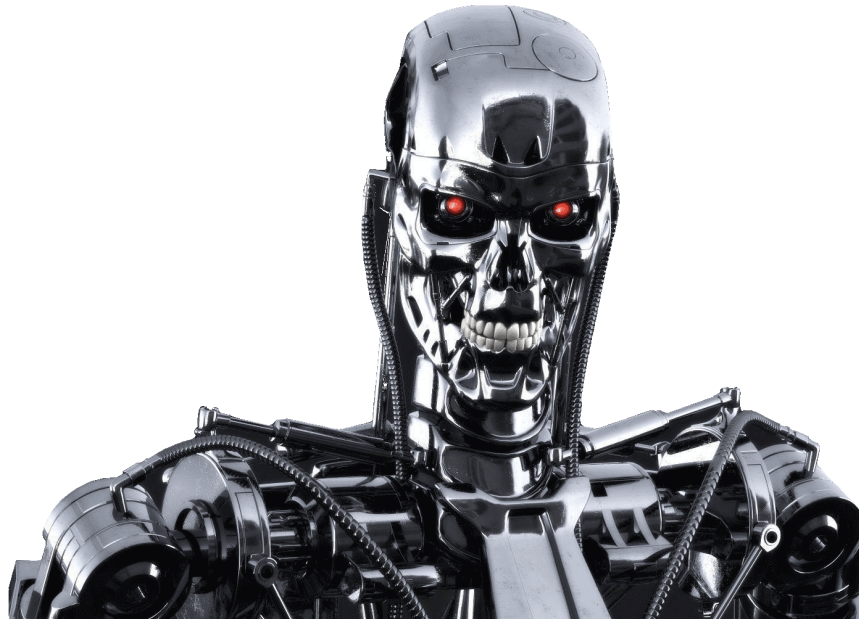


Towards generally intelligent agents?

- Artificial general intelligence
- AIXI
- Artificial life



*From technological breakthroughs...*



... to press coverage.

SingularityHub

TOPICS IN PROGRESS EXPERTS EVENTS

## Will Artificial Intelligence Become Conscious?

By Richard Kiser · Oct 19, 2017 · 42 min

Forget about today's robots. In tomorrow's advanced artificial intelligence, built on the emerging ability of machines to think for themselves, building an artificial mind might be a good idea. Consciousness: A machine that is aware of itself and its surroundings, and that could take its own decisions, receive amounts of data in real time. It could be used on dangerous missions, and placed in conflict. In addition to doing people's work, it might be able to look after us. It could—and even keep humans company when other people aren't nearby.

Don't miss a thing. Sign up for our newsletter. It's free. We'll send you our newsletter. It's free. We'll send you our newsletter. It's free.

Sign up

EXPRESS

HOME NEWS SHOWS & TV SPORT COMMENT FINANCE TRAVEL ENTERTAINMENT LIFE & STYLE

NEWS

### Rise of the machines: Super intelligent robots could 'spell the end of the human race'

Months away from becoming conscious, fast against their makers and overthrow humanity, machines warn

By NAVEEN KHANNA  
Illustration by THE NEW YORK TIMES MAGAZINE/DAVID SHAPIRO

Latest videos

- Robotics Will Soon Be Able to Think for Themselves
- AI Could Be the End of the World
- AI Could Be the End of the World
- AI Could Be the End of the World
- AI Could Be the End of the World
- AI Could Be the End of the World
- AI Could Be the End of the World
- AI Could Be the End of the World

AI

## 'KILLER ROBOTS' WILL START SLAUGHTERING PEOPLE IF THEY'RE NOT BANNED SOON, AI EXPERT WARNS

'These will be weapons of mass destruction'

FRANCE 24

La 4e

Si nous ne faisons rien, l'intelligence artificielle nous écabouillera dans 30 ans

01:26



## Artificial narrow intelligence

Today's artificial intelligence remains **narrow**:

- AI systems often reach super-human level performance, ... but only at **very specific problems**!
- They **do not generalize** to the real world nor to arbitrary tasks.

# The case of AlphaGo

Convenient properties of the game of Go:

- Deterministic (no noise in the game).
- Fully observed (each player has complete information)
- Discrete action space (finite number of actions possible)
- Perfect simulator (the effect of any action is known exactly)
- Short episodes (200 actions per game)
- Clear and fast evaluation (as stated by Go rules)
- Huge dataset available (games)





Can we run AlphaGo on a robot?

# AGI

Artificial general intelligence, or **AGI**, is the intelligence of a machine that could successfully perform any intellectual task that a human being can perform.

The scientific community agrees that AGI would be required to do the following:

- reason, use strategy, solve puzzle, plan,
- make judgments under uncertainty,
- represent knowledge, including commonsense knowledge,
- improve and learn new skills,
- communicate in natural language,
- be creative,
- integrate all these skills towards common goals.

This is similar to our definition of **thinking rationally**, but applied broadly to any set of tasks.



## Roads towards AGI

Several working hypothesis:

1. Learning (supervised, unsupervised, reinforcement)
2. AIXI
3. Artificial life

... or probably something else?

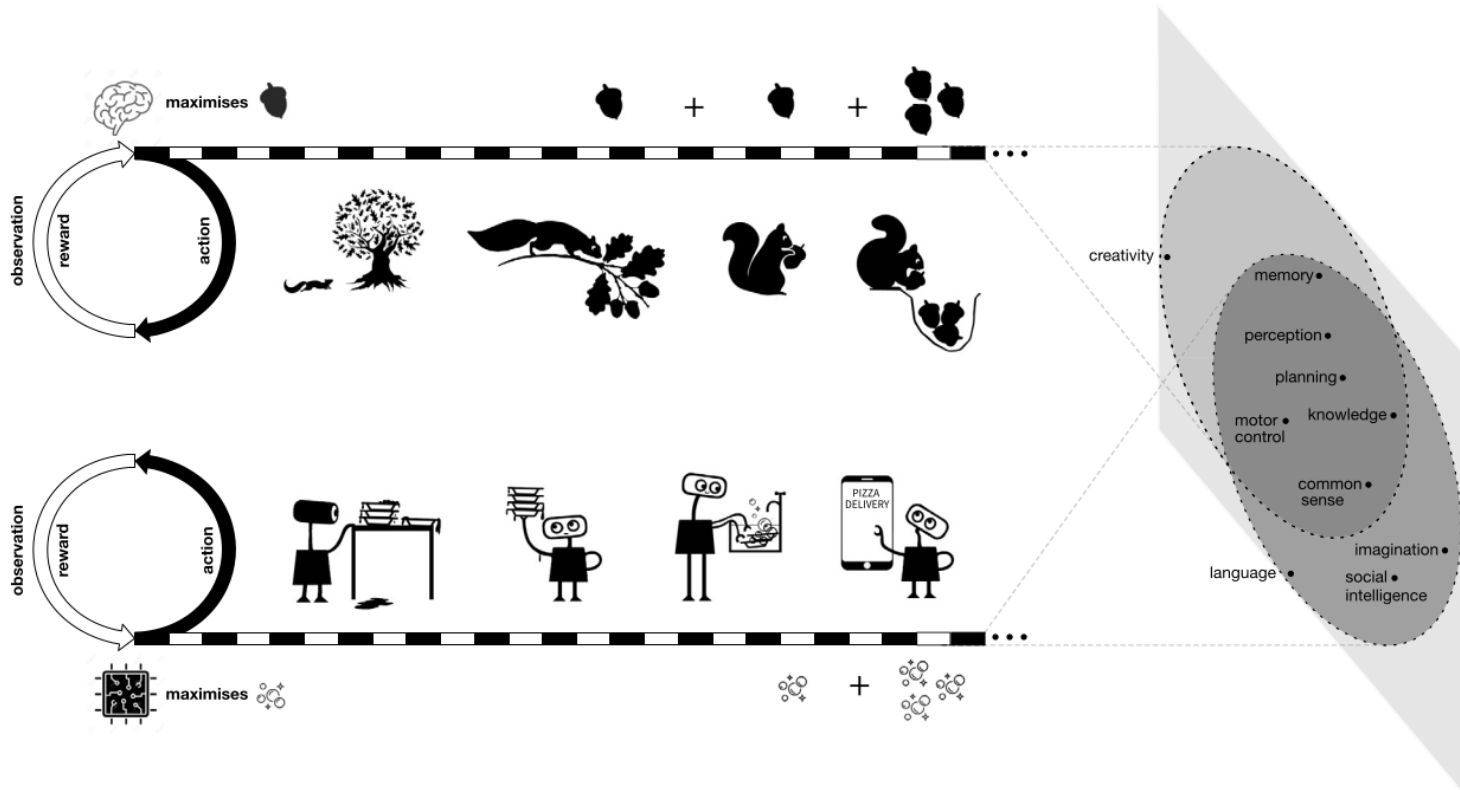
# Learning

# Reward is Enough

## Abstract

In this paper we hypothesise that the objective of maximising reward is enough to drive behaviour that exhibits most if not all attributes of intelligence that are studied in natural and artificial intelligence, including knowledge, learning, perception, social intelligence, language and generalisation. This is in contrast to the view that specialised problem formulations are needed for each attribute of intelligence, based on other signals or objectives. The reward-is-enough hypothesis suggests that agents with powerful reinforcement learning algorithms when placed in rich environments with simple rewards could develop the kind of broad, multi-attribute intelligence that constitutes an artificial general intelligence.

David Silver et al, 2021.



**Fig. 1.** The *reward-is-enough* hypothesis postulates that intelligence, and its associated abilities, can be understood as subserving the maximisation of reward by an agent acting in its environment. For example, a squirrel acts so as to maximise its consumption of food (top, reward depicted by acorn symbol), or a kitchen robot acts to maximise cleanliness (bottom, reward depicted by bubble symbol). To achieve these goals, complex behaviours are required that exhibit a wide variety of abilities associated with intelligence (depicted on the right as a projection from an agent's stream of experience onto a set of abilities expressed within that experience).



Sparks! Could AI be perceived as creative ?

**SPARKS!**  
Serendipity Forum at CERN



Later bekij...



Delen

**Talk**

Jürgen idhuber

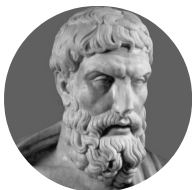
Could AI be perceived as creative? (Jürgen Schmidhuber)

# AIXI

AIXI (Hutter, 2005) is a theoretical mathematical formalism of artificial general intelligence.



Occam: Prefer the simplest consistent hypothesis.



Epicurus: Keep all consistent hypotheses.



$$\text{Bayes: } P(h|d) = \frac{P(d|h)P(h)}{P(d)}$$



Turing: It is possible to invent a single machine which can be used to compute any computable sequence.



Solomonoff: Use computer programs  $\mu$  as hypotheses/environments.

AIXI defines a measure of universal intelligence as

$$\Upsilon(\pi) := \sum_{\mu \in E} 2^{-K(\mu)} V_{\mu}^{\pi}$$

where

- $\Upsilon(\pi)$  formally defines the **universal intelligence** of an agent  $\pi$ .
- $\mu$  is the environment of the agent and  $E$  is the set of all computable reward bounded environments.
- $V_{\mu}^{\pi} = \mathbb{E}[\sum_{i=1}^{\infty} R_i]$  is the expected sum of future rewards when the agent  $\pi$  interacts with environment  $\mu$ .
- $K(\cdot)$  is the Kolmogorov complexity, such that  $2^{-K(\mu)}$  weights the agent's performance in each environment, inversely proportional to its complexity.



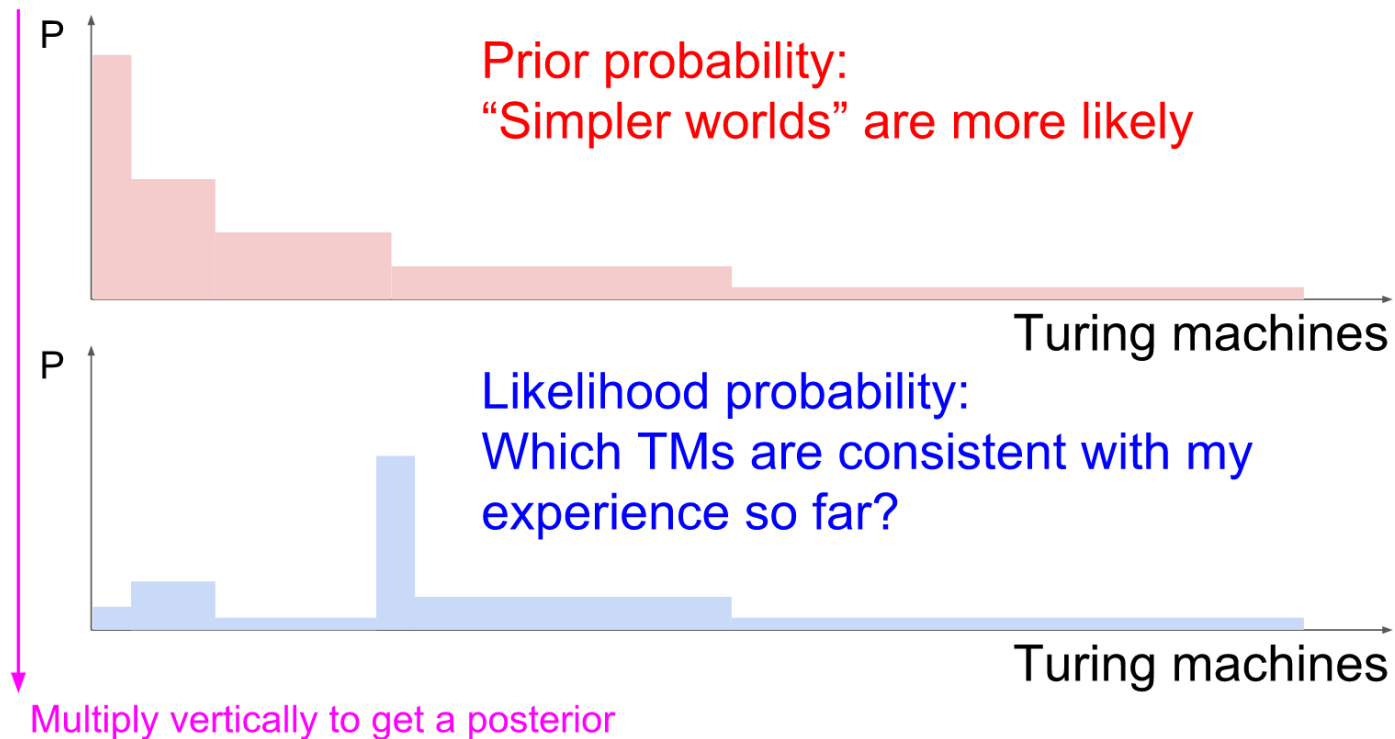
# AIXI

$$\bar{\Upsilon} = \max_{\pi} \Upsilon(\pi) = \Upsilon(\pi^{\text{AIXI}})$$

$\pi^{\text{AIXI}}$  is a **perfect** theoretical agent.

## System identification

- Which Turing machine is the agent in? If it knew, it could plan perfectly.
- Use the [Bayes rule](#) to update the agent beliefs given its experience so far.



## Acting optimally

- The agent always picks the action which has the greatest expected reward.
- For every environment  $\mu \in \mathcal{E}$ , the agent must:
  - Take into account how likely it is that it is facing  $\mu$  given the interaction history so far, and the prior probability of  $\mu$ .
  - Consider all possible future interactions that might occur, assuming optimal future actions.
  - Evaluate how likely they are.
  - Then select the action that maximizes the expected future reward.

$$a_t^{\pi^\xi} := \arg \max_{a_t} \lim_{m \rightarrow \infty} \sum_{x_t} \max_{a_{t+1}} \sum_{x_{t+1}} \cdots \max_{a_m} \sum_{x_m} [\gamma_t r_t + \cdots + \gamma_m r_m] \xi(a_{x_{<t}} a_{x_{t:m}})$$

$$\xi(a_{x_{1:n}}) := \sum_{\nu \in E} 2^{-K(\nu)} \nu(a_{x_{1:n}})$$

(description length of the TM, number of bits)

Complete history of interactions up to this point

$a_{x_{<t}}$

time t

all possible future action-state sequences

time m

Weighted average of the total discounted reward, across all possible Turing Machines.

The weights are [prior] x [likelihood] for each Turing machine.

# AIXI is incomputable

$$a_t^{\pi^\xi} := \arg \max_{a_t} \lim_{m \rightarrow \infty} \sum_{x_t} \max_{a_{t+1}} \sum_{x_{t+1}} \cdots \max_{a_m} \sum_{x_m} [\gamma_t r_t + \cdots + \gamma_m r_m] \xi(\underline{ax}_{<t} \underline{ax}_{t:m})$$

!!!                      !!!

$$\xi(\underline{ax}_{1:n}) := \sum_{\nu \in E} 2^{-K(\nu)} \nu(\underline{ax}_{1:n})$$

## Benefits of AIXI

The AIXI theoretical formalism of AGI provides

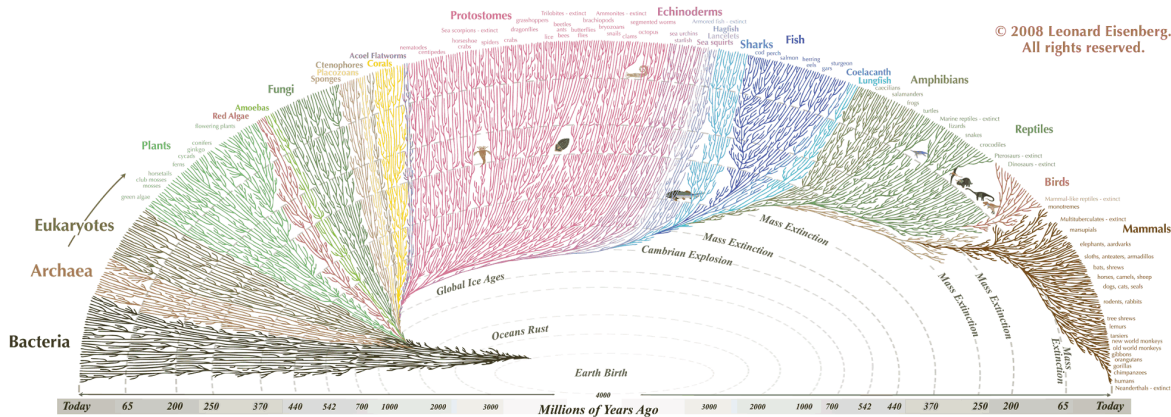
- a high-level **blue-print** or inspiration for design;
- common terminology and goal formulation;
- understand and predict behavior of yet-to-be-built agents;
- appreciation of **fundamental challenges** (e.g., exploration-exploitation);
- **definition/measure** of intelligence.

# Artificial life

# Artificial life

Study of systems related to natural life, its processes and its evolution, through the use of **simulations** with computer models, robotics or biochemistry.

One of its goals is to **synthesize** life in order to understand its origins, development and organization.



All the major and many of the minor living branches of life are shown on this diagram, but only a few of those that have gone extinct are shown. Example: Dinosaurs - extinct

*How did intelligence arise in Nature?*



## Approaches

There are three main kinds of artificial life, named after their approaches:

- Software approaches (soft)
- Hardware approaches (hard)
- Biochemistry approaches (wet)

The field of AI has traditionally used a top down approach. Artificial life generally works from the [bottom up](#).



## Video niet beschikbaar

Deze video bevat content van Discovery Communications, die deze wegens auteursrechtshending heeft geblokkeerd voor je land



Wet artificial life: The line between life and not-life (Martin Hanczyc).

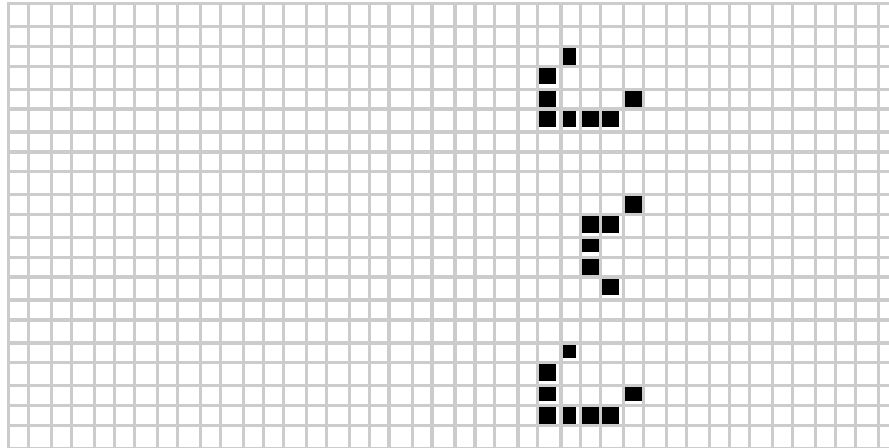
# Evolutionary algorithms

Evolution may **hypothetically** be interpreted as an (unknown) algorithm.

- This algorithm gave rise to AGI (e.g., it induced humans).
- Simulation of the evolutionary process should/could eventually reproduce life and, maybe, intelligence?

## Conway's Game of Life

- Any live cell with two or three live neighbours survives.
- Any dead cell with three live neighbours becomes a live cell.
- All other live cells die in the next generation. Similarly, all other dead cells stay dead.





Let's BUILD a COMPUTER in CONWAY's GAME o...



Later bekij...



Delen

# WHY IS LIFE TURING COMPLETE



Conway's game of life

## Evolutionary algorithms as metaheuristic optimization algorithms

1. Start with a random population of creatures.
2. Repeat until termination:
  1. Each creature is tested for their ability to perform a given task.
  2. Select the fittest creatures for reproduction.
  3. Breed new creatures by combining and mutating the virtual genes of their selected parents.
  4. Replace the least-fit creatures of the population with new creatures.

As this cycle of variation and selection continues, creatures with more and more successful behaviors may **emerge**.



Karl Sims, 1994.



# Learning to Generalize Self-Assembling Agents [...] Generalization w/o Fine-tuning



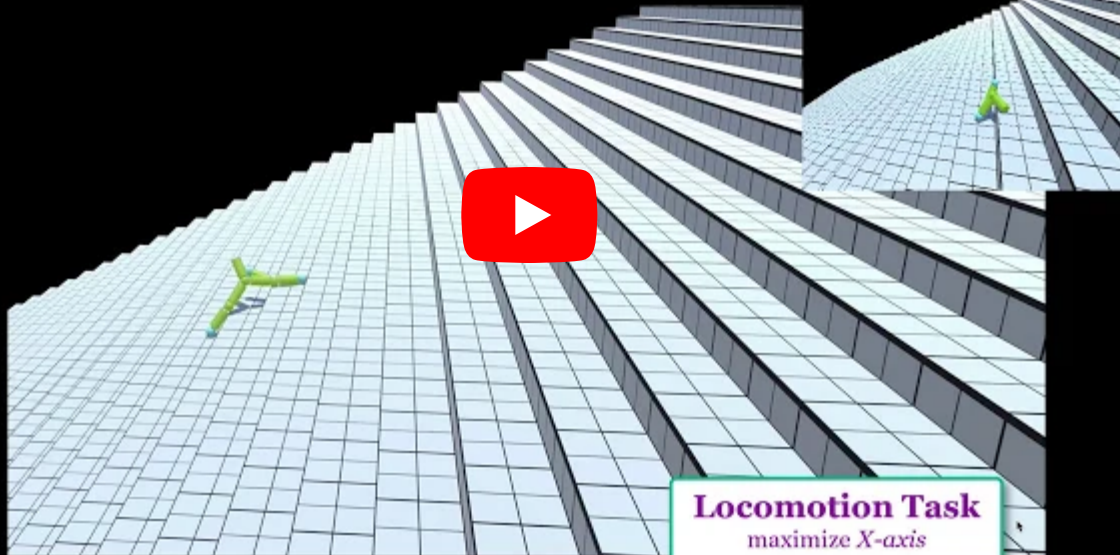
Later bekij...



Delen

terrain with stairs

Vanilla RL



Self-assembling morphologies (Pathak et al, 2019)





LIFE - EP1: Artificial Life



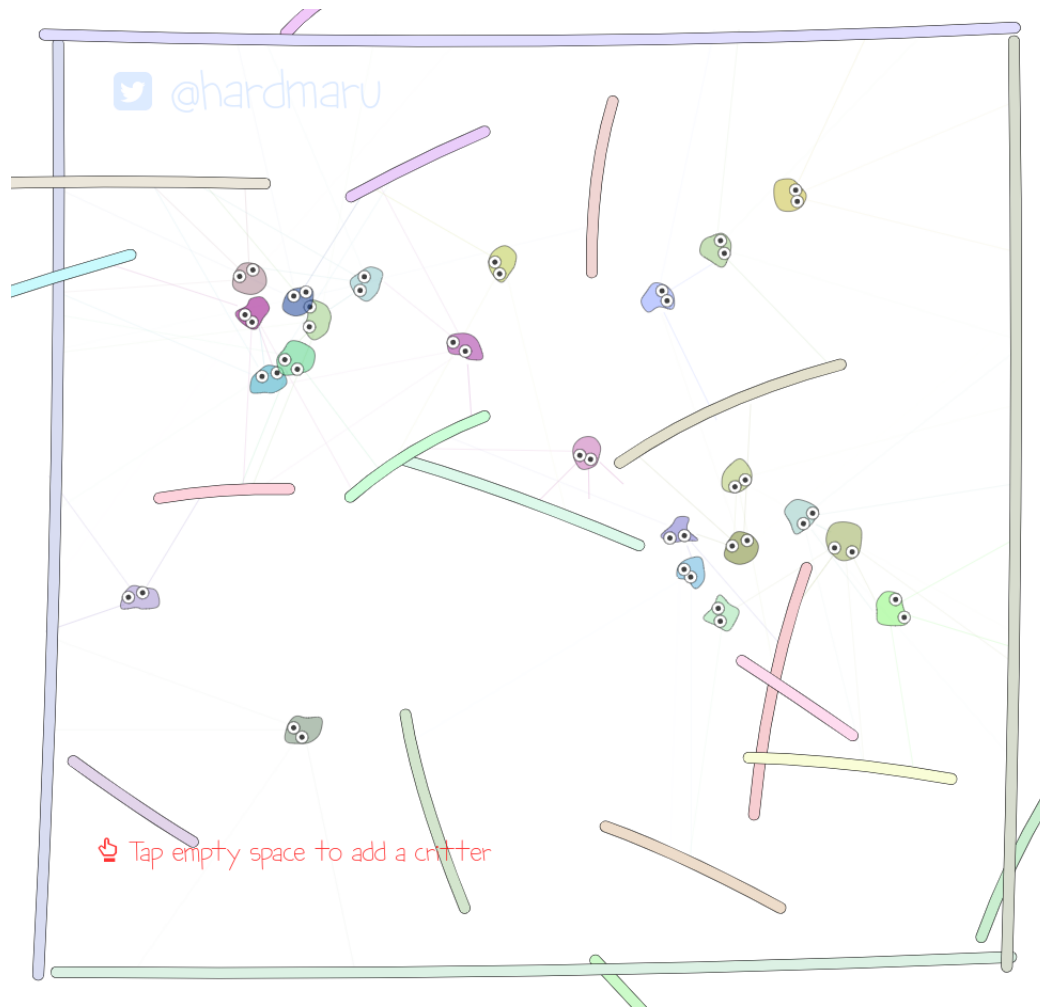
Later bekij...



Delen



Mini-documentary: Artificial life

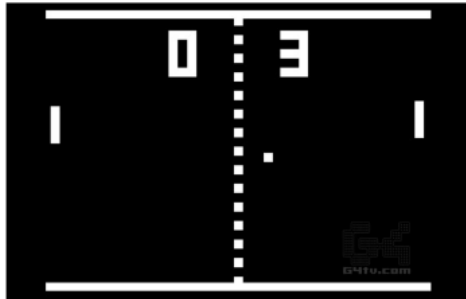


Creatures avoiding planks [demo].

## Environments for AGI?

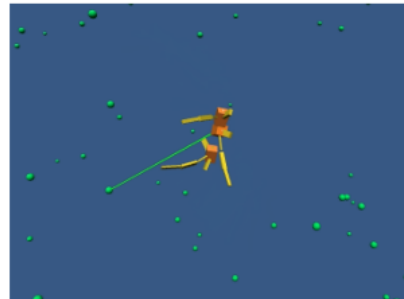
For the emergence of generally intelligent creatures, environments should **incentivize** the emergence of a **cognitive toolkit** (attention, memory, knowledge representation, reasoning, emotions, forward simulation, skill acquisition, ...).

Doing it wrong:



Incentives a lookup table of correct moves.

Doing it right:



Incentivises cognitive tools.

Multi-agent environments are certainly better because of:

- Variety: the environment is parameterized by its agent population. The optimal strategy must be derived dynamically.
- Natural curriculum: the difficulty of the environment is determined by the skill of the other agents.

