

Toward an AI physicist for unsupervised learning

Tailin Wu, Max Tegmark 2018, arxiv:1810.10525

Motivation – warm up

0.0 s

Calculate the 13th root of following number:

70664373816742861022340088302401573757042331707026
32731269721516000395709065419973141914549389684111

→Result: 47941071

→Human world record (2010): 11.8 seconds

→My laptop on average: 15 nanoseconds

Outline

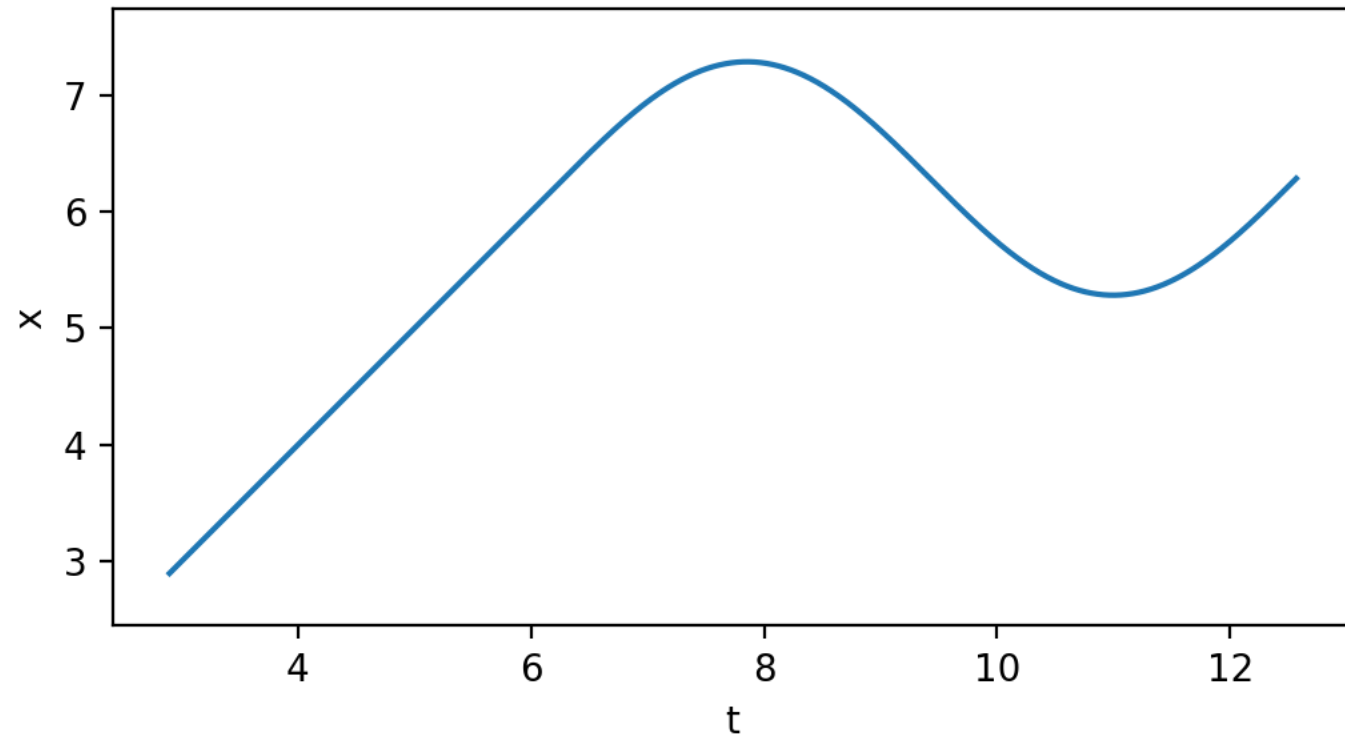
- What is a theory?
- AI physicist architecture
 - Divide and conquer
 - Occam's razor
 - Unification
 - Life long learning
- Experiments
- Conclusion/Discussion

What is a theory?

- Definition:
 - Theory \mathcal{T} is 2-tupel (\mathbf{f}, c)
 - Input $\mathbf{x}_t = (\mathbf{y}_{t-T}, \dots, \mathbf{y}_{t-1})$, $\mathbf{y}_i \in \mathbb{R}^d$
 - \mathbf{f} is prediction function $\rightarrow \mathbf{y}_t$ (3 layer NN with linear activation function)
 - c is domain classifier (3 layer NN with leakyReLU activation function)

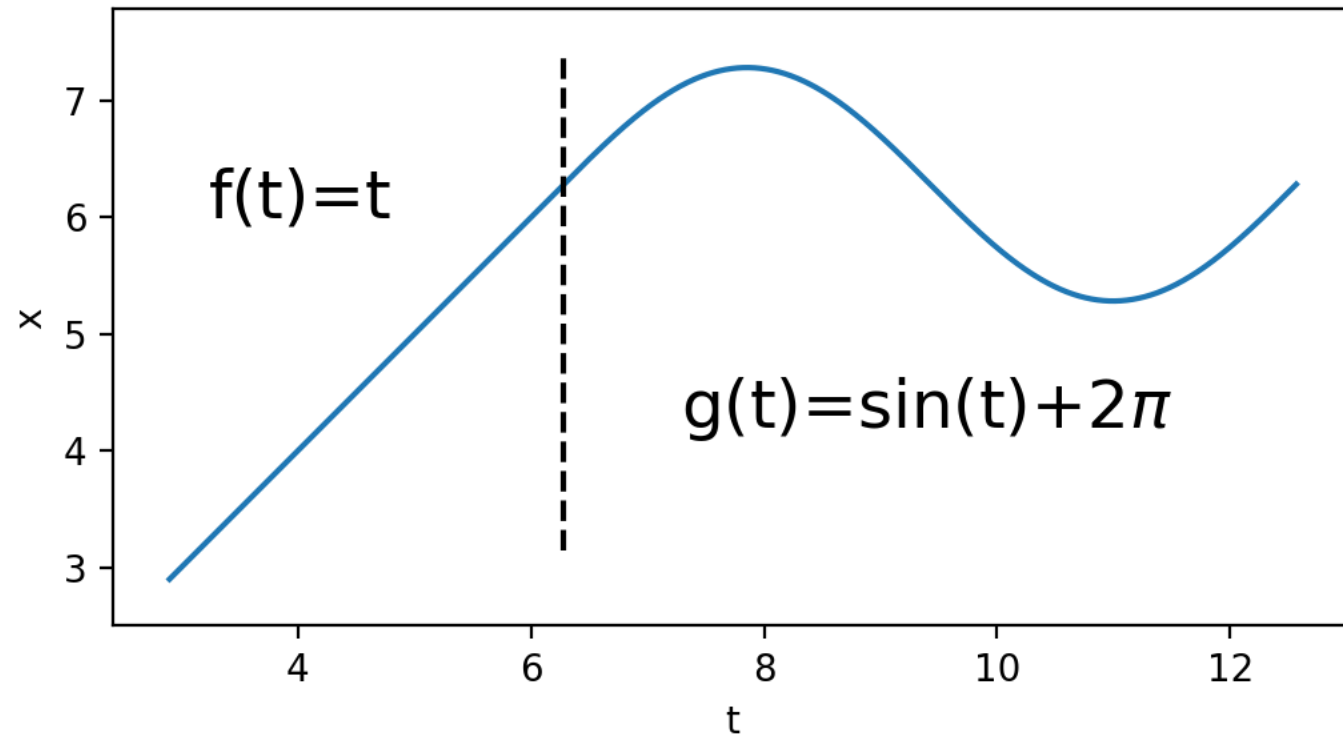
What is a theory?

- Example:



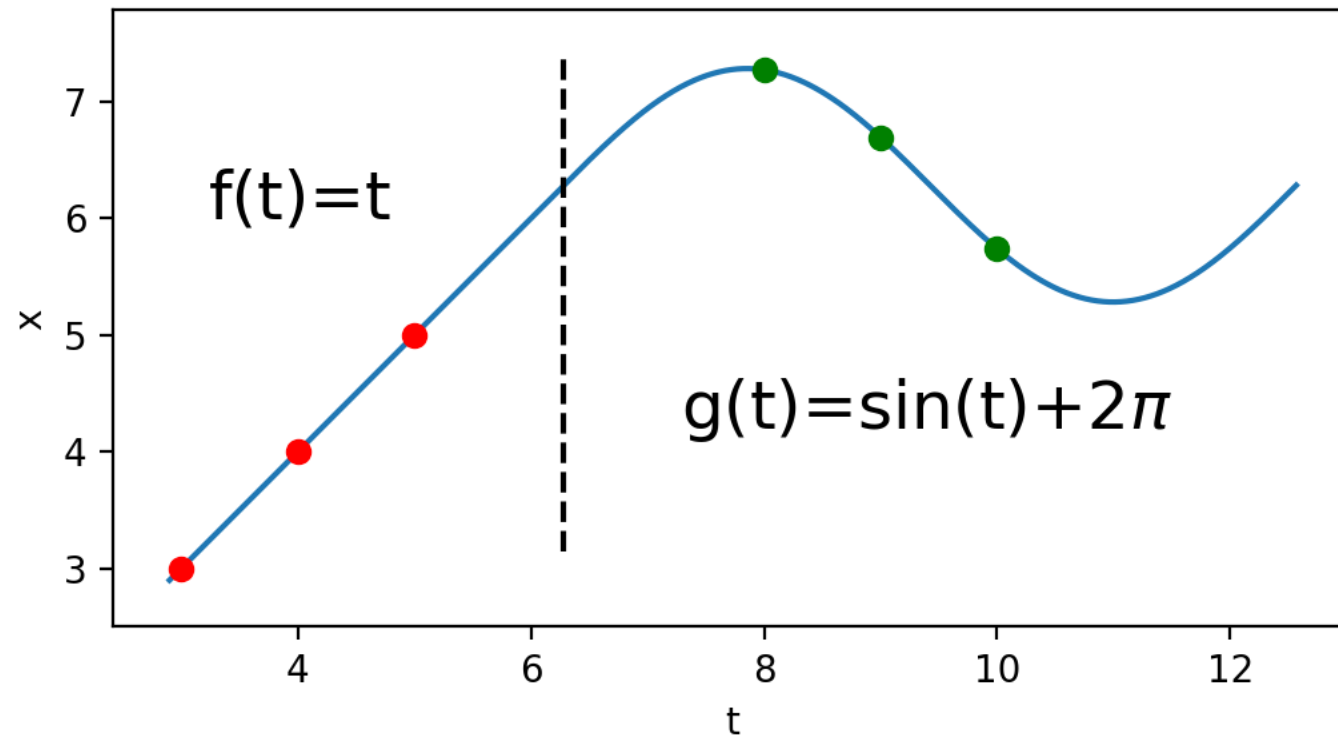
What is a theory?

- Example:



What is a theory?

- Example:



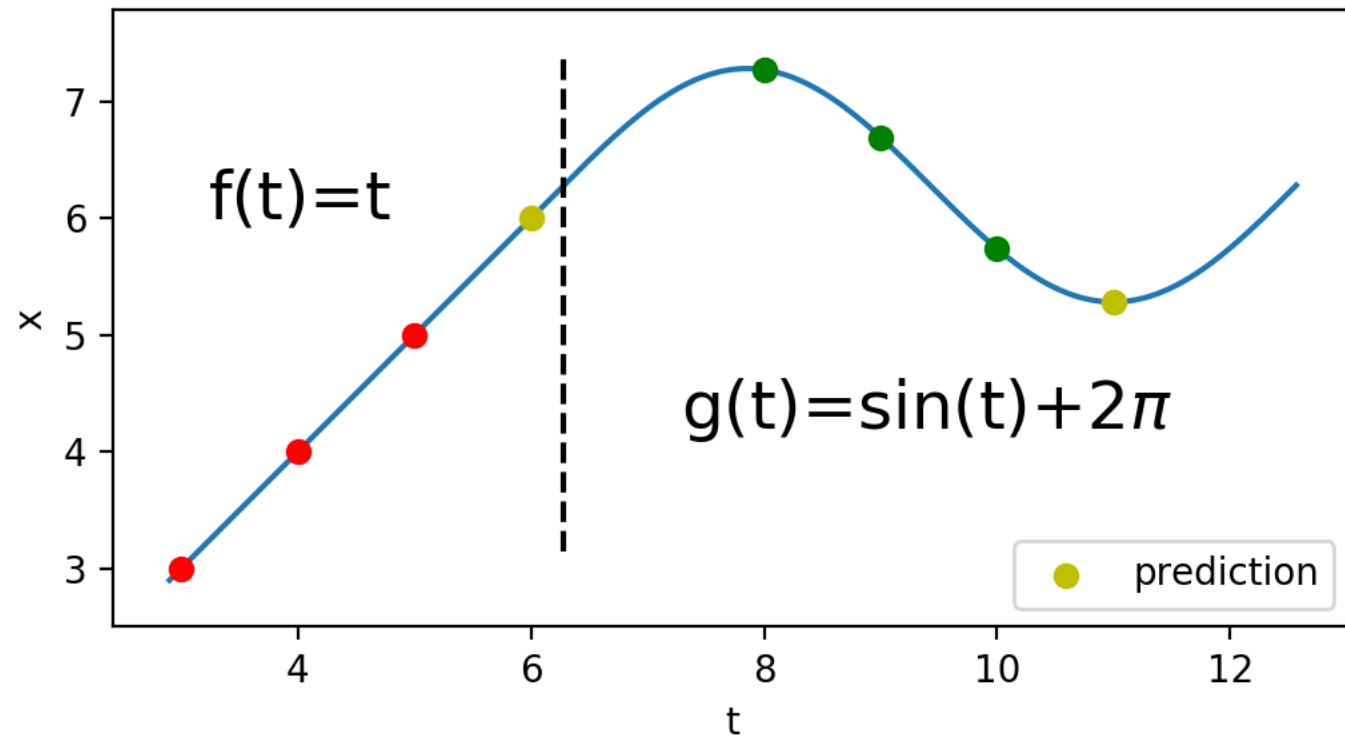
Data:

$$\mathbf{x}_t = (3.0, 4.0, 5.0) \xrightarrow{c} \mathbf{f}(\mathbf{x}_t)$$

$$\mathbf{x}_t = (7.3, 6.7, 5.7) \xrightarrow{c} \mathbf{g}(\mathbf{x}_t)$$

What is a theory?

- Example:



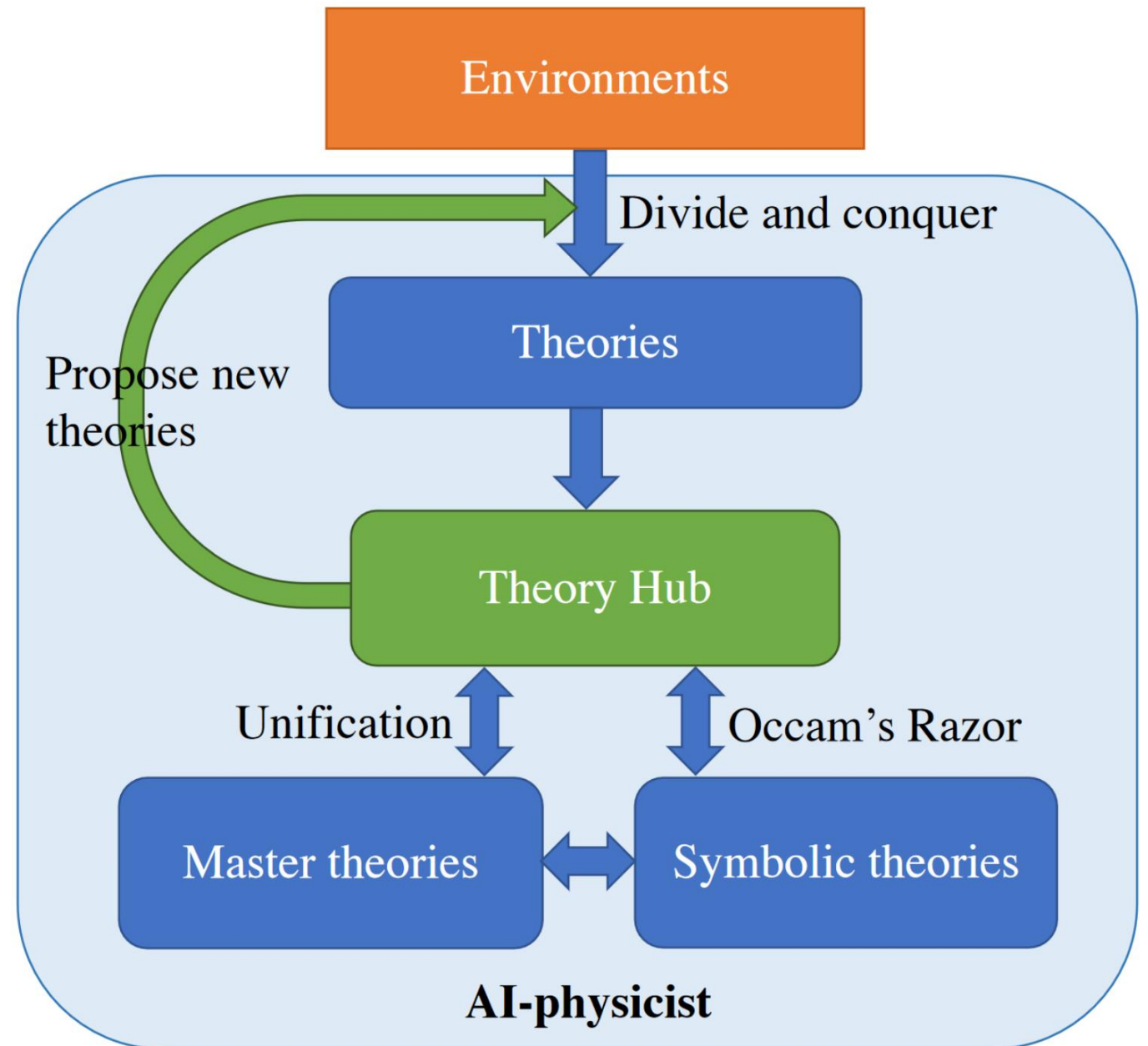
Data:

$$\mathbf{x}_t = (3.0, 4.0, 5.0) \xrightarrow{c} \mathbf{f}(\mathbf{x}_t) = 6.0$$

$$\mathbf{x}_t = (7.3, 6.7, 5.7) \xrightarrow{c} \mathbf{g}(\mathbf{x}_t) = 5.3$$

Architecture

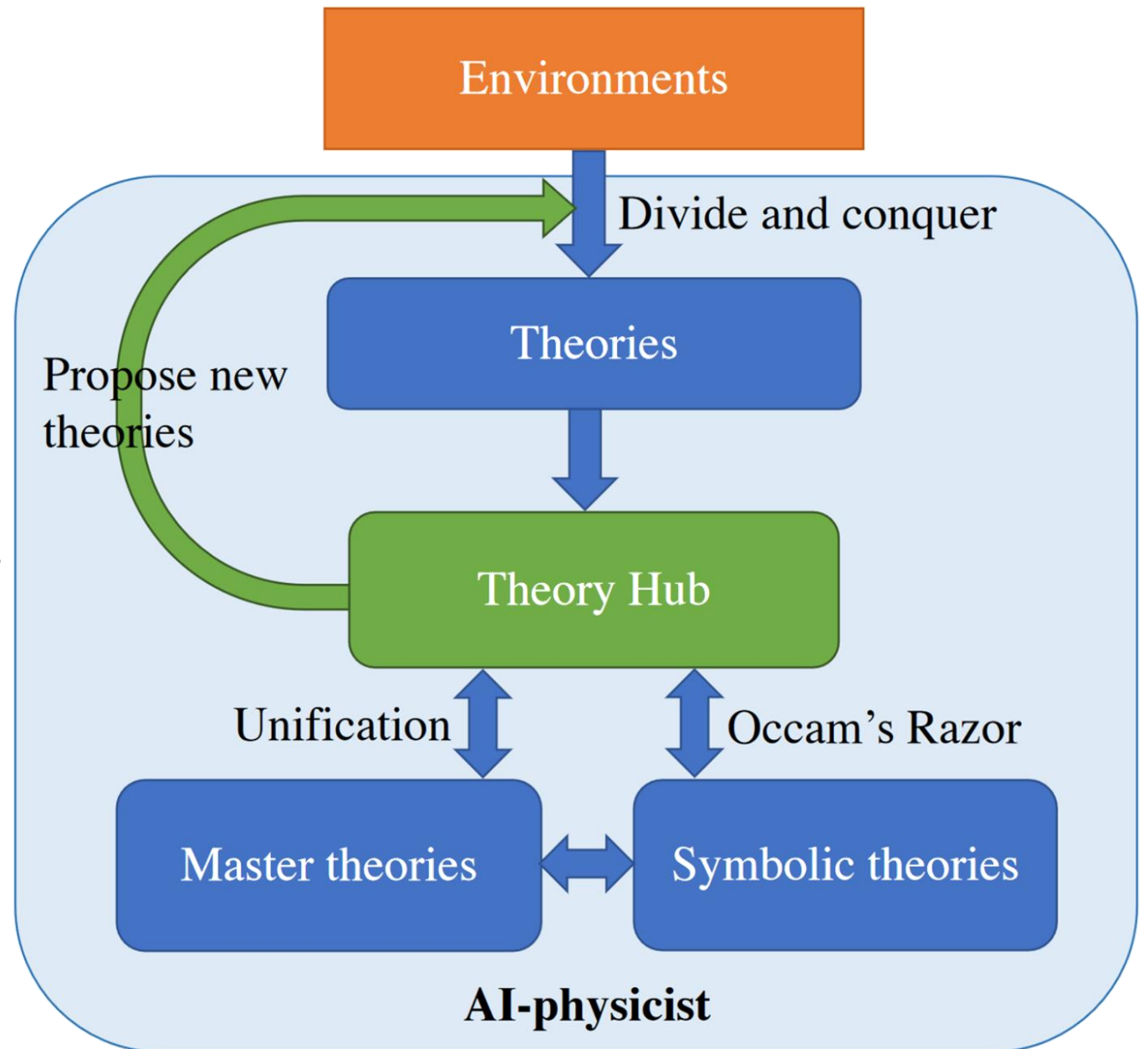
- Divide-and-conquer
- Lifelong learning
- Occam's Razor
- Unification



Source: [1]

Architecture

- AI physicist leaves many small models applicable in different domains
- Not approximate predictions but near exact predictions with complete intelligibility



Source: [1]

Differentiable divide and conquer (DDAC)

- Idea: Not predict everything in one model but only parts of the data
 - DDAC learns prediction functions and corresponding domain classifiers
- Each prediction function will specialize in its domain due to:
 - Generalized mean loss
 - Bisected training

Differentiable divide and conquer (DDAC)

Algorithmic idea:

- Initialize random theories
- Iteratively train prediction functions and domain classifier

- Loss function:
$$\mathcal{L}_\gamma = \sum_t \left(\frac{1}{M} \sum_{i=1}^M \ell[\mathbf{f}_i(\mathbf{x}_t), \mathbf{y}_t]^\gamma \right)^{1/\gamma}$$

- M : Number of theories
- ℓ : description length loss ($\ell[\mathbf{f}_i(\mathbf{x}_t), \mathbf{y}_t] = \frac{1}{2} \log_2 \left(1 + \left(\frac{\mathbf{f}_i(\mathbf{x}_t) - \mathbf{y}_t}{\epsilon} \right)^2 \right)$)
- γ : parameter ($\gamma = -1$)
- \mathbf{f} : prediction function
- \mathbf{x}_t : input data
- \mathbf{y}_t : label

Differentiable divide and conquer (DDAC)

Training:

- $\mathbf{f}_\theta = (\mathbf{f}_1, \dots, \mathbf{f}_M)$, $\mathbf{c}_\phi = (c_1, \dots, c_M)$
- Gradient descent on \mathbf{f}_θ :

$$\mathbf{g}_f \leftarrow \nabla_{\theta} \mathcal{L}[\mathcal{T}, D, \ell]$$

- Gradient descent on \mathbf{c}_ϕ :

$$b_t \leftarrow \arg \min_i \ell[\mathbf{f}_i(\mathbf{x}_t), \mathbf{y}_t] \forall t$$

$$\mathbf{g}_c \leftarrow \nabla_{\phi} \sum_{(\mathbf{x}_t, \cdot) \in D} \text{CrossEntropy}[\text{softmax}(\mathbf{c}_\phi(\mathbf{x}_t)), b_t]$$

Occam's Razor – The simpler the better!

- How to characterize simplicity? → Description length!
- Finding the minimum description length is a hard problem!
- Therefore Heuristic for description length:

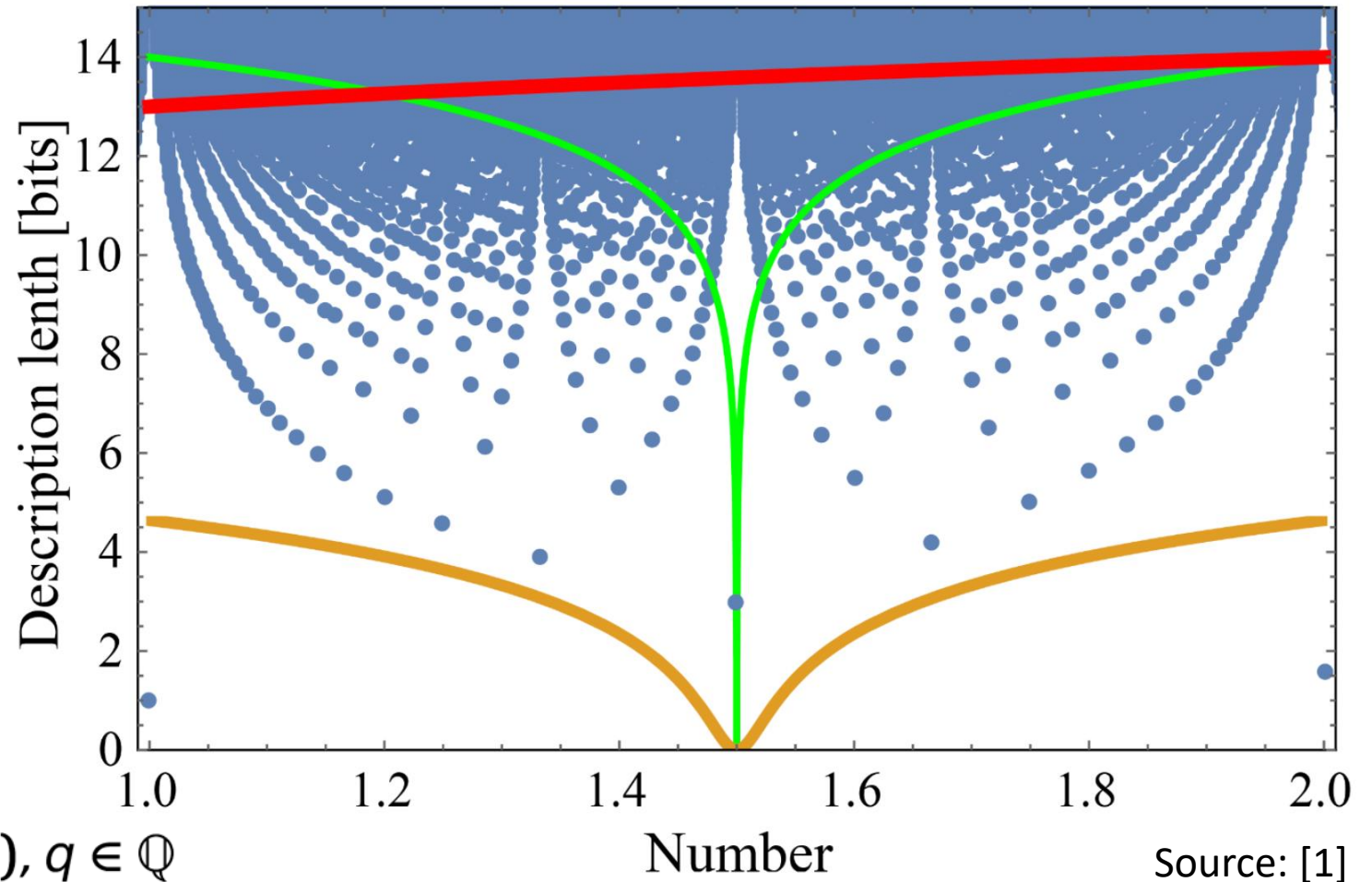
$$DL(n) = \log_2(n), n \in \mathbb{N}$$

$$DL(m) = \log_2(1 + |m|), m \in \mathbb{Z}$$

$$DL(q = \frac{m}{n}) = \log_2((1 + |m|)n), q \in \mathbb{Q}$$

$$DL(r) = \log_+ \left(\frac{r}{\epsilon} \right), r \in \mathbb{R} \text{ and } \log_+(x) = \frac{1}{2} \log_2(1 + x^2)$$

Occam's Razor



$$DL(q = \frac{m}{n}) = \log_2((1 + m)n), q \in \mathbb{Q}$$

$$DL(r) = \log_+ \left(\frac{r}{\epsilon} \right), r \in \mathbb{R} \text{ and } \log_+(x) = \frac{1}{2} \log_2(1 + x^2) \quad \text{Precision: } \epsilon = 2^{-14} \approx 6 * 10^{-5}$$

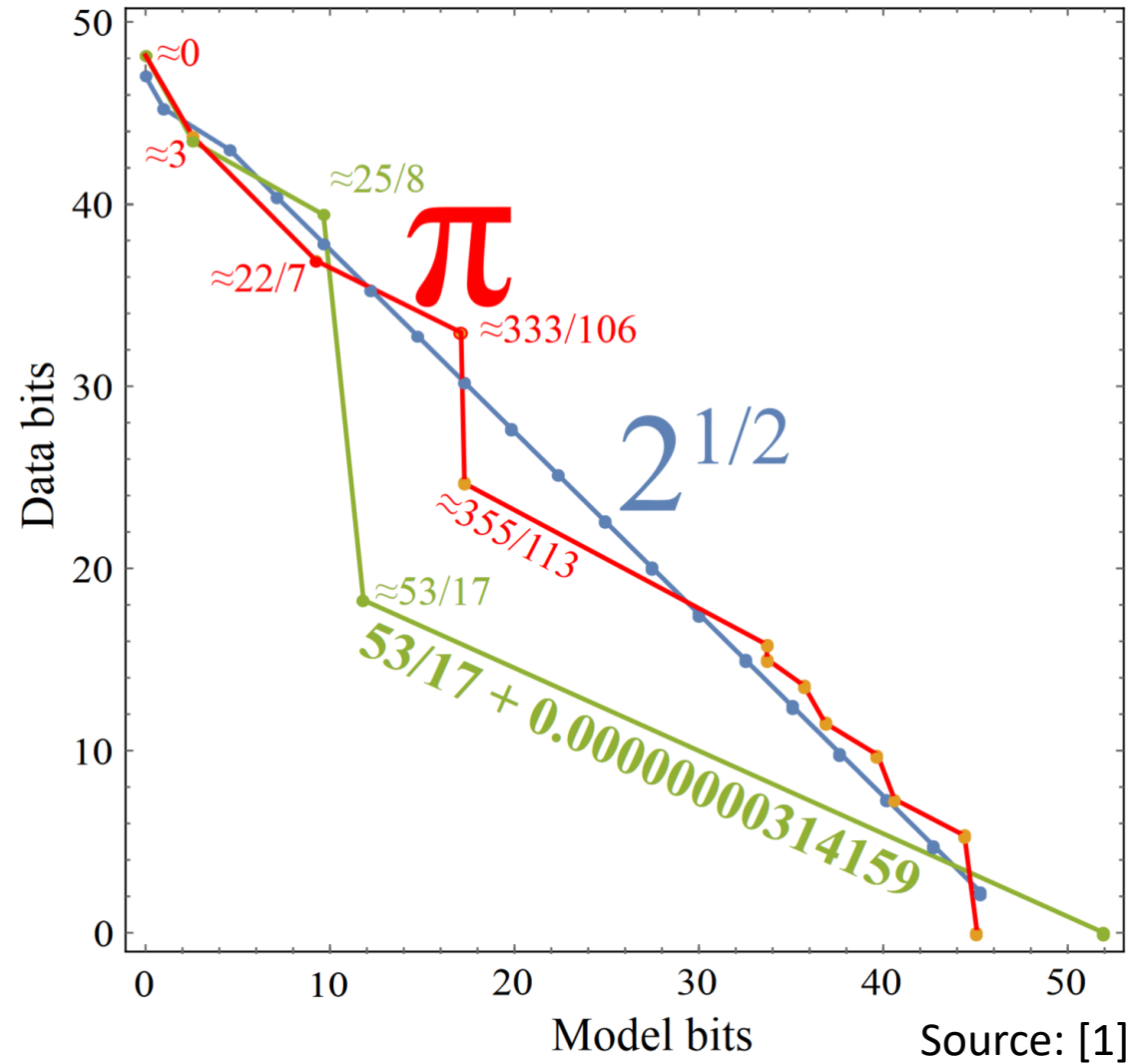
Occam's Razor

- The total description length (DL):

$$DL(\mathcal{T}, D) = DL(\mathcal{T}) + \sum_t DL(\mathbf{u}_t) \quad \mathbf{u}_t = |\hat{\mathbf{y}}_t - \mathbf{y}_t|$$

- First term is the DL of the parameters of the model
- Second term is the DL of the prediction error
- Algorithmically:
 1. Set all parameters to real numbers to minimize $\sum_t DL(\mathbf{u}_t)$
 2. Finding close rational numbers with continued fraction expansion to minimize $DL(\mathcal{T}, D)$

Occam's Razor



Unification

- Goal: Finding underlying similarities between theories and unify them
 - Master theory $\mathcal{T} = (\mathbf{f}_{\mathbf{p}}, \cdot)$, varying the parameter vector $\mathbf{p} \in \mathbb{R}^n$ can generate a continuum of theories.
- Algorithmic idea:
 1. Description length of every prediction function (as symbolic function)
 2. Clustering on the theories
 3. Finding similarities and variations between the symbolic functions in every cluster

Unification – The algorithm

Algorithm 4 AI Physicist: Theory Unification

Require Hub: theory hub

Require C : initial number of clusters

1: for (f_i, c_i) in Hub.all-symbolic-theories do:

2: $dl^{(i)} \leftarrow DL(f_i)$

3: end for

Preparation

4: $\{S_k\} \leftarrow \text{Cluster } \{f_i\}$ into C clusters based on $dl^{(i)}$

Clustering

5: for S_k in $\{S_k\}$ do:

6: $(g_{i_k}, h_{i_k}) \leftarrow \text{Canonicalize}(f_{i_k}), \forall f_{i_k} \in S_k$

7: $h_k^* \leftarrow \text{Mode of } \{h_{i_k} | f_{i_k} \in S_k\}$.

8: $G_k \leftarrow \{g_{i_k} | h_{i_k} = h_k^*\}$

9: $g_{p_k} \leftarrow \text{Traverse all } g_{i_k} \in G_k$ with synchronized steps,

replacing the coefficient by a p_{j_k} when not all coefficients at the same position are identical.

10: $f_{p_k} \leftarrow \text{toPlainForm}(g_{p_k})$

11: end for

Generalization

12: $\mathcal{T} \leftarrow \{(f_{p_k}, \cdot)\}, k = 1, 2, \dots, C$

13: $\mathcal{T} \leftarrow \text{MergeSameForm}(\mathcal{T})$

14: return \mathcal{T}

subroutine Canonicalize(f_i):

s1: $g_i \leftarrow \text{ToTreeForm}(f_i)$

s2: $h_i \leftarrow \text{Replace all non-input coefficient by a symbol } s$

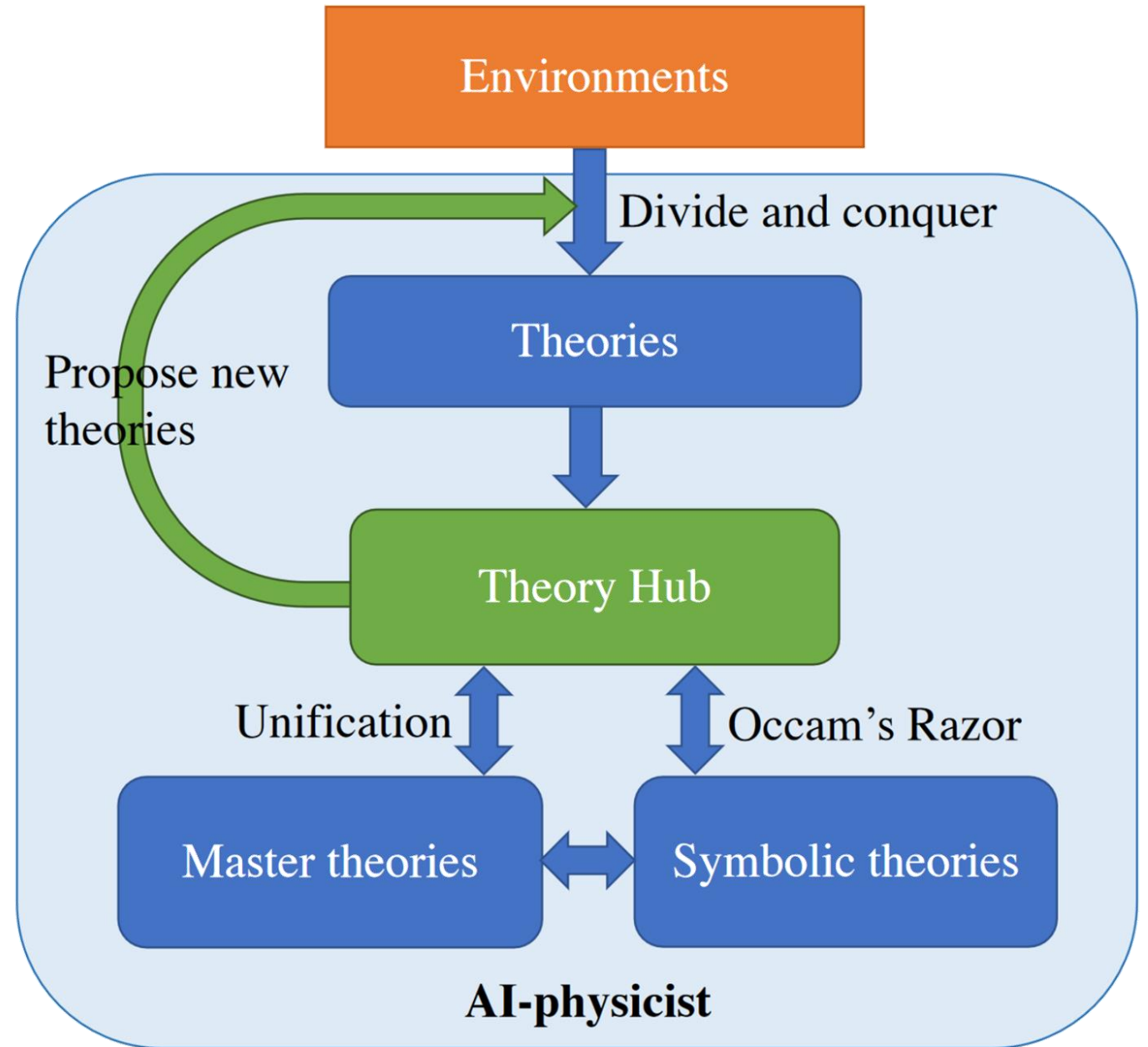
return (g_i, h_i)

Transformation

Example on Blackboard ←

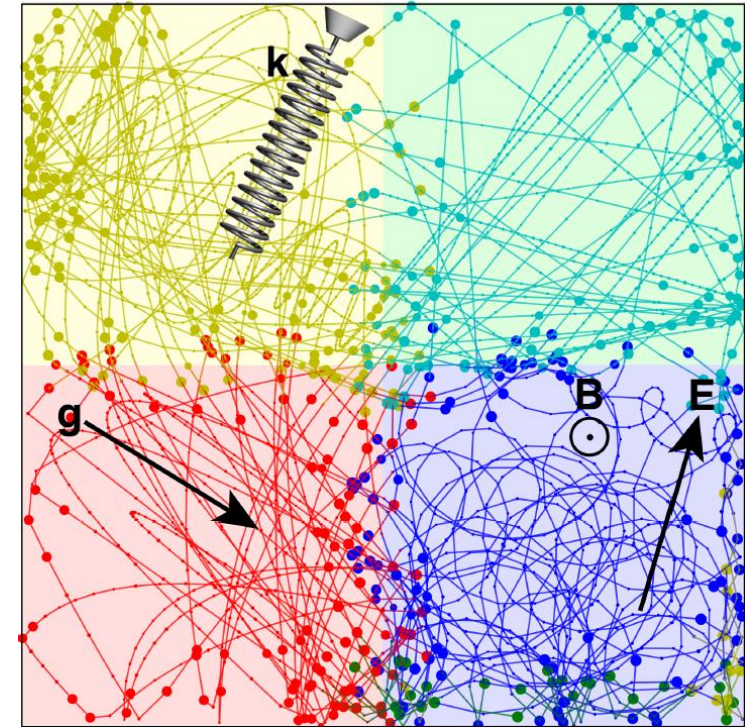
Life long learning

- Idea: Past experiences/knowledge give us the ability to model new environment faster
- Represented in the architecture through the theory hub



The experiments

- AI physicist was tested in two randomized environments:
 1. Mystery world
$$V \propto (ax + by + c)^n$$
 1. Gravity (n=1)
 2. E-field (and optionally uniform B-field) (n=1)
 3. Springs/Hooke's law (n=2)
 4. Bounce boundaries (n= ∞)
 2. Charged double pendulum in two adjacent E-fields



Own visualization, data from [2]

The experiments – results

Benchmark	Baseline	Newborn	AI Physicist
\log_{10} mean-squared error	-3.89	-13.95	-13.88

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \hat{\mathbf{y}}_i)^2$$

Source: [1]

The experiments – results

Benchmark	Baseline	Newborn	AI Physicist
\log_{10} mean-squared error	-3.89	-13.95	-13.88
Classification accuracy	67.56%	100.00%	100.00%

Source: [1]

The experiments – results

Benchmark	Baseline	Newborn	AI Physicist
\log_{10} mean-squared error	-3.89	-13.95	-13.88
Classification accuracy	67.56%	100.00%	100.00%
Fraction of worlds solved	0.00%	90.00%	92.50%

- A domain is solved then:
 - Any rational number must be calculated exactly
 - Any irrational number in the theory must be recovered by an accuracy of 10^{-4}

Source: [1]

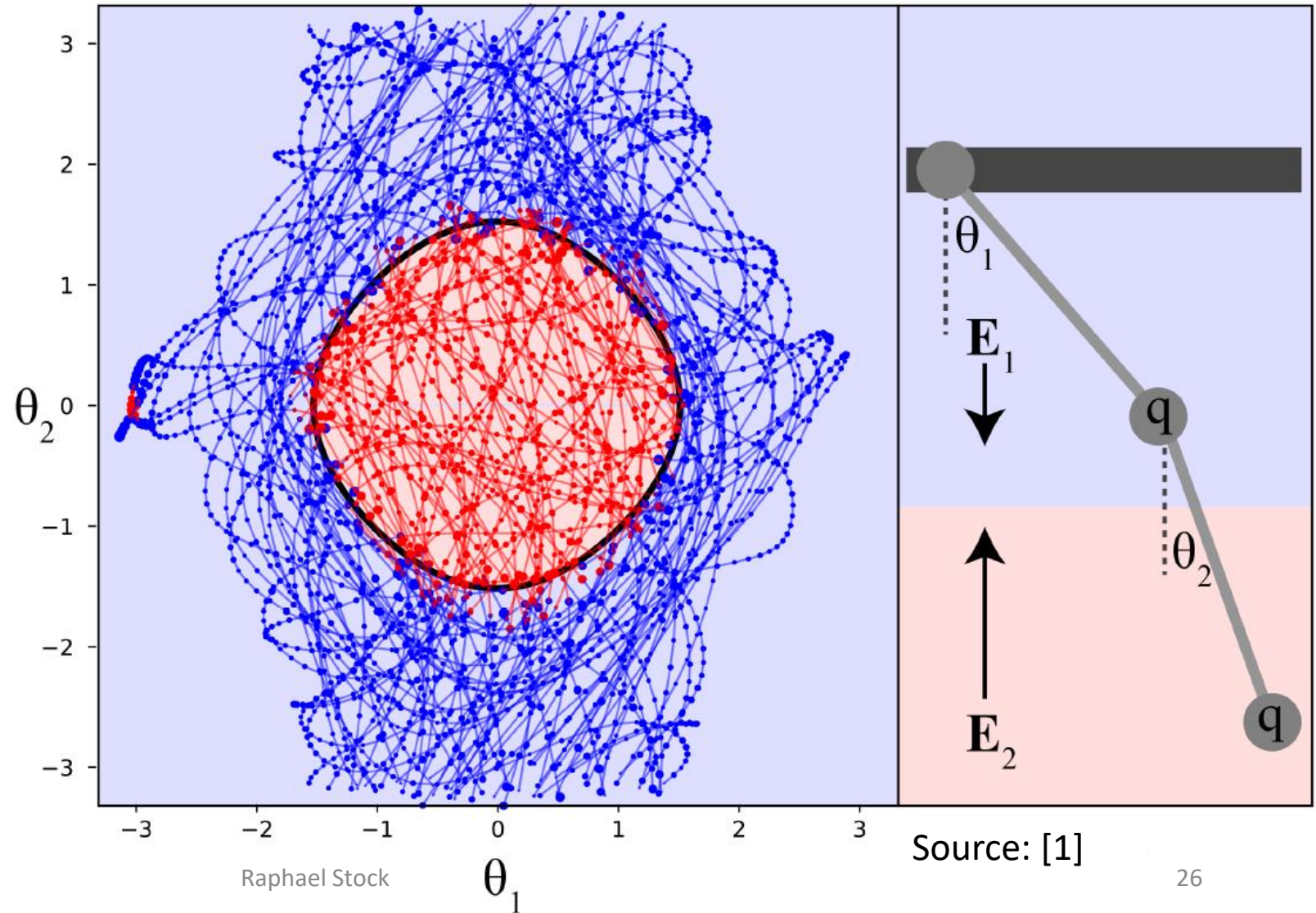
The experiments – results

Benchmark	Baseline	Newborn	AI Physicist
\log_{10} mean-squared error	-3.89	-13.95	-13.88
Classification accuracy	67.56%	100.00%	100.00%
Fraction of worlds solved	0.00%	90.00%	92.50%
Description length for \mathbf{f}	11,338.7	198.9	198.9
Epochs until 10^{-2} MSE	95	83	15
Epochs until 10^{-4} MSE	6925	330	45
Epochs until 10^{-6} MSE	∞	5403	3895
Epochs until 10^{-8} MSE	∞	6590	5100

Source: [1]

The experiments – Charged double pendulum

- No double-pendulum was solved exactly
- Domain classification accuracies:
 - Baseline: 76.9%
 - Newborn: 96.5%



Summary & Conclusion

- Approach towards AI physicist for finding EOM's, whose architecture is orientated on principles a real physicist uses
- The agent works well on simple “mystery world”
- Agent has troubles for harder problems

Discussion

- Classification accuracy is whitewashed by only considering domain centres
- Questionable benchmark with baseline neural network compensating the more complex AI physicist only by double the neurons
- AI physicist obviously outperforms on solved domains and description length because of architecture
- Mystery worlds are limited in their appearance

Thank you for your Attention!

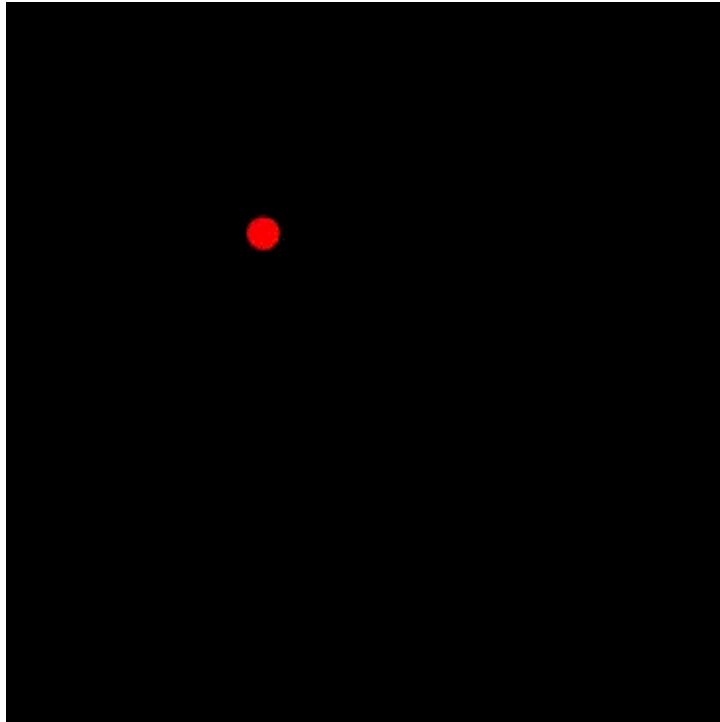
Sources

[1] Tailin Wu, Max Tegmark. Toward an AI Physicist for Unsupervised Learning

[2] <https://space.mit.edu/home/tegmark/aiphysicist.html>

[3] <https://www.sueddeutsche.de/panorama/weltrekord-im-kopfrechnen-in-11-8-sekunden-zur-13-wurzel-1.664320>

Further questions?



Source: [2]

