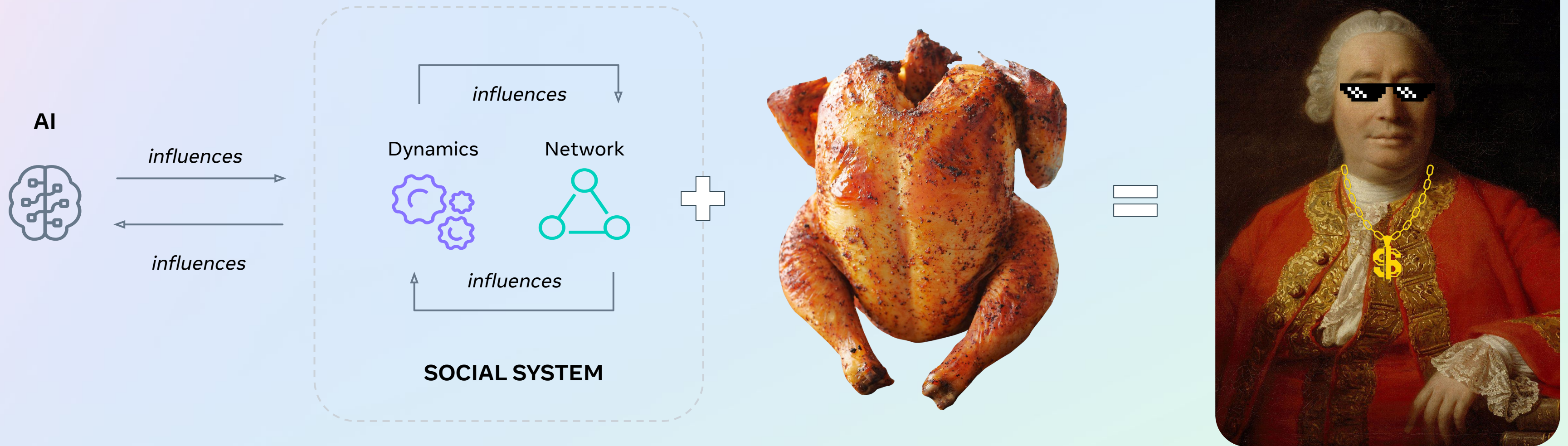# Agenda for today

*Convince you that this equation holds and is relevant to your work*

1) We need to establish **new foundations for modern AI** to work as intended

2) This requires to understand AI interacts with the **social systems** in which they operate

3) This affects **all areas of AI,** not only fairness, responsible AI, etc.

*via insights from complex social systems* _

can we **understand** ...

Are our **benchmarks** give insights into the intended tasks or do they project a false image of quality?

Will naïve **scaling** solve all our problems?

...ed to ...mpressive ... ways at the same t...

Tell me if this is an Iris or not.

Answer any question truthfully about any object in the known universe.

# AI Paradigm Shift_
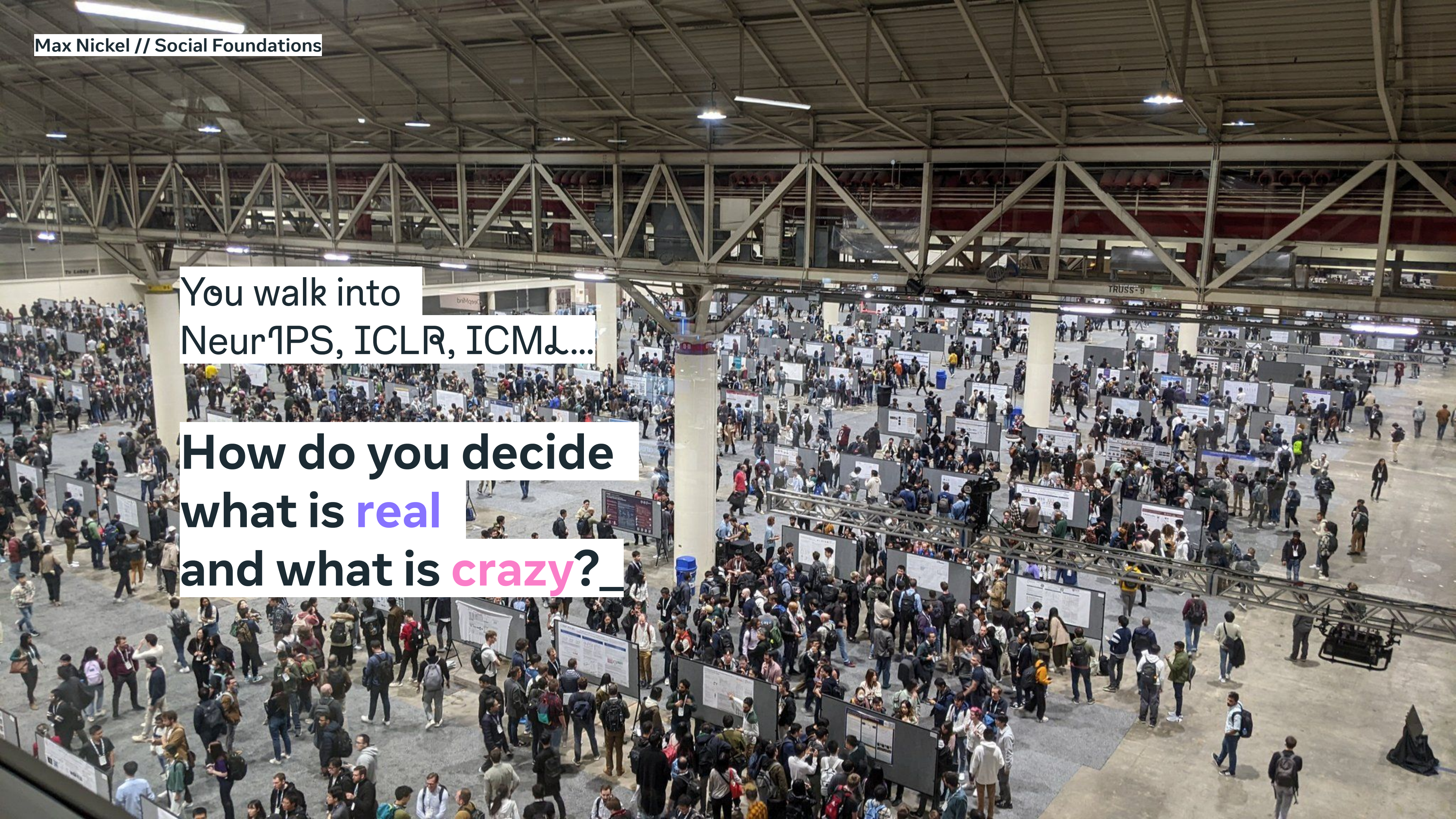
Our theory (read justification)
is built for this case

Does it still hold for what we are doing now?
Easy: Obviously not...
**Harder: CAN IT EVER** be valid?

?

# AI Paradigm
# Shift_

*Pivoine de la Chine*

*Paeonia*

*P.J. Redouté _ 39.*

*Victor*

You walk into
NeurIPS, ICLR, ICML...
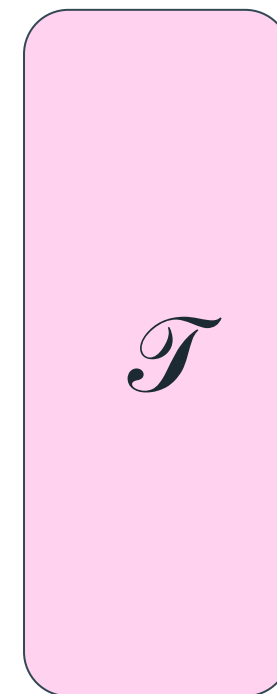
How do you decide
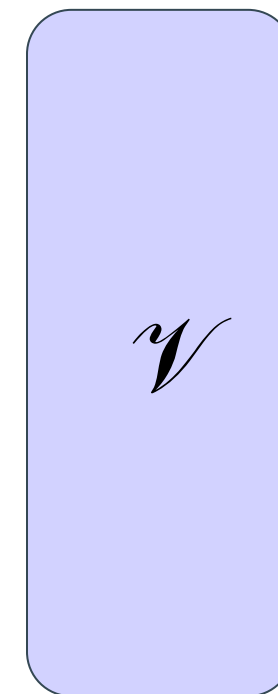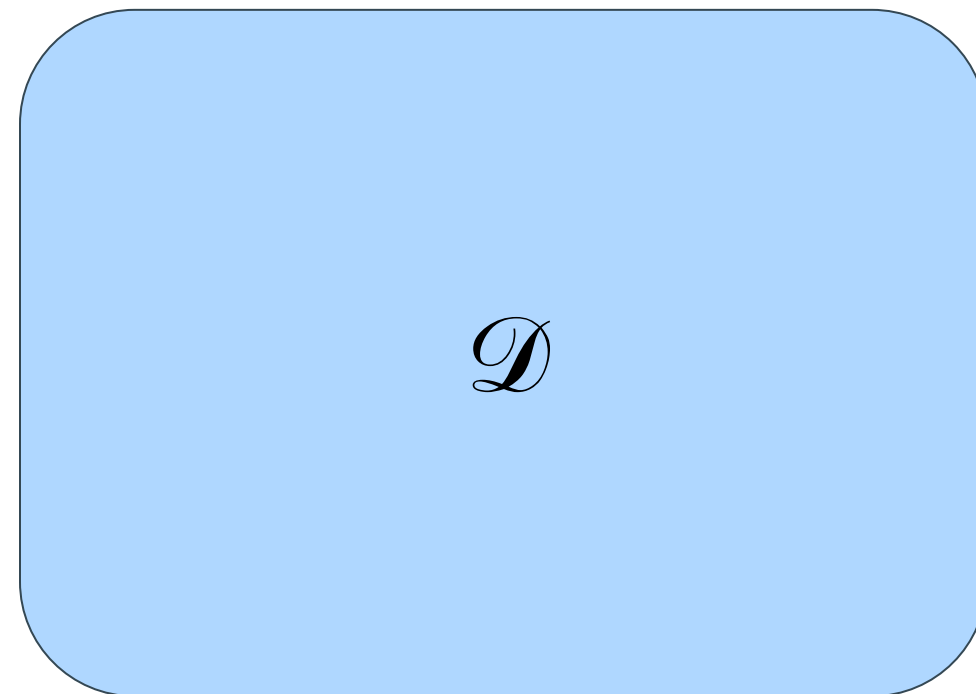what is **real**
and what is **crazy?**_

# Train-Test Paradigm

THE dominant / only
approach to model
validation in modern ML.

- ○ Training set $\mathscr{D}$
- ○ Validation set $\mathscr{V}$
- ○ Test set $\mathscr{T}$

Rapid model validation via
the train-test paradigm
has been a key driver for
the breathtaking progress
in machine learning and AI
(e.g. see *Bottou 2015*).

**The Design and Analysis of Pattern
Recognition Experiments**

By W. H. HIGHLEYMAN

(Manuscript received March 2, 1961)

This is the only thing we care about.

$$\mathscr{D}$$

$$\mathscr{V} \quad \mathscr{T}$$

# **Train-Test** Paradigm

The Design and Analysis of Pattern Recognition Experiments

By W. H. HIGHLEYMAN

(Manuscript received March 2, 1961)

Rapid model validation via the train-test paradigm has been a key driver for the breathtaking progress in machine learning and AI (e.g. see *Bottou 2015*).

- ○ Estimated model $h$
- ○ True world $f$
- ○ Loss function $\ell$
- ○ Target distribution $T$
- ○ Test set $\mathscr{T}$

We usually focus on this part…

This is the only thing we care about.

$$\mathscr{D} \qquad \mathscr{V} \qquad \frac{1}{m}\sum_{x\in\mathcal{T}}\ell(h(x), f(x)) \; \approx \qquad \mathbb{E}_{X\sim\mathsf{T}}\left[\ell(h(X), f(X))\right]$$

because we assume this part is fine.

# **Train-Test** Paradigm

The Design and Analysis of Pattern
Recognition Experiments

By W. H. HIGHLEYMAN
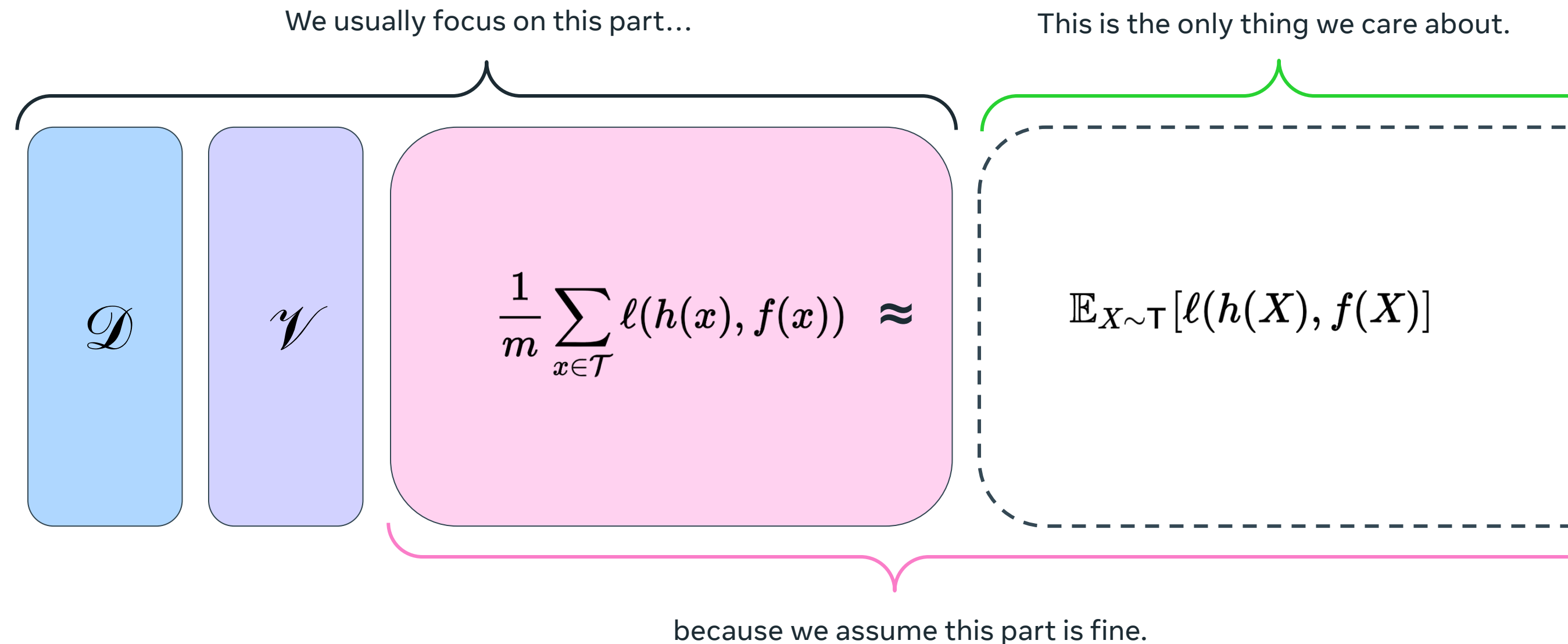
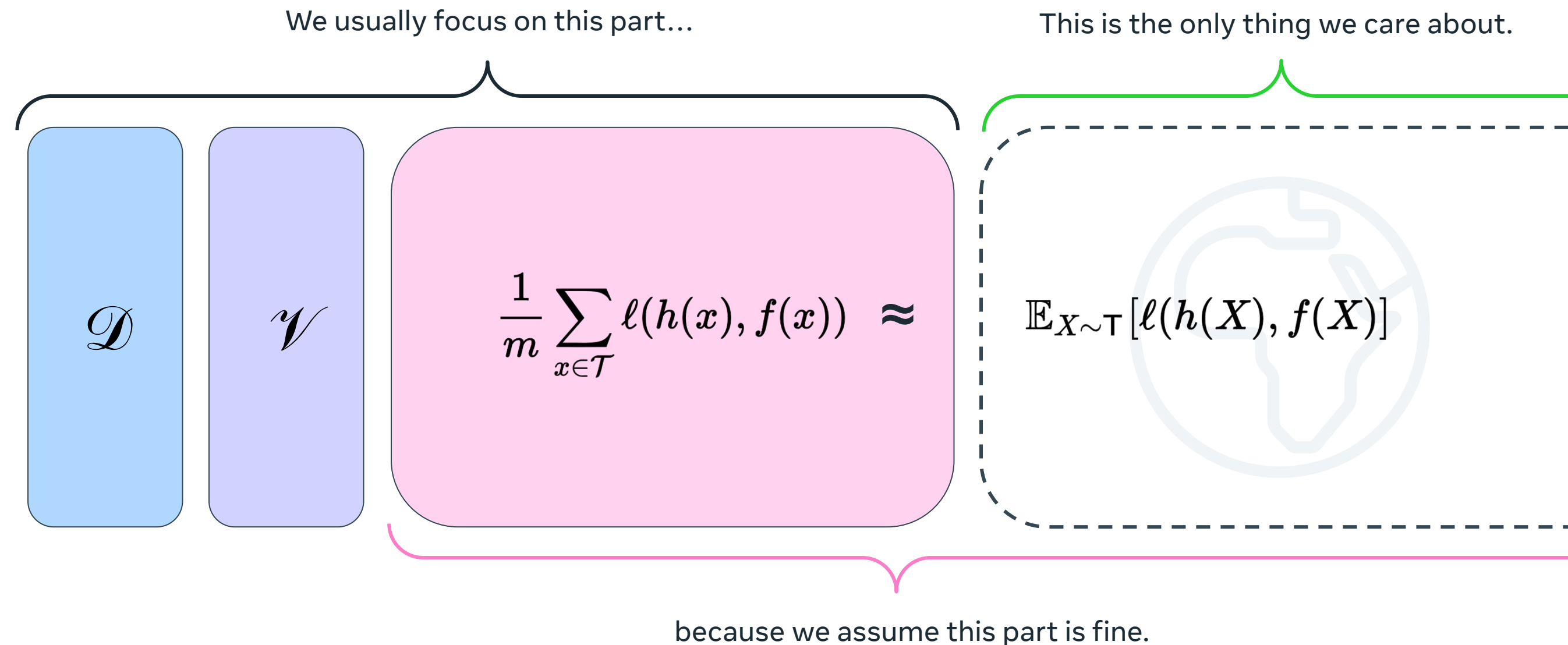(Manuscript received March 2, 1961)

Rapid model validation via
the train-test paradigm
has been a key driver for
the breathtaking progress
in machine learning and AI
(e.g. see *Bottou 2015*).

○ Estimated model $h$
○ True world $f$
○ Loss function $\ell$
○ Target distribution $T$
○ Test set $\mathscr{T}$

We usually focus on this part…

This is the only thing we care about.

$$\mathscr{D} \qquad \mathscr{V} \qquad \frac{1}{m}\sum_{x\in\mathcal{T}}\ell(h(x), f(x)) \ \approx \ \mathbb{E}_{X\sim\mathsf{T}}[\ell(h(X), f(X)]$$

because we assume this part is fine.
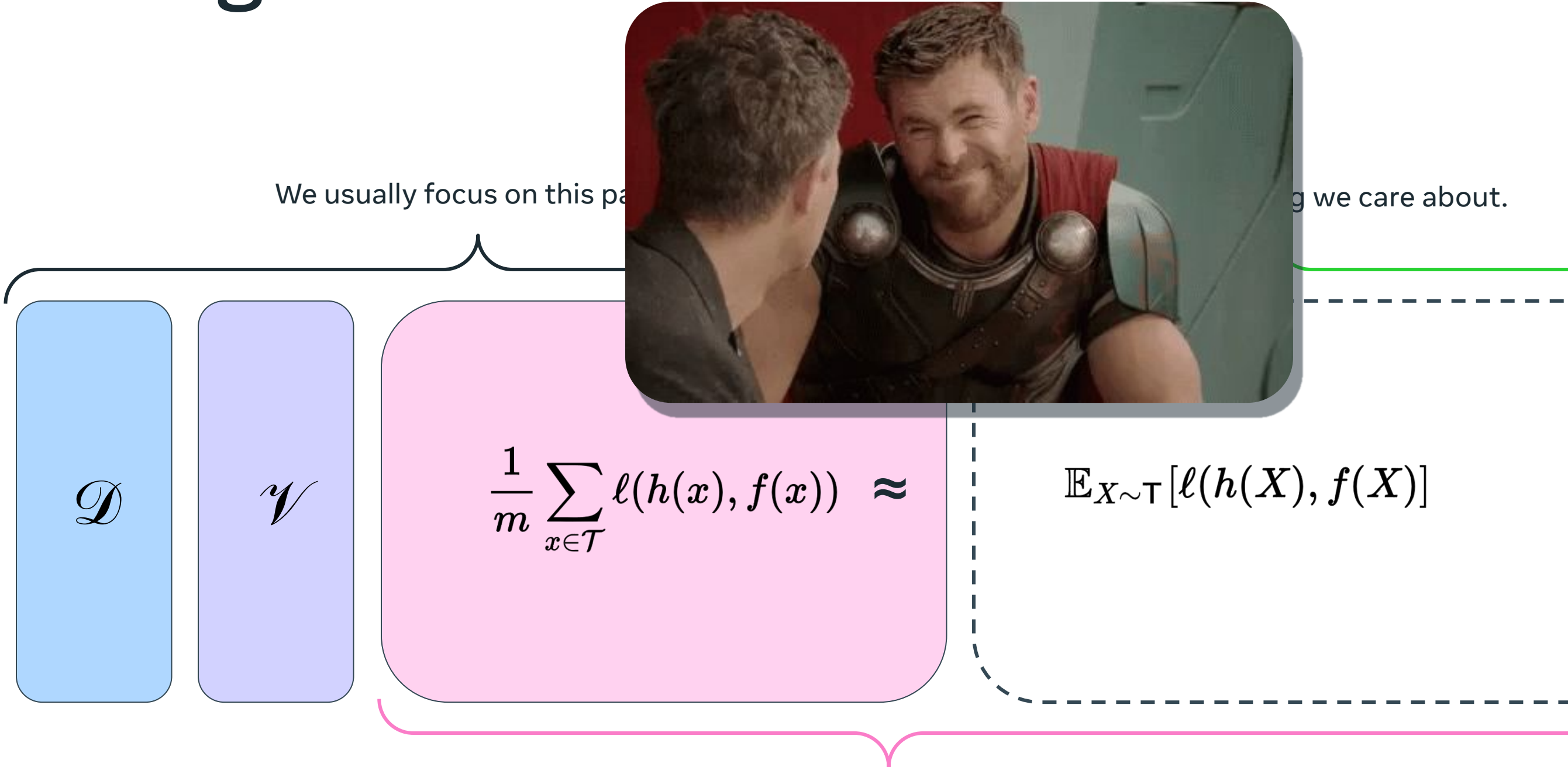
# Train-Test Paradigm

Rapid model validation via the train-test paradigm has been a key driver for the breathtaking progress in machine learning and AI (e.g. see *Bottou 2015*).

- ○ Estimated model $h$
- ○ True world $f$
- ○ Loss function $\ell$
- ○ Target distribution $T$
- ○ Test set $\mathscr{T}$

We usually focus on this pa...                                        ...g we care about.

$$\mathscr{D}$$

$$\mathscr{V}$$

$$\frac{1}{m}\sum_{x\in\mathcal{T}}\ell(h(x),f(x)) \;\approx$$

$$\mathbb{E}_{X\sim\mathsf{T}}[\ell(h(X),f(X))]$$

because we assume this part is fine.

# Train-Test Paradigm

Rapid model validation via the train-test paradigm has been a key driver for the breathtaking progress in machine learning and AI (e.g. see *Bottou 2015*).

- Estimated model $h$
- True world $f$
- Loss function $\ell$
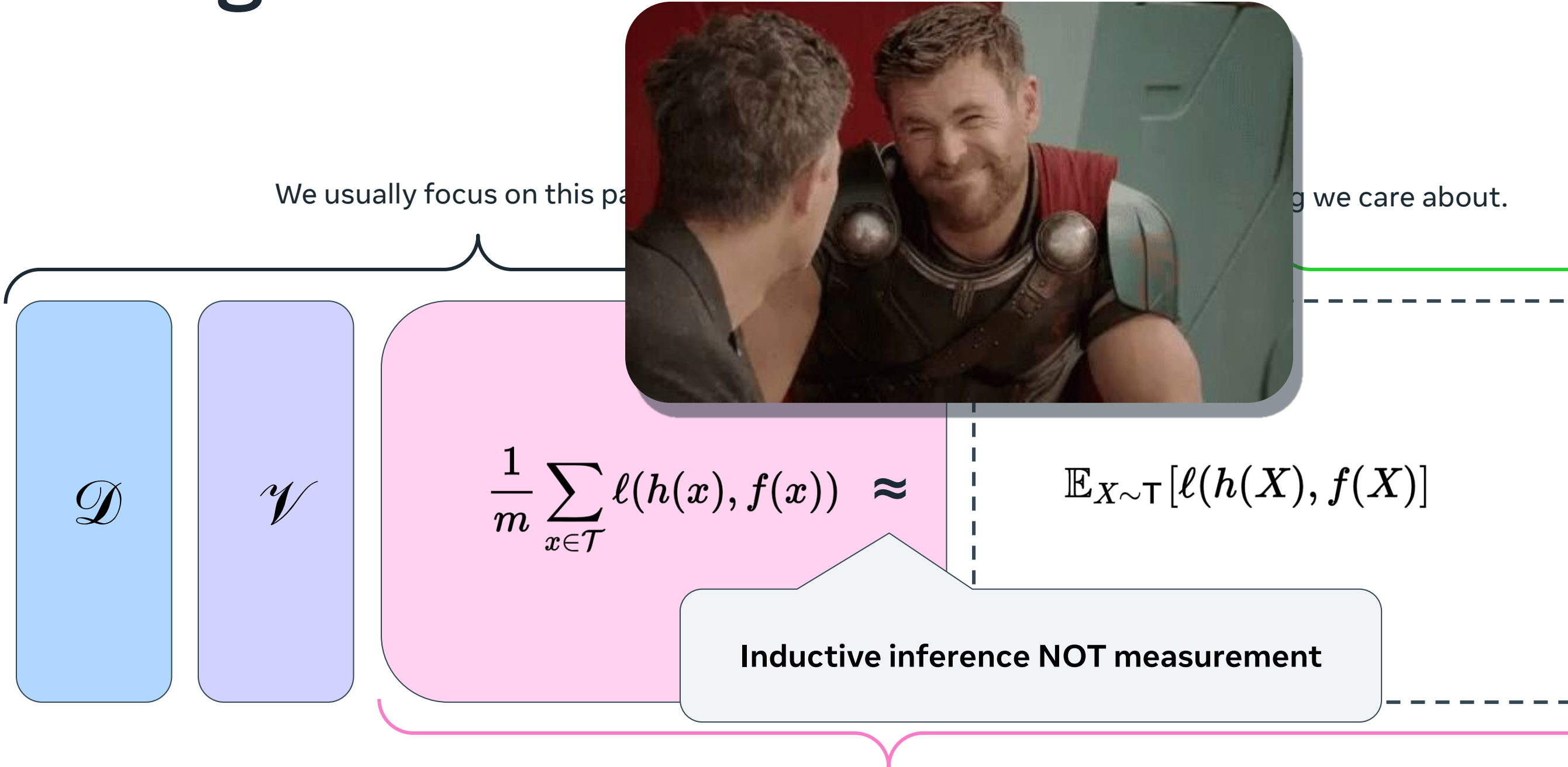- Target distribution $T$
- Test set $\mathscr{T}$

We usually focus on this pa... ...g we care about.

$$\mathscr{D} \quad \mathscr{V} \quad \frac{1}{m}\sum_{x\in\mathcal{T}}\ell(h(x), f(x)) \approx \quad \mathbb{E}_{X\sim\mathsf{T}}\left[\ell(h(X), f(X))\right]$$

**Inductive inference NOT measurement**

because we assume this part is fine.

# Is **Induction** Possible?

Fundamental question in science, dating back at least to Hume's **problem of induction** *(1739)*:

"*even after the observation of the frequent or constant conjunction of objects, we have no reason to draw any inference concerning any object beyond those of which we have had experience*"

In ML, Wolpert's **No-Free-Lunch** theorem *(1996)* established formally that statistical learning/prediction is **impossible without making assumptions** about the world.

# *Staying close to*
# **Non-Uniformity** of Nature

"Domestic animals expect food when they see the person who usually feeds them. We know that all these rather crude expectations of uniformity are liable to be misleading. The man who has fed the chicken every day throughout its life at last wrings its neck instead, showing that **more refined views as to the uniformity of nature would have been useful to the chicken**."

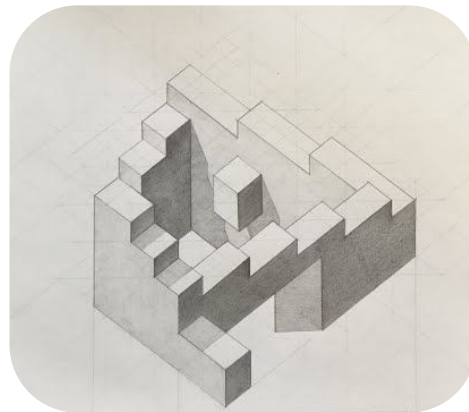Bertrand Russell — *The Problems of Philosophy*

# Ontological **Parsimony**

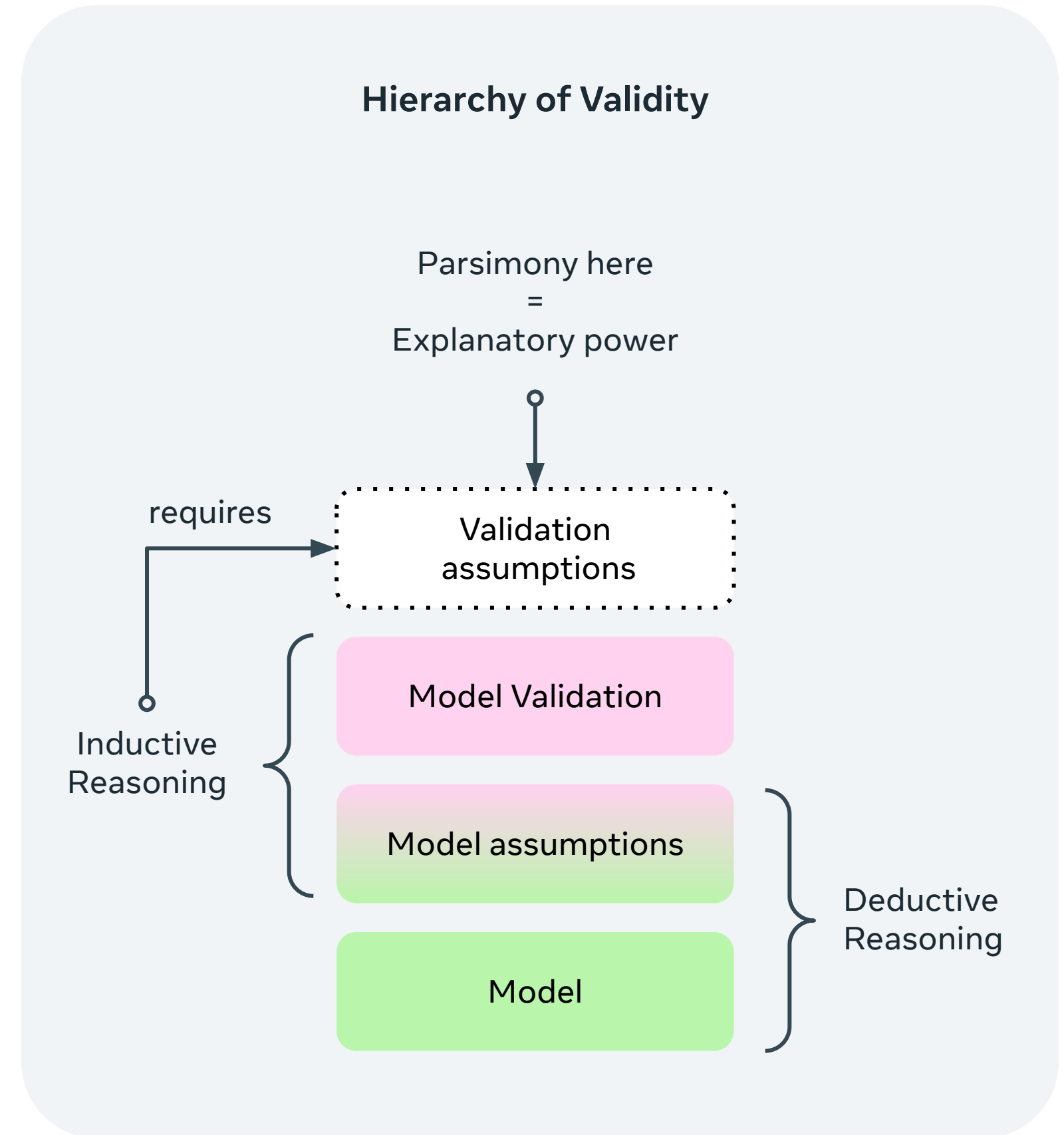Hume's argument, or to say it w/ W. v. O. Quine (1969)

*"The Humean condition is the human condition."*

But we can ask: Are there some **reasonable assumptions** that we might be willing to make that can **neutralize the problem of induction**?

- ○ We can validate model assumptions via model validation.
- ○ We never can validate assumptions necessary to ensure the validity of the model validation itself w/o **circular reasoning or infinite regress.**

Preferring **parsimonious** hypotheses is rational — they have greater **explanatory power** than less parsimonious alternatives. *(Baker, 2003)*

**Hierarchy of Validity**

Parsimony here
=
Explanatory power

requires → Validation assumptions

Inductive Reasoning

Model Validation

Model assumptions

Model

Deductive Reasoning

# David Hume *hates* this one simple trick

True risk $L_{fh}^{\mathsf{T}}$  Empirical risk $\theta$

$$\underbrace{\mathbb{E}_{X \sim \mathsf{T}}[\ell(h(X), f(X))]}_{} \approx \underbrace{\frac{1}{m} \sum_{x \in \mathcal{T}} \ell(h(x), f(x))}_{}$$
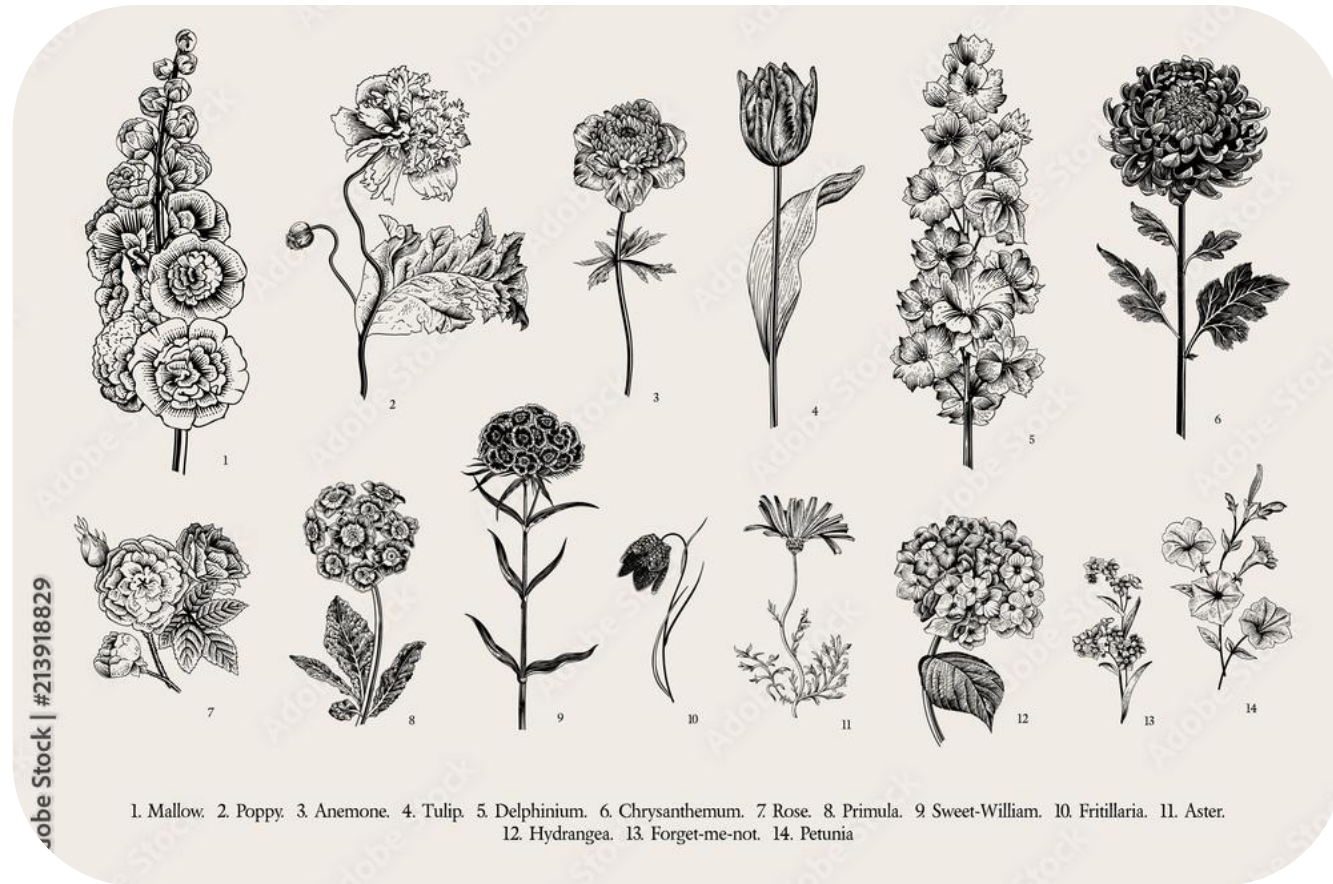
**I.I.D. assumption**
enables straightforward proof via Hoeffding's inequality

$$\mathbb{P}_{\mathcal{T} \sim \mathsf{T}^m} \left( |\theta - L_{fh}^{\mathsf{T}}| \leq \sqrt{\frac{\log(2/\delta)}{2m}} \right) \geq 1 - \delta$$

**(ε,δ)–guarantee**

With probability larger than 1 - $\delta$, the error will be smaller than $\varepsilon$
- ○  $\varepsilon$ = accuracy parameter
- ○  $\delta$ = confidence parameter

Data generating system
(Delivery service)

# What we can **justify**

- ○ **Actively collected data** to satisfy IID assumption
- ○ Corresponds to target distribution
- ○ Very costly, does not scale to large data sets
- ○ **Scope: closed** domain

# What we are **doing**

- ● **Passively collected data** from *some* data generating system
- ● Non-IID, does not correspond to target distribution
- ● Cheap, easy to scale when access to the system
- ● Provides massive datasets <u>required</u> for modern AI
- ● **Scope: open** domain

# Complex Social Systems

Our **data generating systems** are _

- ○ Internet platforms (recommendations) _
- ○ The internet (reasoning & QA) _
- ○ Human knowledge (reasoning & QA) _
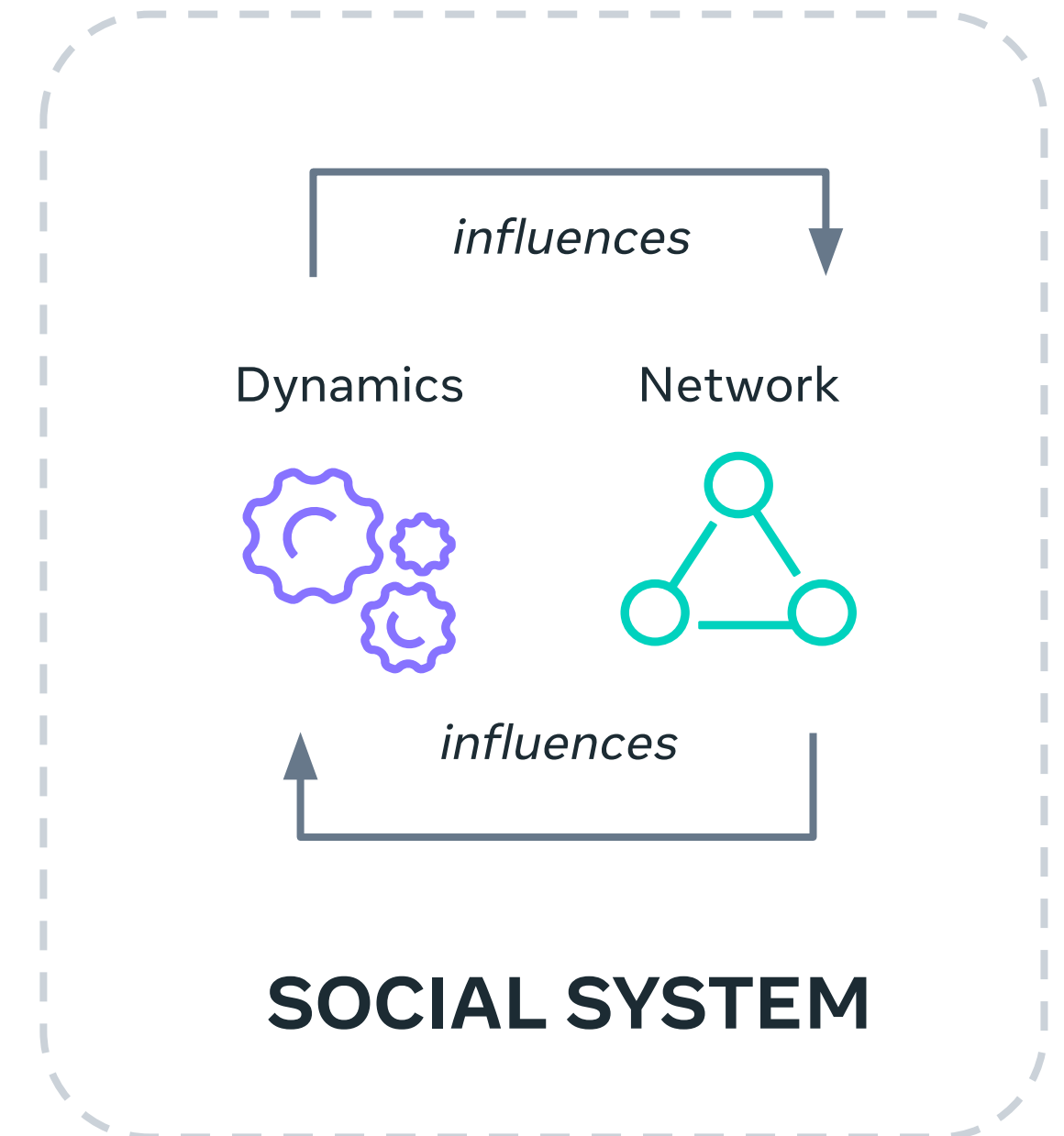
# **Social** Systems
# ≈
# **Complex** Systems

Often, we can understand social systems as complex systems, i.e., as systems with

- **Interactions** of their parts
- Internal **dynamics**
- **Non-linearities** and chaotic behavior
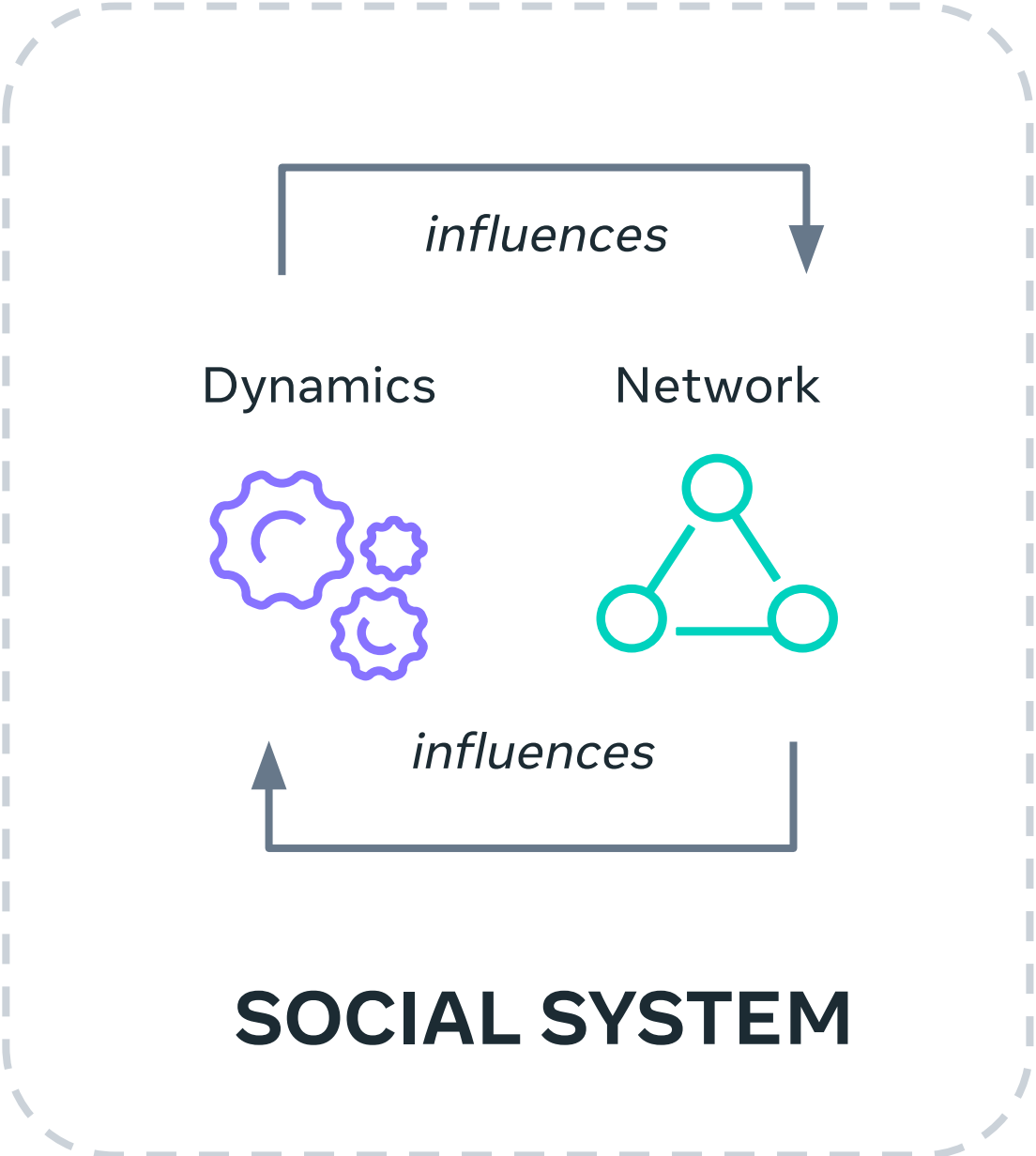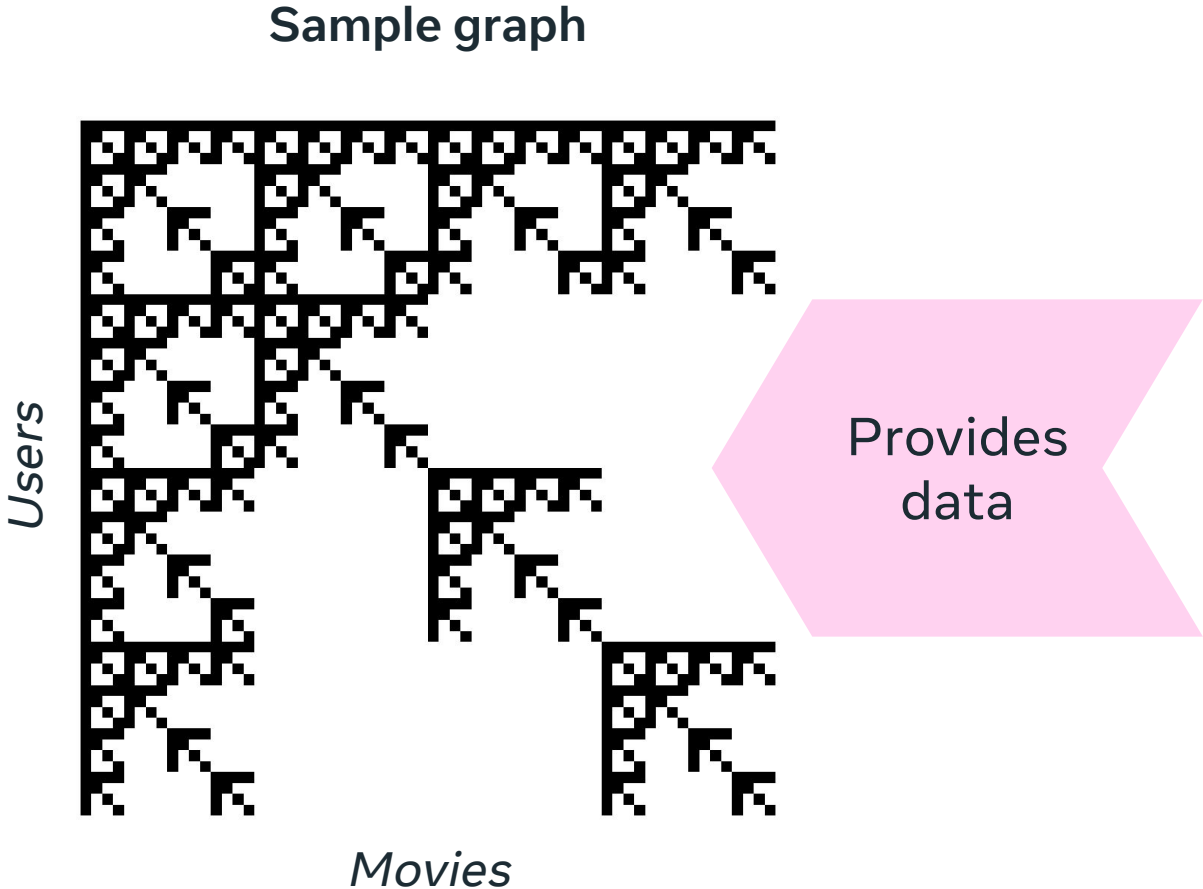- Memory and **feedback**
- **Emergent** properties and behavior

**AI**

Provides data

*influences*

Dynamics          Network
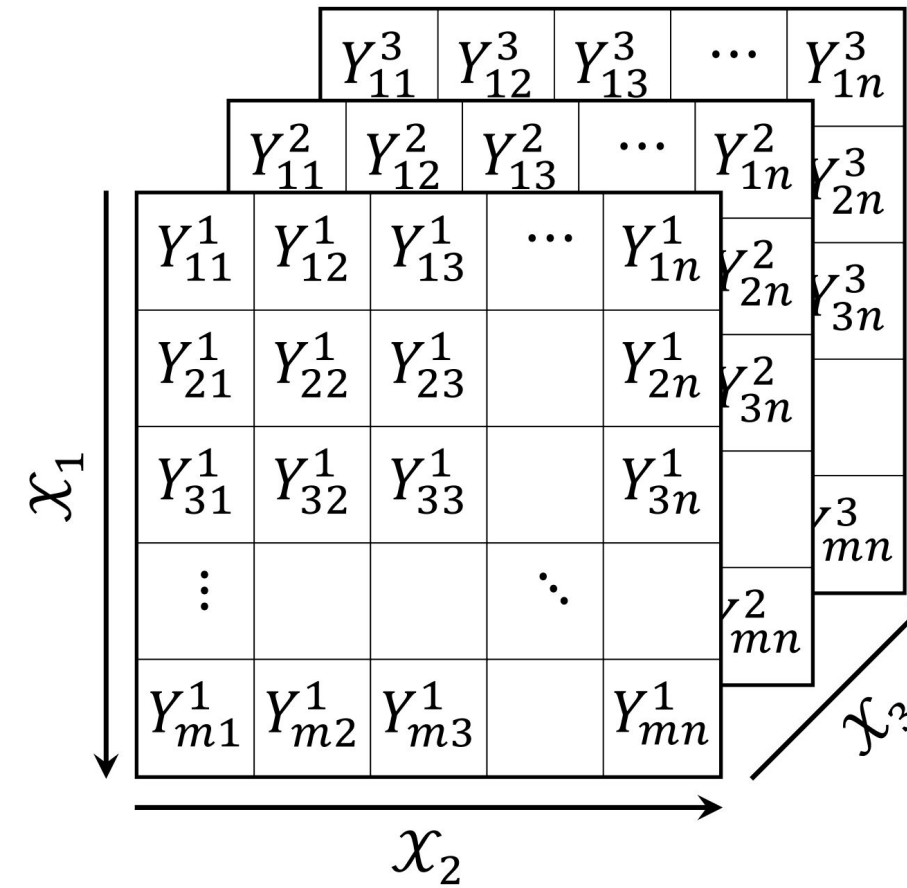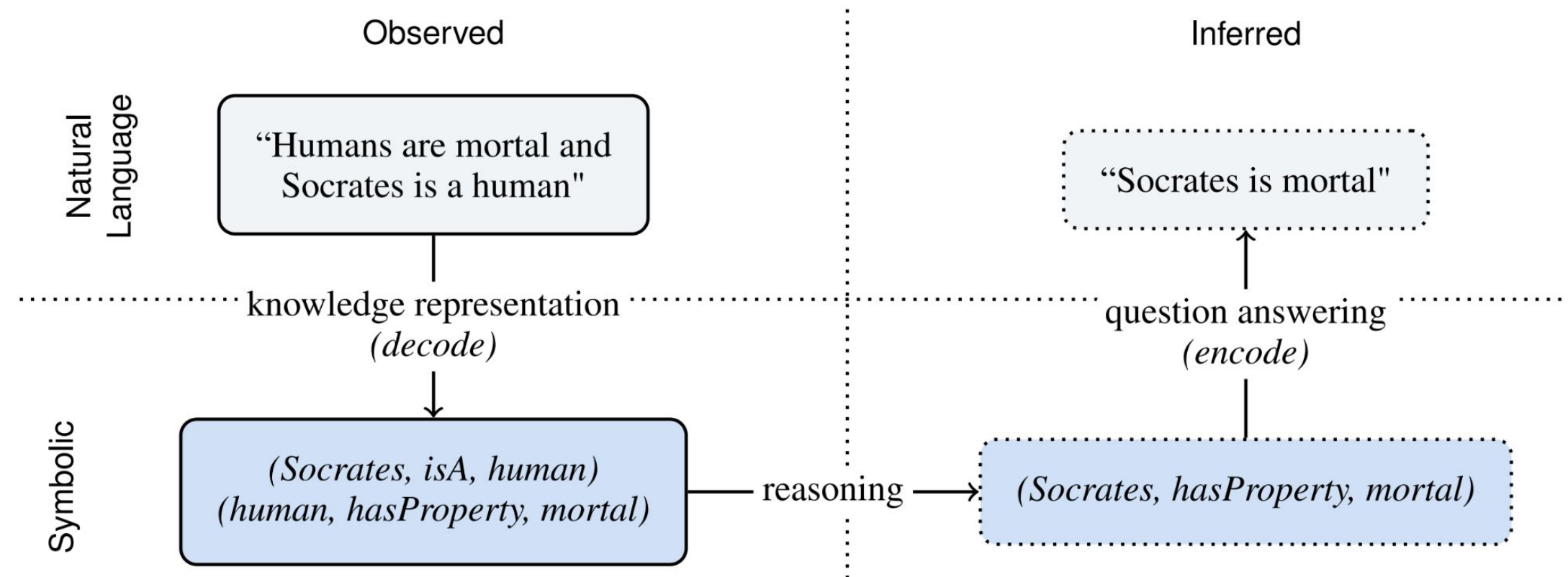
*influences*

**SOCIAL SYSTEM**

# Modeling passive data collection

Formalize data collection via **sample graphs**

**Edge** in a sample graph denotes an observed data point (noise free)
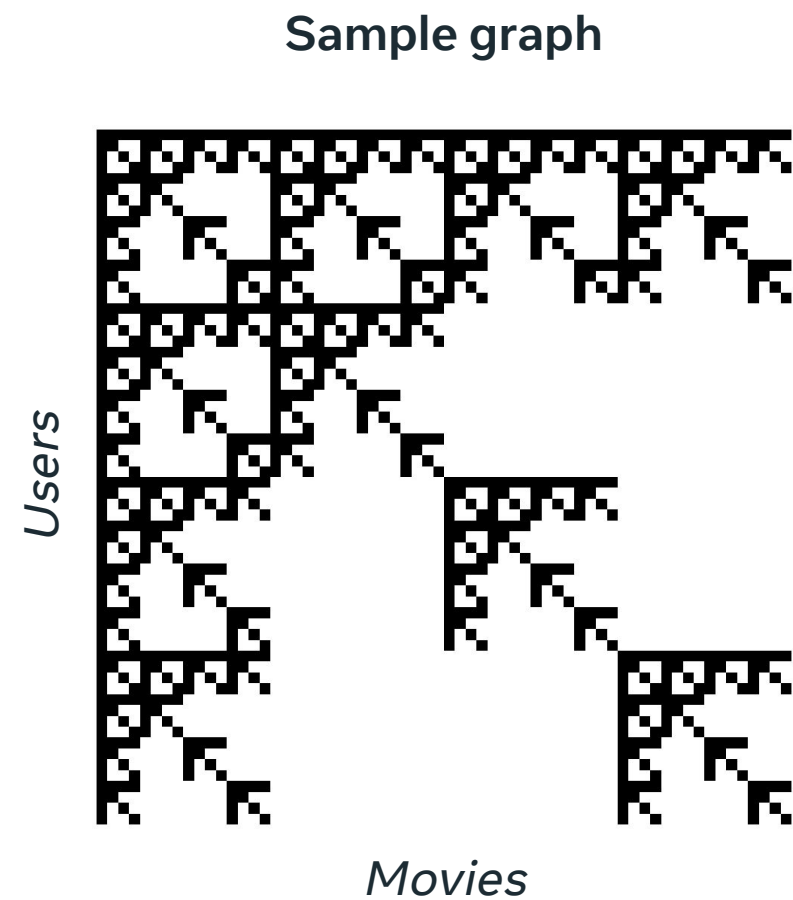
**Sample graph**

*Users*

*Movies*



Provides data

influences

Dynamics          Network

influences

**SOCIAL SYSTEM**

# Modeling passive data collection



$$\mathcal{X}_1, \mathcal{X}_2, \mathcal{X}_3$$

$$
\begin{array}{|c|c|c|c|c|}
\hline
Y^3_{11} & Y^3_{12} & Y^3_{13} & \cdots & Y^3_{1n} \\
\hline
\end{array}
$$

Formalize data collection via **sample graphs**

**Edge** in a sample graph denotes an observed data point (noise free)



Observed

Natural Language

"Humans are mortal and Socrates is a human"

Inferred

"Socrates is mortal"

knowledge representation
*(decode)*

question answering
*(encode)*

Symbolic

(Socrates, isA, human)
(human, hasProperty, mortal)

→ reasoning →

(Socrates, hasProperty, mortal)

# Modeling passive data collection

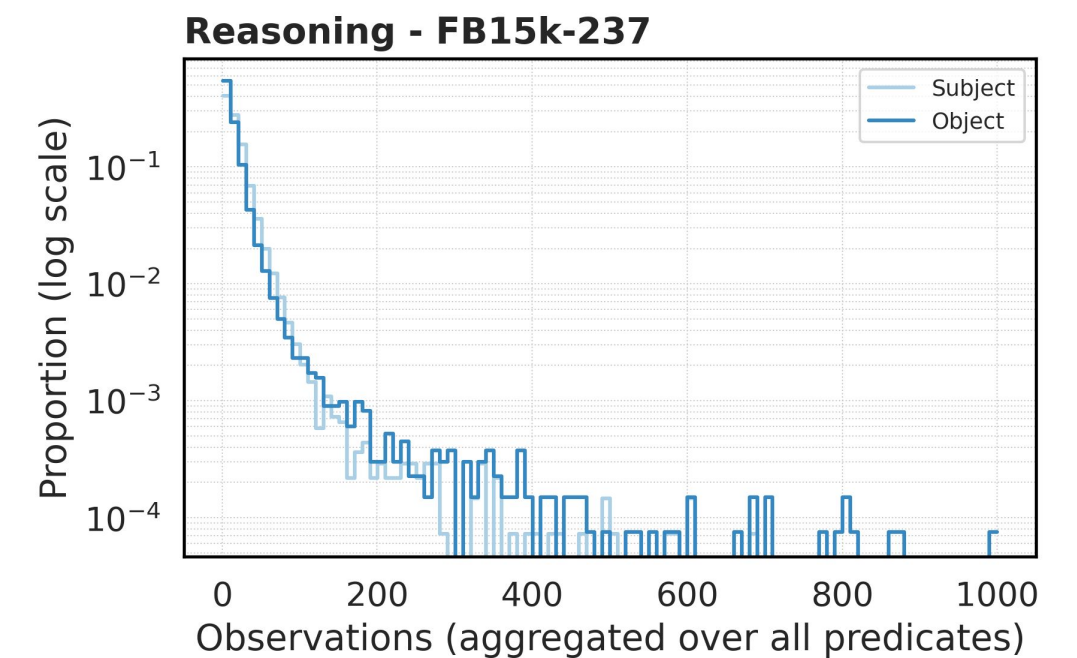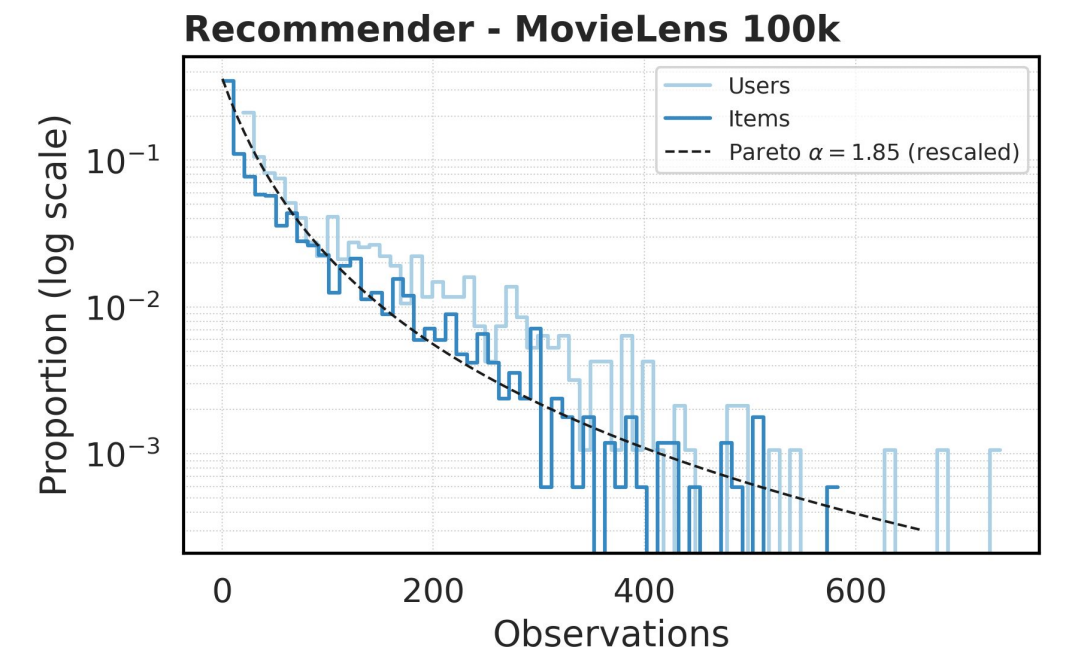What kind of data do our sample generating systems generate? Since they are social systems, they usually exhibit

i) **Sample bias**
ii) **Heavy-tailed / Power-law distributed observations**

$$\mathbb{P}(K_1 > k) = u_1(k)k^{-\alpha_1} \quad \text{and} \quad \mathbb{P}(K_2 > k) = u_2(k)k^{-\alpha_2}$$

caused by well-documented processes such as popularity bias, homophily, feedback loops, etc

**Sample graph**



*Users*

*Movies*

Has Structure

**Recommender - MovieLens 100k**



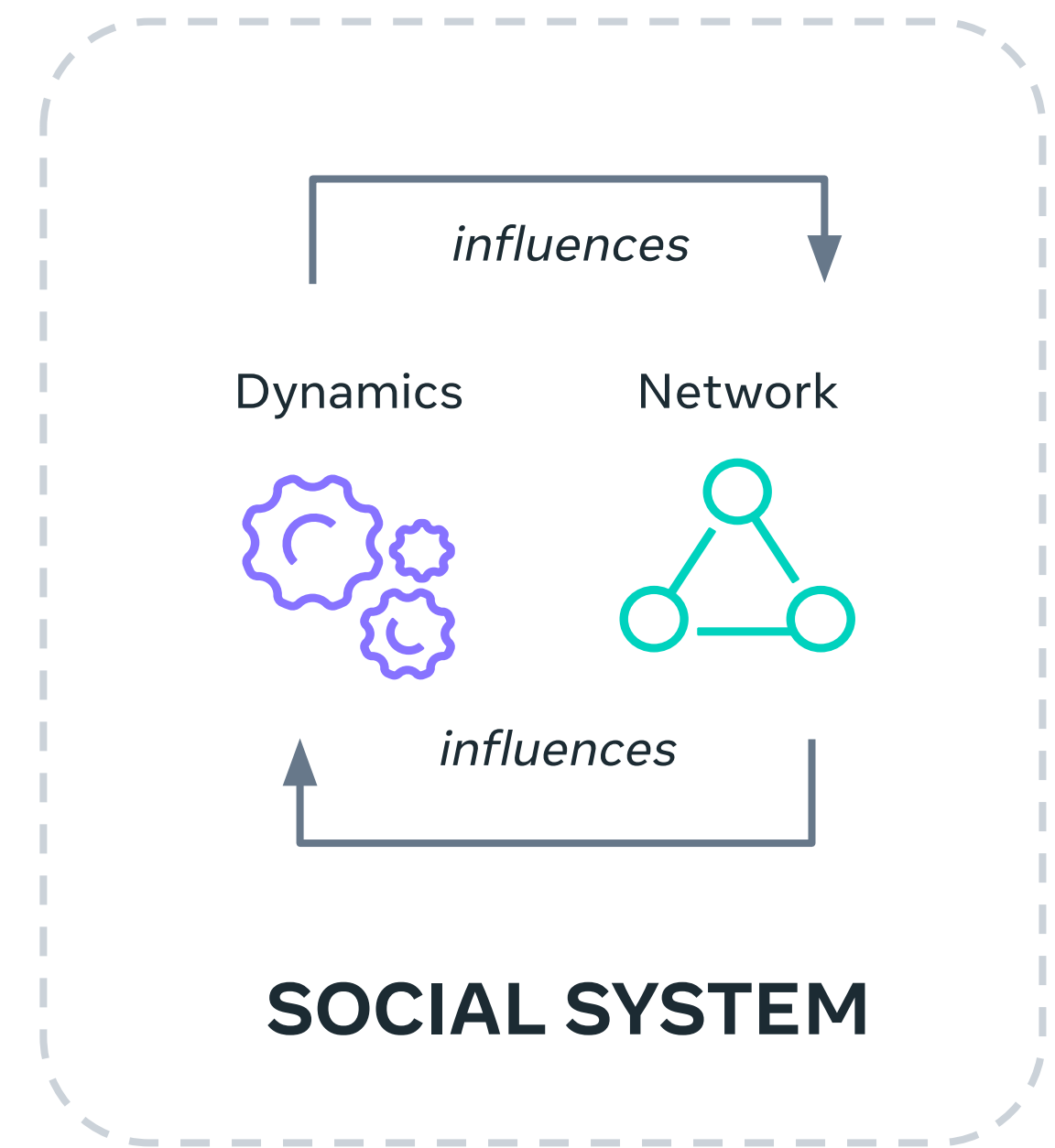**Reasoning - FB15k-237**

# Modeling passive data collection

Formalize data collection via **sample graphs**

**Edge** in a sample graph denotes an observed data point (noise free)

**Sample graph**



**SOCIAL SYSTEM**

Dynamics      Network

*influences*

*influences*

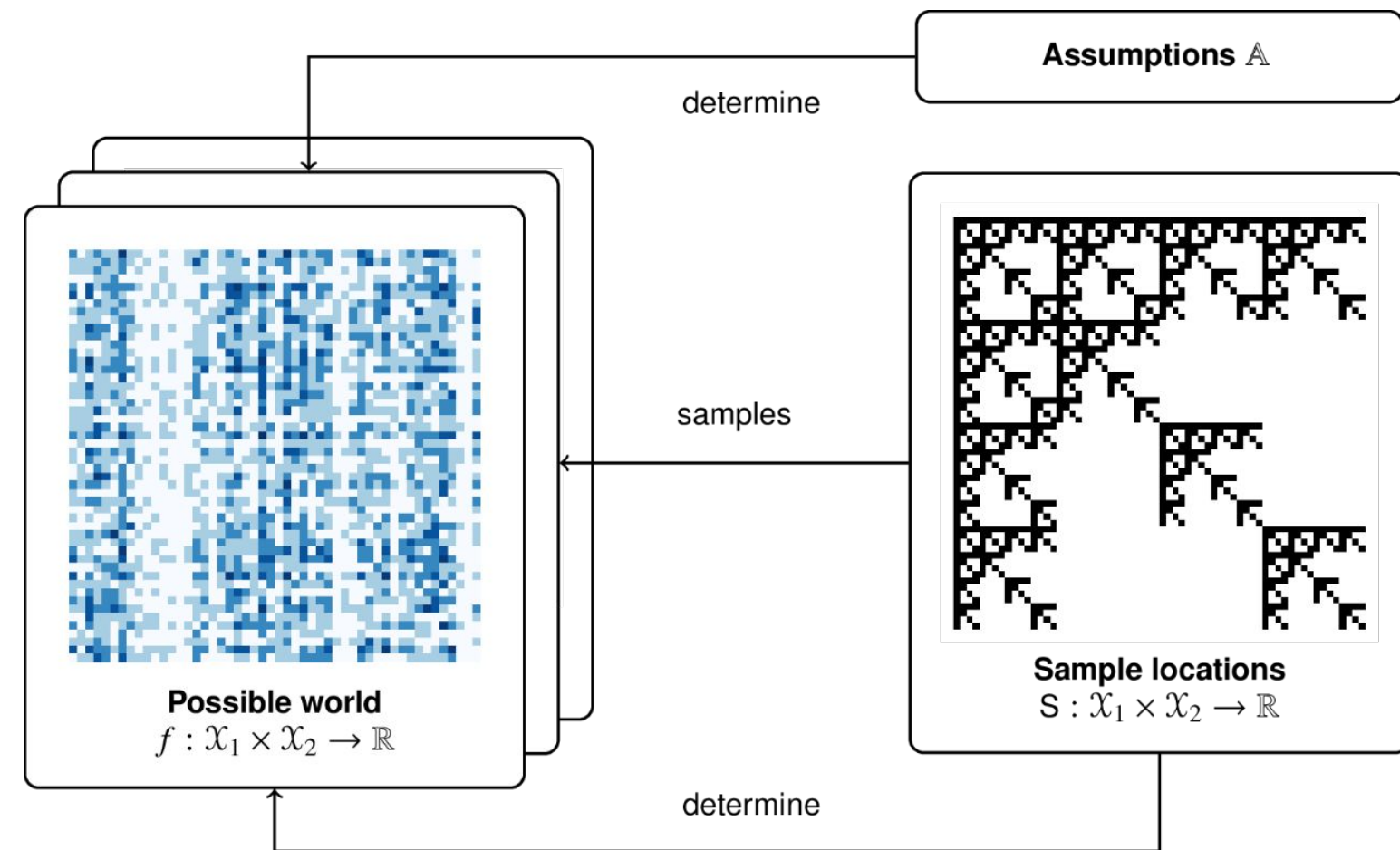| | **Domain** $\mathcal{X}$ | **Possible world** $f$ | **Sample distribution** S | **Target distribution** T |
|---|---|---|---|---|
| Recommender systems | $\mathcal{U} \times \mathcal{I}$ | User preferences | Probability of user interacting with item, heavy-tailed in $\mathcal{U}$ and $\mathcal{I}$ | Uniform, $p_T(u, i) = 1/|\mathcal{U} \times \mathcal{I}|$ |
| Symbolic reasoning | $\mathcal{S} \times \mathcal{P} \times \mathcal{O}$ | Truth value of factoids | Probability of observing factoid, heavy-tailed in $\mathcal{S}$, $\mathcal{P}$, and $\mathcal{O}$ | Uniform, $p_T(s, p, o) = 1/|\mathcal{S} \times \mathcal{P} \times \mathcal{O}|$ |

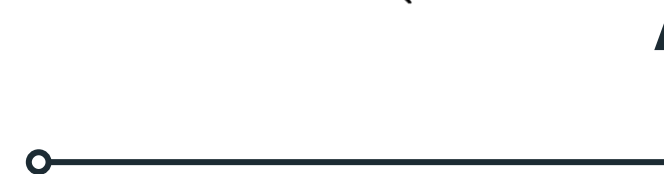# Possible Worlds Semantics

We want to evaluate how well an

○ estimated model $h$

○ approximates the true world $f$

**Observations + assumptions** define
possible worlds that are consistent with both.

**Test validity**: can we bound the error of of *a risk estimator* $\theta$
compared to the true risk over all possible worlds?



Assumptions $\mathbb{A}$

determine

Possible world
$f : \mathcal{X}_1 \times \mathcal{X}_2 \to \mathbb{R}$

samples

Sample locations
$\mathsf{S} : \mathcal{X}_1 \times \mathcal{X}_2 \to \mathbb{R}$

determine

$$\mathbb{P}_{f \sim \mathsf{F}} \left( |\theta - L_{fh}^{\mathsf{T}}| \leq \epsilon \right) \geq 1 - \delta$$

# *A rigorous impossibility result*
# No Free Delivery Service *(Nickel, 2024)*

**Theorem 1** (Informal). For passively collected data in complex social systems the train-test paradigm cannot be valid under ontological minimality for the vast majority of the system. This includes widely employed variants of recommender systems and question answering via LLMs.

**Theorem 1** (Test validity in complex social systems). *Let $(\mathbb{A}, \mathcal{D}, \mathsf{T}, \mathsf{F})$ be identical to lemma 2. Furthermore, let $\mathcal{S} \sim \mathsf{S}^m$ where $\mathsf{S}$ follows power-law distributions such that the degrees of $x \in \mathcal{X}_i$ in the sample graph $\mathcal{S}$ are drawn i.i.d. from a regularly-varying power-law distribution $\mathbb{P}(\deg(x) > k) = u(k)k^{-\alpha_i}$. Furthermore, let $n_i = |\mathcal{X}_i|$ be the size of domain $\mathcal{X}_i$. Then, the number $V_i$ of nodes in $\mathcal{X}_i$ for which test validity holds decreases with a power-law decay in $\mathrm{rank}(f) = k$, i.e,*

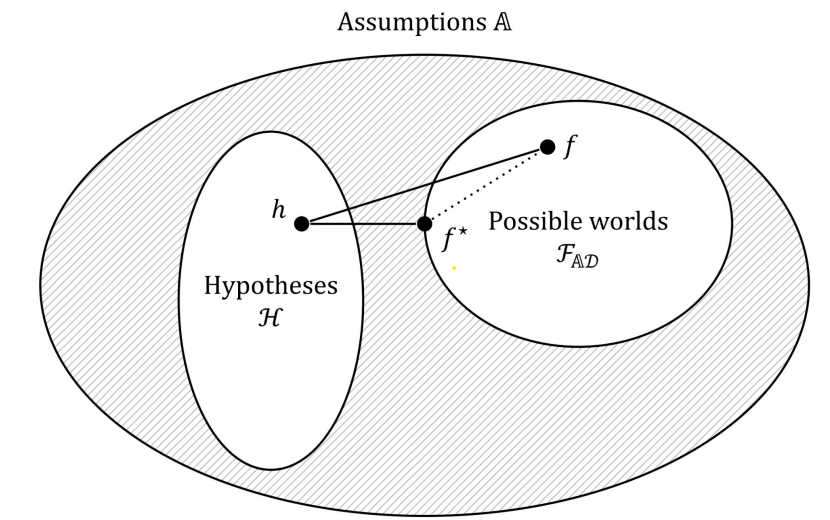$$\mathbb{E}[V_i] \leq n_i u(k) k^{-\alpha_i}.$$

# **Proof Sketch**

Assumptions

$$\mathbb{A} = \{f \mid \operatorname{rank}(f) \leq k\}$$

Test validity

$$\mathbb{P}_{f \sim \mathsf{F}}\left(|\theta - L_{fh}^{\mathsf{T}}| \leq \epsilon\right) \geq 1 - \delta$$

$\leq$

Necessary conditions

$$\mathbb{P}_{f \sim \mathsf{F}}\left(L_{fh}^{\mathsf{T}} \leq \epsilon + \theta\right) \geq 1 - \delta$$

$\leq$

$$\mathbb{P}_{f \sim \mathsf{F}}\left(L_{ff^{\star}}^{\mathsf{T}} \leq \epsilon + \theta\right) \geq 1 - \delta$$

Grounding in
○ $\mathscr{D}$ ~ complex social system
○ Ontological parsimony

$$\nexists \epsilon : \mathbb{P}_{f \sim \mathsf{U}}\left(L_{ff^{\star}} \leq \epsilon\right) > 0$$

$$f^{\star} = \arg\min_{f} L_{fh}^{\mathsf{T}} : L_{ff^{\star}}^{\mathsf{T}} \leq L_{fh}^{\mathsf{T}}$$

$\mathcal{F}_{\mathbb{A}\mathscr{D}}$ is a vector space if $k$-connectivity of $\mathscr{D}$ is smaller than *rank(f)*
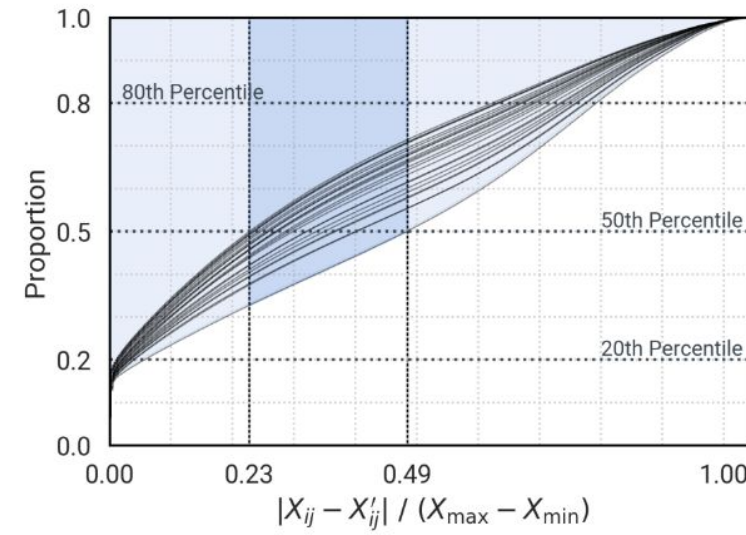

Assumptions $\mathbb{A}$

# MovieLens (100k)
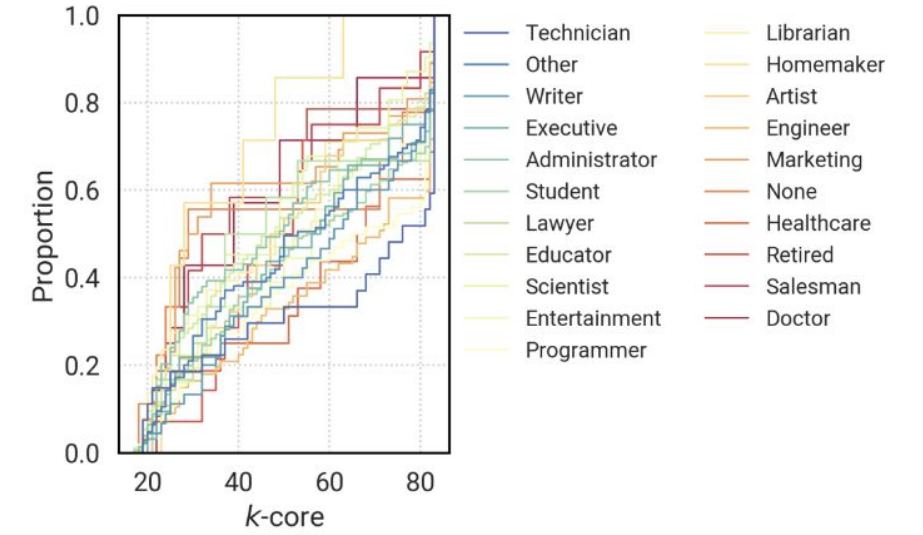
THE recommender systems benchmark since 1998

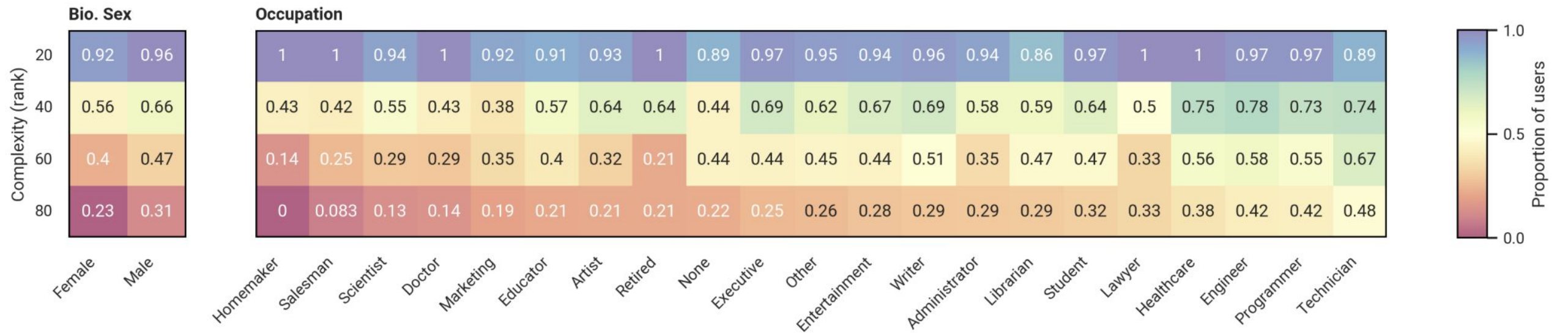Yet, it's invalid for the evaluation tasks that we (typically) use it for... 😬



**(a)** eCDF of Maximum NAE

**(b)** eCDF of Pairwise NAE
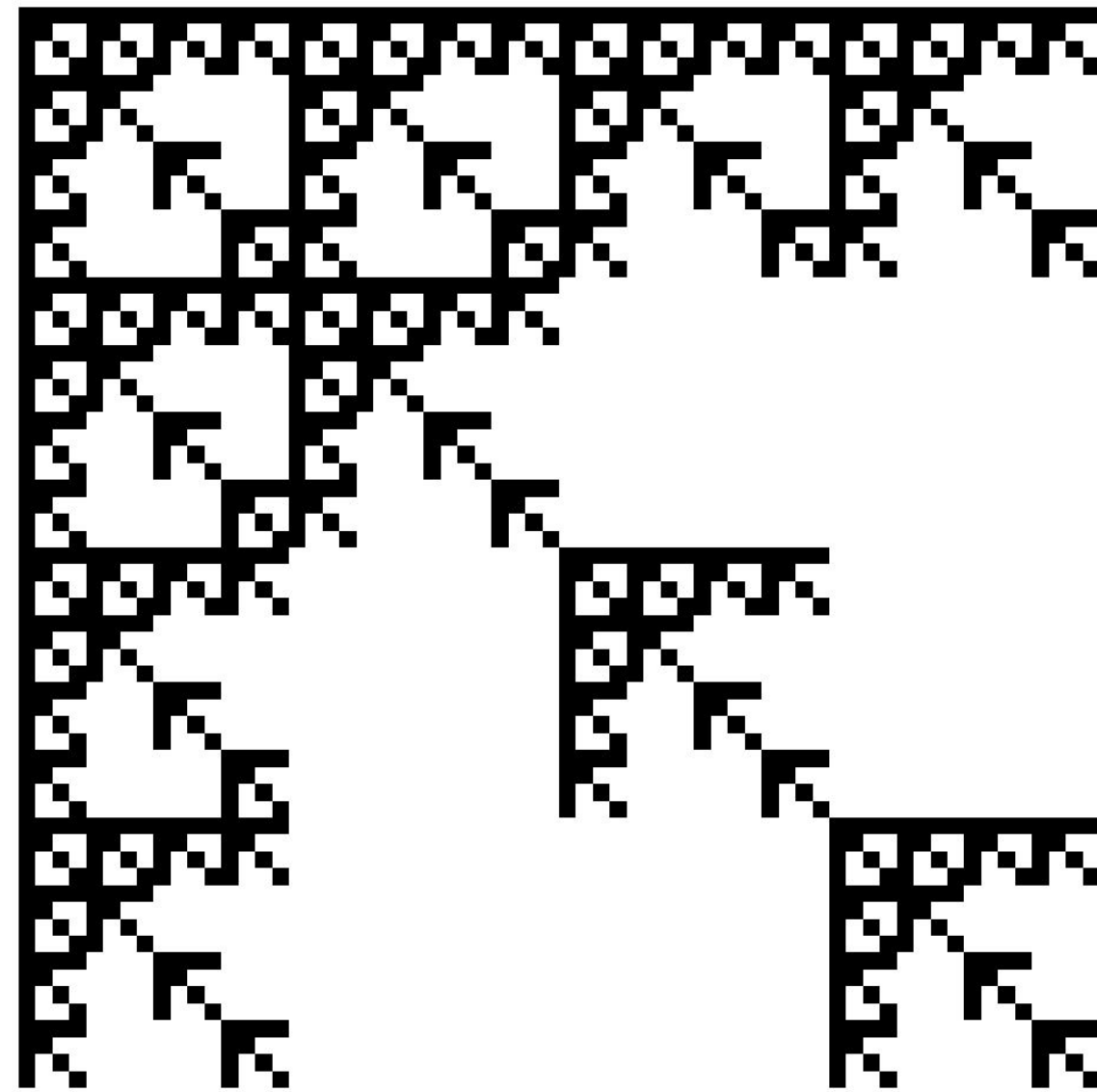
**(c)** $k$-core per occupation

**(d)** Test-validity per demographic group and model complexity

# **Naive scaling** and **manual benchmarks** won't fix it…

**Corollary 4** (Informal). Naïve scaling and selective benchmarks are prohibitively inefficient to address theorem 1 and therefore not suited to attain test validity in complex social systems.

| $\alpha$ | $x_{\min}$ | $|\mathcal{X}|$ | **Scaling**<br>Samples needed to increase k-core of random node | **Benchmarks**<br>Nodes with less than 100 observations |
|---|---|---|---|---|
| 2.5 | 5 | $10^7$ | $\mathbb{E}_{i\sim\mathsf{U}}[T_i] \geq (|\mathcal{X}|/2)^{\alpha+1}/(\alpha x_{\min}^{\alpha}) \quad = \quad 2\cdot 10^{21}$ | $\mathbb{E}[N] = |\mathcal{X}|(1-(x_{\min}/x)^{\alpha}) \quad > \quad 9.9\cdot 10^6$ |

| $\alpha$ | $x_{\min}$ | $|\mathcal{X}|$ | **Book Crossing** (Ziegler et al., 2005)<br>Fraction of users with large enough degrees such that train-test measures and inferences are valid |
|---|---|---|---|
| 2.38 | 8 | $10^5$ | Rank 8: 100%,   Rank 10: 58.8%,   Rank 20: 11.3%,   Rank 100: 0.2% |

# No need for an existential crisis 😱 advances are real **but realism is needed**
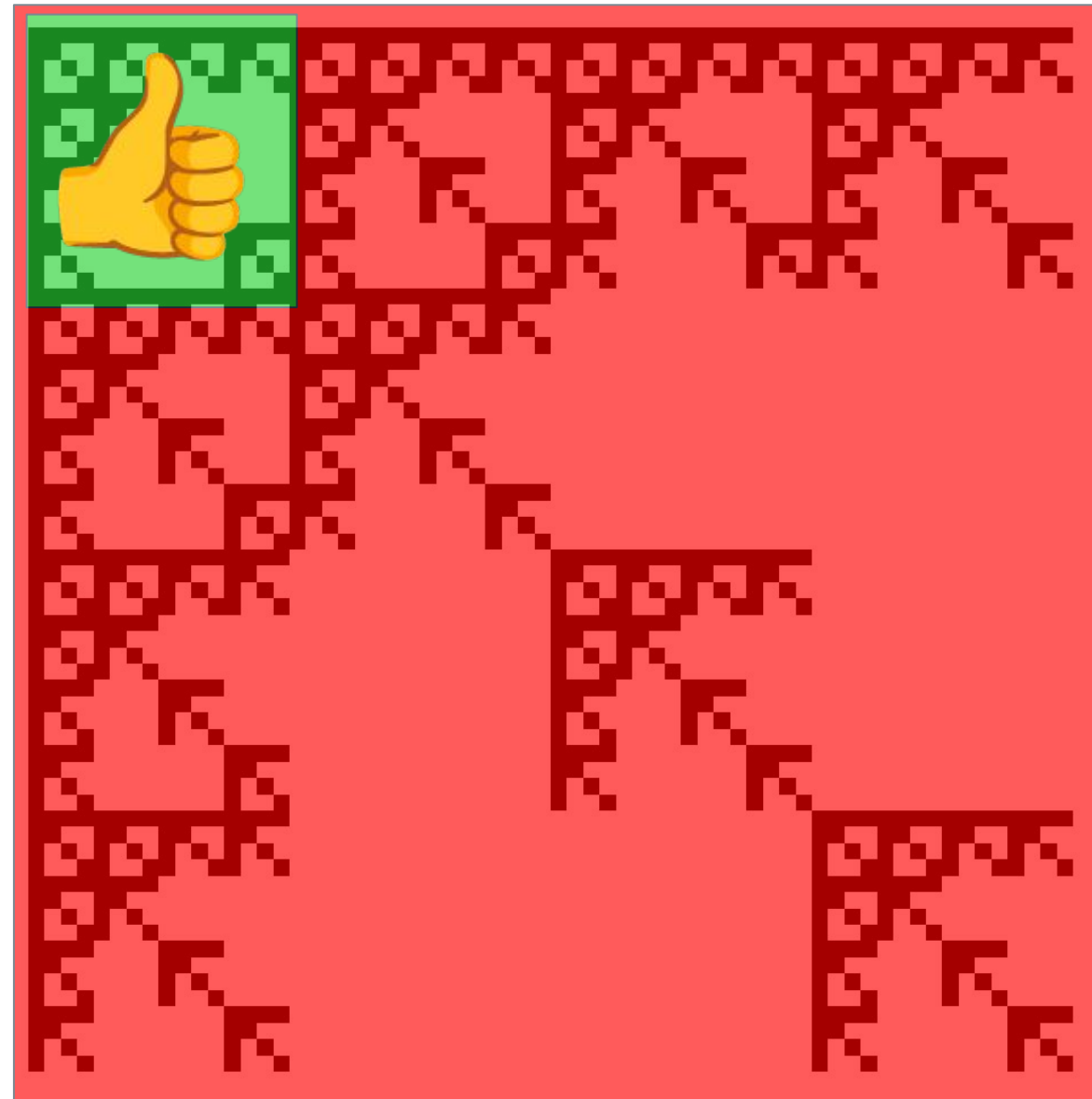


**Sample Graph**

👍 **Test Valid**
k-core(Sample Graph) ≥ Complexity of world

👎 **Test Invalid**
k-core(Sample Graph) < Complexity of world

# No need for an existential crisis 😱 advances are real **but realism is needed**



**Sample Graph**

**Test Valid**
k-core(Sample Graph) ≥ Complexity of world

**Test Invalid**
k-core(Sample Graph) < Complexity of world

# **Hume** is back…
# back again!

Hume's **problem of induction** is now back in a slightly different form and **renders the train-test paradigm ineffective** for our **current data collection practices.**
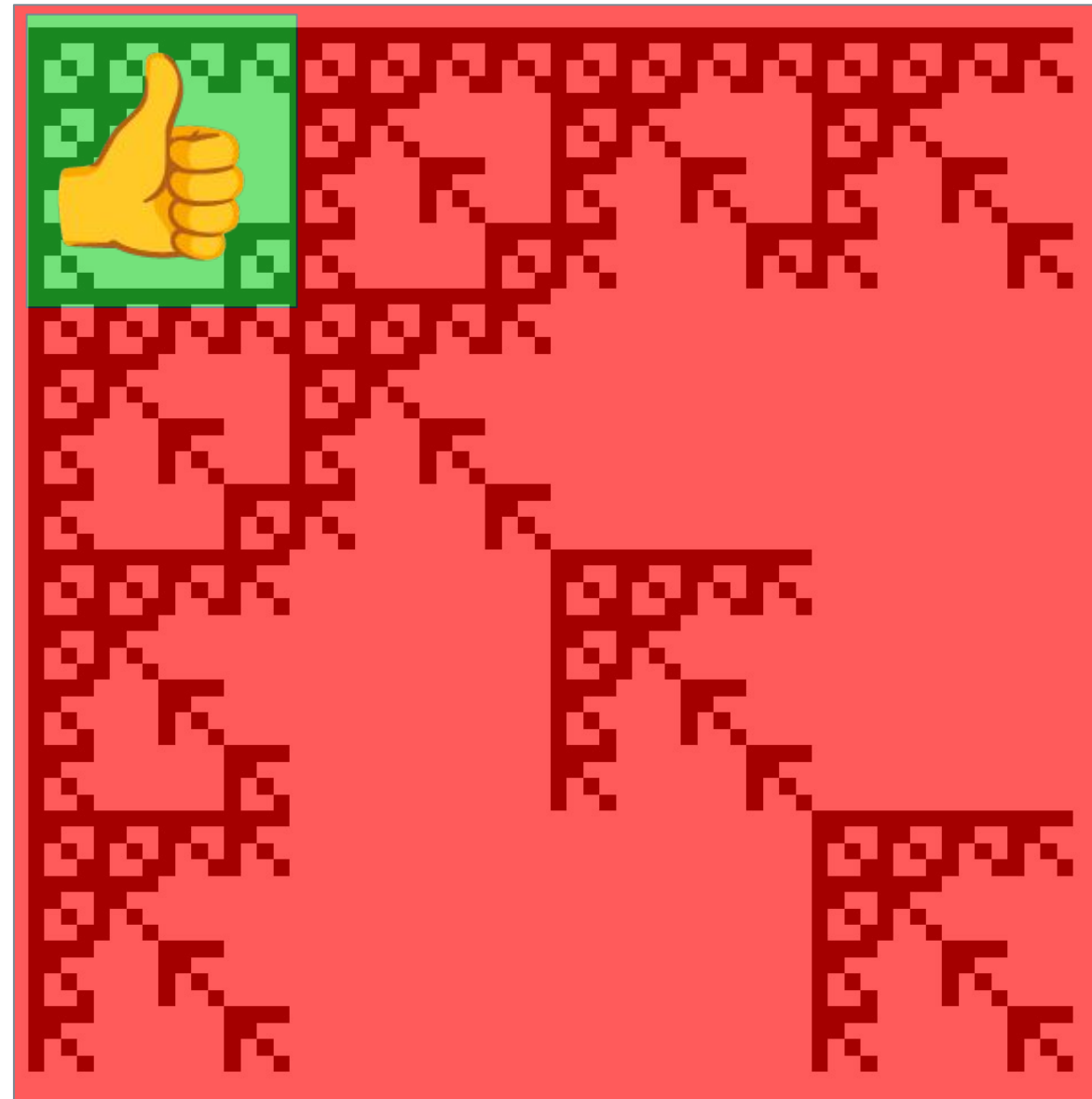
*Solving many complex AI tasks will not come for free through scaling or for cheap through extrapolating from small-scale benchmarks.*

There is an **inherent trade-off between data quality, quantity, and task complexity.** If we want to avoid asking AI systems to solve simpler tasks (e.g., non-out-of-distribution or smaller scale), **new data curation efforts** are needed.

# *Future work*
# Provably fair cooperative data collection



Sample Graph

**Test Valid**
k-core(Sample Graph) ≥ Complexity of world

- Increase size of the green area.

- We know where to collect data via k-core condition!

- Number of test-valid data points is a **supermodular** function with regard to datasets

- Shapley value!

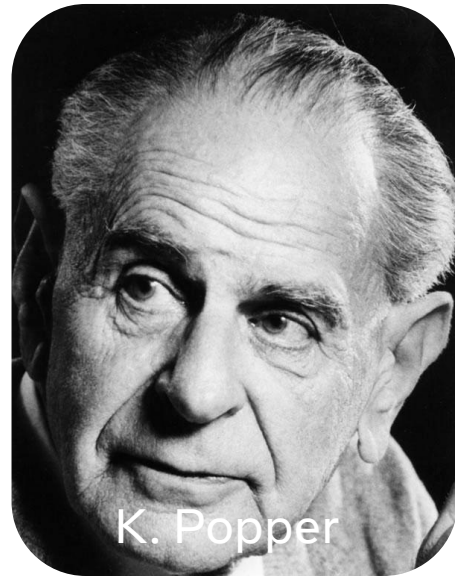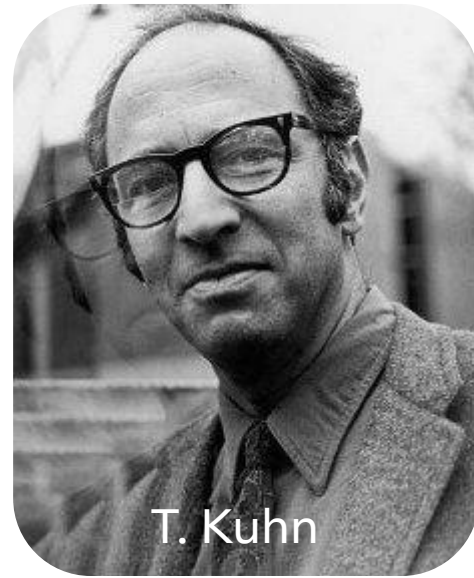- Scalable at the level of organizations, e.g., in **open science**

AI

Machine learning

~~Statistics~~ is the **science of induction**

K. Popper


T. Kuhn


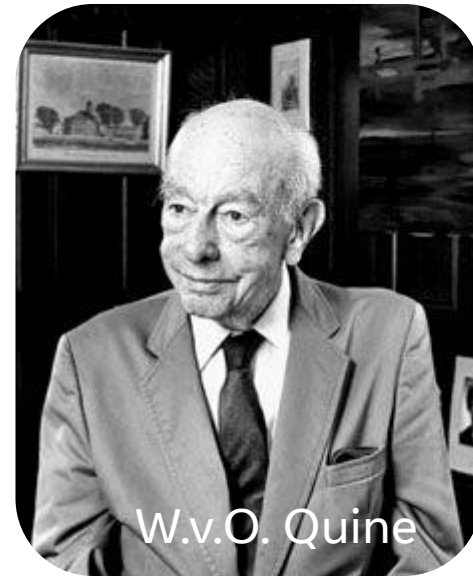W.v.O. Quine


Alchemy, Phrenology, Astrology, … Us?

**Demarcation**

**Demarcation**

# Science

Trying to provide a definite answer to what science is, is a good way to get your philosopher friend upset.

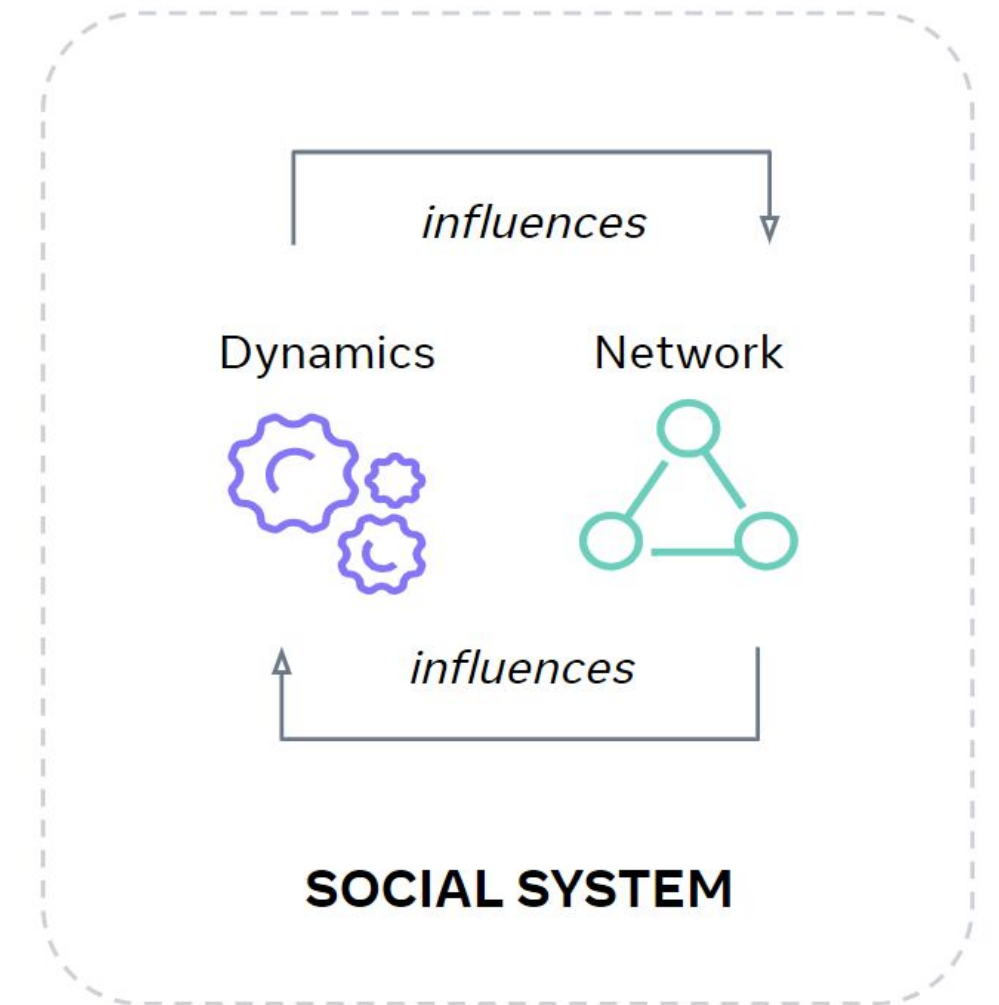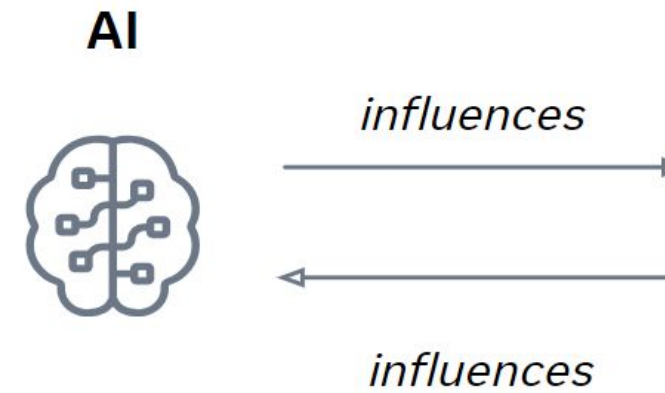The ultimate issue is **"how to determine which beliefs are epistemically [justified]"** *(Fuller 1985)*.

# Pseudo-Science

What exactly constitutes pseudo-science is not clear either, but roughly, it amounts to *(Hansson 1996):*

1) it is **not scientific**, and
2) its major proponents try to **create the impression that it is scientific**.

# No Free Delivery Service
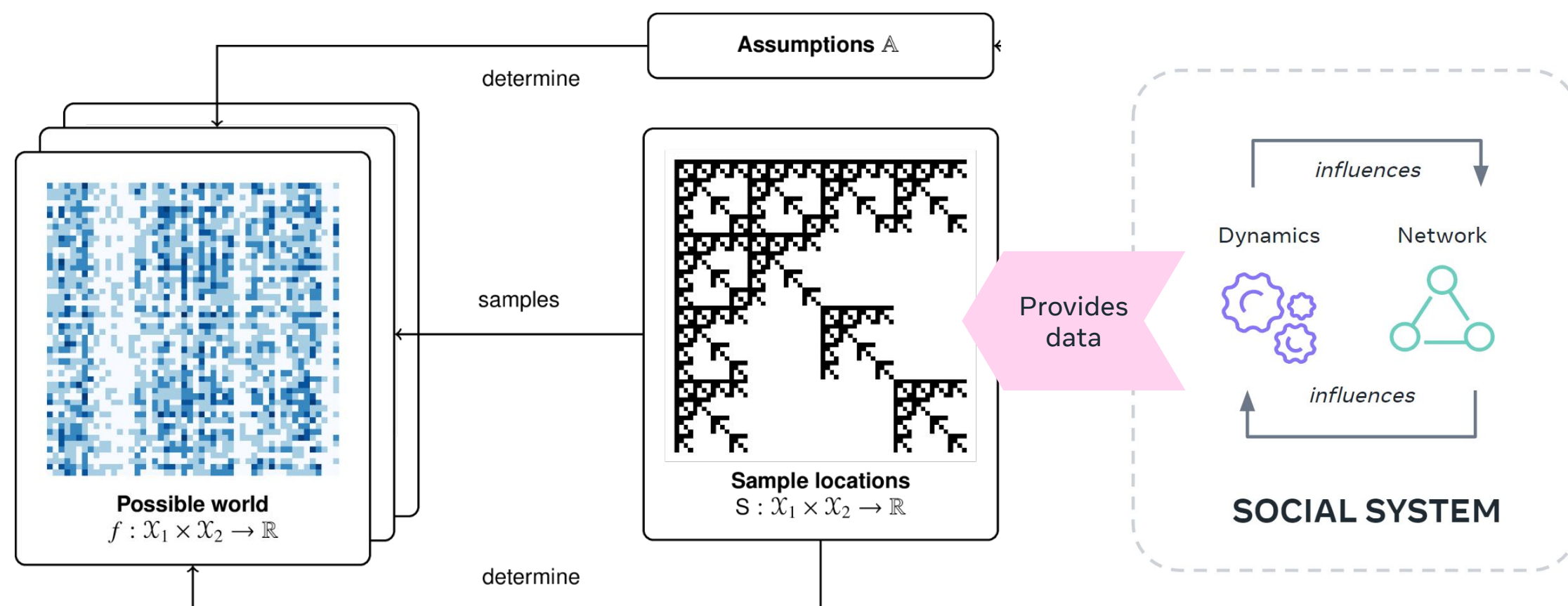
https://arxiv.org/abs/2411.13653

**Theorem 1** (Informal) For passively collected data in complex social systems, the train–test paradigm <u>cannot</u> be valid under ontological parsimony for the vast majority of the system. This incluses widely employed variants of recommender systems and QA via LLMs.
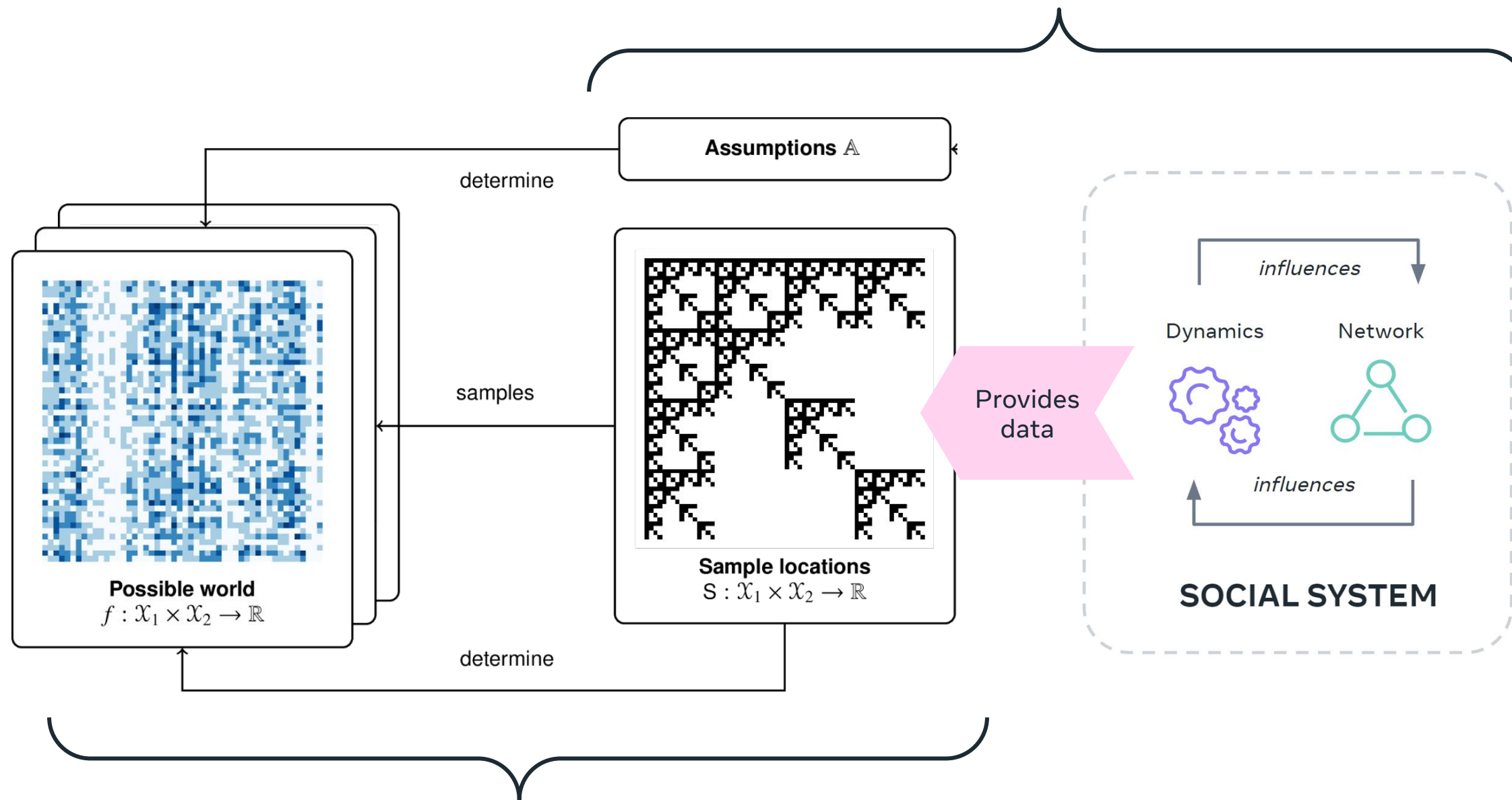
**How it works:**

○ Observations + assumptions define **possible worlds** that are consistent with both.
○ The error of of *any risk estimator θ* cannot be bounded whp over these possible worlds due to the structure of the data generating system.

$$\mathbb{P}_{f\sim\mathsf{F}}\left(|\theta - L_{fh}^{\mathsf{T}}| \le \epsilon\right) \ge 1-\delta$$

**Formalization of data collection** for validation of AGI tasks

Assumptions $\mathbb{A}$

determine

**Possible world**
$f : \mathcal{X}_1 \times \mathcal{X}_2 \to \mathbb{R}$

samples

**Sample locations**
$S : \mathcal{X}_1 \times \mathcal{X}_2 \to \mathbb{R}$

determine

Provides data

influences

Dynamics

Network

influences

**SOCIAL SYSTEM**

**Test validity**: Can the test error be informative about the true generalization error?

$$\mathbb{P}_{f \sim \mathsf{F}} \left( |\theta - L_{fh}^{\mathsf{T}}| \leq \epsilon \right) \geq 1 - \delta$$

∞ Meta