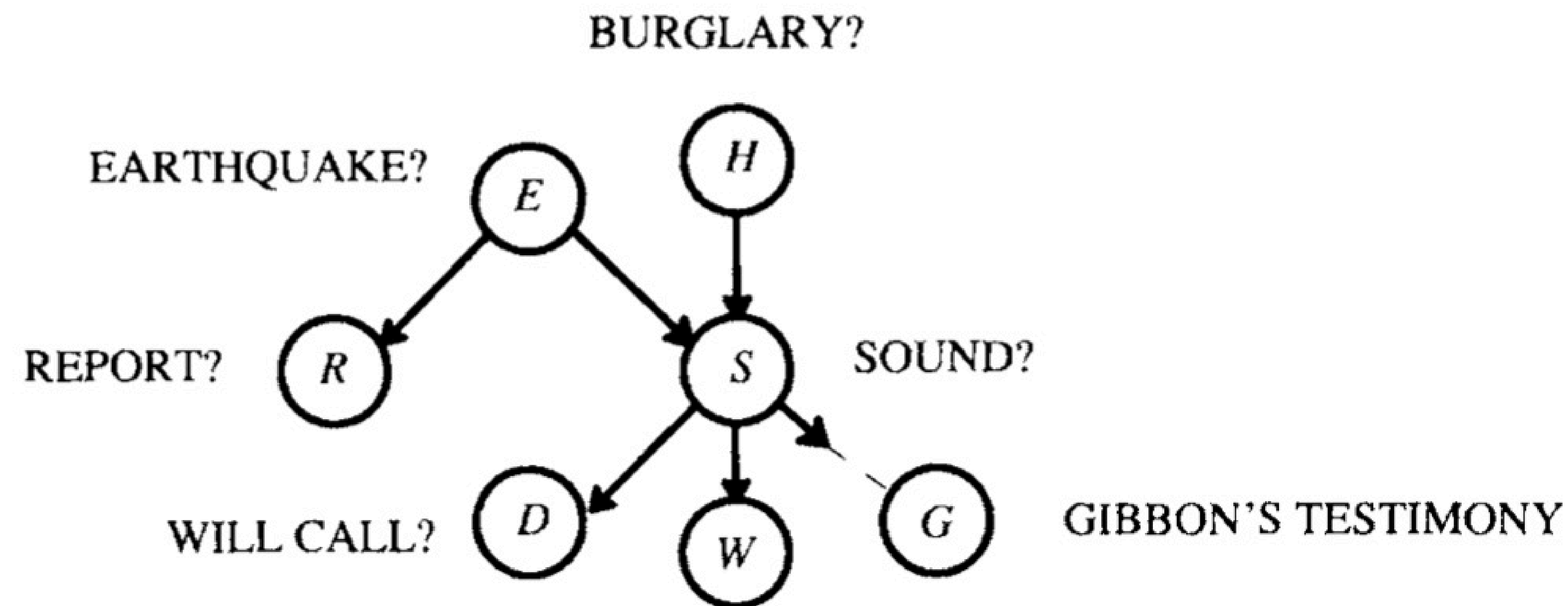


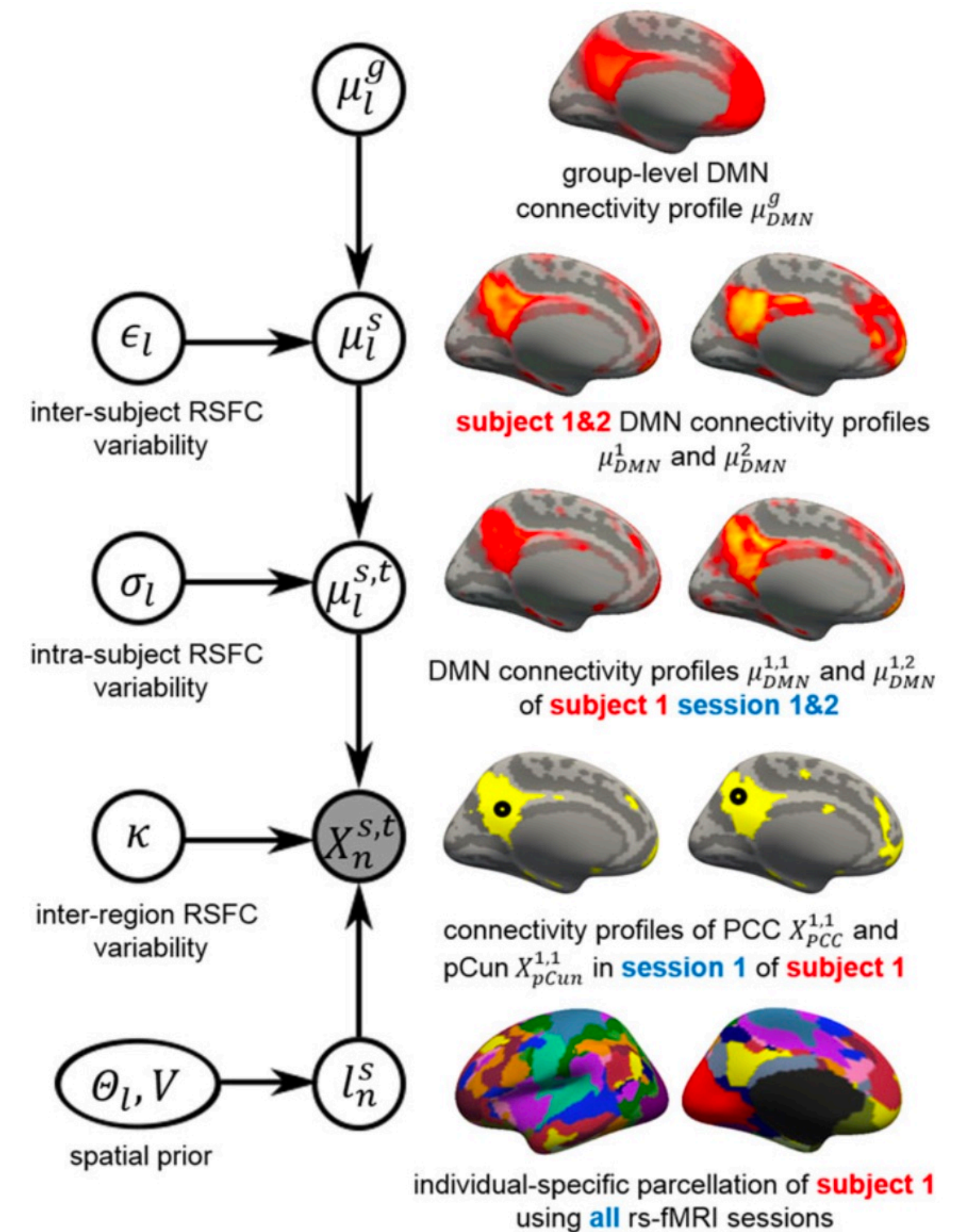
# Graphical Models and Simulation-Based Inference

Graphical Models: Discrete Inference and Learning

# Introduction to DAG and their relationship with Probability Functions (Pearl)

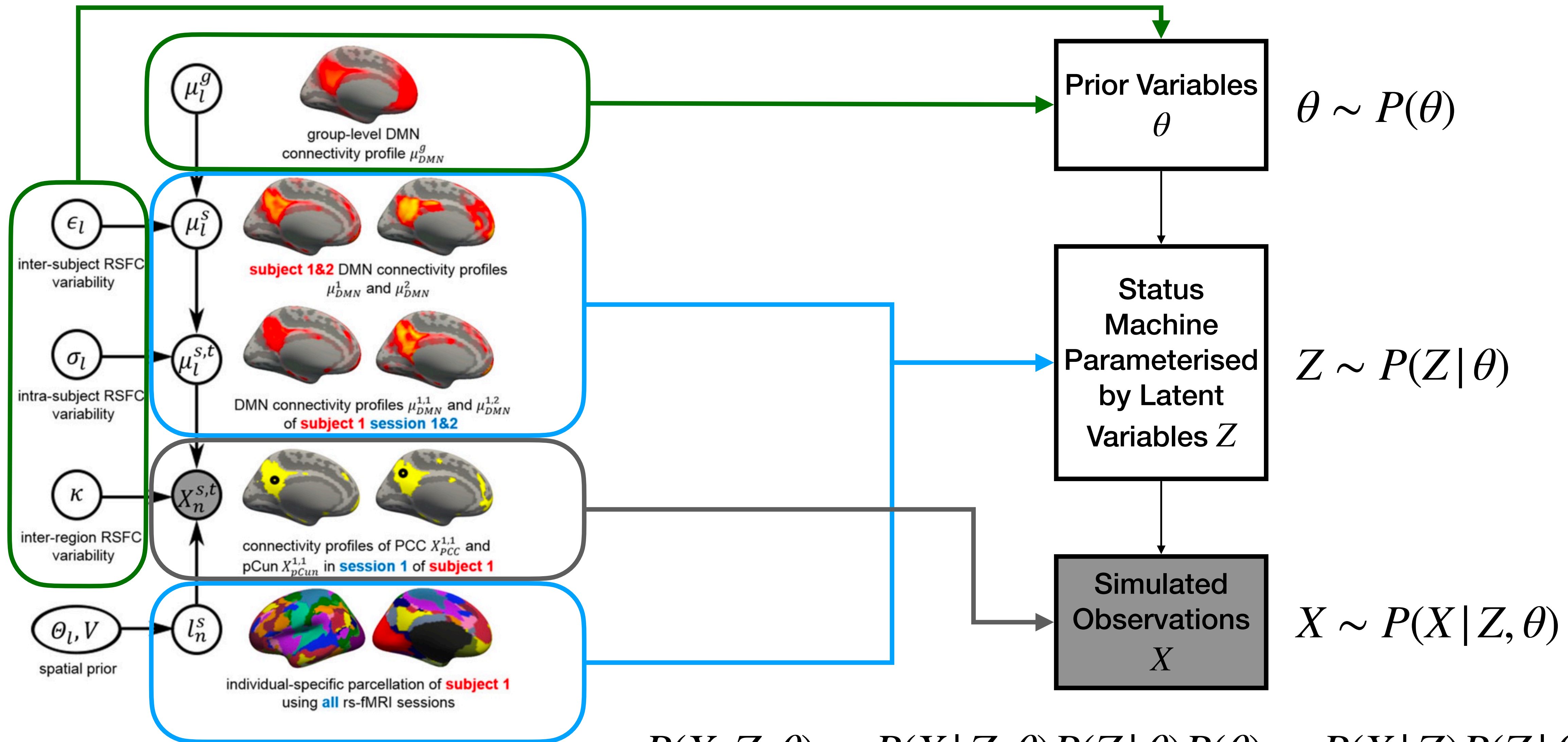


[Pearl 1987]



[Kong et al 2019]

# Graphical Models and Simulation Systems



$$P(X, Z, \theta) \stackrel{3}{=} P(X | Z, \theta)P(Z | \theta)P(\theta) = P(X | Z)P(Z | \theta)P(\theta)$$

# General Inference Notation

$\theta$ : parameters     $X$ : observations

$$P(\theta | X) = \frac{\overset{\text{Likelihood}}{P(X | \theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Evidence}}{P(X)}}$$

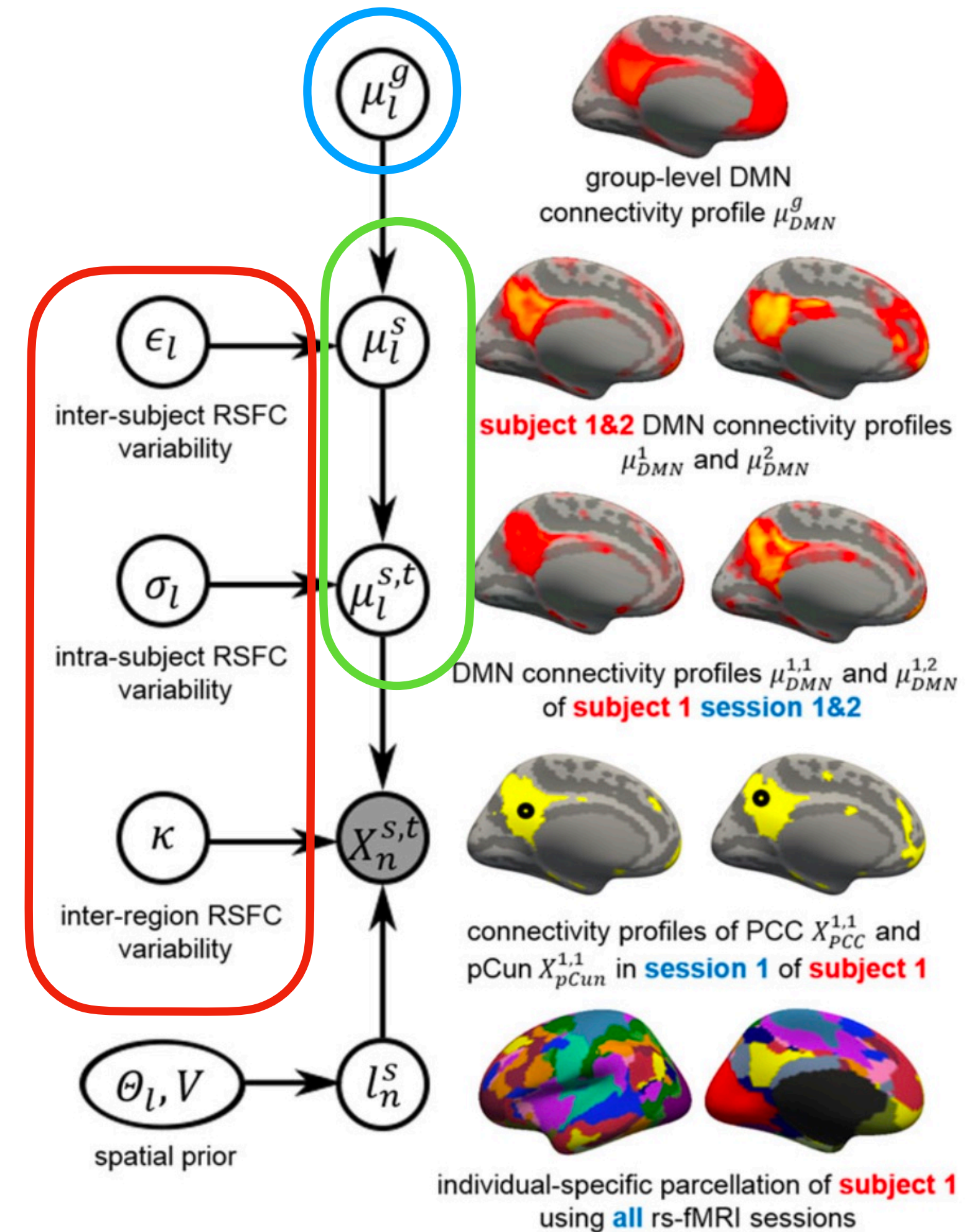
Posterior

$Z$ : latent random variables

$$P(\theta | X) = \frac{\mathbb{E}_Z [P(X | Z, \theta)] P(\theta)}{P(X)}$$

$\eta$ : nuisance random variables

$$P(\theta | X) = \frac{\mathbb{E}_\eta [P(X | \theta, \eta)] P(\theta)}{P(X)}$$



# General Inference Notation

$\theta$ : parameters     $X$ : observations

$$P(\theta | X) = \frac{\overset{\text{Likelihood}}{P(X | \theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Evidence}}{P(X)}}$$

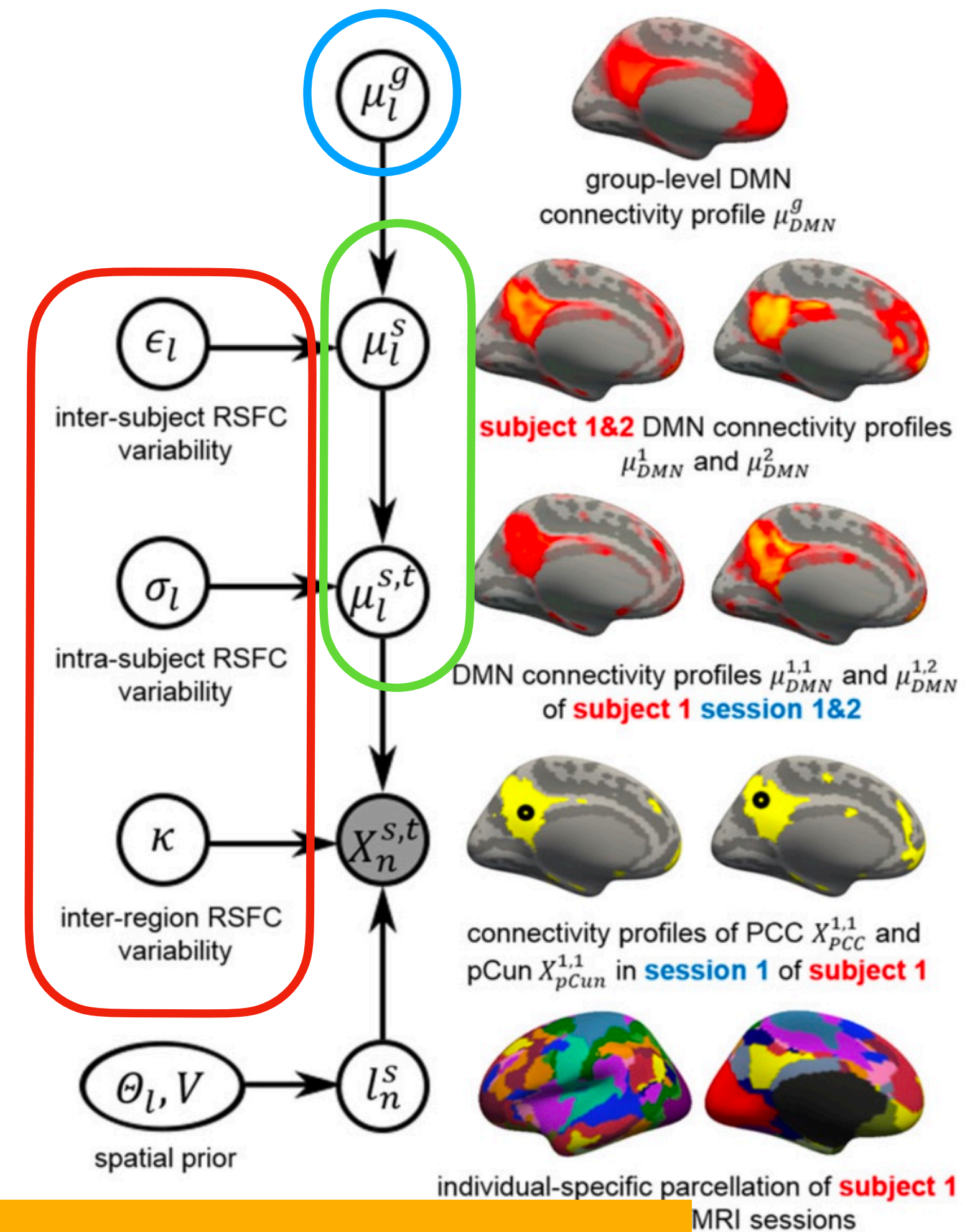
Posterior

$Z$ : latent random variables

$$P(\theta | X) = \frac{\mathbb{E}_Z [P(X | Z, \theta)] P(\theta)}{P(X)}$$

$\eta$ : nuisance random variables

$$P(\theta | X) = \frac{\mathbb{E}_\eta [P(X | \theta, \eta)] P(\theta)}{P(X)}$$



Intractable in general:  
full likelihood impossible to evaluated or computation cost is extremely high

# Likelihood computation is hard: Enter Mechanistic, Example Models Galton Board

$\theta$ : parameters     $X$ : observations

$$P(\theta | X) = \frac{\overset{\text{Likelihood}}{P(X | \theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Evidence}}{P(X)}}$$

Posterior

$Z$ : latent random variables

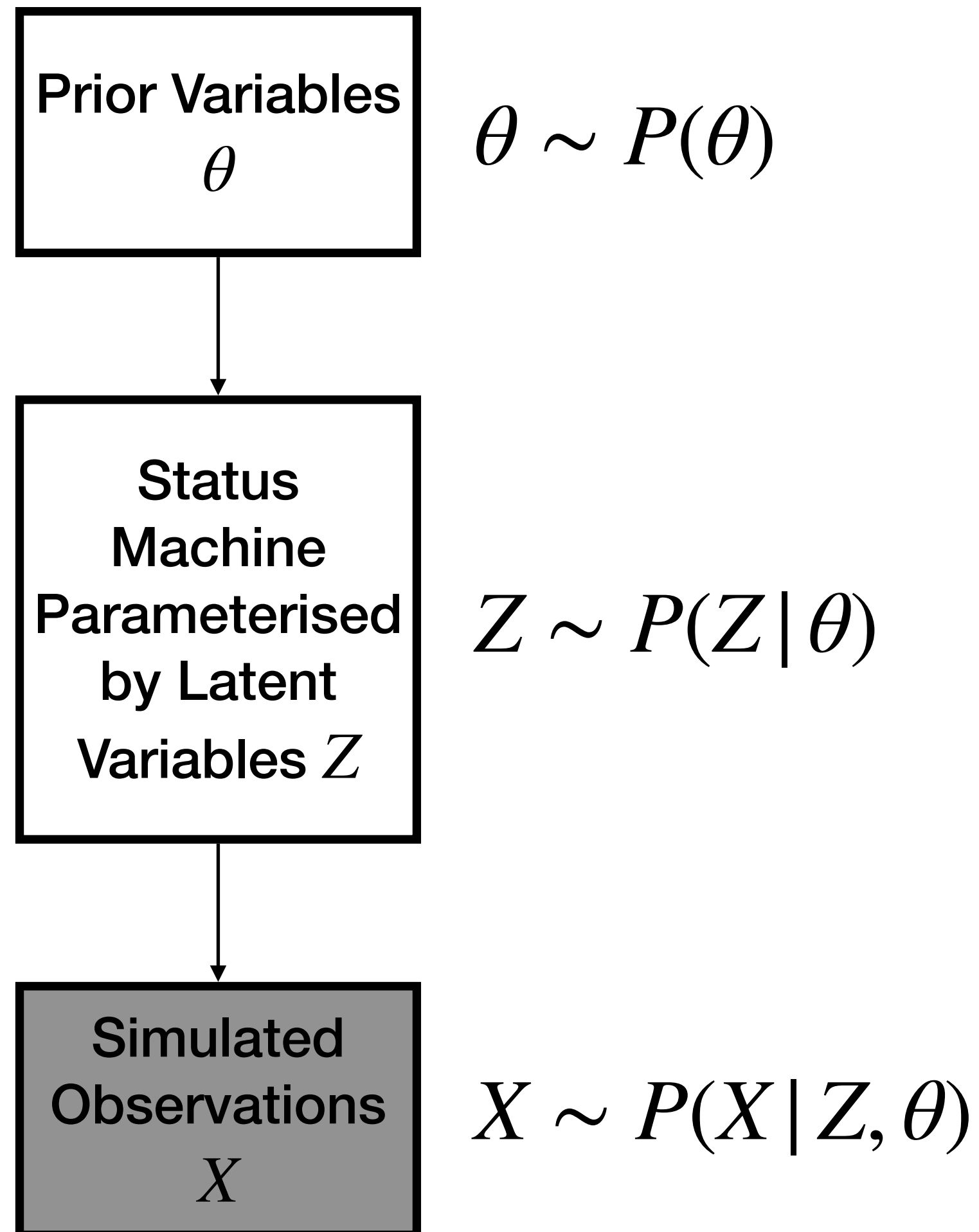
$$P(\theta | X) = \frac{\mathbb{E}_Z [P(X | \mathbf{Z}, \theta)] P(\theta)}{P(X)}$$

$\eta$ : nuisance random variables

$$P(\theta | X) = \frac{\mathbb{E}_\eta [P(X | \theta, \eta)] P(\theta)}{P(X)}$$



# Simulation-Based Inference



- Inference is defined as finding the  $\theta$  that could be at the origin of an observation  $X$ . Specifically computing  $P(\theta | X) = \mathbb{E}_Z[P(\theta, Z | X)]$
- For this, we use Bayes 
$$P(\theta, Z | X) = \frac{P(X | Z, \theta)P(Z, \theta)}{P(X)}$$
, nonetheless the likelihood  $P(X | Z, \theta)$  is often unknown or intractable.
- Hence simulation-based inference either approximates or eliminates the need for an explicit likelihood by simulating observations.

# Simulation-Based Inference: Neural Network Approximations

Prior Variables  
 $\theta$

$$\theta \sim P(\theta)$$

Status  
Machine  
Parameterised  
by Latent  
Variables  $Z$

$$Z \sim P(Z | \theta)$$

Simulated  
Observations  
 $X$

$$X \sim P(X | Z, \theta)$$

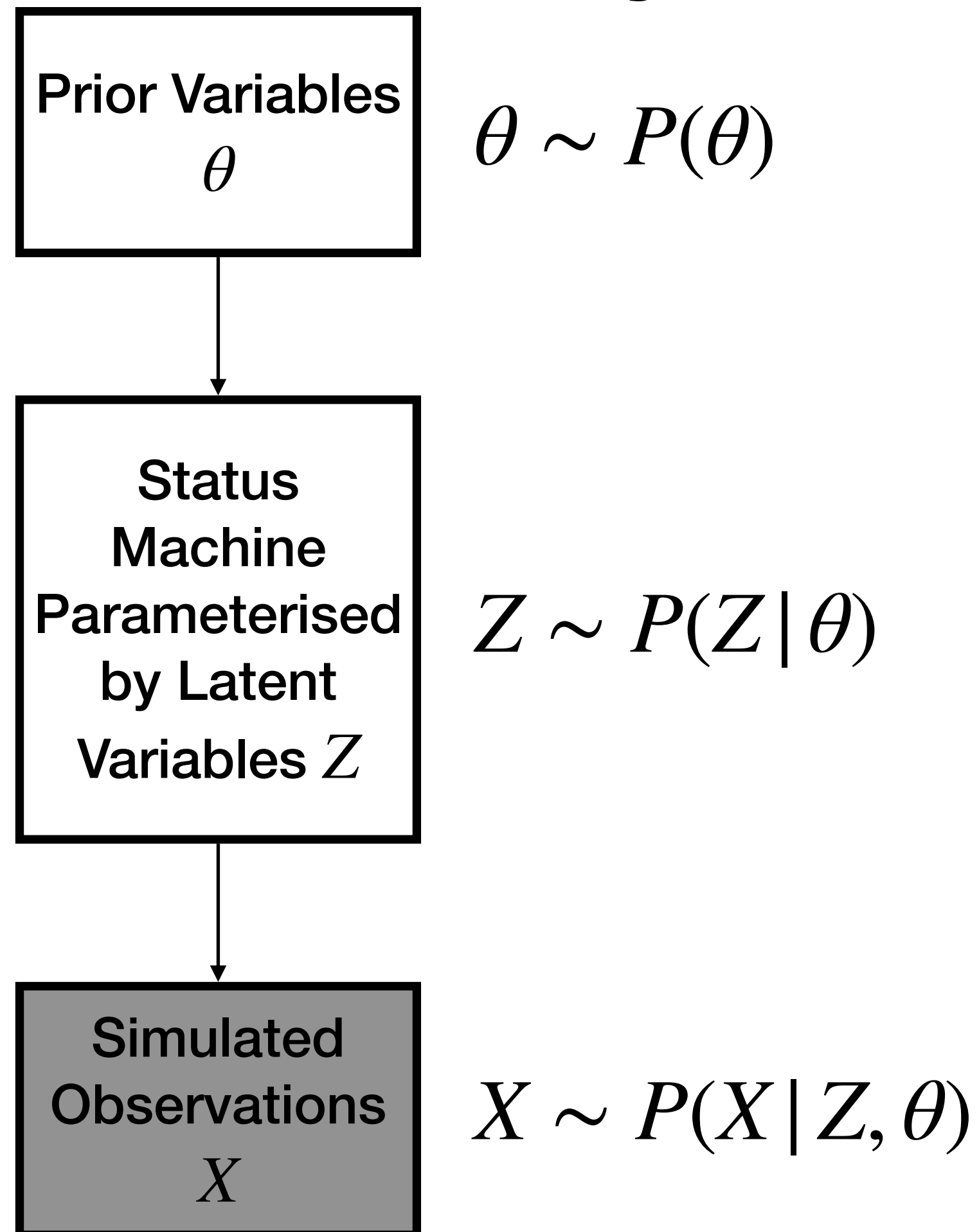
$$P(\theta | X) = \frac{\overset{\text{Likelihood}}{P(X | \theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Evidence}}{P(X)}}$$

$\theta$ : parameters    $X$ : observations

Posterior

- $P(\theta | X)$  approximated through “Neural Posterior” estimators
- $P(X | \theta)$  approximated through “Neural Likelihood” estimators
- $\frac{P(X | \theta)}{P(X)}$  approximated through the “Neural ratio” estimators

# Simulation-Based Inference: Why now it works (Cranmer et al 2019)

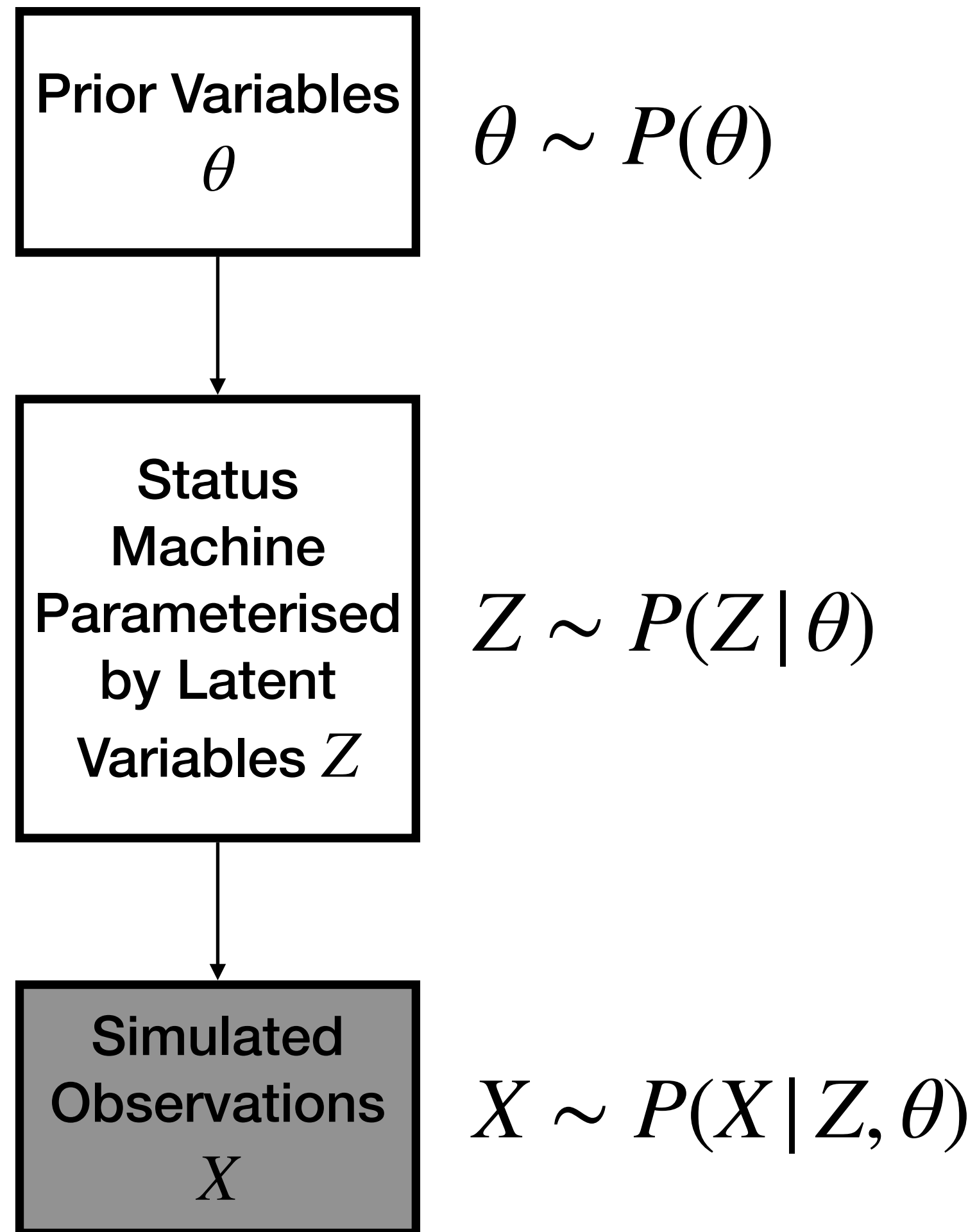


$$P(\theta | X) = \frac{P(X | \theta)P(\theta)}{P(X)}$$

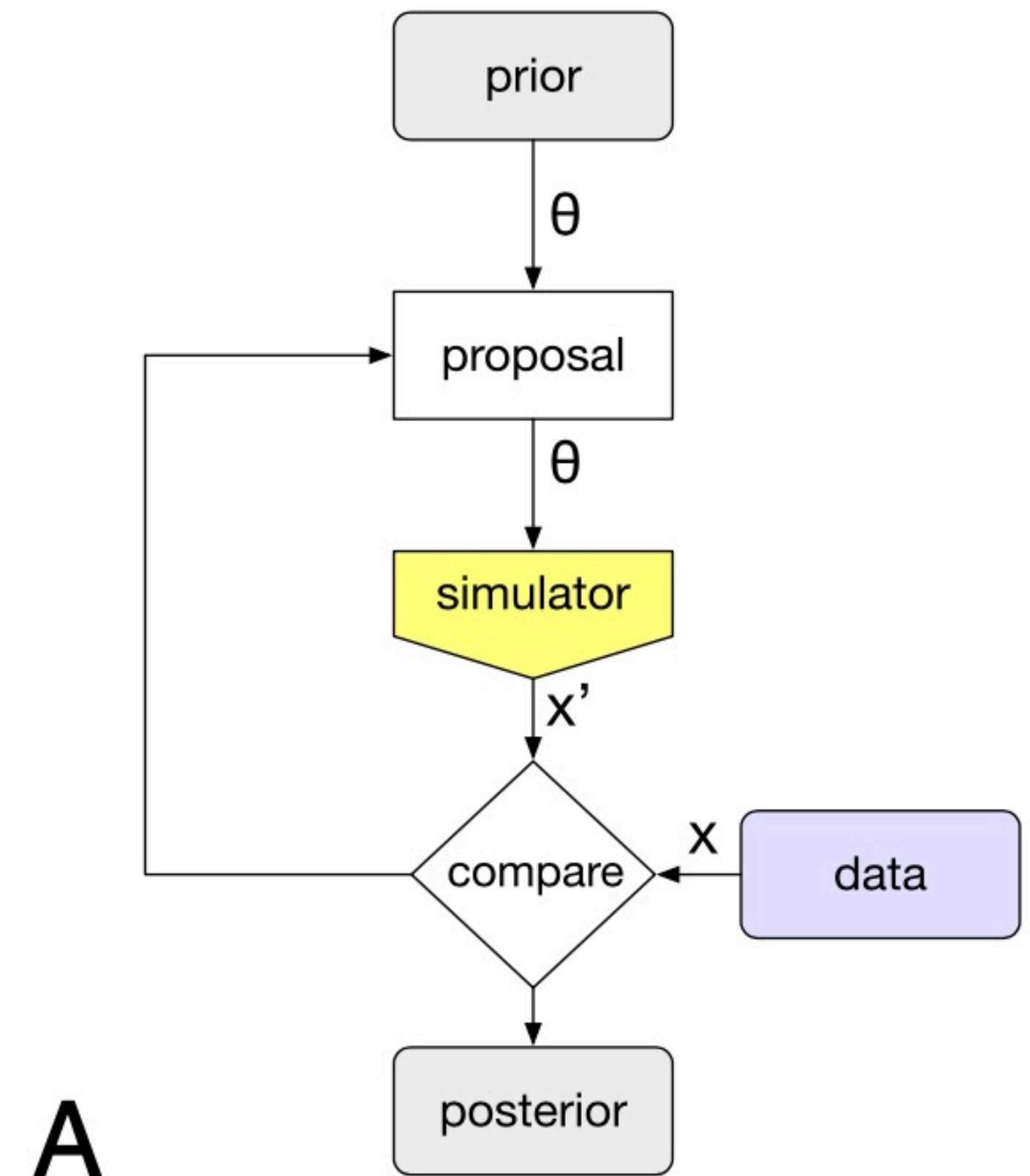
$\theta$ : parameters     $X$ : observations  
 Posterior                      Likelihood    Prior  
 Evidence

- Novel ML-based approaches allow us to massively generate simulated observations
- Autodifferentiation and neural network approaches are great non-linear function estimators
- Active learning can help improving sampling efficiency much better than Markov Chains

# Simulation-Based Inference: Approximate Bayesian MC

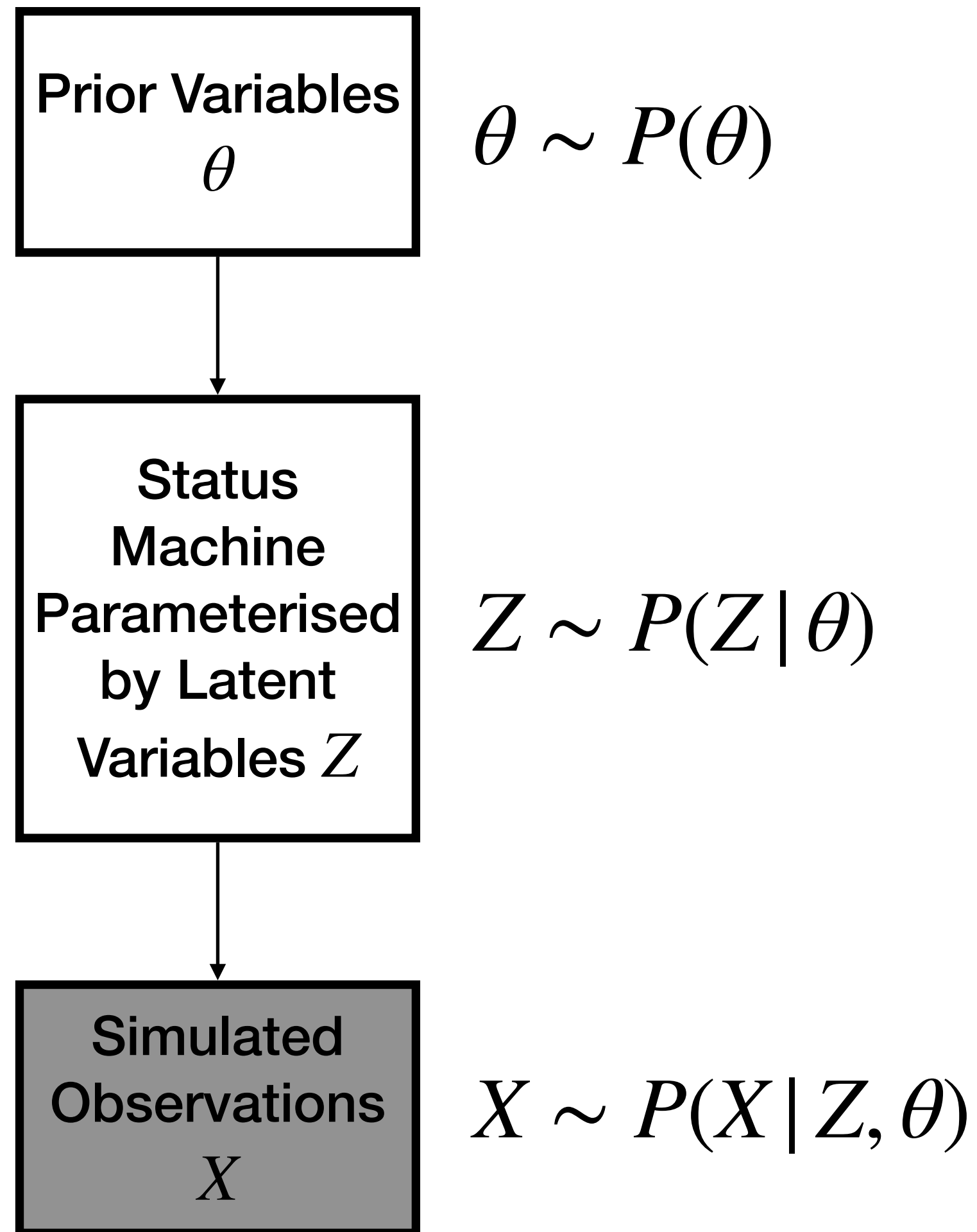


Approximate Bayesian Computation  
with Monte Carlo sampling

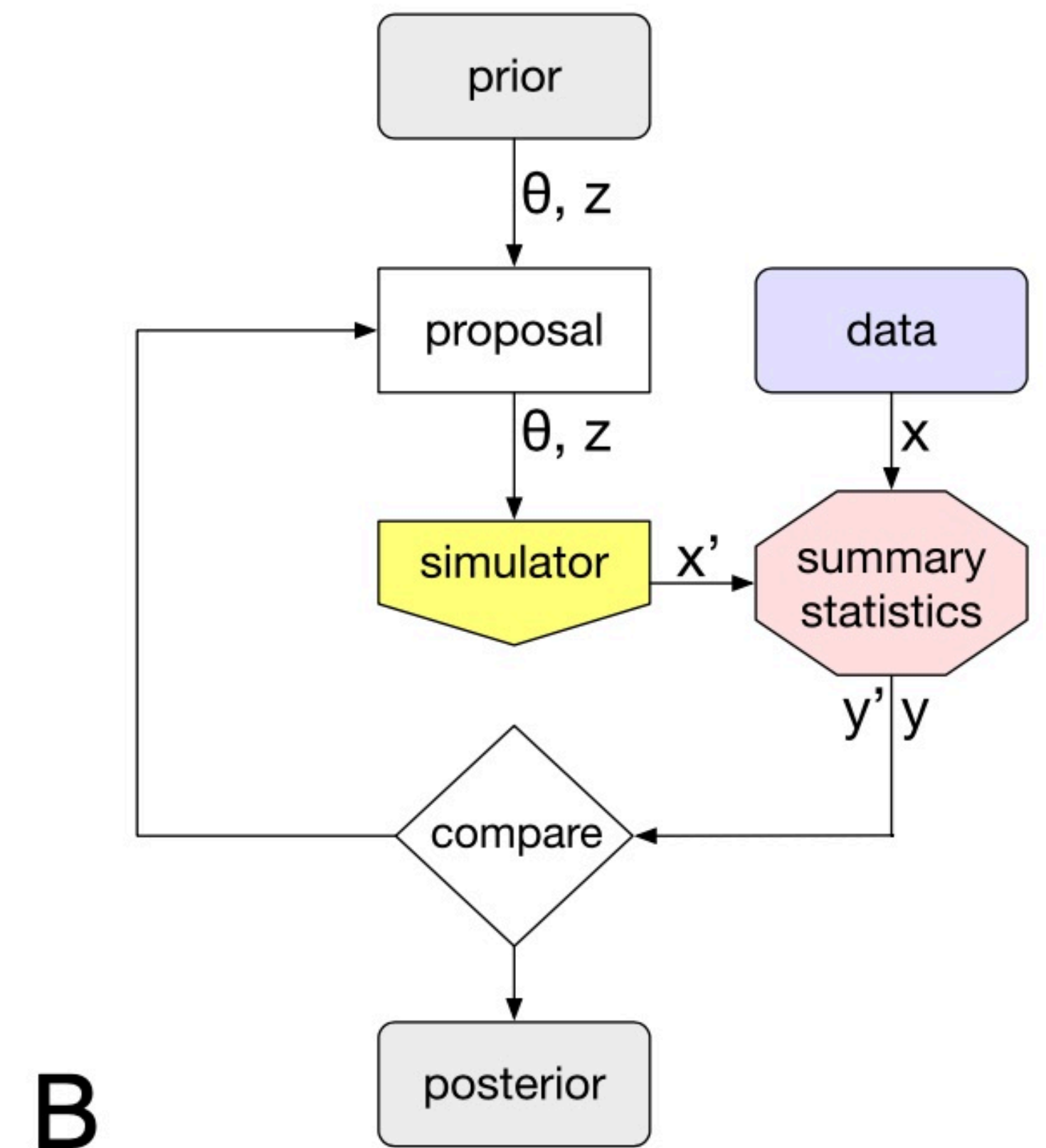


(Cranmer et al 2019)

# Simulation-Based Inference: Approximate Bayesian MC

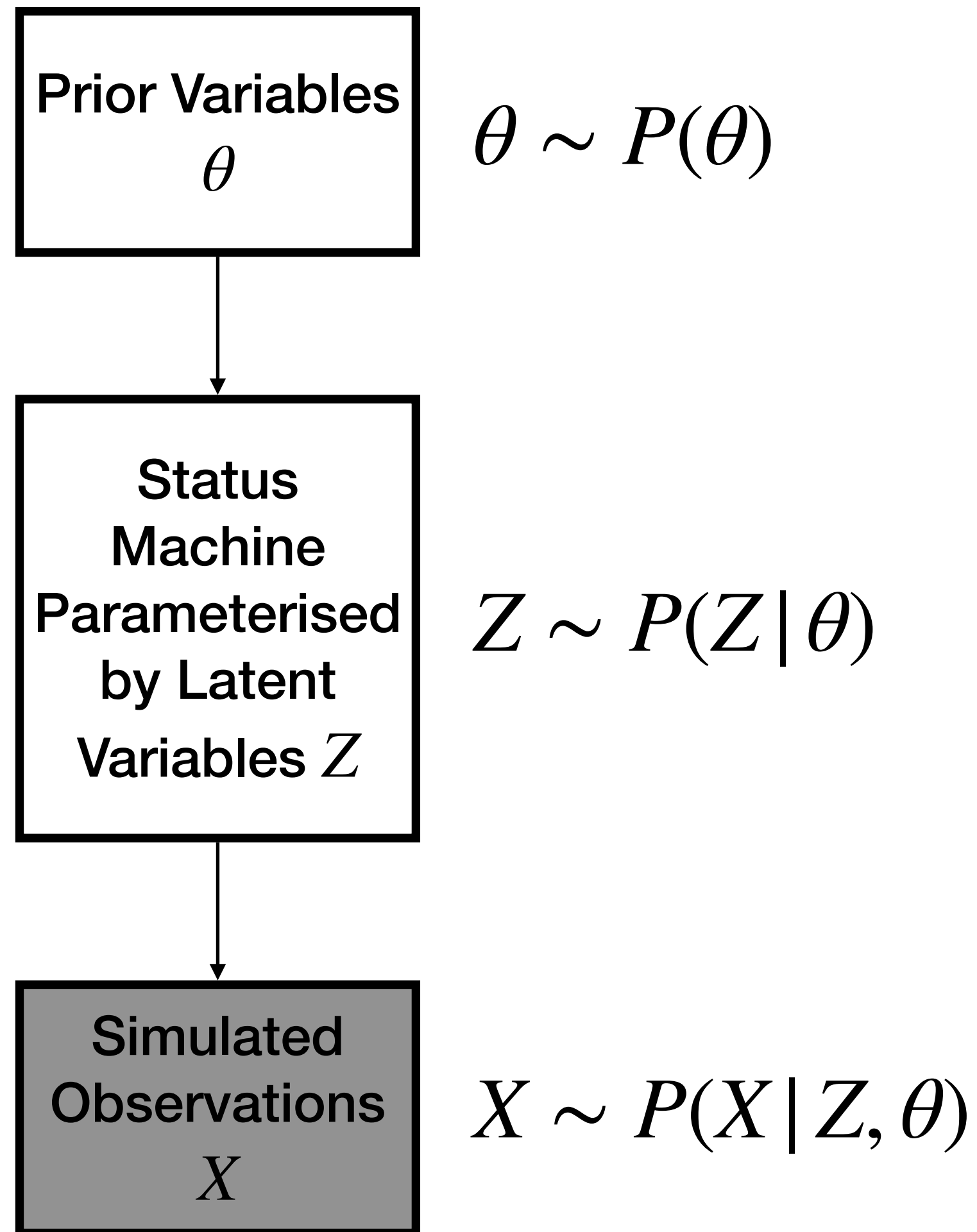


Approximate Bayesian Computation  
with learned summary statistics

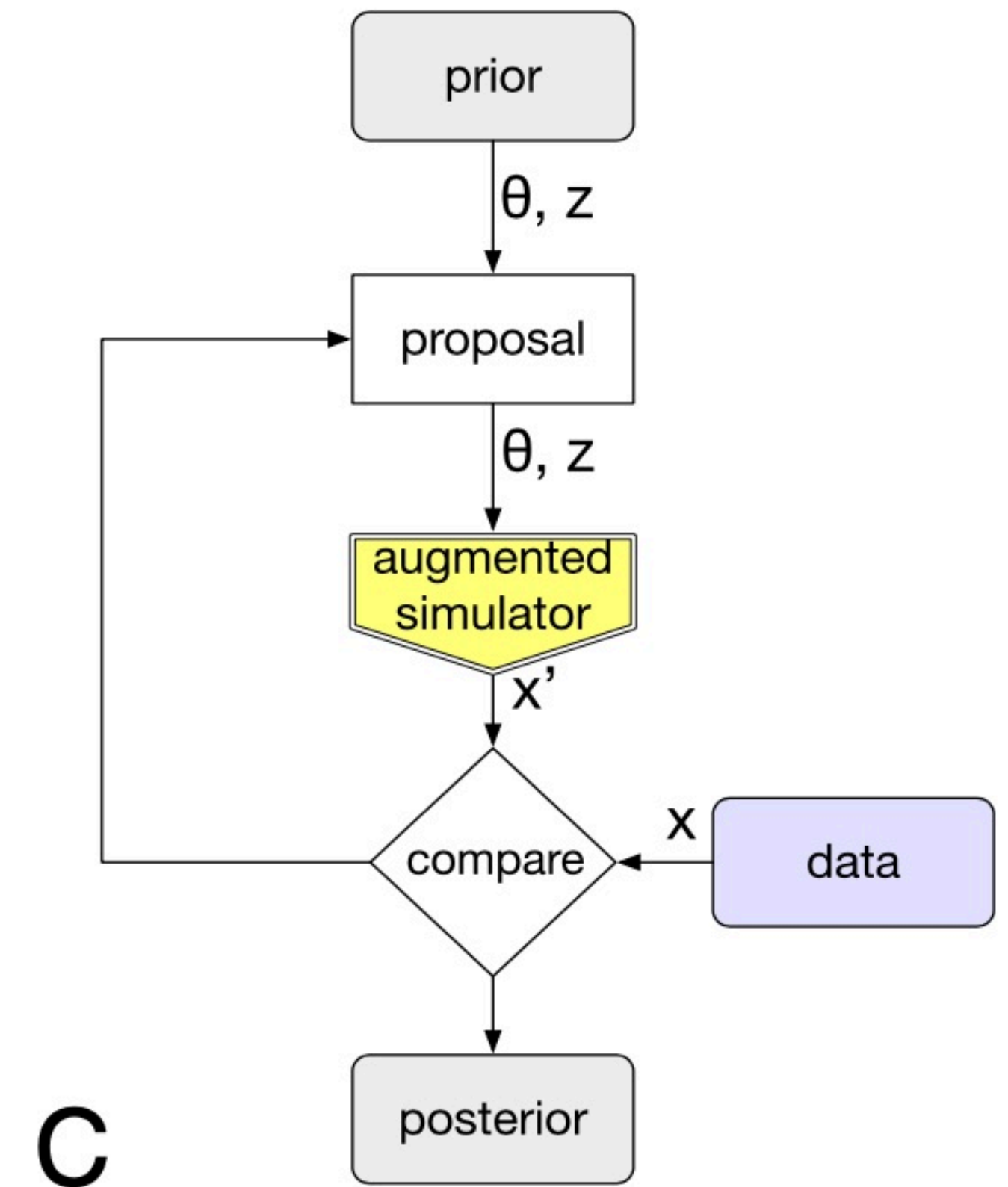


(Cranmer et al 2019)

# Simulation-Based Inference: Approximate Bayesian MC



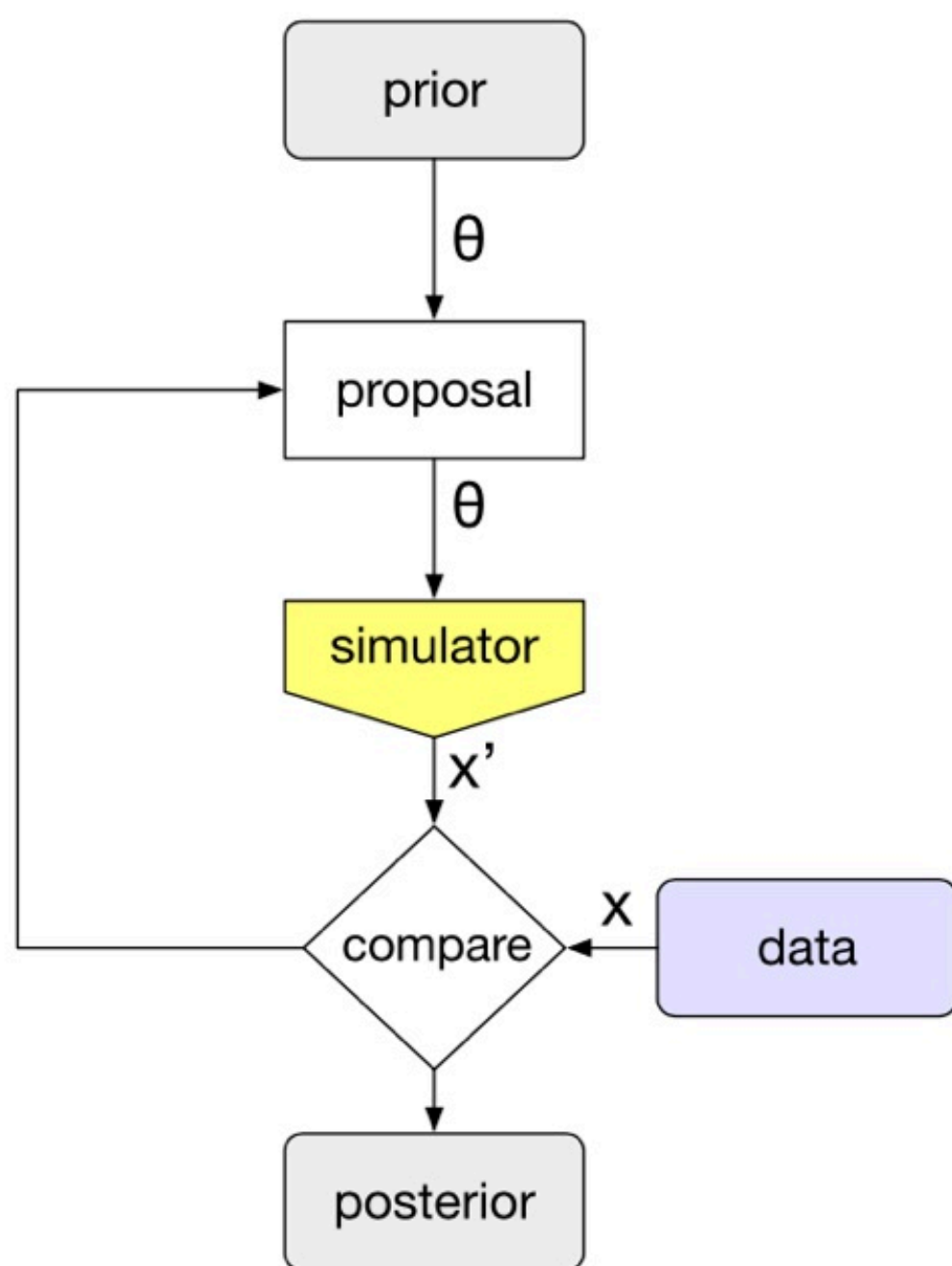
Probabilistic Programming  
with Monte Carlo sampling



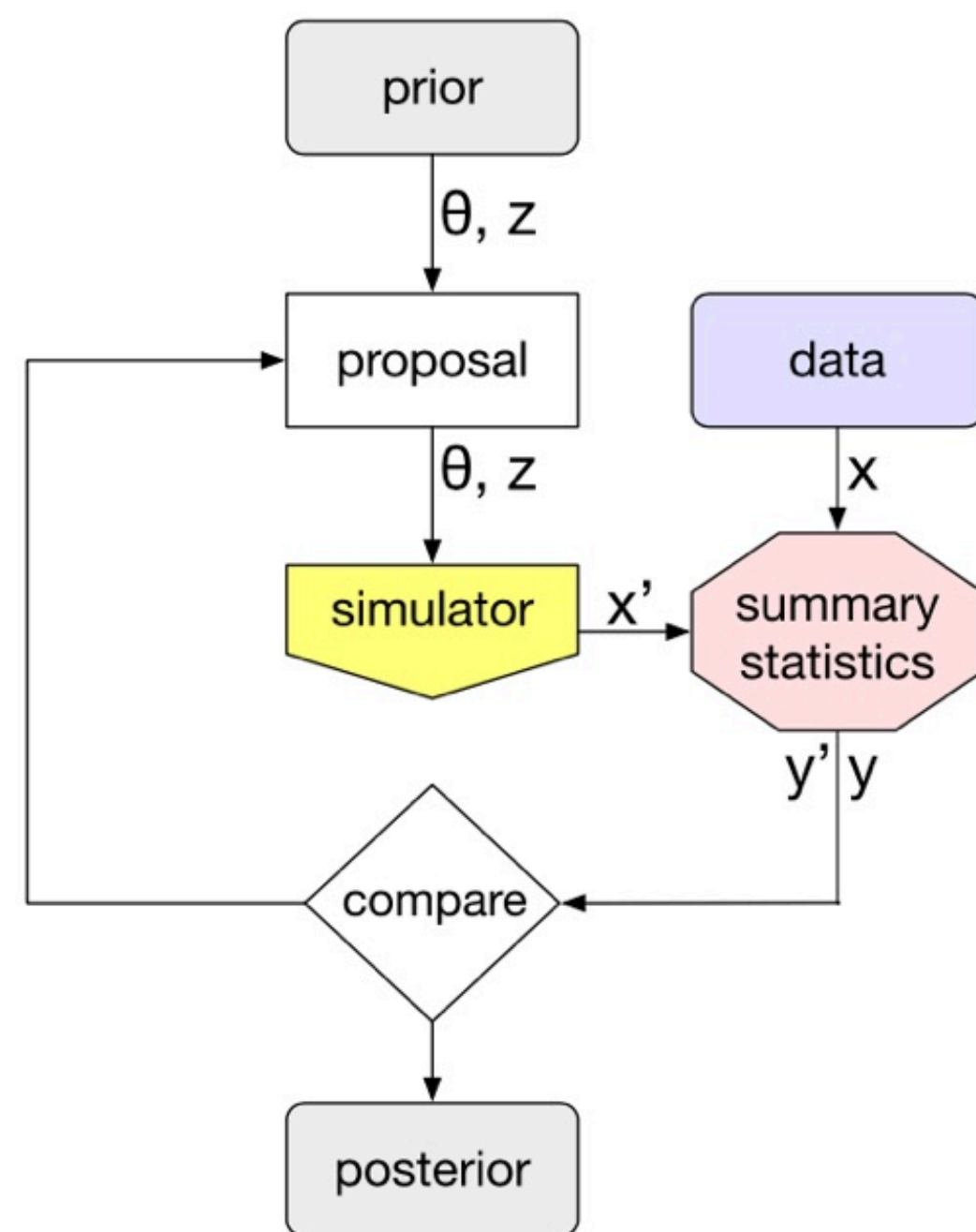
(Cranmer et al 2019)

# Simulation-Based Inference: Approximate Bayesian MC

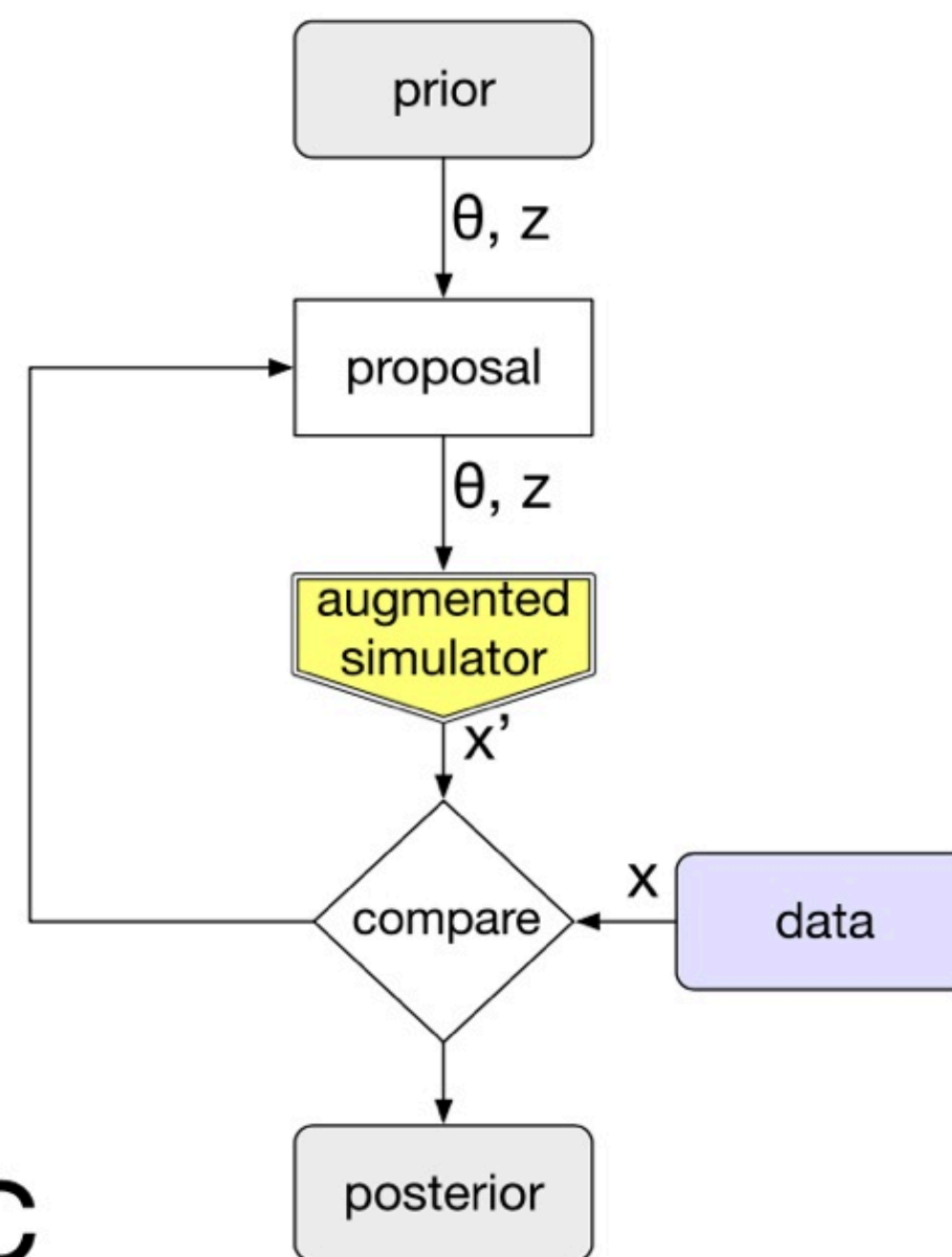
Approximate Bayesian Computation  
with Monte Carlo sampling



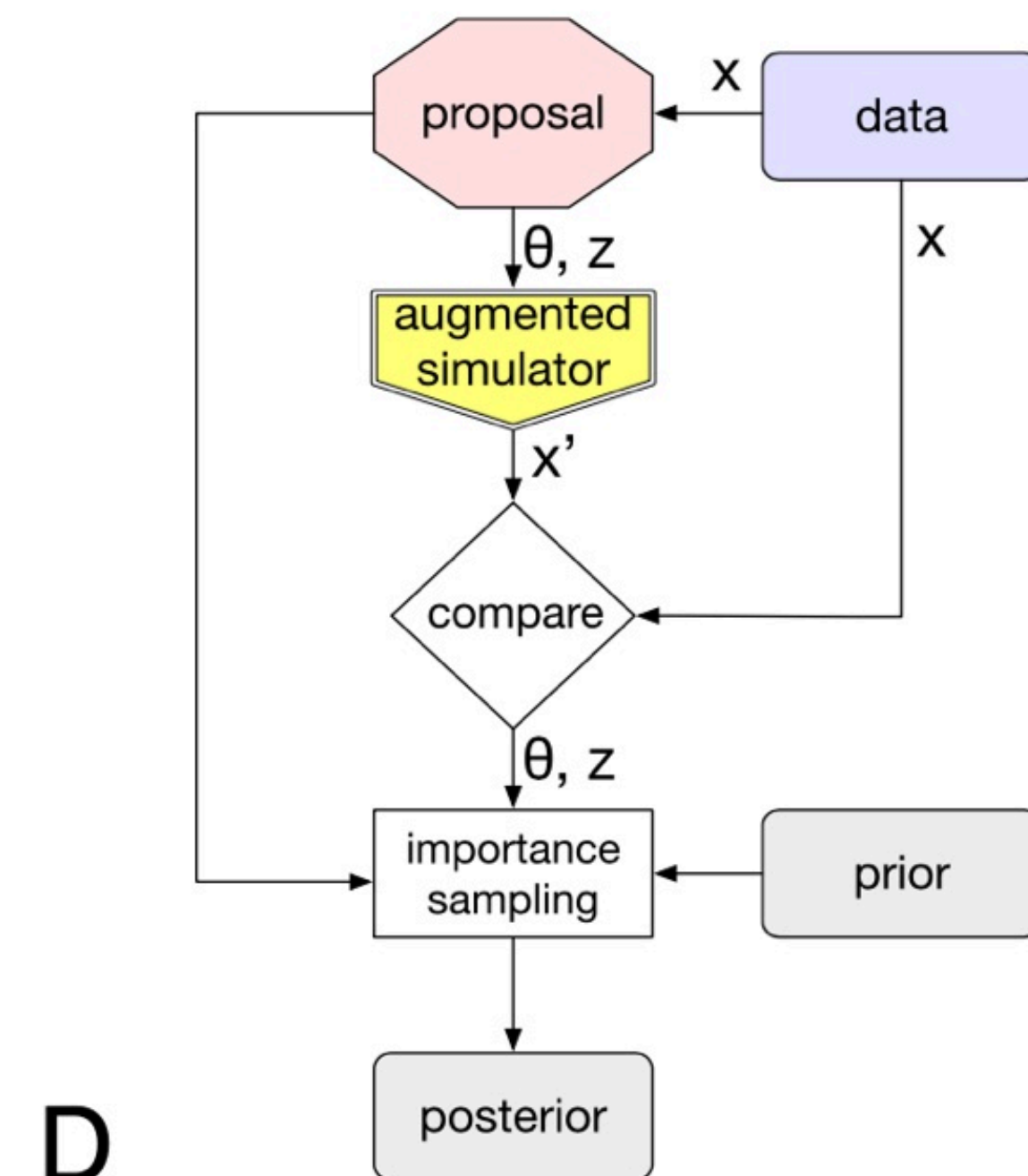
Approximate Bayesian Computation  
with learned summary statistics



Probabilistic Programming  
with Monte Carlo sampling



Probabilistic Programming  
with Inference Compilation



A

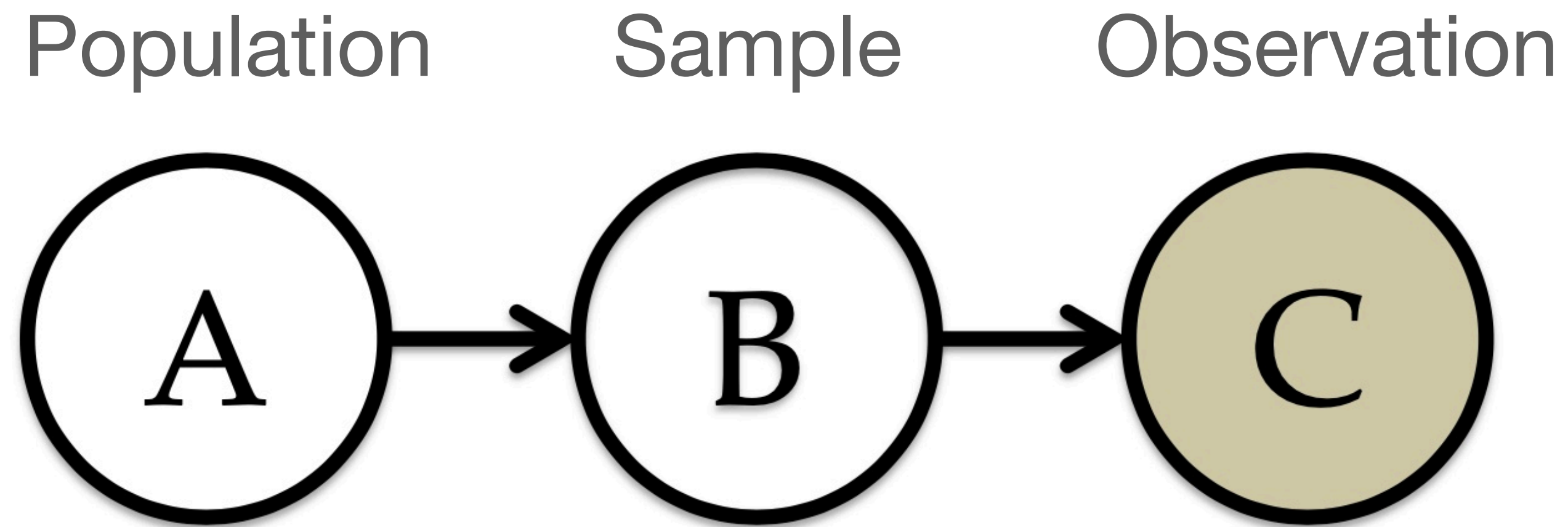
B

C

D

(Cranmer et al 2019)

# Simulation-Based Inference: Amortization

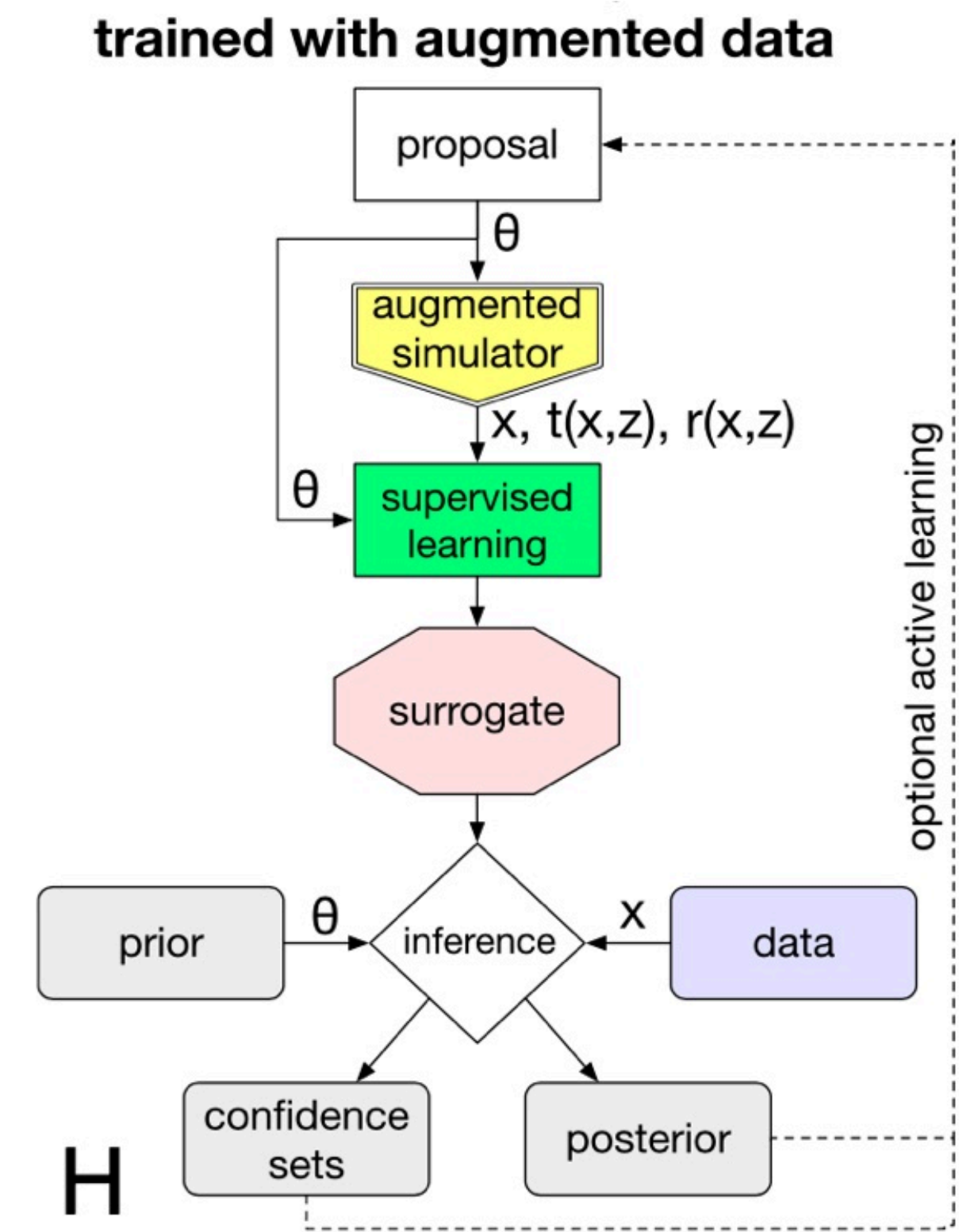
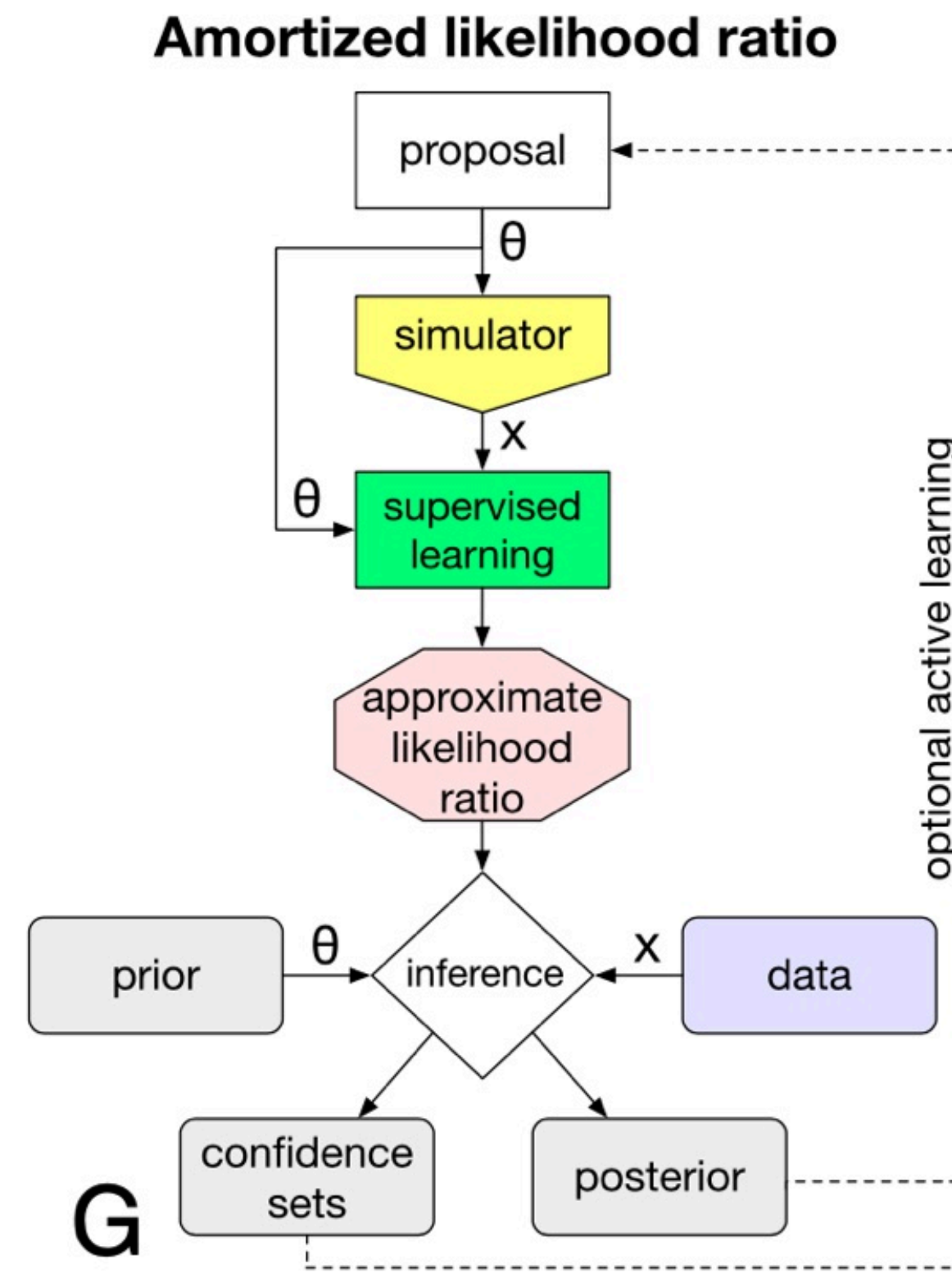
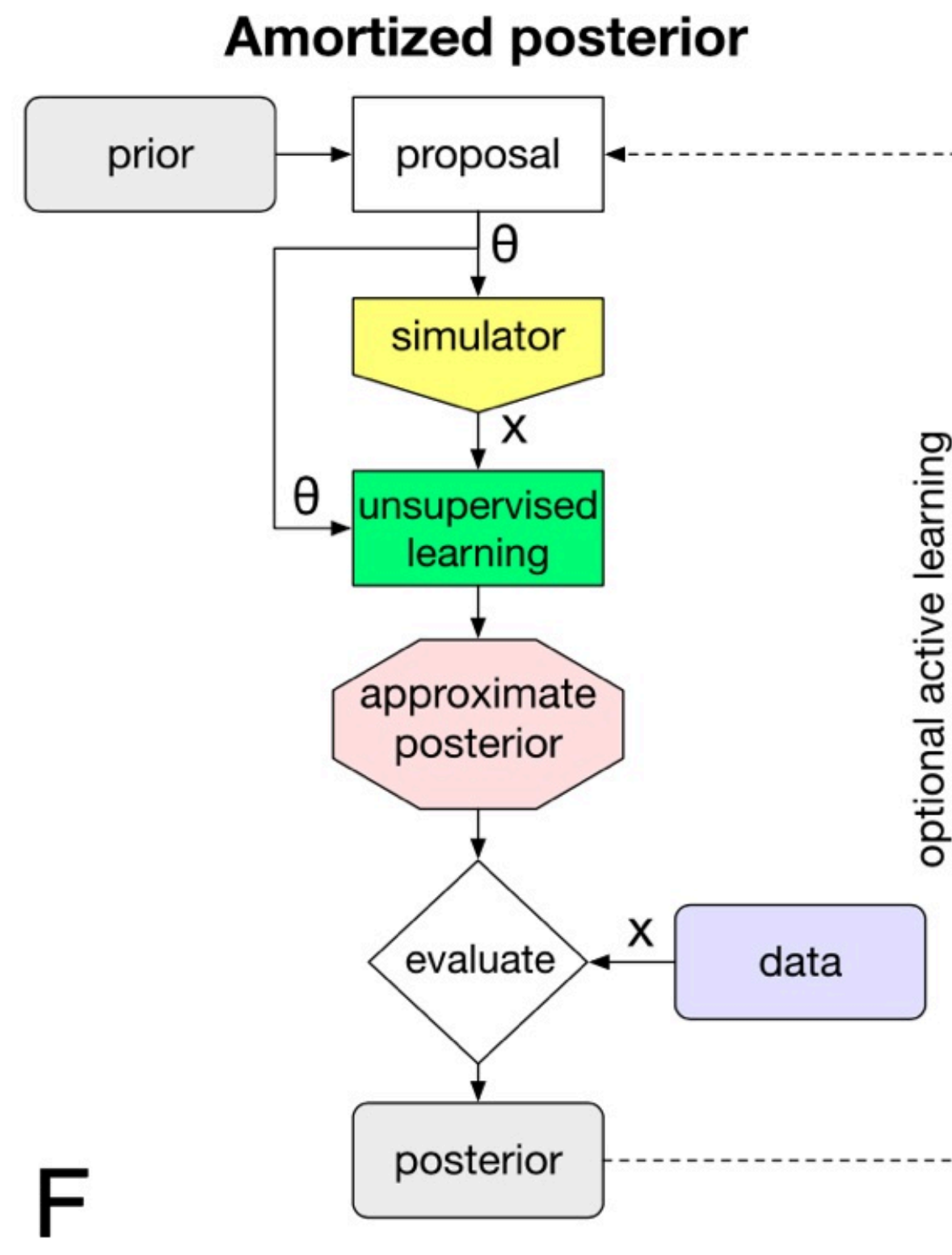
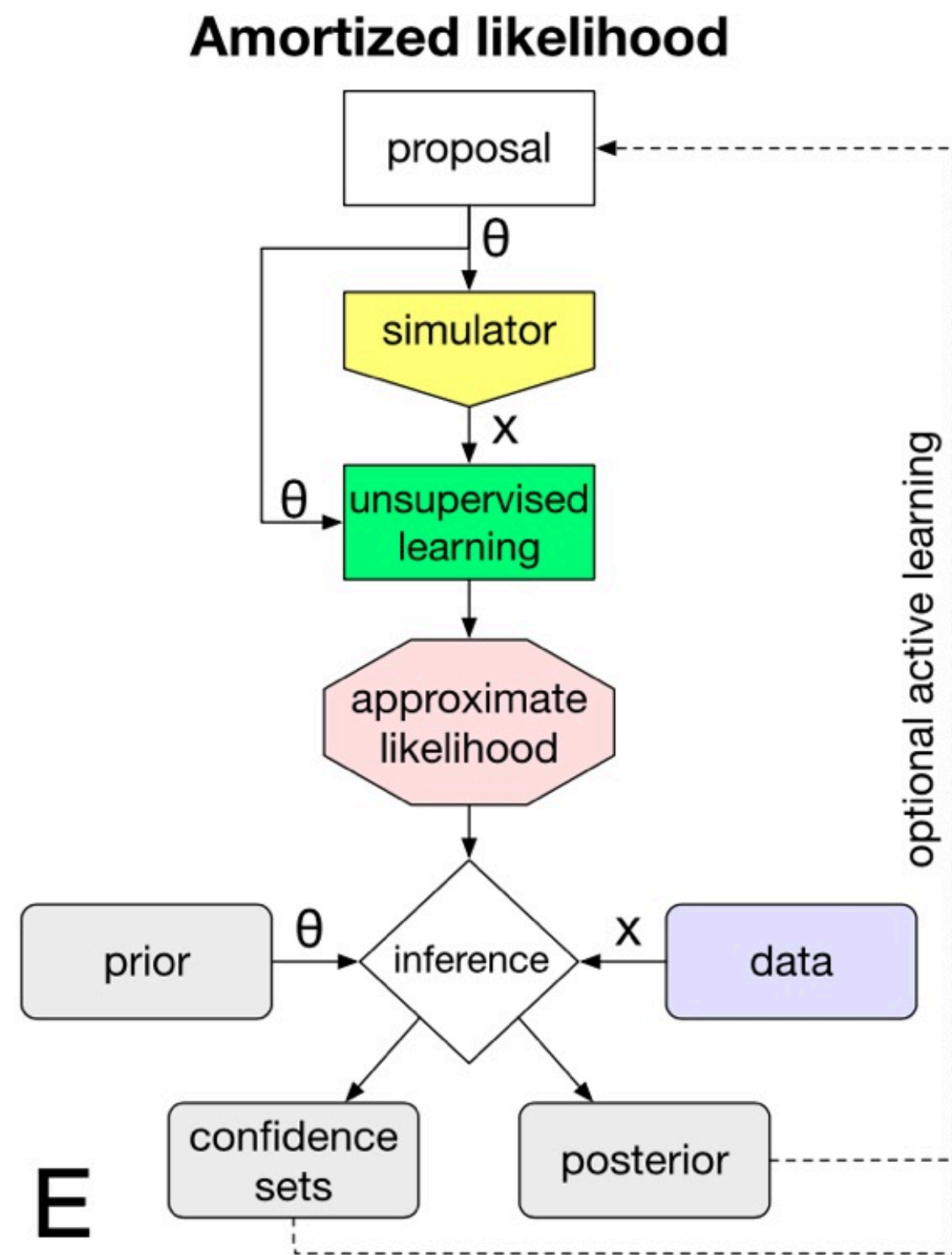


**Query 1:**  $P(B|C) = P(C|B)P(B)/P(C)$

**Query 2:**  $P(A|C) = \sum_B P(A|B)P(B|C)$

(Gershman et al 2014)

# Simulation-Based Inference: Amortisation Techniques



(Cranmer et al 2019)

# Simulation-Based Inference: Neural Network Approximations

Prior Variables  
 $\theta$

$$\theta \sim P(\theta)$$

Status  
Machine  
Parameterised  
by Latent  
Variables  $Z$

$$Z \sim P(Z | \theta)$$

Simulated  
Observations  
 $X$

$$X \sim P(X | Z, \theta)$$

$$P(\theta | X) = \frac{\overset{\text{Likelihood}}{P(X | \theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Evidence}}{P(X)}}$$

$\theta$ : parameters    $X$ : observations

- $P(\theta | X)$  approximated through “Neural Posterior” estimators
- $P(X | \theta)$  approximated through “Neural Likelihood” estimators
- $\frac{P(X | \theta)}{P(X)}$  approximated through the “Neural ratio” estimators

# Simulation-Based Inference: Neural Network Approximations Through Stochastic Flows

$$P(\theta | X) = \frac{\overset{\text{Likelihood}}{P(X | \theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Evidence}}{P(X)}}$$

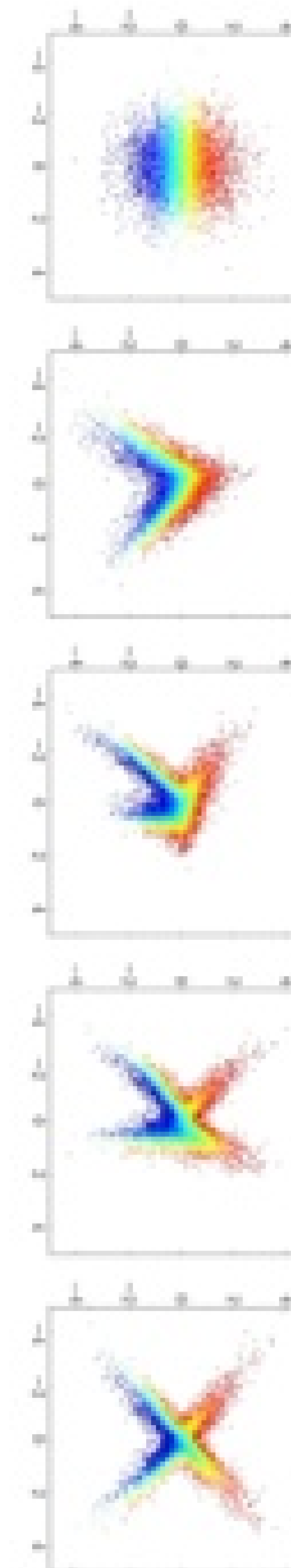
$\theta$ : parameters     $X$ : observations

Posterior

$$f(X, \theta) = N_{\mu, \Sigma}(\phi(X, \theta)) |J_{\phi}(X, \theta)|$$

$f$  the Neural estimator and  $\phi$  the stochastic flow

- $P(\theta | X)$  approximated through “Neural Posterior” estimators
- $P(X | \theta)$  approximated through “Neural Likelihood” estimators
- $\frac{P(X | \theta)}{P(X)}$  approximated through the “Neural ratio” estimators



# Simulation-Based Inference: Automatic Posterior Transformation (Greenberg et al 2019)

$$P(\theta | X) = \frac{\overset{\text{Likelihood}}{P(X | \theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Evidence}}{P(X)}}$$

$\theta$ : parameters    $X$ : observations

- $P(\theta | X)$  approximated through “Neural posterior” by a flow  $Q_{F(x_0, \phi)}(\theta)$

- Loss function:

$$\tilde{q}_{x, \phi}(\theta) = q_{F(x, \phi)}(\theta) \frac{\tilde{p}(\theta)}{p(\theta)} \frac{1}{Z(x, \phi)}, \quad (2)$$

- Where a proposal posterior is

$$\tilde{p}(\theta | x) = p(\theta | x) \frac{\tilde{p}(\theta) p(x)}{p(\theta) \tilde{p}(x)}$$

---

## Algorithm 1 APT with per-round proposal updates

---

**Input:** simulator with (implicit) density  $p(x|\theta)$ , data  $x_o$ , prior  $p(\theta)$ , density family  $q_\psi$ , neural network  $F(x, \phi)$ , simulations per round  $N$ , number of rounds  $R$ .

```

 $\tilde{p}_1(\theta) := p(\theta)$ 
for  $r = 1$  to  $R$  do
  for  $j = 1$  to  $N$  do
    Sample  $\theta_{r,j} \sim \tilde{p}_r(\theta)$ 
    Simulate  $x_{r,j} \sim p(x|\theta_{r,j})$ 
  end for
   $\phi \leftarrow \underset{\phi}{\operatorname{argmin}} \sum_{i=1}^r \sum_{j=1}^N -\log \tilde{q}_{x_{i,j}, \phi}(\theta_{i,j})$     using (2)
   $\tilde{p}_{r+1}(\theta) := q_{F(x_o, \phi)}(\theta)$ 
end for
return  $q_{F(x_o, \phi)}(\theta)$ 

```

---

# Simulation-Based Inference: Sequential Neural Likelihood (Papamakarios et al 2019)

$$P(\theta | X) = \frac{\overset{\text{Likelihood}}{P(X | \theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Evidence}}{P(X)}}$$

$\theta$ : parameters     $X$ : observations

Posterior

- $P(X | \theta)$  approximated through “Neural likelihood” by a flow  $Q_\phi(X | \theta)$

---

## Algorithm 1: Sequential Neural Likelihood (SNL)

---

**Input** : observed data  $\mathbf{x}_o$ , estimator  $q_\phi(\mathbf{x} | \theta)$ , number of rounds  $R$ , simulations per round  $N$

**Output**: approximate posterior  $\hat{p}(\theta | \mathbf{x}_o)$

set  $\hat{p}_0(\theta | \mathbf{x}_o) = p(\theta)$  and  $\mathcal{D} = \{\}$

**for**  $r = 1 : R$  **do**

**for**  $n = 1 : N$  **do**

sample  $\theta_n \sim \hat{p}_{r-1}(\theta | \mathbf{x}_o)$  with MCMC

simulate  $\mathbf{x}_n \sim p(\mathbf{x} | \theta_n)$

add  $(\theta_n, \mathbf{x}_n)$  into  $\mathcal{D}$

(re-)train  $q_\phi(\mathbf{x} | \theta)$  on  $\mathcal{D}$  and set

$\hat{p}_r(\theta | \mathbf{x}_o) \propto q_\phi(\mathbf{x}_o | \theta) p(\theta)$

**return**  $\hat{p}_R(\theta | \mathbf{x}_o)$

---

# Simulation-Based Inference: Neural Ratio (Hermans et al 2020)

$$P(\theta | X) = \frac{\overset{\text{Likelihood}}{P(X | \theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Evidence}}{P(X)}}$$

$\theta$ : parameters    $X$ : observations

Posterior

- $P(X | \theta)/P(X)$  approximated through “Neural ratio” by a flow  $d_\phi(X | \theta)$

---

## Algorithm 1 Optimization of $d_\phi(\mathbf{x}, \theta)$ .

---

*Inputs:*                      Criterion  $\ell$  (e.g., BCE)  
                                      Implicit generative model  $p(\mathbf{x} | \theta)$   
                                      Prior  $p(\theta)$

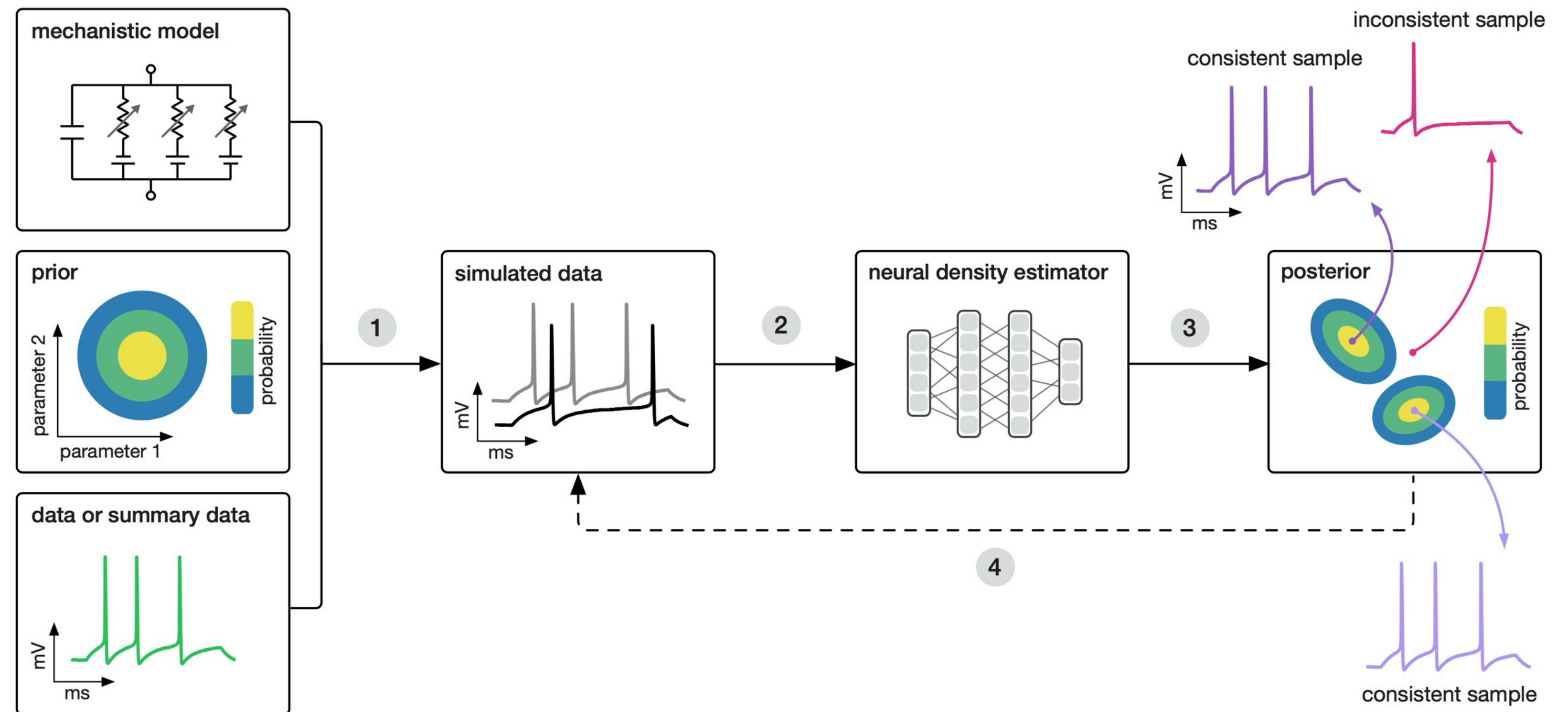
*Outputs:*                     Parameterized classifier  $d_\phi(\mathbf{x}, \theta)$

*Hyperparameters:*        Batch-size  $M$

- 1: **while not converged do**
  - 2:     **Sample**  $\theta \leftarrow \{\theta_m \sim p(\theta)\}_{m=1}^M$
  - 3:     **Sample**  $\theta' \leftarrow \{\theta'_m \sim p(\theta)\}_{m=1}^M$
  - 4:     **Simulate**  $\mathbf{x} \leftarrow \{\mathbf{x}_m \sim p(\mathbf{x} | \theta_m)\}_{m=1}^M$
  - 5:      $\mathcal{L} \leftarrow \ell(d_\phi(\mathbf{x}, \theta), 1) + \ell(d_\phi(\mathbf{x}, \theta'), 0)$
  - 6:      $\phi \leftarrow \text{OPTIMIZER}(\phi, \nabla_\phi \mathcal{L})$
  - 7: **end while**
  - 8: **return**  $d_\phi$
-

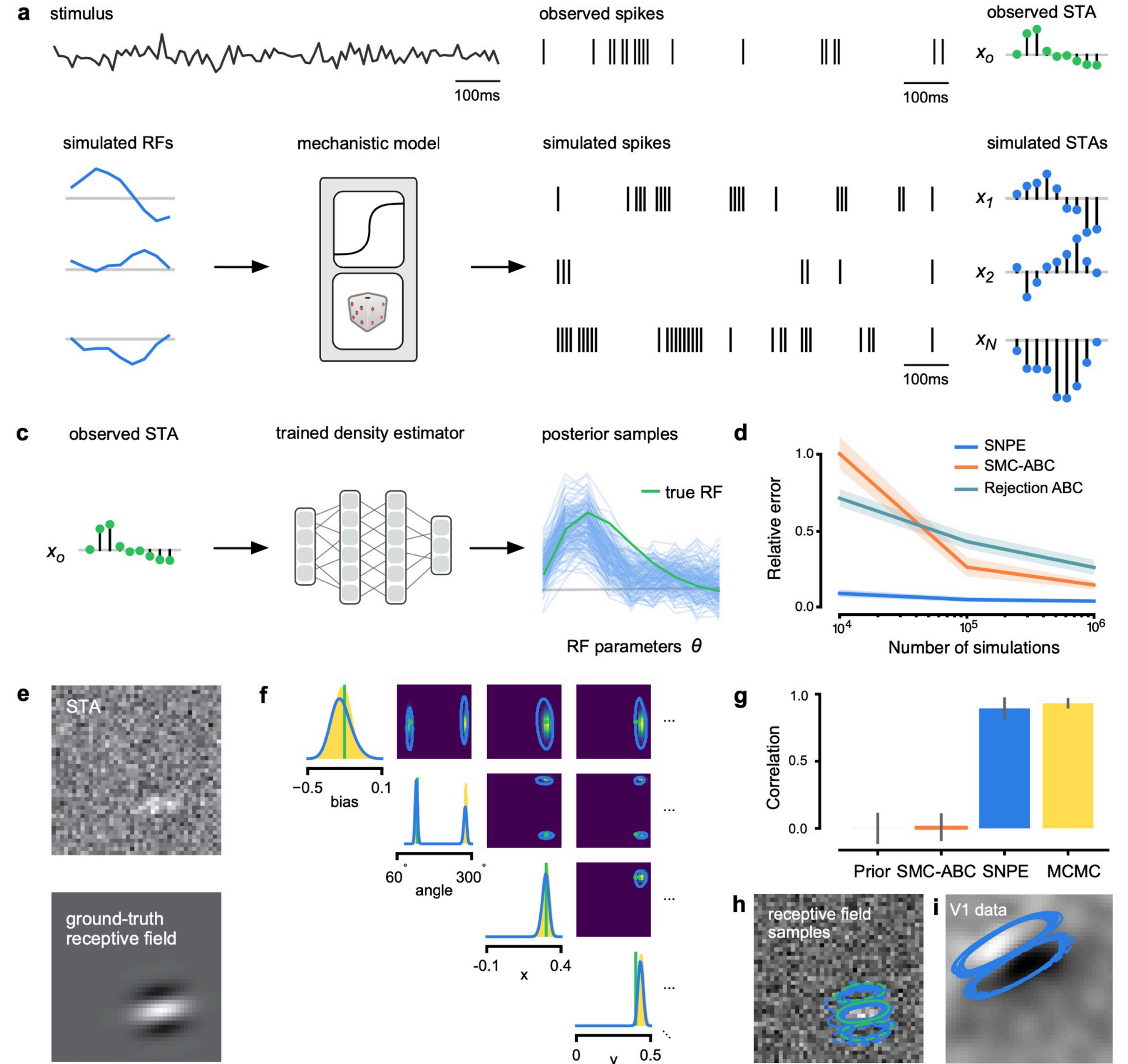
# Training deep neural density estimators to identify mechanistic models of neural dynamics

Pedro J Gonçalves<sup>1,2†\*</sup>, Jan-Matthis Lueckmann<sup>1,2†\*</sup>, Michael Deistler<sup>1,3†\*</sup>, Marcel Nonnenmacher<sup>1,2,4</sup>, Kaan Öcal<sup>2,5</sup>, Giacomo Bassetto<sup>1,2</sup>, Chaitanya Chintaluri<sup>6,7</sup>, William F Podlaski<sup>6</sup>, Sara A Haddad<sup>8</sup>, Tim P Vogels<sup>6,7</sup>, David S Greenberg<sup>1,4</sup>, Jakob H Macke<sup>1,2,3,9\*</sup>



# Training deep neural density estimators to identify mechanistic models of neural dynamics

Pedro J Gonçalves<sup>1,2†\*</sup>, Jan-Matthis Lueckmann<sup>1,2†\*</sup>, Michael Deistler<sup>1,3†\*</sup>, Marcel Nonnenmacher<sup>1,2,4</sup>, Kaan Öcal<sup>2,5</sup>, Giacomo Bassetto<sup>1,2</sup>, Chaitanya Chintaluri<sup>6,7</sup>, William F Podlaski<sup>6</sup>, Sara A Haddad<sup>8</sup>, Tim P Vogels<sup>6,7</sup>, David S Greenberg<sup>1,4</sup>, Jakob H Macke<sup>1,2,3,9\*</sup>



---

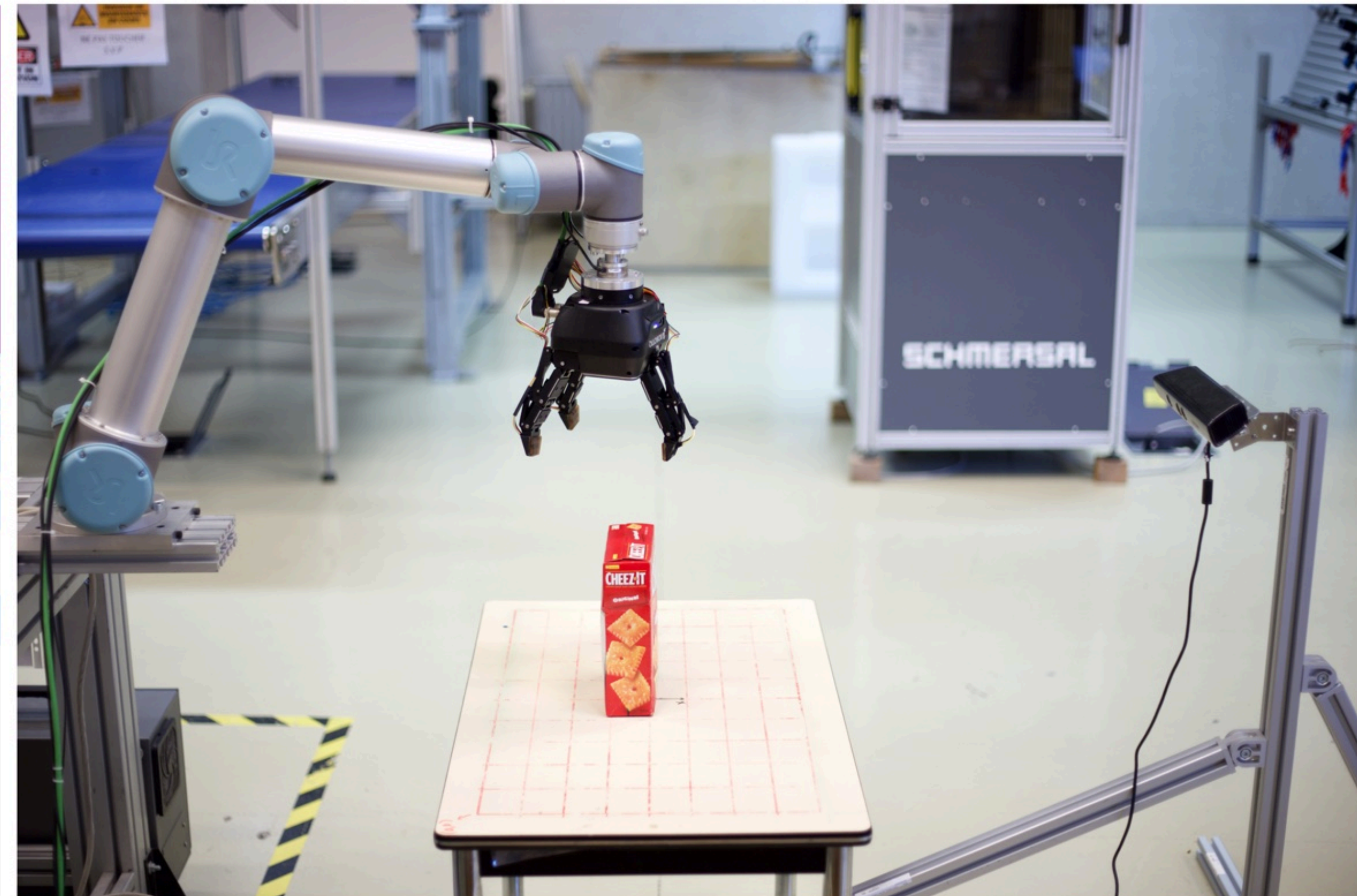
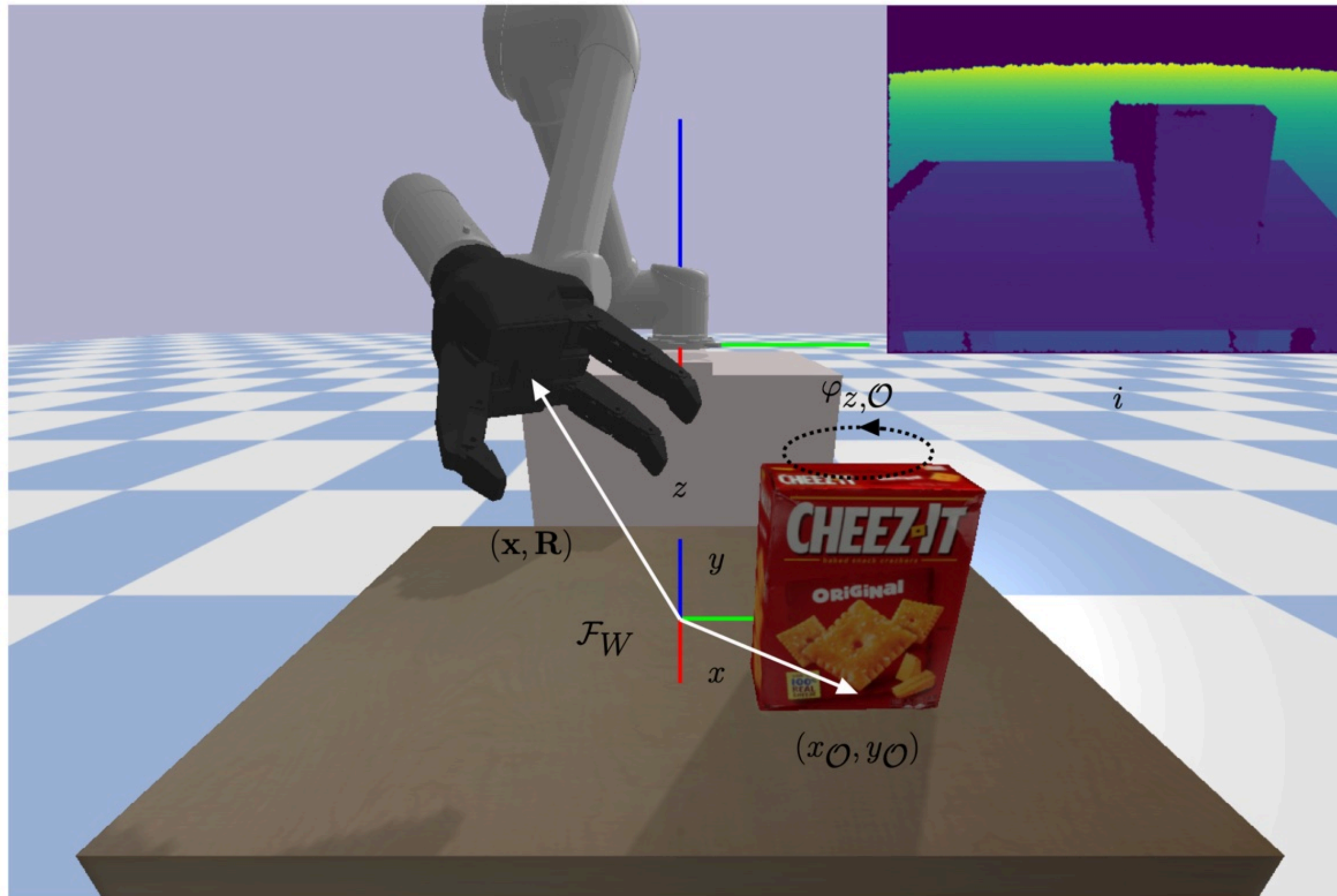
# SIMULATION-BASED BAYESIAN INFERENCE FOR MULTI-FINGERED ROBOTIC GRASPING

---

**Norman Marlier**  
University of Liège  
norman.marlier@uliege.be

**Olivier Bruls**  
University of Liège  
o.bruls@uliege.be

**Gilles Louppe**  
University of Liège  
g.louppe@uliege.be

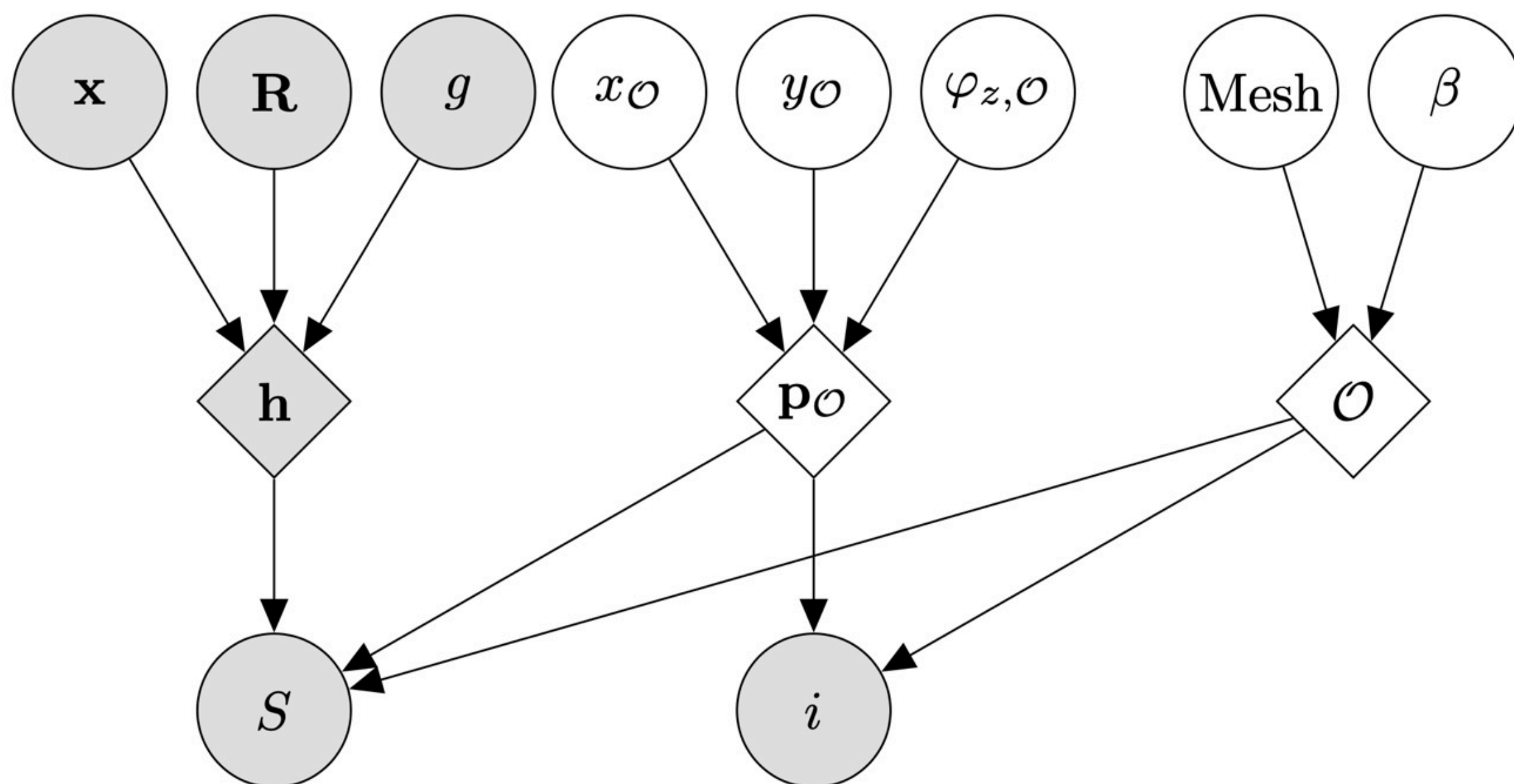


# SIMULATION-BASED BAYESIAN INFERENCE FOR MULTI-FINGERED ROBOTIC GRASPING

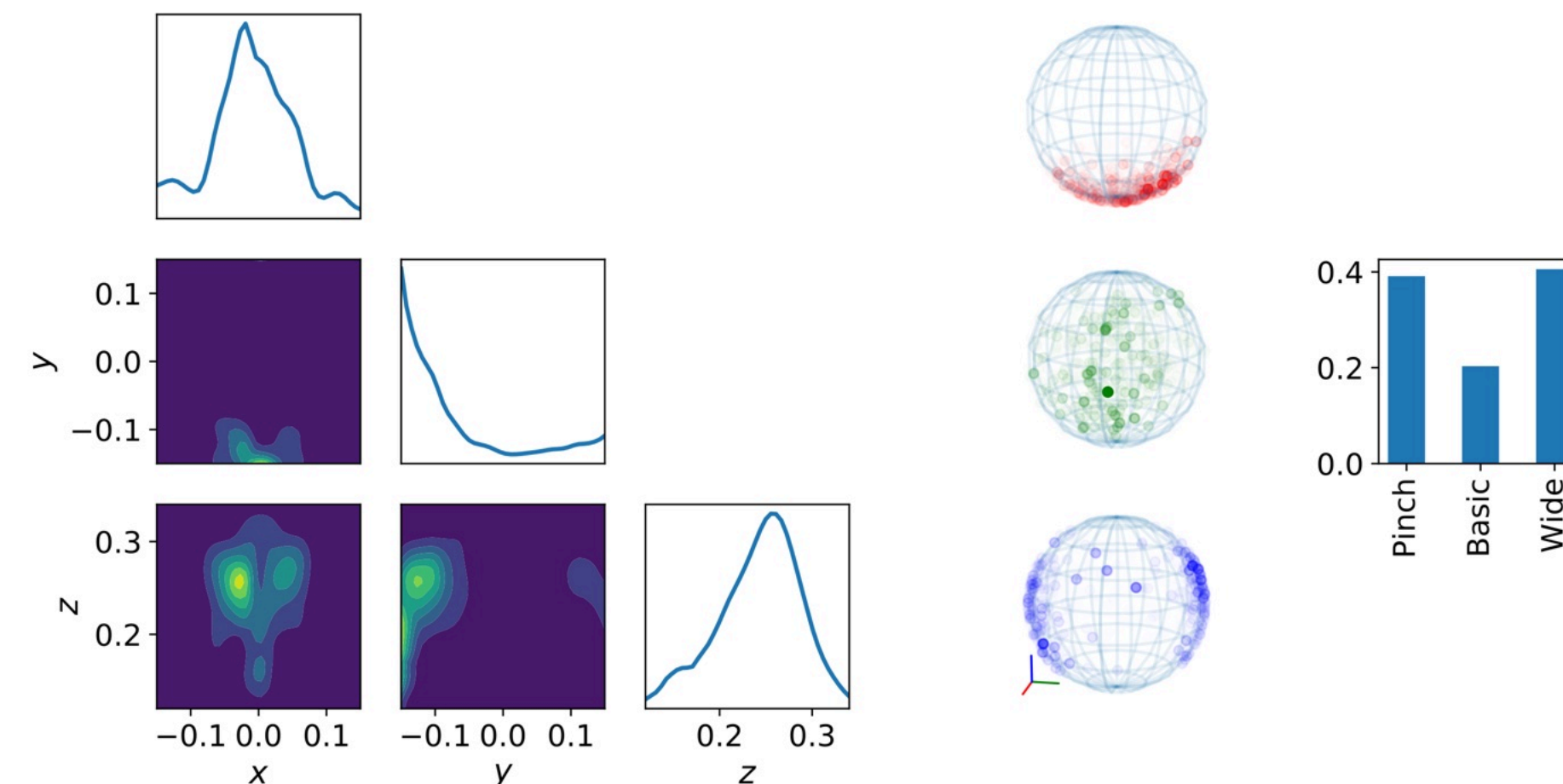
**Norman Marlier**  
University of Liège  
norman.marlier@uliege.be

**Olivier Brùls**  
University of Liège  
o.bruls@uliege.be

**Gilles Louppe**  
University of Liège  
g.louppe@uliege.be



(a)



Variable	Prior
$x$	$\text{uniform}(-0.15, 0.15)$
$y$	$\text{uniform}(-0.15, 0.15)$
$z$	$\text{uniform}(0.12, 0.34)$
$\mathbf{R}$	mixture of power spherical( $\mu_i, \kappa$ )
$g$	categorical( $\{\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\}$ )
$x_{\mathcal{O}}$	$\text{uniform}(-0.05, 0.05)$
$y_{\mathcal{O}}$	$\text{uniform}(-0.05, 0.05)$
$\varphi_{z,\mathcal{O}}$	$\text{uniform}(-\pi, \pi)$
Mesh	uniform in the set of objects
$\beta$	$\text{uniform}(0.9, 1.1)$

(b)

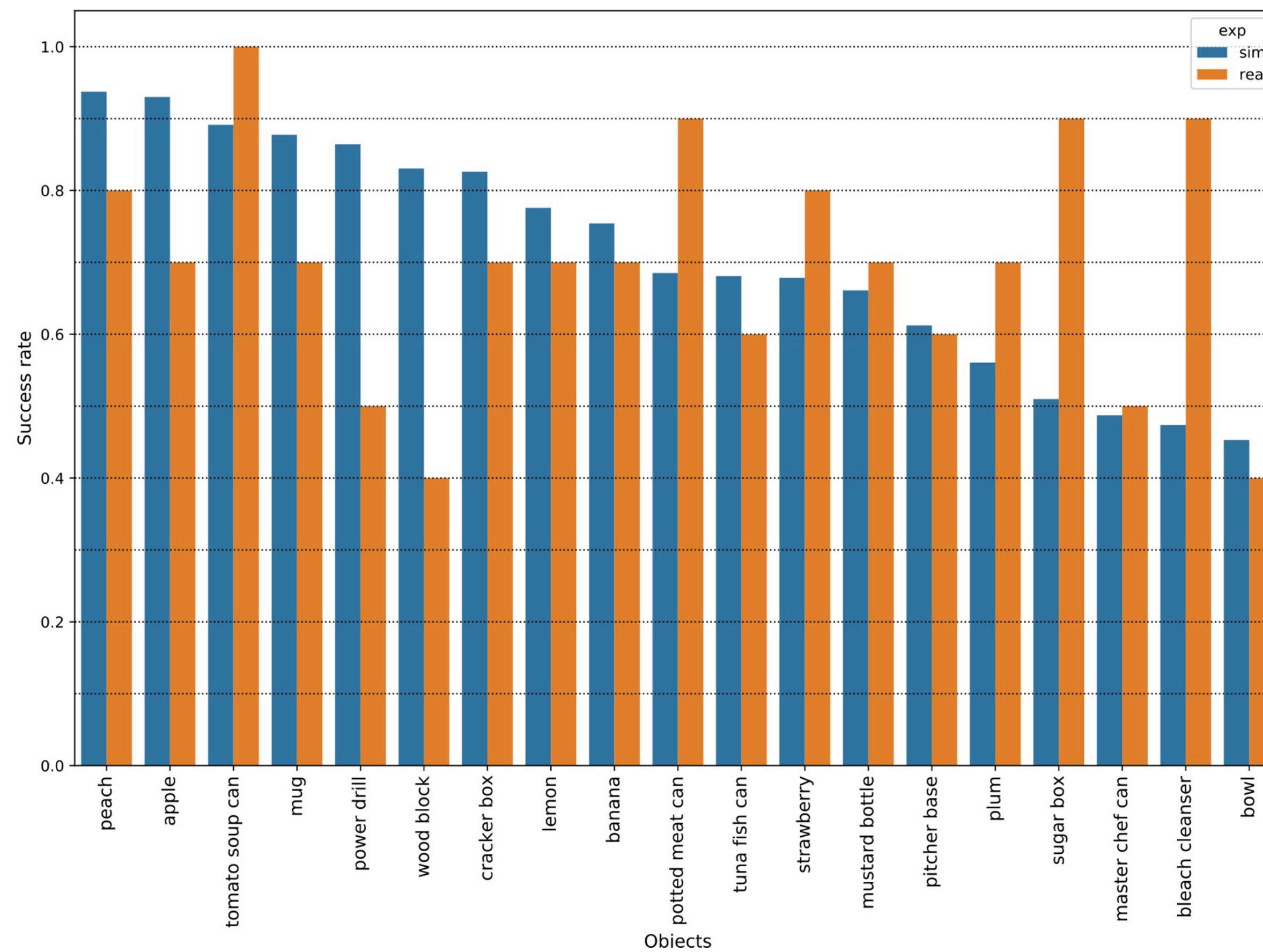
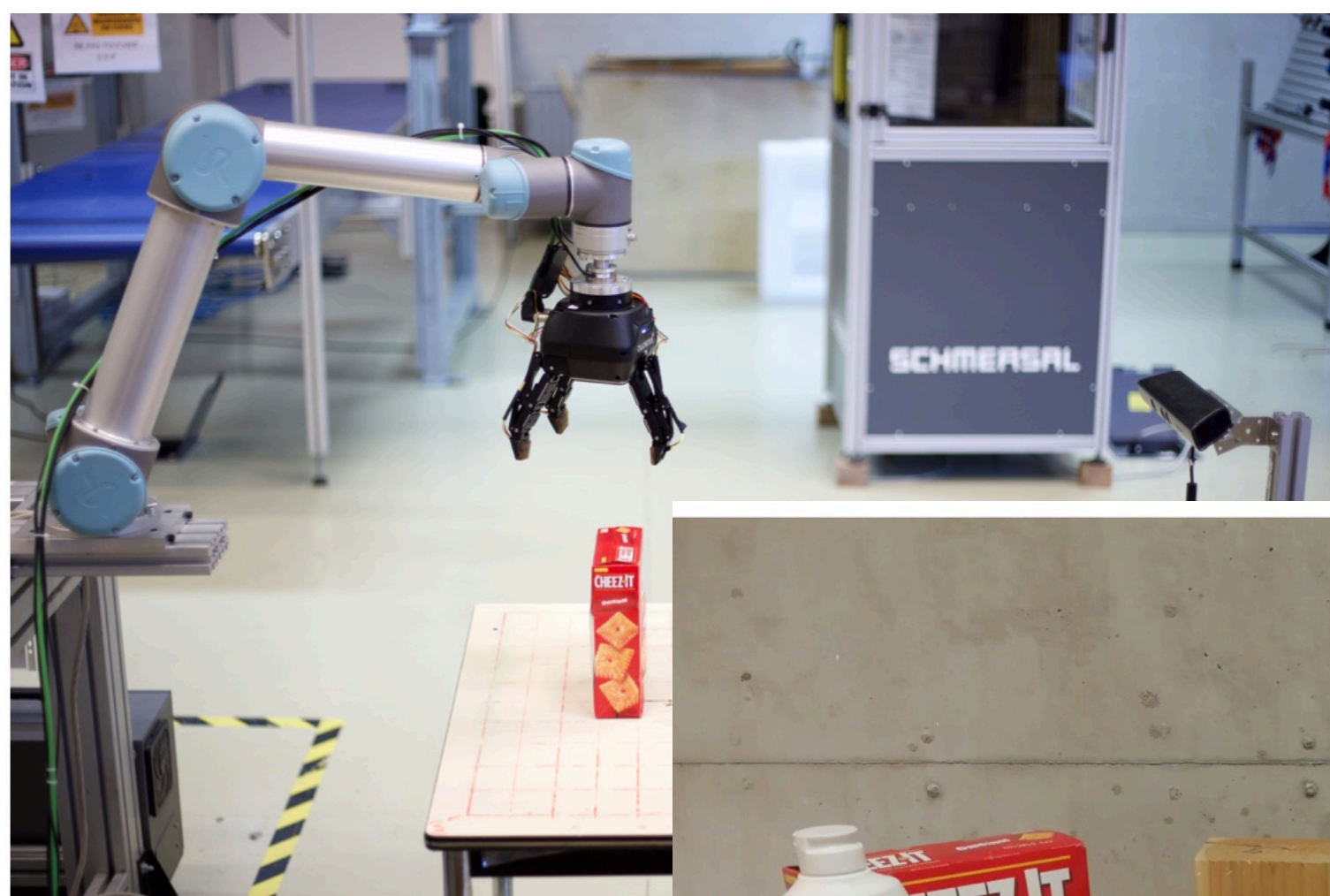
Figure 2: (a) Probabilistic graphical model of the environment. Gray nodes correspond to observed variables and white nodes to unobserved variables. (b) Prior distributions.

# SIMULATION-BASED BAYESIAN INFERENCE FOR MULTI-FINGERED ROBOTIC GRASPING

**Norman Marlier**  
University of Liège  
norman.marlier@uliege.be

**Olivier Brûls**  
University of Liège  
o.bruls@uliege.be

**Gilles Louppe**  
University of Liège  
g.louppe@uliege.be



# Current Uses of Simulation-Based Inference

## Tabular Foundation Models

nature

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [articles](#) > article

Article | [Open access](#) | Published: 08 January 2025

### Accurate predictions on small data with a tabular foundation model

[Noah Hollmann](#) ✉, [Samuel Müller](#) ✉, [Lennart Purucker](#), [Arjun Krishnakumar](#), [Max Körfer](#), [Shi Bin Hoo](#),

[Robin Tibor Schirrmeyer](#) & [Frank Hutter](#) ✉

*Nature* **637**, 319–326 (2025) | [Cite this article](#)

441k Accesses | 475 Citations | 518 Altmetric | [Metrics](#)

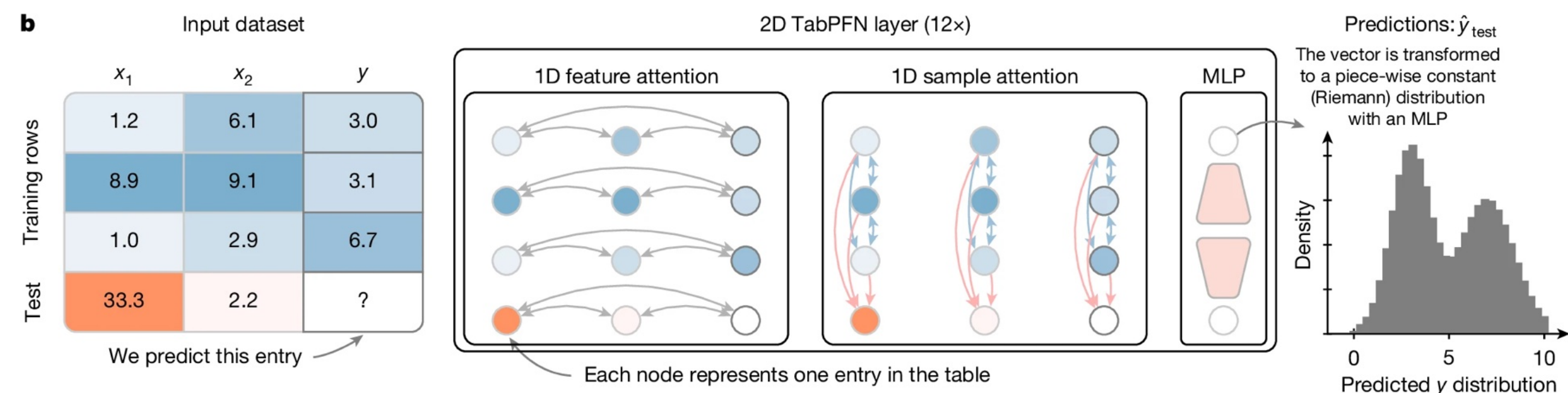
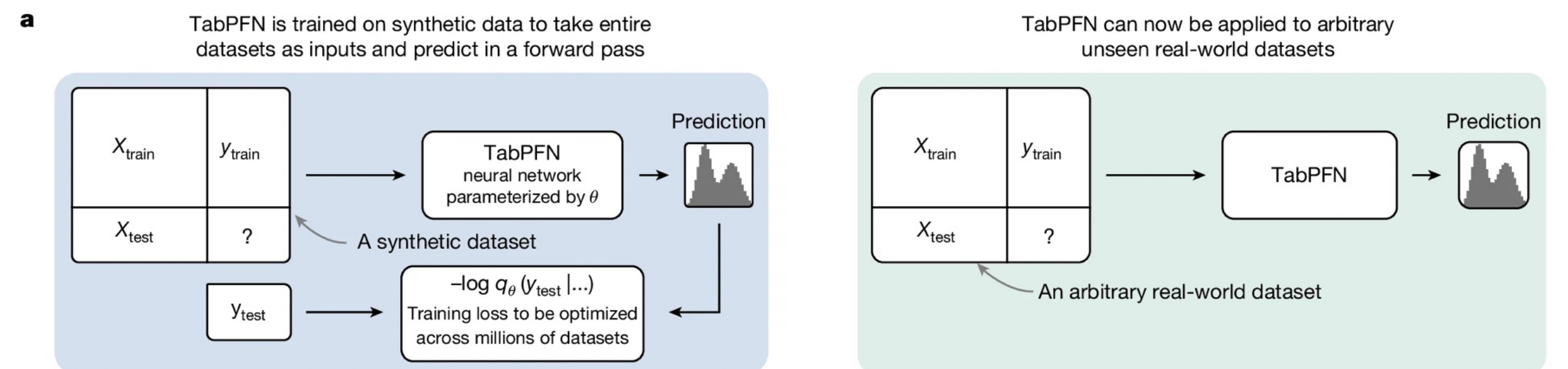
arXiv > cs > arXiv:2207.01848

Computer Science > Machine Learning

[Submitted on 5 Jul 2022 (v1), last revised 16 Sep 2023 (this version, v6)]

### TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second

Noah Hollmann, Samuel Müller, Katharina Eggenberger, Frank Hutter



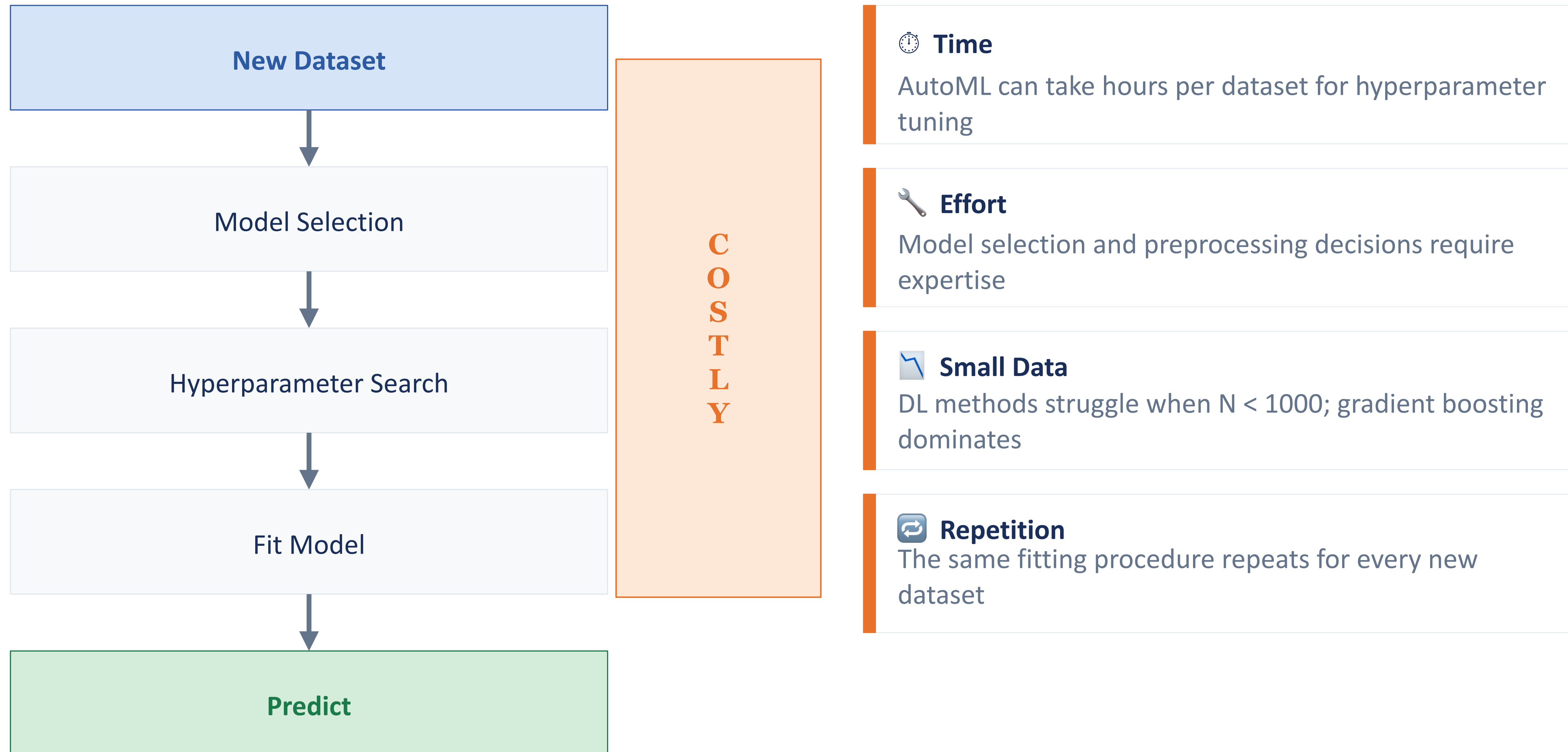
# Current Uses of Simulation-Based Inference

## Tabular Foundation Models

- The problem:  
Why tabular classification is hard; fitting a new model per dataset is expensive and brittle.
- PFN Framework  
Prior-Data Fitted Networks: approximating the Bayesian Posterior Predictive Distribution offline.
- Graphical Model Prior  
Structural Causal Models (SCMs) and BNNs as data generators — the backbone of TabPFN's world model.
- Simulation-Based Inference  
Generating millions of synthetic datasets from the prior; training the transformer to invert the simulator.
- In-Context Learning  
How a single transformer forward pass replaces dataset-specific fitting at inference time.
- Results & Insights  
What the causal/SCM prior buys us; performance vs. boosted trees and AutoML.

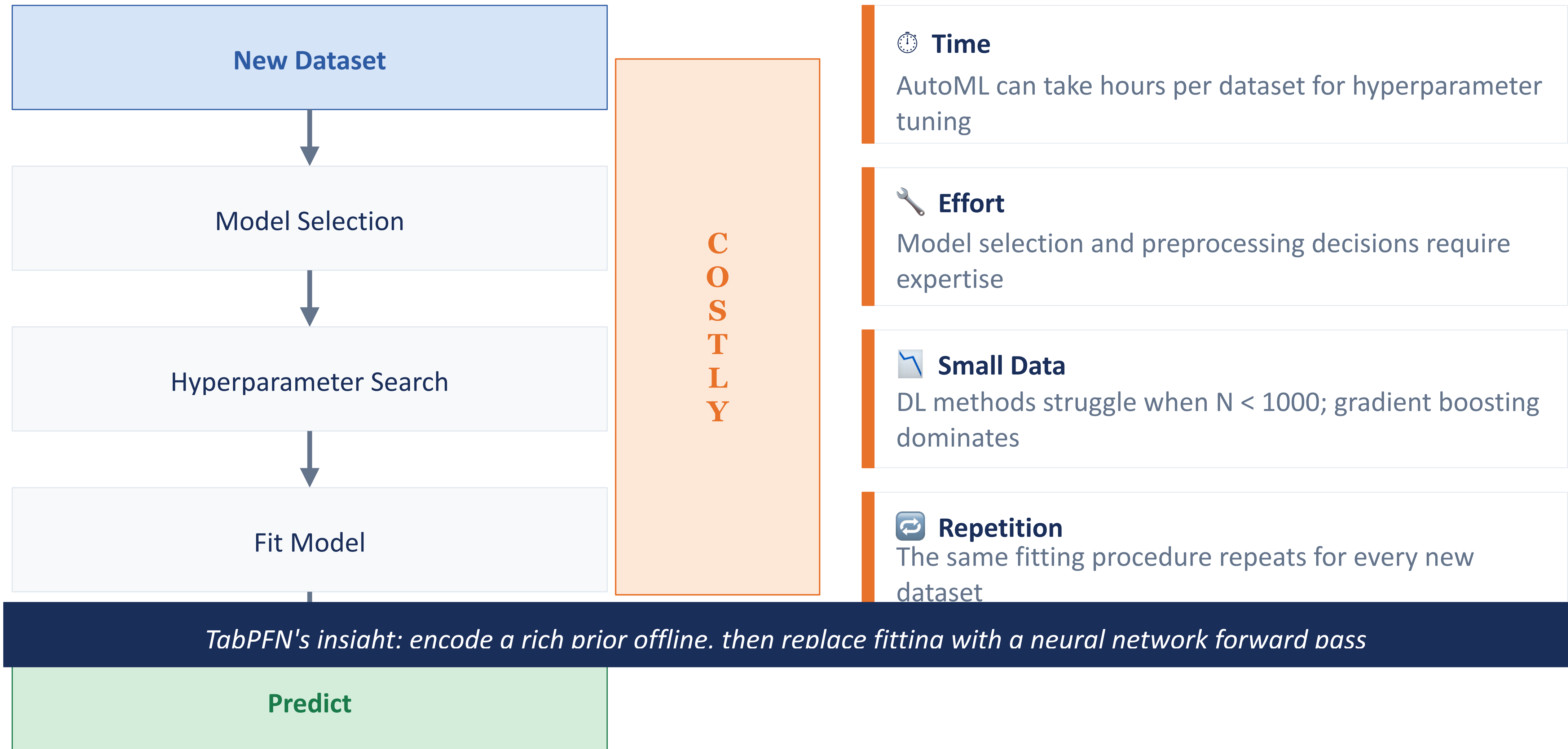
# The Problem to be Solved

*Each new dataset requires fitting from scratch — expensive and hard to tune*



# The Problem to be Solved

*Each new dataset requires fitting from scratch — expensive and hard to tune*



# Prior-Data Fitted Networks: Amortized Bayesian Inference

*Learn to approximate the Posterior Predictive Distribution (PPD) across a whole prior*

$$p(y^* \mid x^*, D_{\text{train}}) = \int p(y^* \mid x^*, \phi) \cdot p(\phi \mid D_{\text{train}}) d\phi$$

*PPD: target posterior*

*likelihood*

*parameter posterior*

*marginalise over  $\phi$*

## Classical Bayes

1. Fix prior  $p(\phi)$
2. Observe  $D_{\text{train}}$
3. Compute posterior  $p(\phi \mid D_{\text{train}})$
4. Marginalize analytically / MCMC
5.  $\rightarrow$  prediction for  $x^*$

*Expensive per dataset*

## PFN (TabPFN) Approach

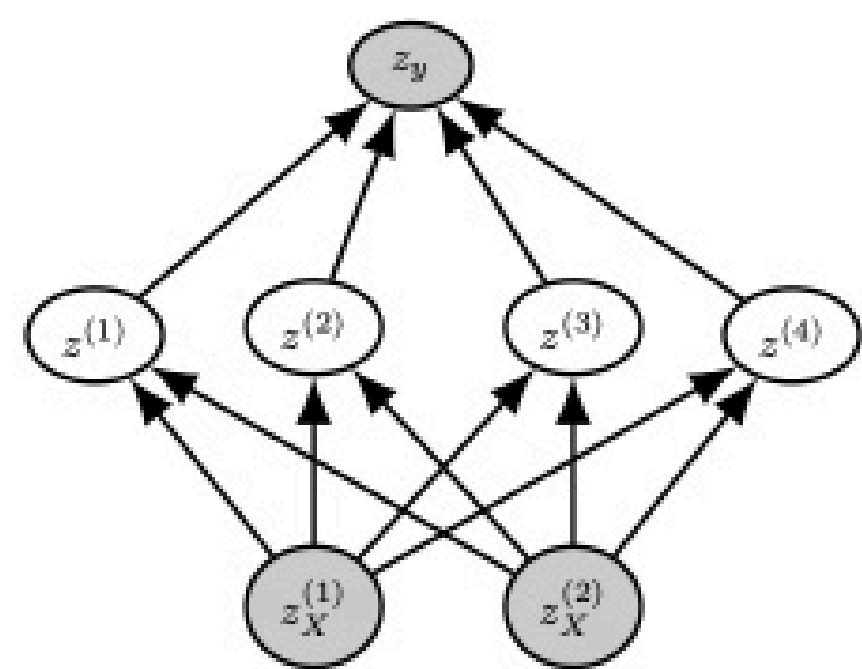
1. Design prior  $p(\phi)$  as a simulator
2. Sample datasets  $D \sim p(D \mid \phi), \phi \sim p(\phi)$
3. Train transformer:  $(D_{\text{train}}, x^*) \rightarrow \hat{y}$
4. Amortizes the integral over all  $\phi$
5.  $\rightarrow$  prediction = one forward pass

*Fast at inference; expensive only once offline*

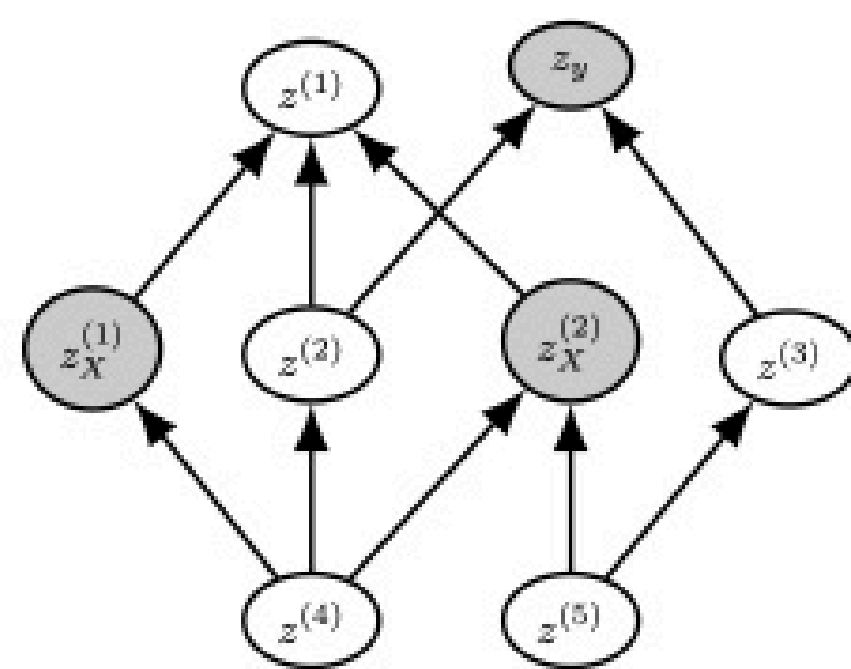
VS

*Key insight: the transformer weight vector  $\vartheta \equiv$  an amortized representation of the full posterior over the prior's hypothesis space*

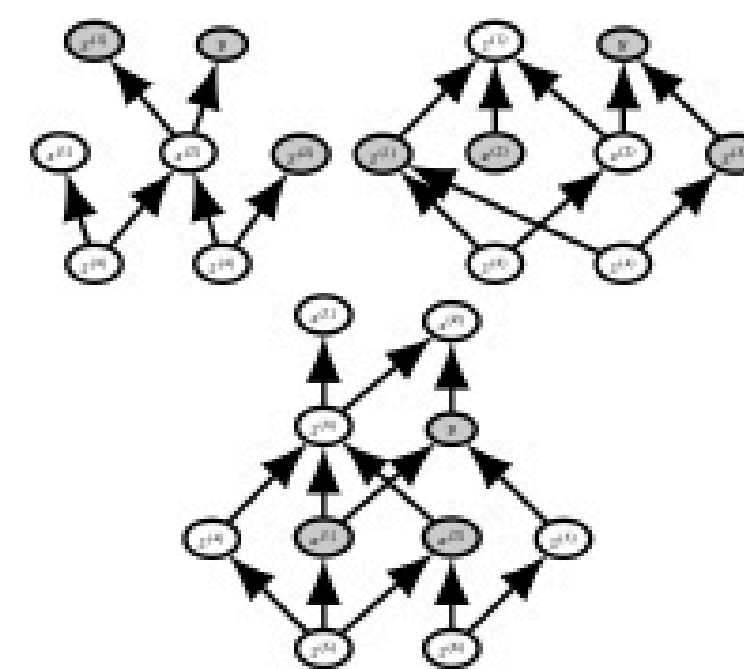
# Prior-Data Fitted Networks: Amortized Bayesian Inference



(a) A BNN



(b) An SCM

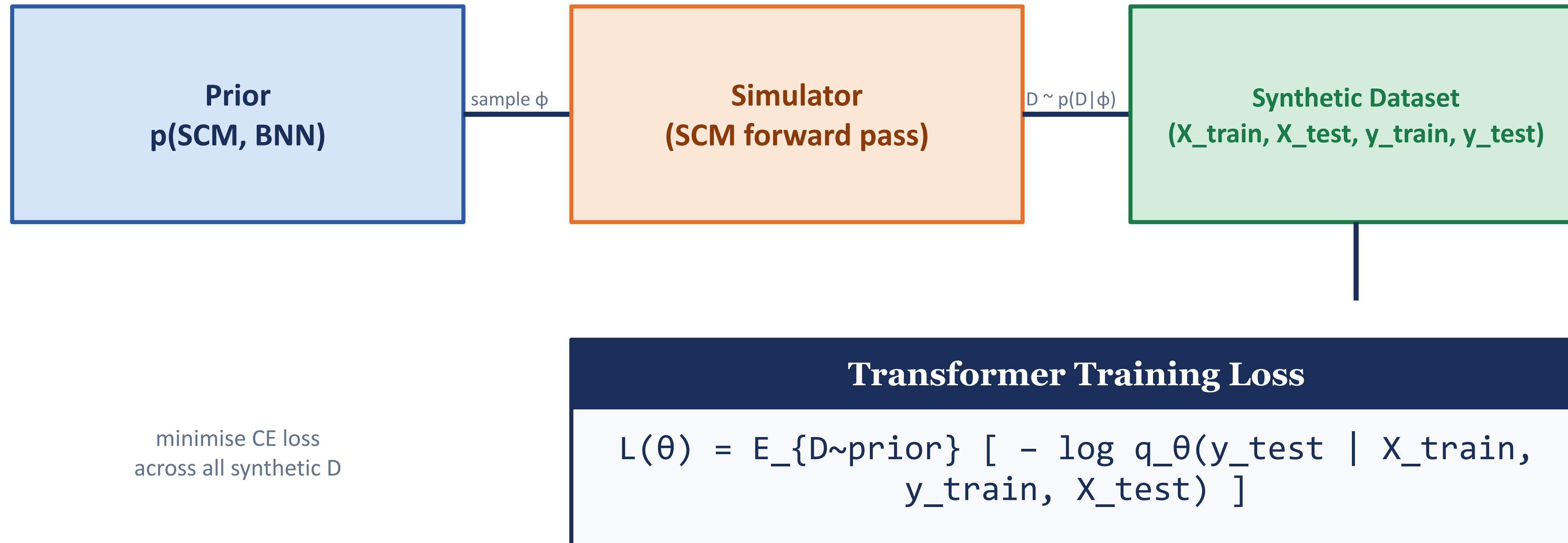


(c) SCMs sampled from the prior

Figure 2: Overview of graphs generating data in our prior. Inputs  $x$  are mapped to the output  $y$  through unobserved nodes  $z$ . Plots based on Müller et al. (2022).

# Simulation-Based Inference (SBI): The Training Paradigm

*TabPFN is trained like a neural posterior estimator — but applied to supervised prediction tasks*



TabPFN vs. Classical SBI		
	Classical SBI (e.g. NPE/SNPE)	TabPFN
<b>Goal</b>	Estimate $p(\theta   x_{\text{obs}})$	Estimate $p(y^*   x^*, D_{\text{train}})$
<b>Simulator</b>	Scientific forward model	SCM/BNN generative prior
<b>Observed data</b>	Summary statistics of real experiment	Labeled training set $D_{\text{train}}$
<b>Amortization</b>	Over observations $x_{\text{obs}}$	Over all datasets in prior

# Inside the Prior: SCMs + Bayesian Neural Networks

Two complementary mechanisms generate the diversity of synthetic datasets

## Structural Causal Model (SCM) Path

- Sample random DAG  $G \sim \text{Erdős-Rényi}(N, P)$
- Each edge: random nonlinear function  $f_{i,j}$   
(MLP, linear, multiplicative, step, ...)
- Root nodes seeded with noise  $\epsilon \sim \text{Uniform/Gaussian}/\dots$
- Forward-propagate all  $N$  samples through  $G$
- Designate random subset of nodes as features  $X$
- Remaining node(s) become target  $Y$
- Implicit causal structure:  $X \perp\!\!\!\perp Y \mid \text{pa}(Y)$

## Bayesian Neural Network (BNN) Path

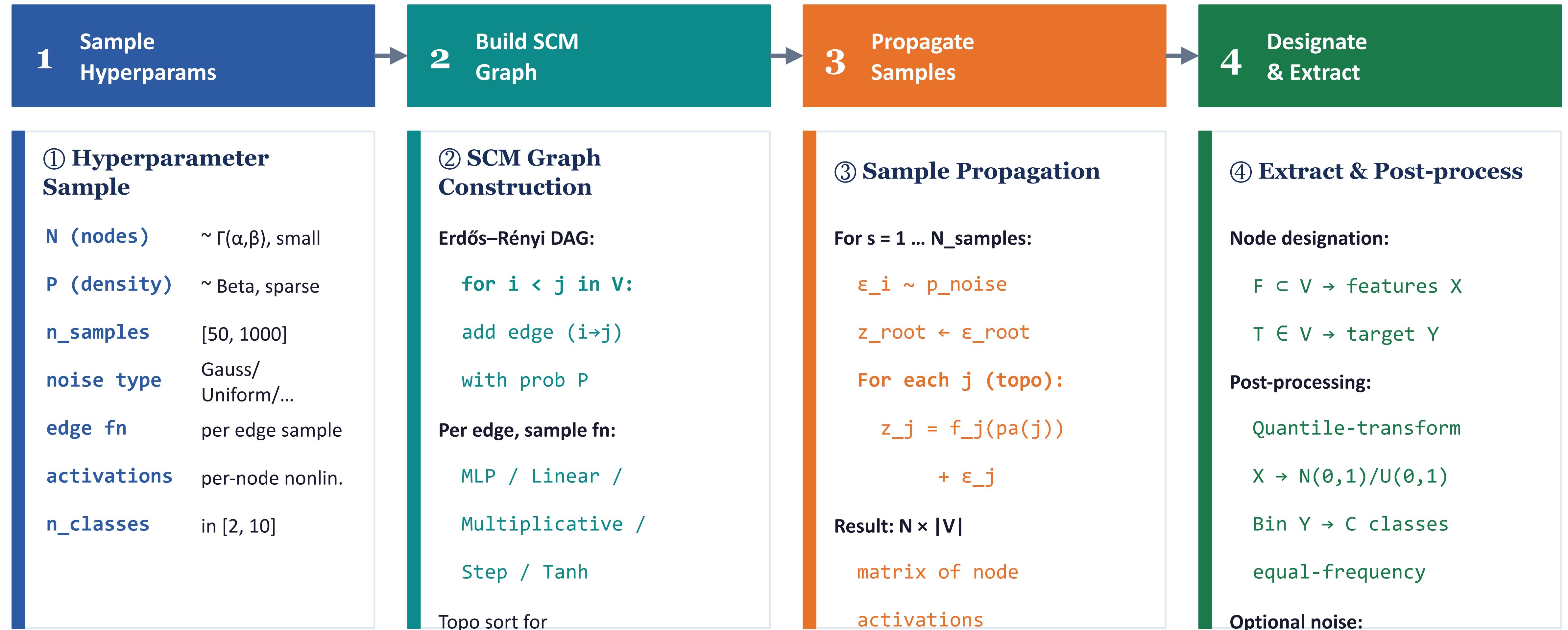
- Sample a random MLP architecture  
(depth, width, activation drawn from priors)
- Sample weights  $W \sim N(0, \sigma^2)$  per layer
- Generate inputs  $X \sim \text{some distribution}$
- Compute  $f_W(X) = \text{regression output}$
- Complementary coverage: smooth function classes  
that may not arise naturally from SCMs
- Adds regularity and interpolation bias to prior

+

**Hyperprior key detail:** Almost all prior hyperparameters (graph size, edge density, activation, noise scale...) are themselves drawn from distributions rather than fixed — this creates enormous diversity and prevents overfitting to any single data-generating mechanism.

# Forward Simulation in Detail: How One Synthetic Dataset is Generated

Four precise steps from raw hyperparameters to a labelled  $(X, y)$  dataset ready for transformer training



**Training batch:** 512 synthetic datasets  $\times$  18 000 batches  $\approx$  9 M unique prior samples seen by the transformer during offline training

# The Edge Function Zoo: Diversity in the Simulator

Each directed edge in the SCM is independently assigned one of several random nonlinear mapping types

## MLP (Neural Network)

$$f(x) = W_L \sigma(\dots \sigma(W_1 x + b_1)\dots) + b_L$$

*depth*  $\sim U\{1,3\}$ ; *width*  $\sim U\{4,64\}$ ;  $W \sim N(0,1)$ ;  $act \in \{ReLU, Tanh, GELU\}$

→ Captures arbitrary smooth nonlinear relationships between parent and child nodes.

## Linear / Affine

$$f(x) = a \cdot x + b, \quad a, b \sim N(0, 1)$$

*Slope and intercept sampled independently per edge; ensures linear subspace covered.*

→ Guarantees linear functions are well-represented — prevents prior from being purely nonlinear.

## Multiplicative (Interaction)

$$f(x_1, x_2) = x_1 \cdot x_2 \quad (\text{multi-parent nodes})$$

*Applied when a node has  $\geq 2$  parents; creates non-additive feature interaction terms.*

→ Generates interaction effects ( $x_1 \cdot x_2$  type) that are common in real-world tabular datasets.

## Step / Threshold Function

$$f(x) = 1[x > \tau], \quad \tau \sim \text{Uniform}(x_{\min}, x_{\max})$$

*Threshold  $\tau$  drawn from the data range; creates piecewise-constant decision boundaries.*

→ Mimics decision-rule-like structure naturally arising in categorical / rule-based processes.

## Smooth Nonlinear (Tanh / Sigmoid)

$$f(x) = \tanh(a \cdot x + b) \quad \text{or} \quad \sigma(a \cdot x + b)$$

*Gain  $a \sim U(0.5, 5)$ ; bias  $b \sim N(0, 1)$ ; saturating output keeps activations bounded.*

→ Saturating functions produce bounded outputs typical of physiological or economic signals.

## Noise (Root Nodes — Exogenous)

$$z_{\text{root}} = \varepsilon, \quad \varepsilon \sim p_{\text{noise}}(\cdot)$$

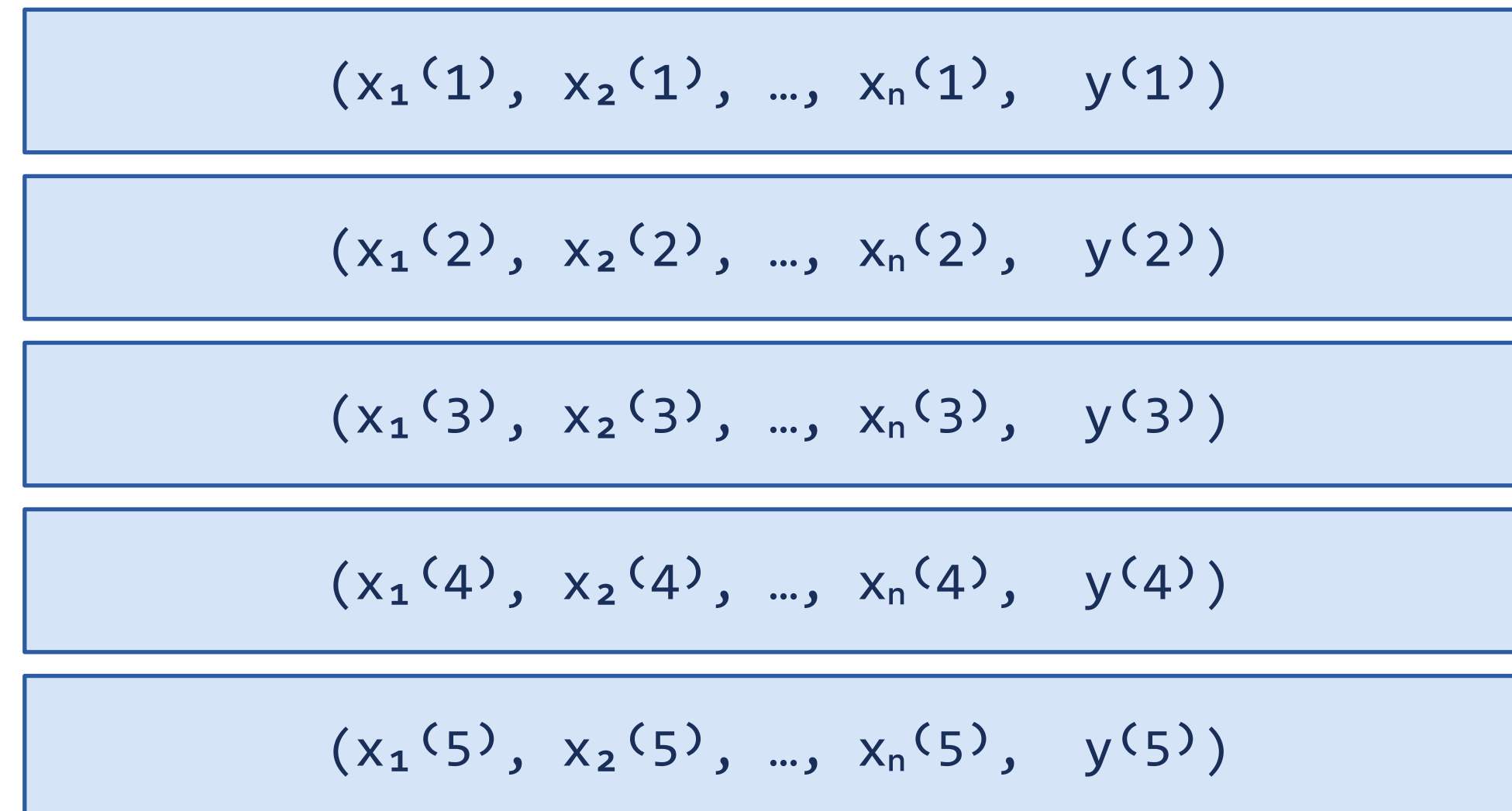
*$p_{\text{noise}} \in \{N(0,1), \text{Uniform}[0,1], \text{Beta}(\alpha,\beta), \text{Laplace}(0,b), \dots\}$  sampled per node.*

→ Varied marginal distributions give diverse feature marginals across synthetic datasets.

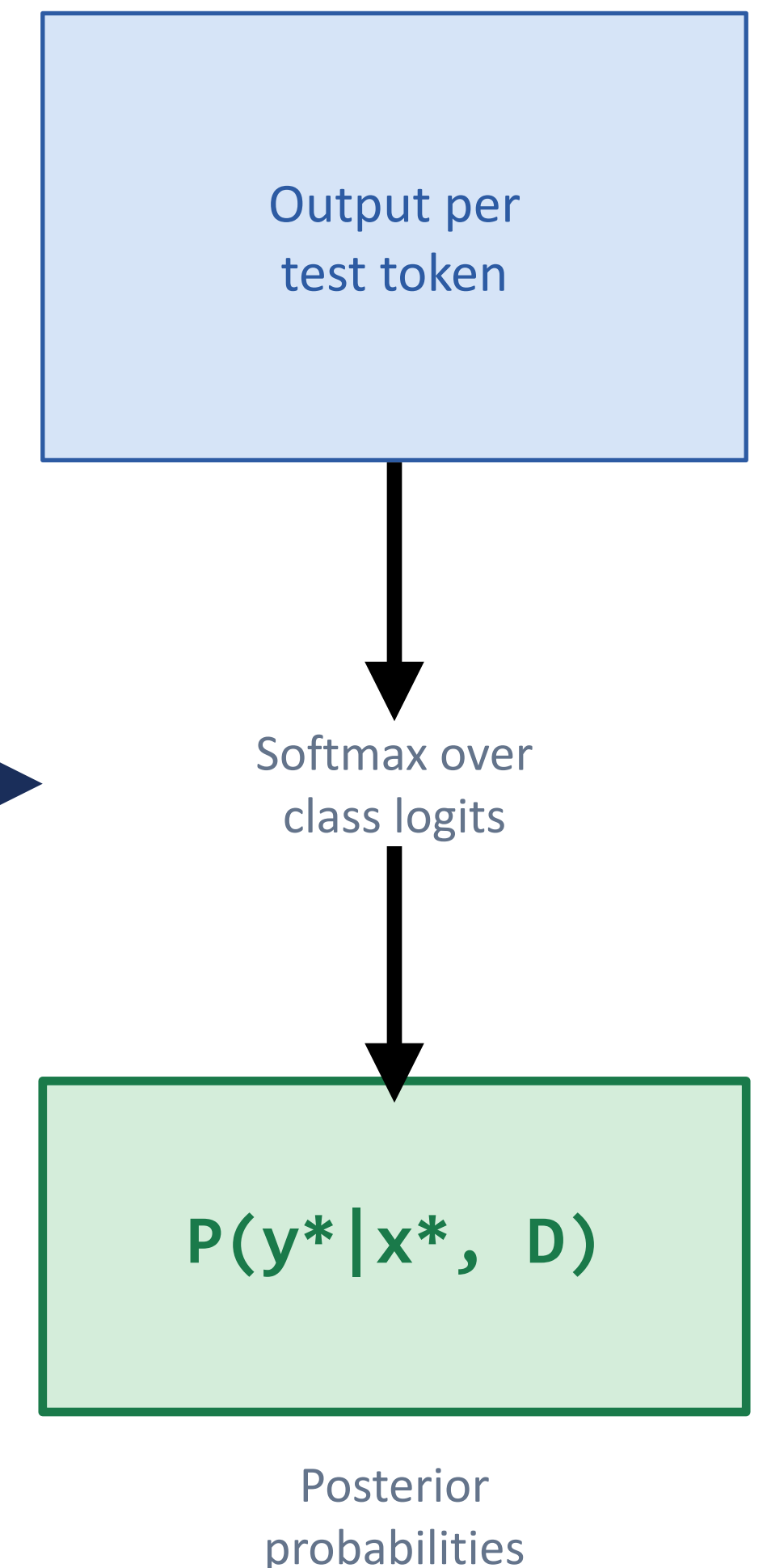
# Architecture: In-Context Learning via Transformer

The network processes training + test samples jointly in one forward pass

## Training set $D_{\text{train}}$



## Test set $X_{\text{test}}$



Permutation invariant over training samples

All test predictions in one pass

No positional encoding (set not sequence)

Row & feature attention combined

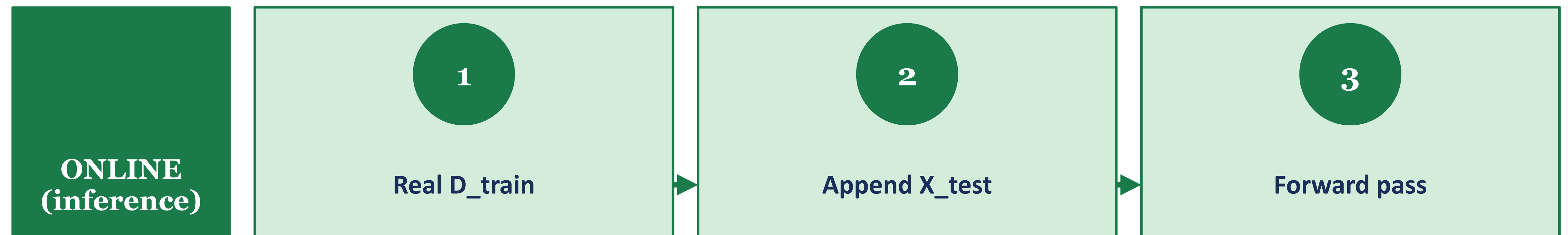
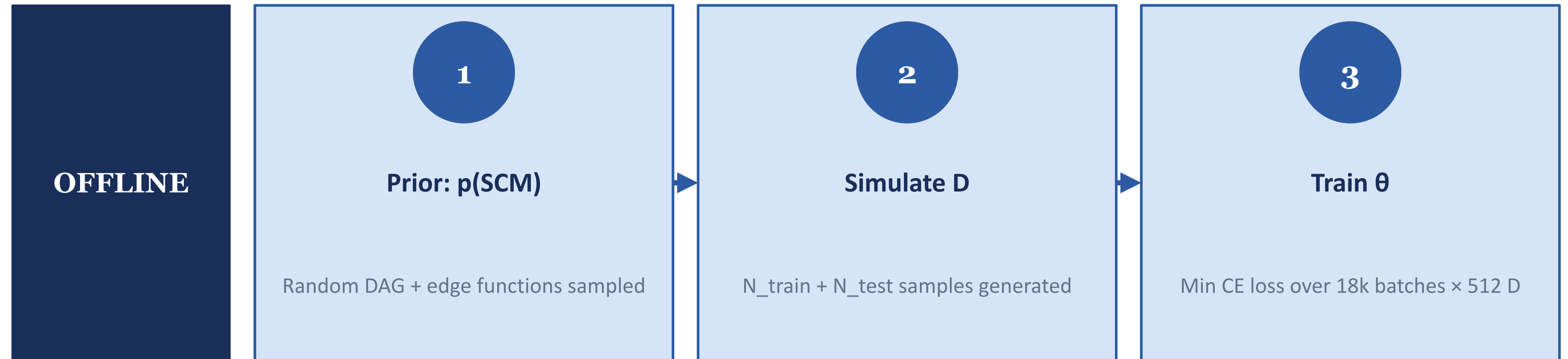
# Full Pipeline: From Causal Prior to Prediction

End-to-end view connecting the graphical model, simulation, training and inference



# Full Pipeline: From Causal Prior to Prediction

End-to-end view connecting the graphical model, simulation, training and inference



The transformer weights  $\theta$  are the amortized posterior — they implicitly encode  $p(\phi | D)$  averaged over the entire SCM prior

# Inductive Bias: Simplicity & Causality in the Prior

*Two principled choices that help TabPFN generalise to real-world tabular data*

## Simplicity Bias (Occam's Razor)

- Prior places higher weight on simpler SCM structures: smaller DAGs, fewer hidden nodes, sparser edges.
- Motivated by Occam's razor and MDL principle: simpler explanations are preferred in the absence of evidence.
- Implemented by skewing hyperprior distributions toward smaller N (nodes) and sparser P (edges).
- Hypothesis: real-world tabular datasets arise from relatively simple underlying mechanisms.
- This prevents the posterior from collapsing to overly complex, overfit explanations.

## Causal (SCM) Inductive Bias

- Using SCMs as data generators builds in an assumption that features and labels have a joint causal history.
- Pearl's Ladder of Causation: TabPFN operates at 'rung 1.5' — association-based but causally-structured.
- TabPFN does NOT perform causal discovery. It bypasses explicit graph recovery and predicts Y directly.
- Benefits: handles spurious correlations, distribution shifts from observational data more gracefully.
- Empirical check: TabPFN's predictions correlate with simple SCM-generated hypotheses (cf. Fig. 8 in paper).

*Both biases are baked into the prior distribution — not into the transformer architecture itself, making them transparent and modifiable.*

# Results: TabPFN vs. GBDT and AutoML

OpenML-CC18 benchmark — 18 datasets,  $\leq 1000$  training samples,  $\leq 100$  numerical features

**230x**

speedup vs AutoML  
(CPU)

**5700x**

speedup vs AutoML  
(GPU)

**< 1s**

total inference  
time

**#1**

on 18 small  
datasets

Method	Relative Accuracy	Train + Predict Time	Notes
<b>TabPFN</b>	<b>Ref (SOTA)</b>	<b>&lt; 1 s</b>	<b>No HPO needed</b>
XGBoost (tuned)	~-2% ROC	minutes	Hyperparameter search required
AutoML (AutoSklearn)	comparable	hours	Full pipeline search
Random Forest (def)	< TabPFN	seconds	No tuning but lower accuracy
KNN	< TabPFN	ms	Simple baseline

TabPFN also validates on 67 additional OpenML datasets — consistent advantage at small-data regime.