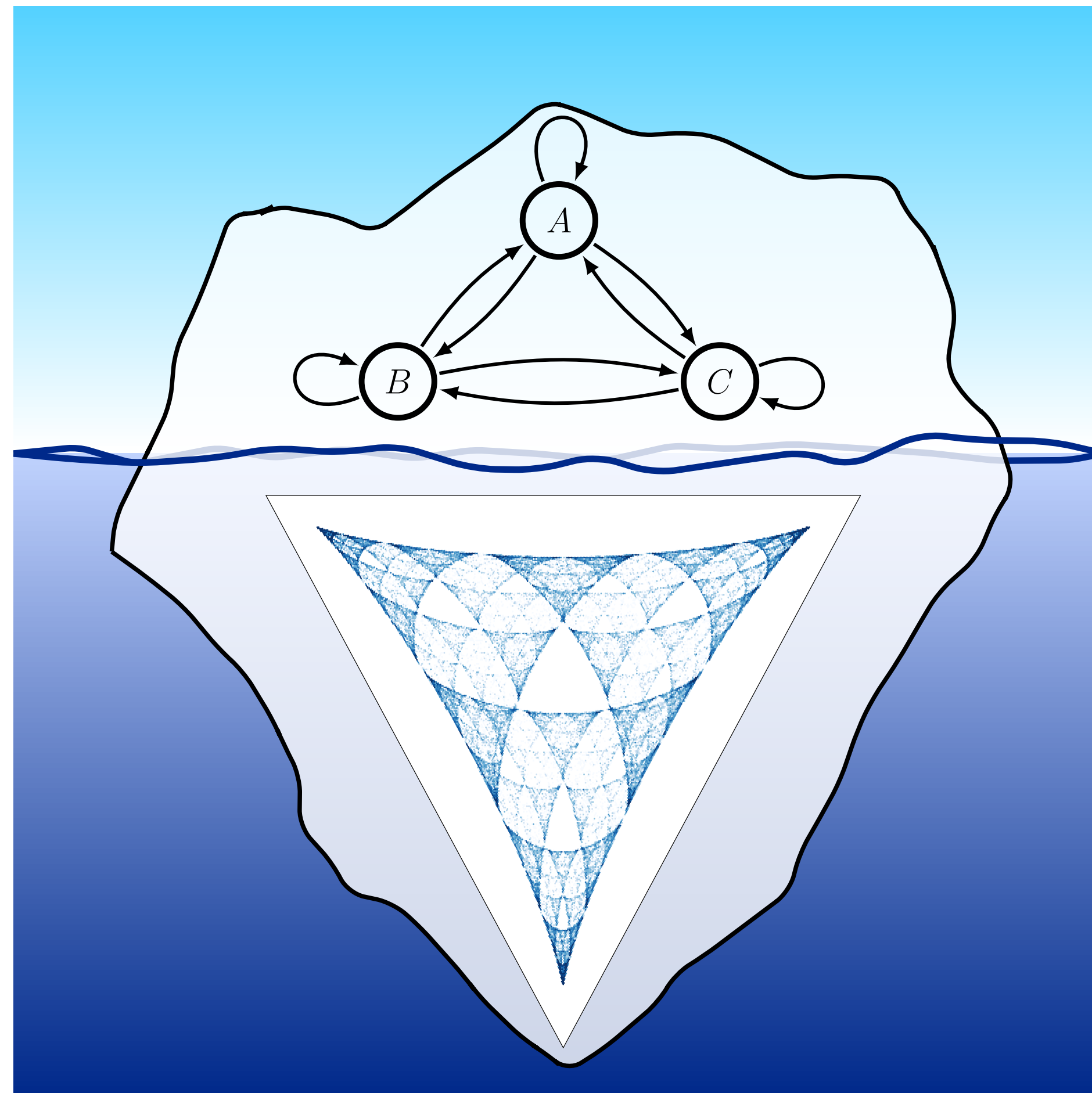


Finite and Infinite Models: Optimal Prediction of Hidden Markov Processes



Alexandra M. Jurgens

INRIA Sud Ouest - Bordeaux

Information Theory as a Bridge Across the
Geosciences and Modeling Sciences

11/09/2023



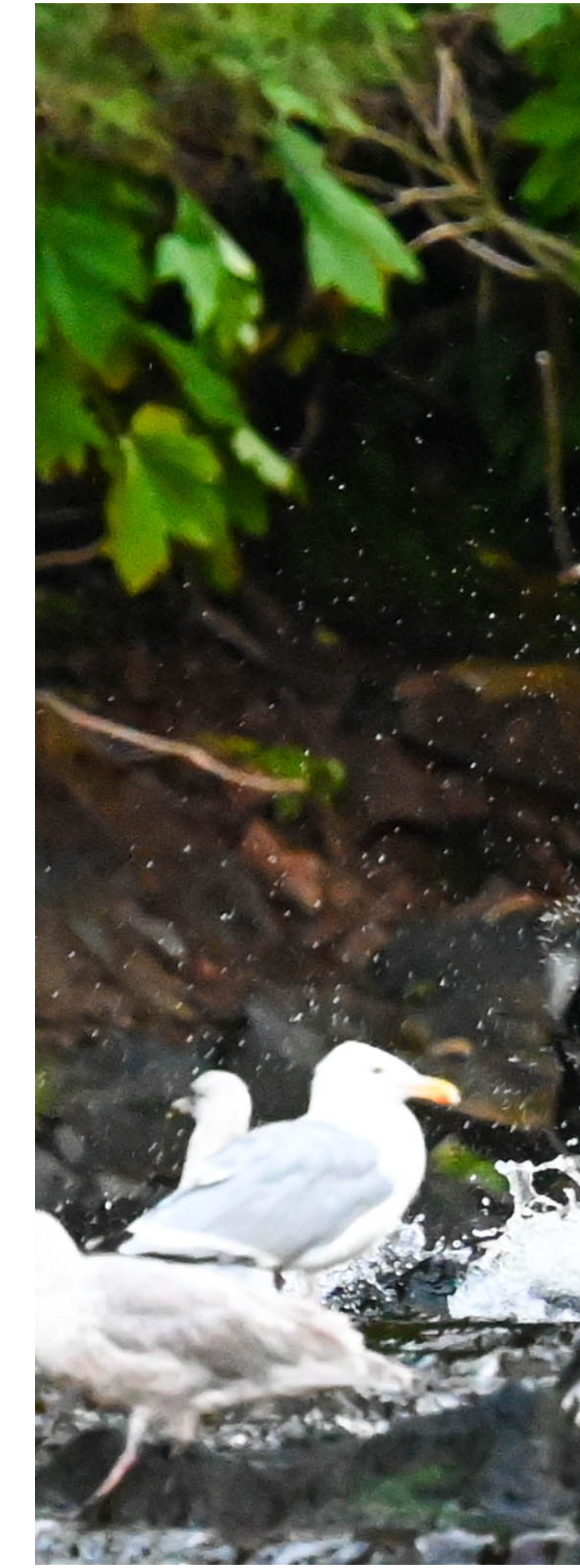
Inria

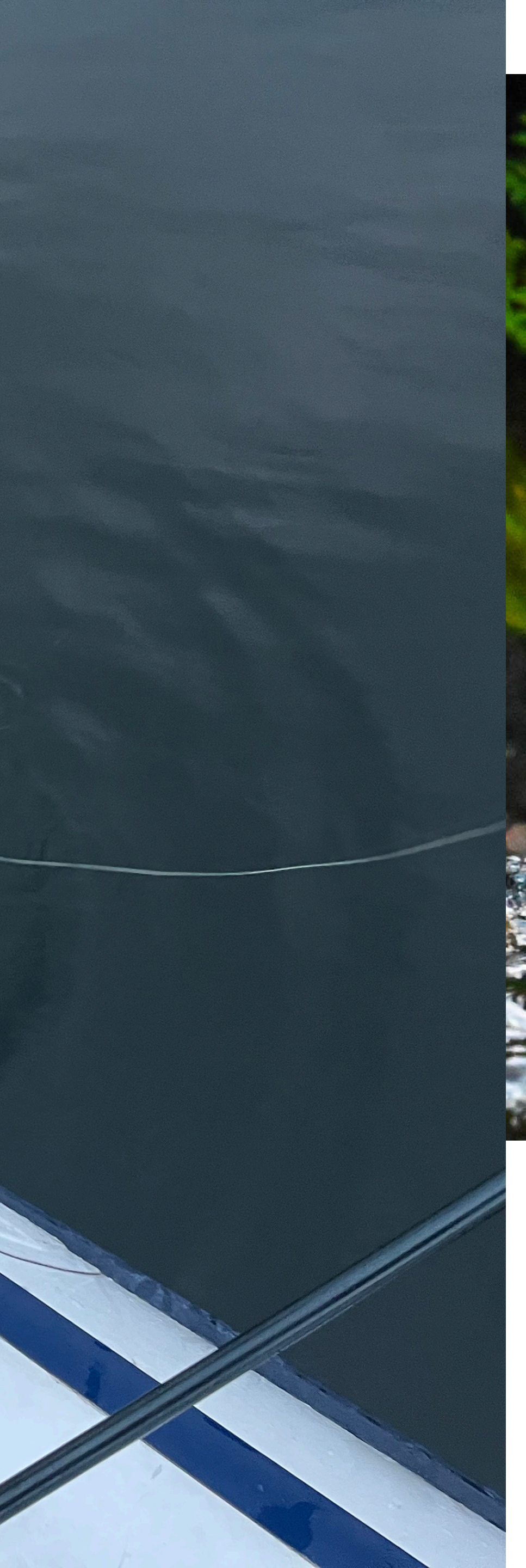












review article

Simple mathematical models with very complicated dynamics

Robert M. May*

First-order difference equations arise in many contexts in the biological, economic and social sciences. Such equations, even though simple and deterministic, can exhibit a surprising array of dynamical behaviour, from stable points, to a bifurcating hierarchy of stable cycles, to apparently random fluctuations. There are consequently many fascinating problems, some concerned with delicate mathematical aspects of the fine structure of the trajectories, and some concerned with the practical implications and applications. This is an interpretive review of them.

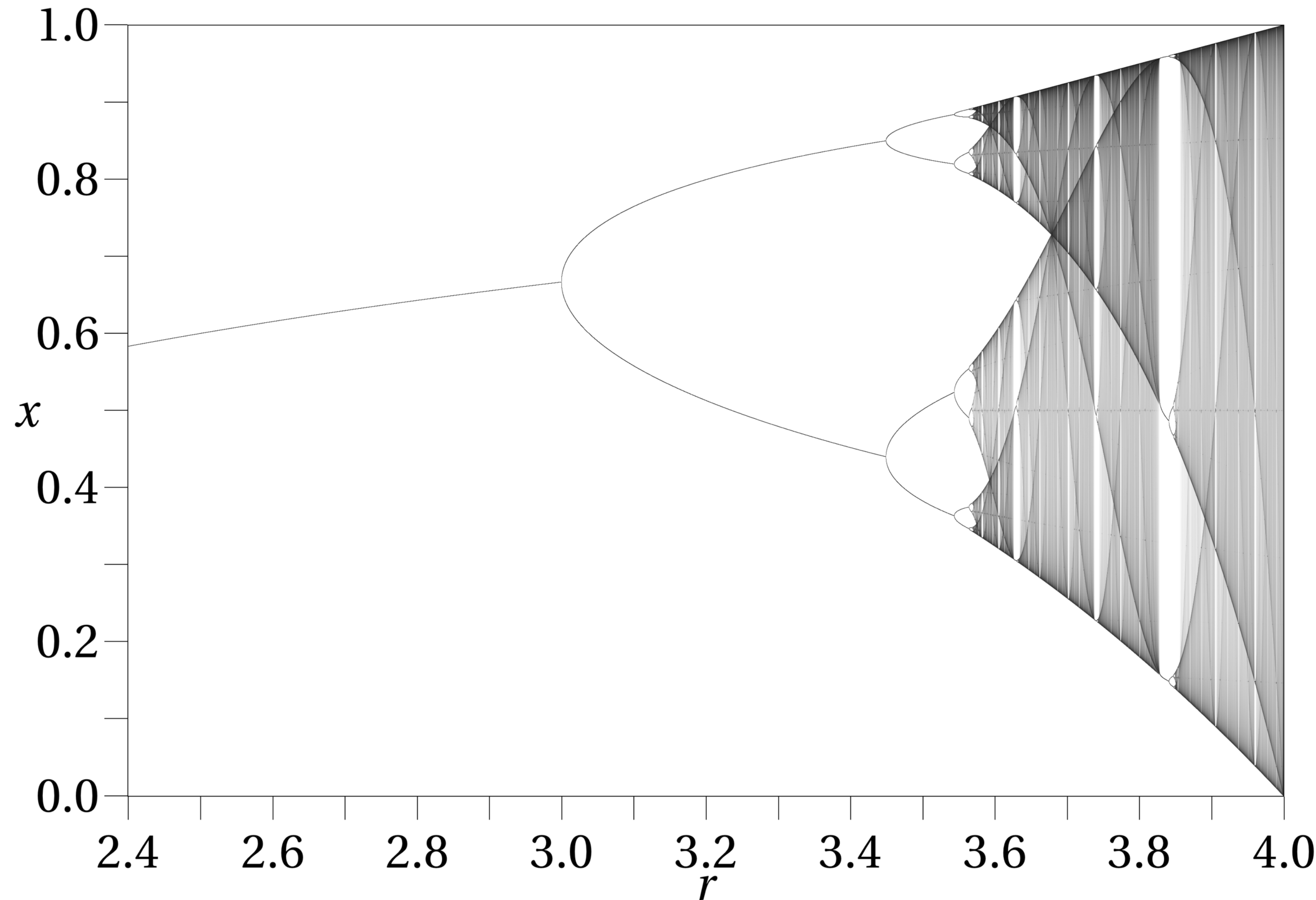
review article

Simple mathematical models with very complicated dynamics

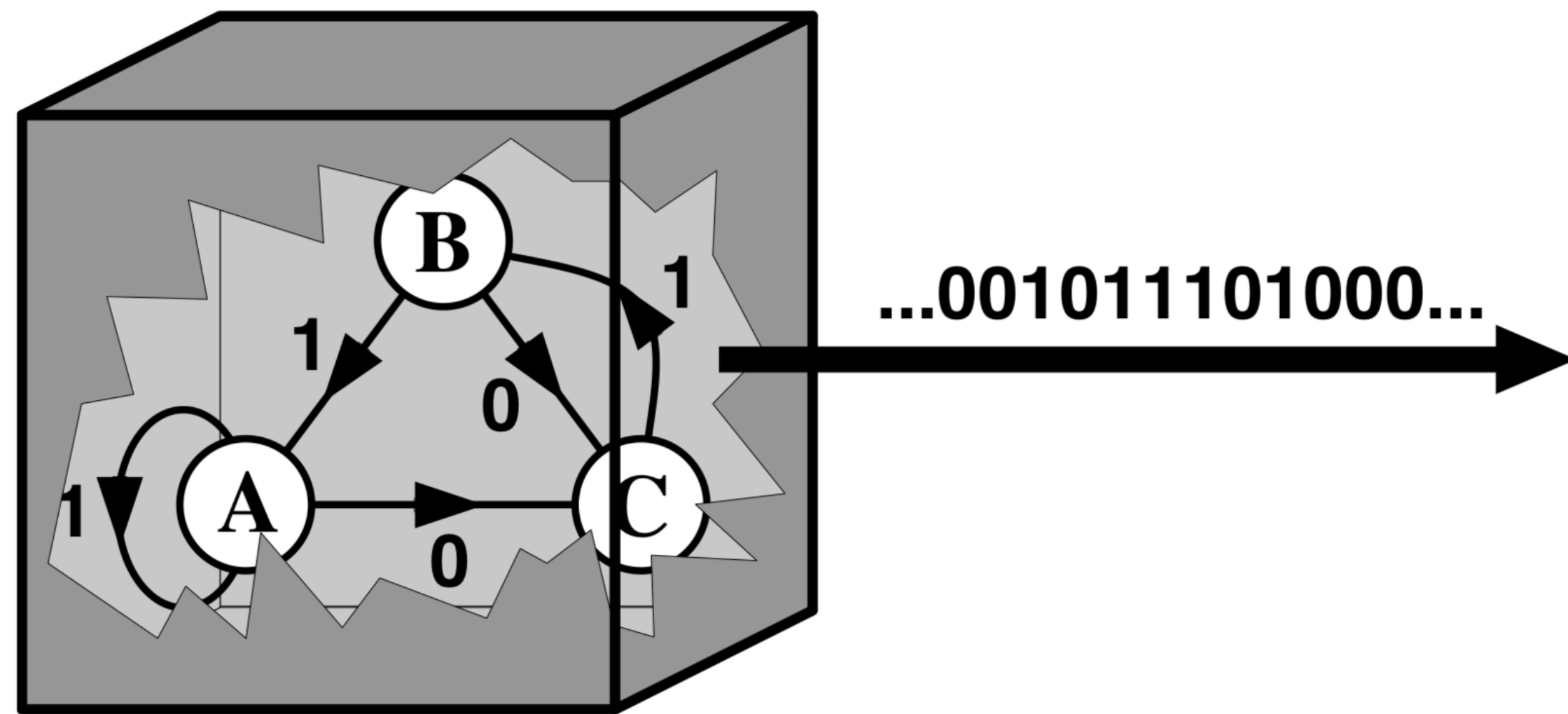
Robert M. May*

First-order difference equations arise in many contexts in the biological, economic and social sciences. Such equations, even though simple and deterministic, can exhibit a surprising array of dynamical behaviour, from stable points, to a bifurcating hierarchy of stable cycles, to apparently random fluctuations. There are consequently many fascinating problems, some concerned with delicate mathematical aspects of the fine structure of the trajectories, and some concerned with the practical implications and applications. This is an interpretive review of them.

$$x_{n+1} = rx_n (1 - x_n)$$



Object of Study: Processes



System

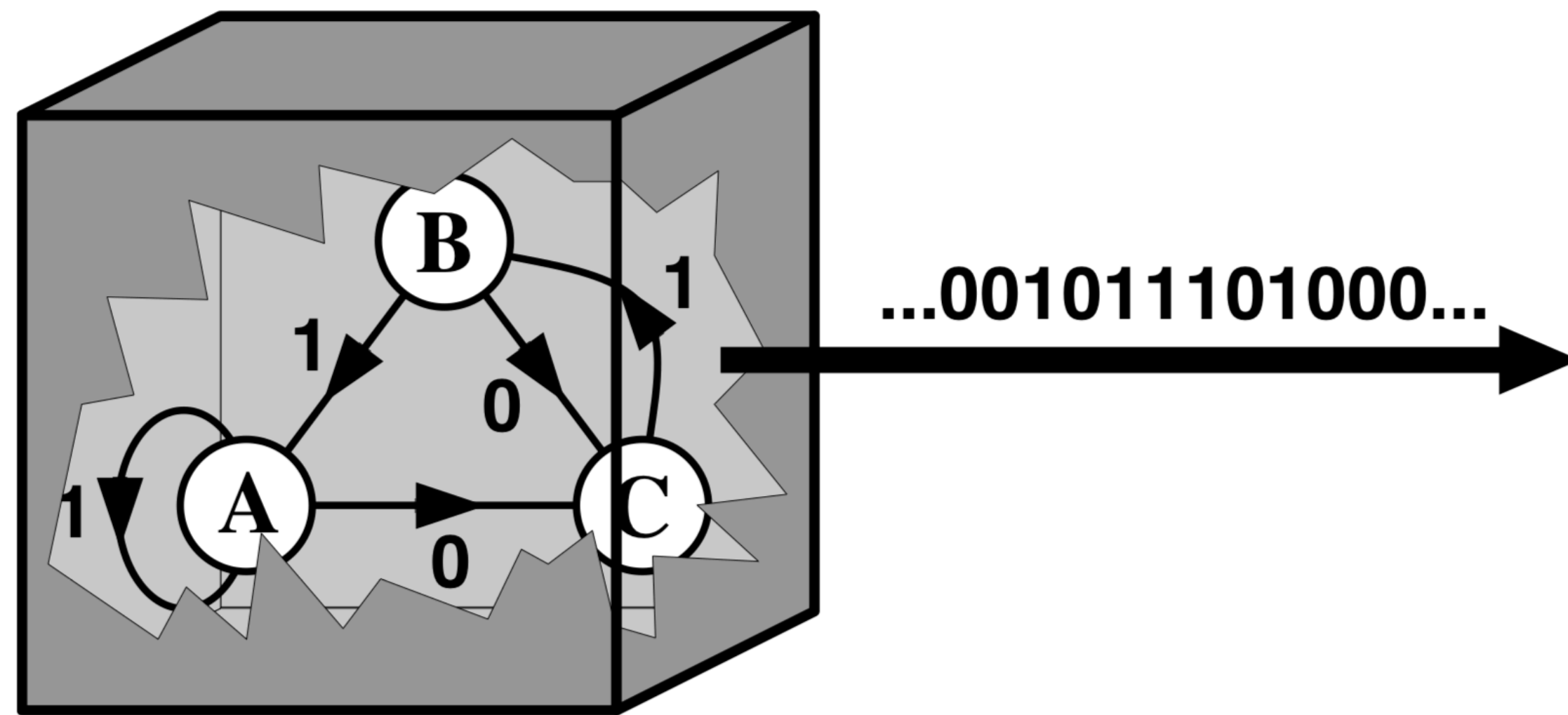
Instrument

Process

The random variables X_i may take on values in alphabet A :

$$\vec{X} = \dots X_{-1} X_0 X_1 \dots$$

Object of Study: Processes



System

Instrument

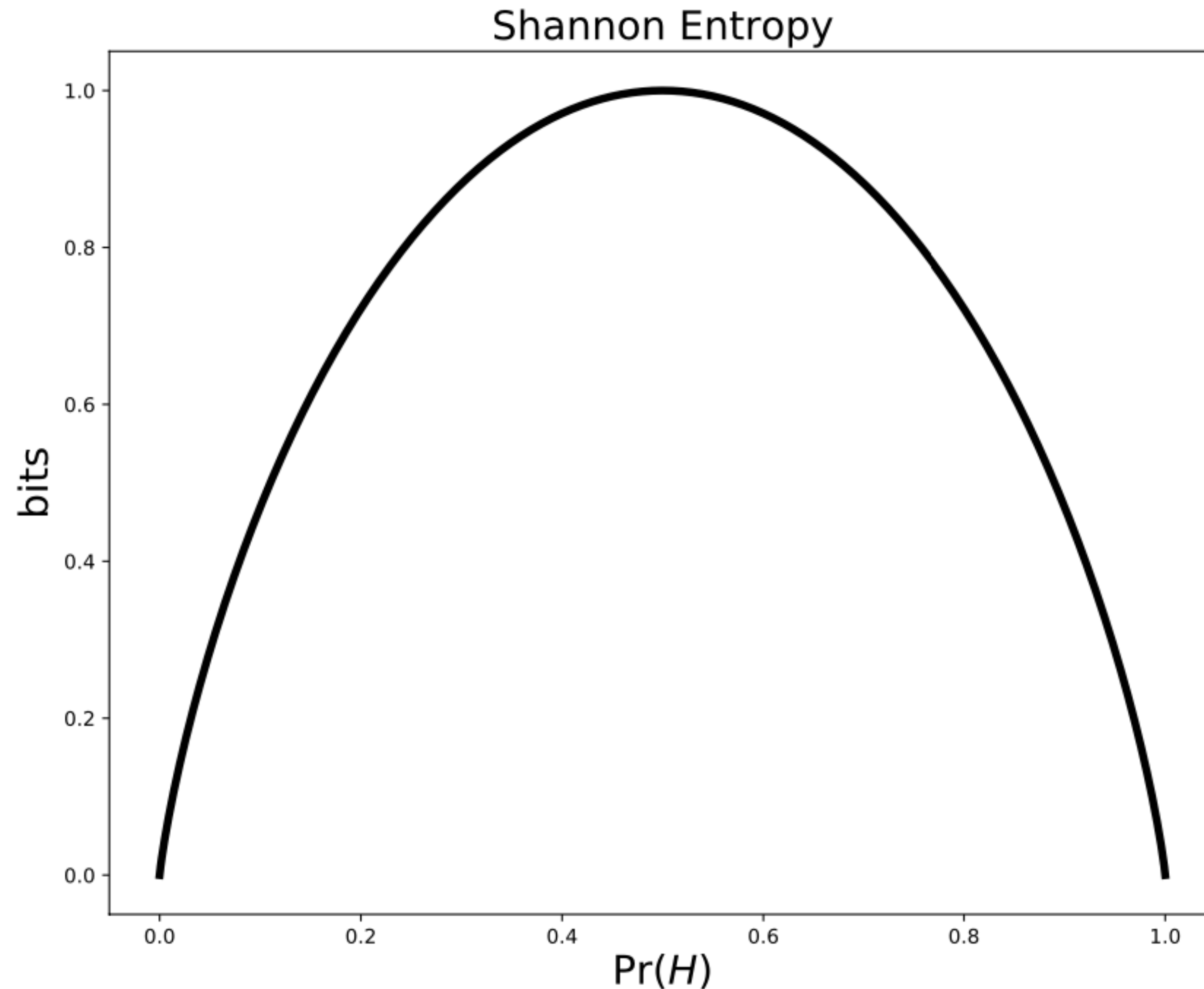
Process

Let $A = \{0, 1\}$. Then a *realization* of the process is written

$$\begin{aligned} \overleftrightarrow{x} &= \dots x_{-1} x_0 x_1 \dots \\ &= \dots 011 \dots \end{aligned}$$

A process P is defined as the probability distribution over bi-infinite strings.

Information Theory



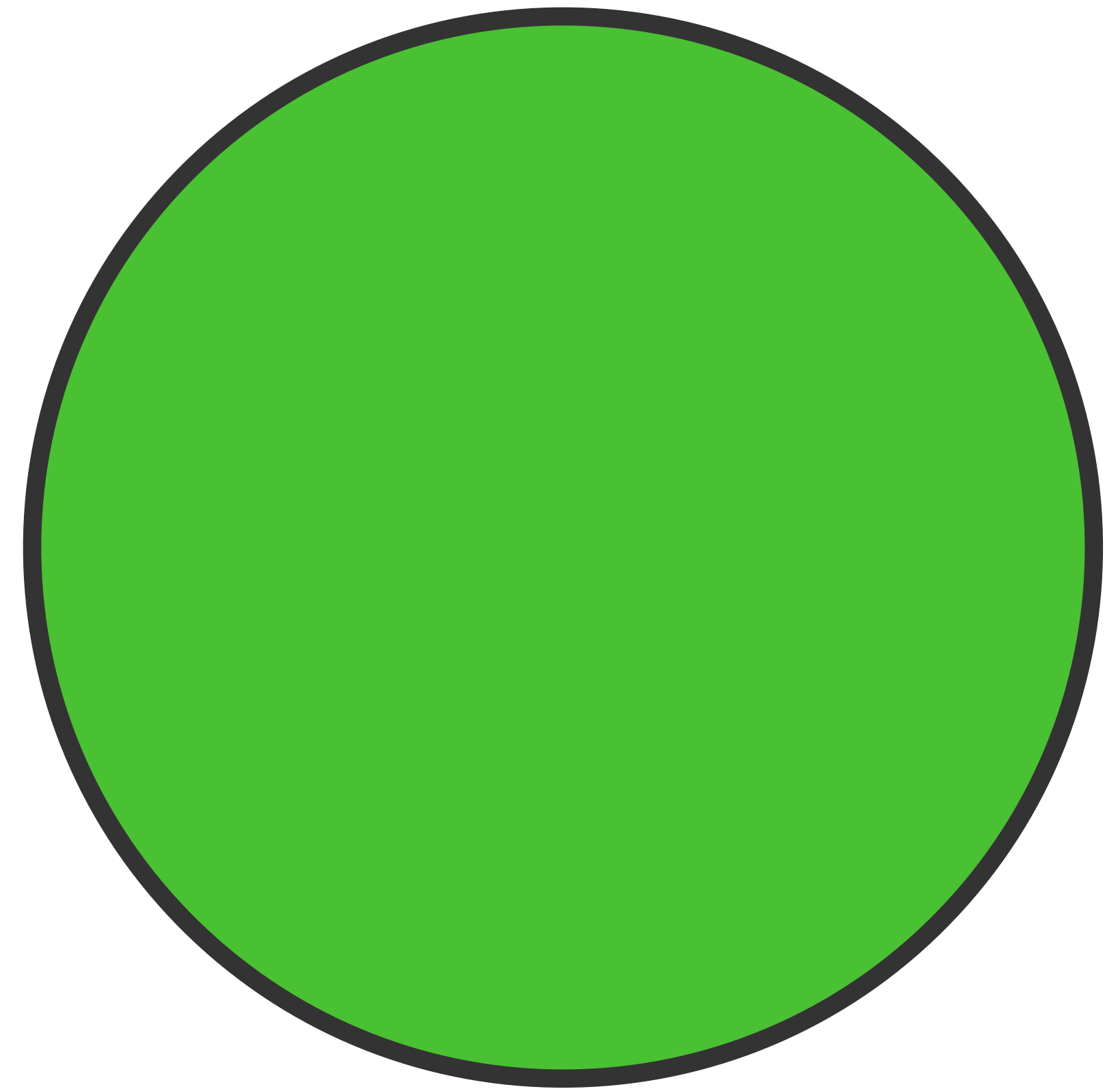
The *Shannon entropy* over a random variable is defined:

$$H[X] = - \sum_{x \in A} \Pr(X = x) \log_2 \Pr(X = x)$$

Information Theory

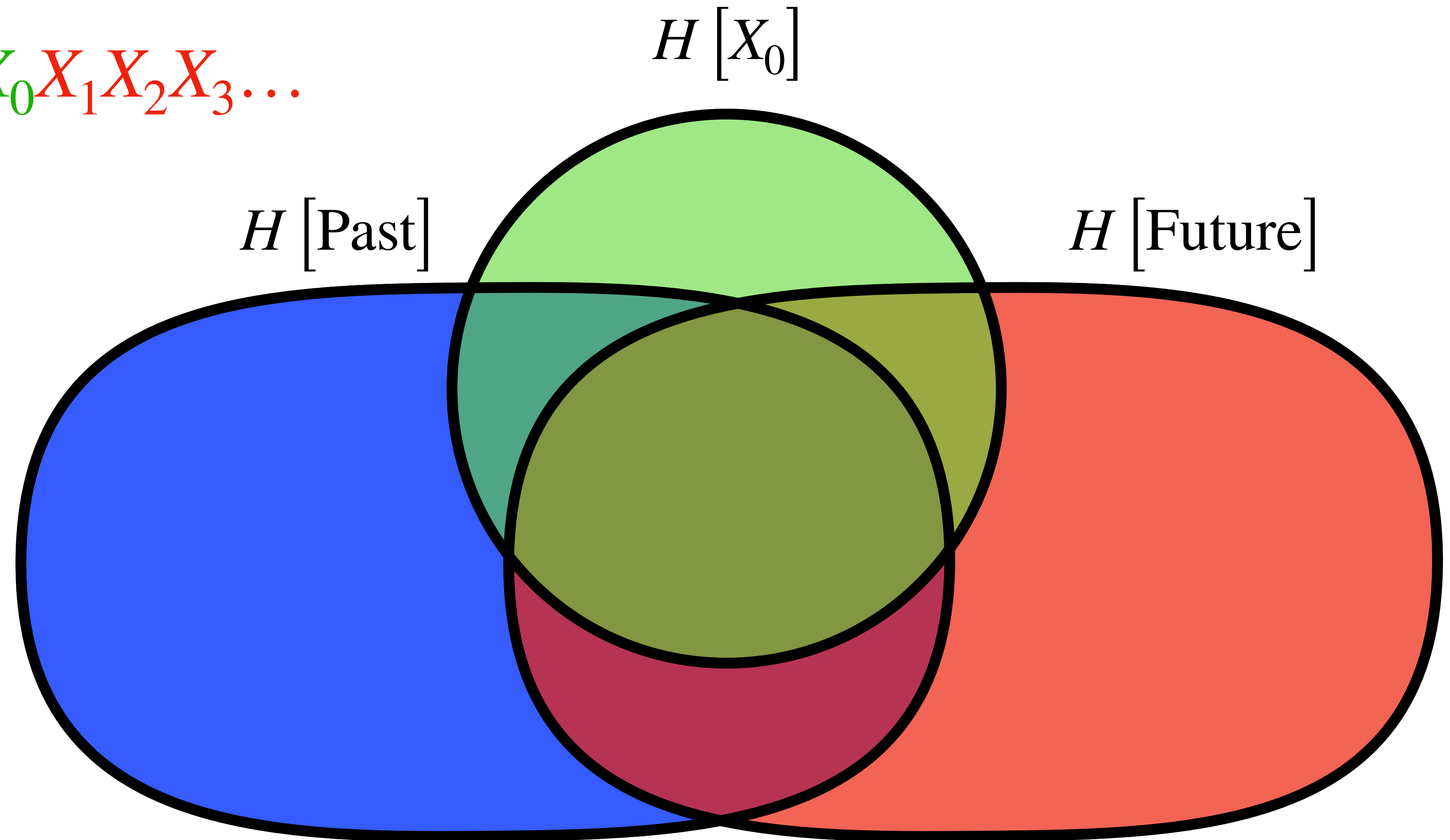
$H [X_0]$

$$H [X] = - \sum_{x \in A} \Pr (X = x) \log_2 \Pr (X = x)$$



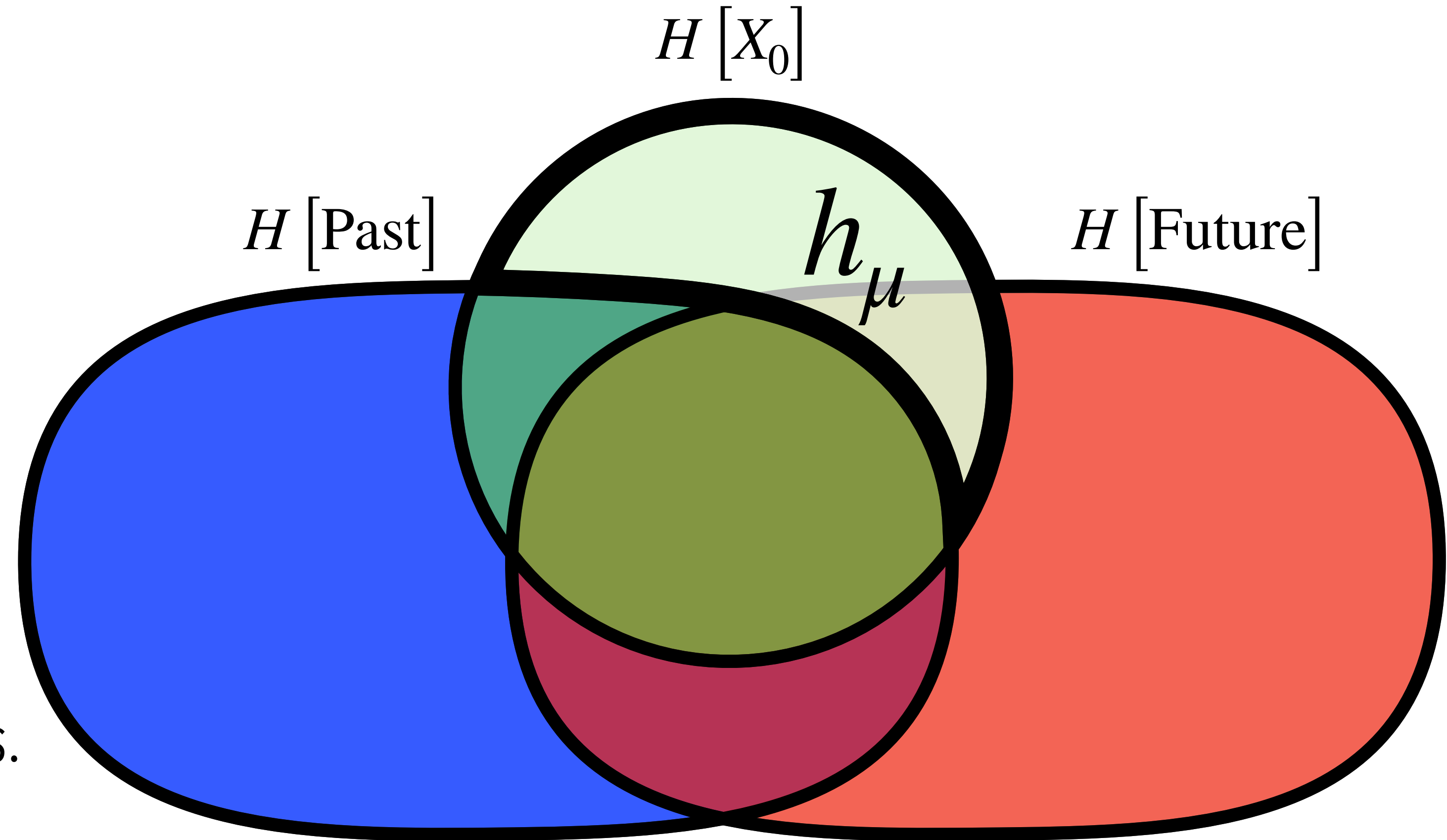
Information Theory

$$\vec{X} = \dots X_{-3} X_{-2} X_{-1} X_0 X_1 X_2 X_3 \dots$$



Entropy Rate

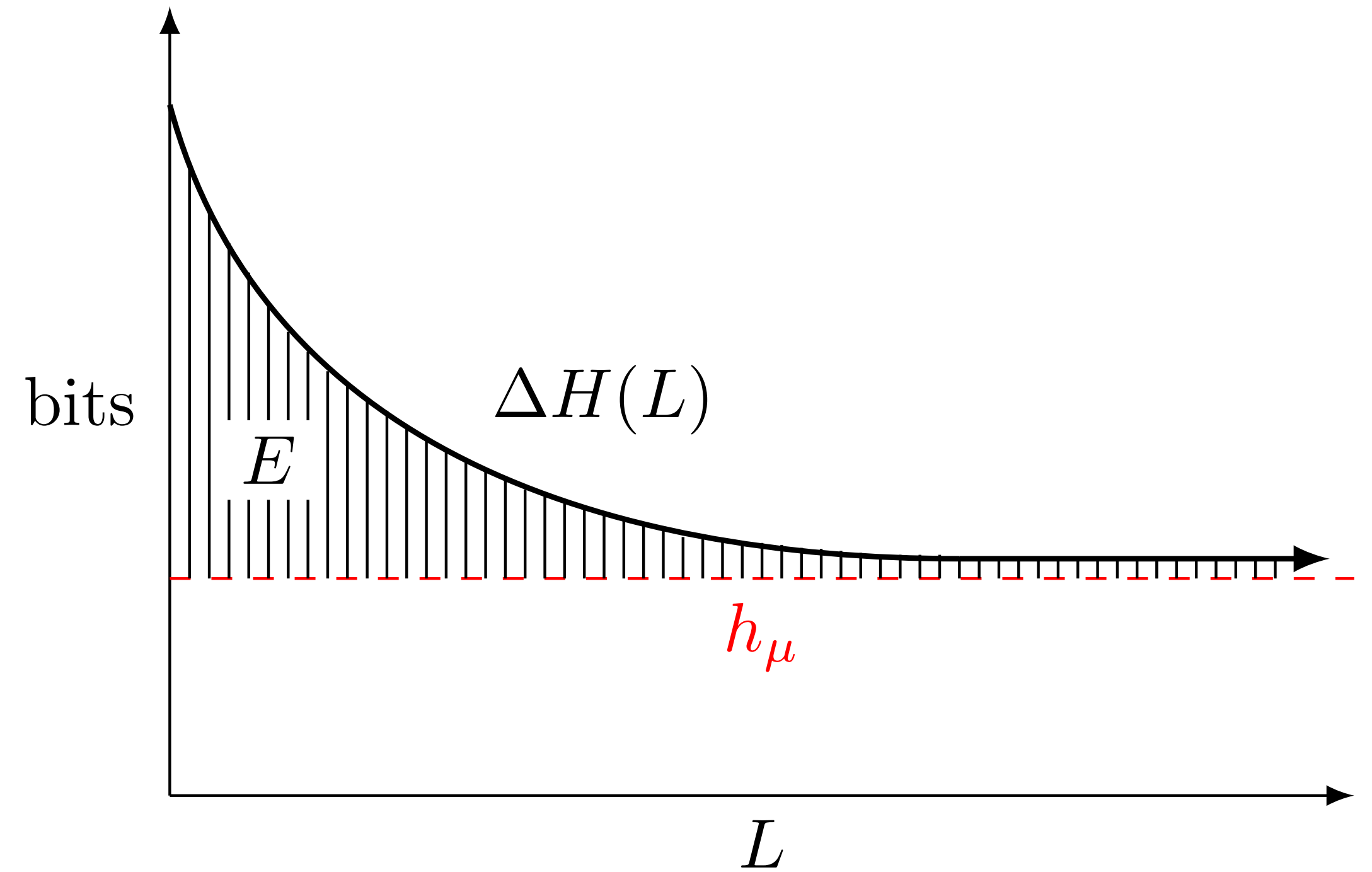
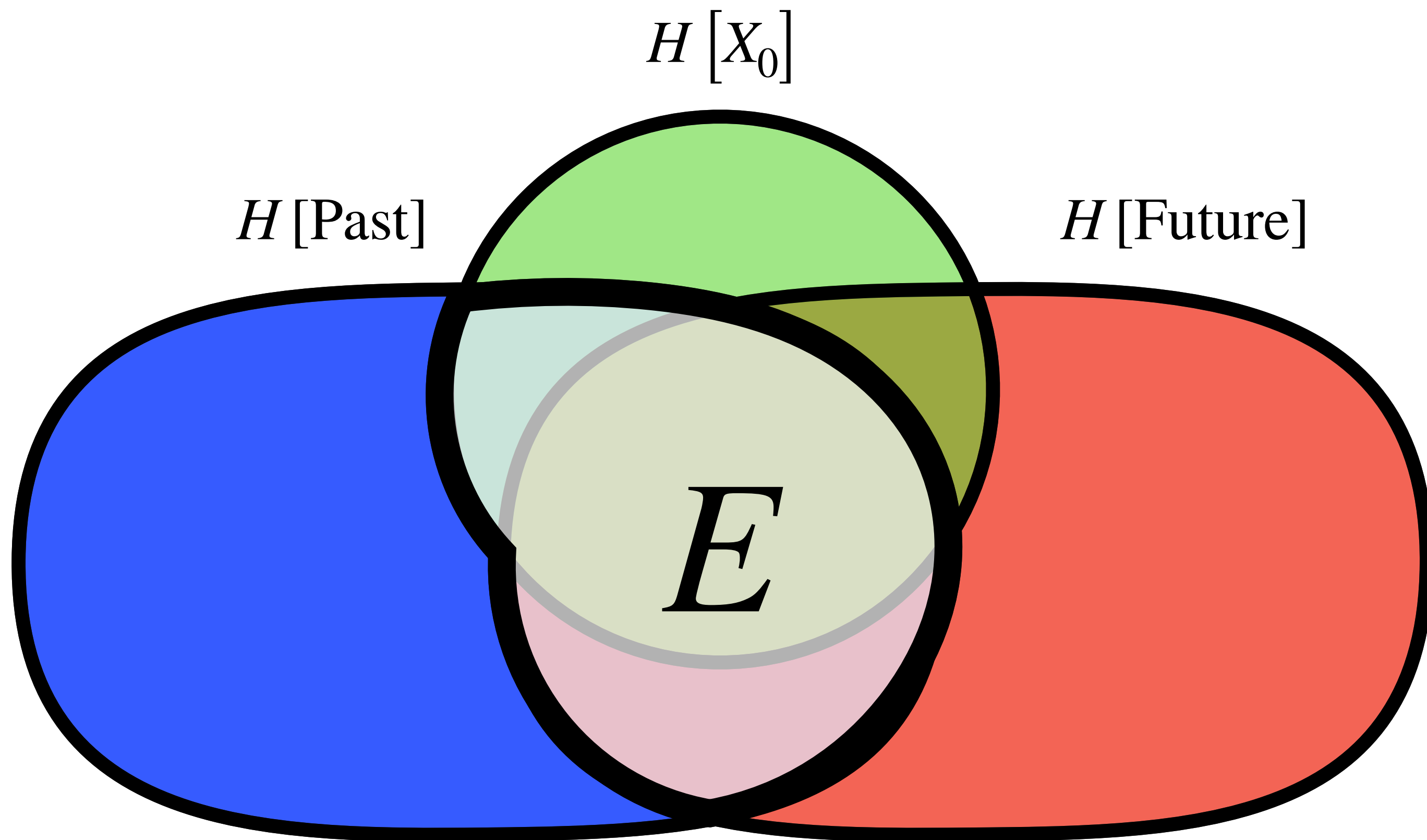
$$h_\mu = \lim_{L \rightarrow \infty} H[X_0 | \text{Past}]$$
$$= \lim_{L \rightarrow \infty} \left(H[X_0, \text{Past}] - H[\text{Past}] \right)$$



h_μ is the *irreducible randomness* of a process.

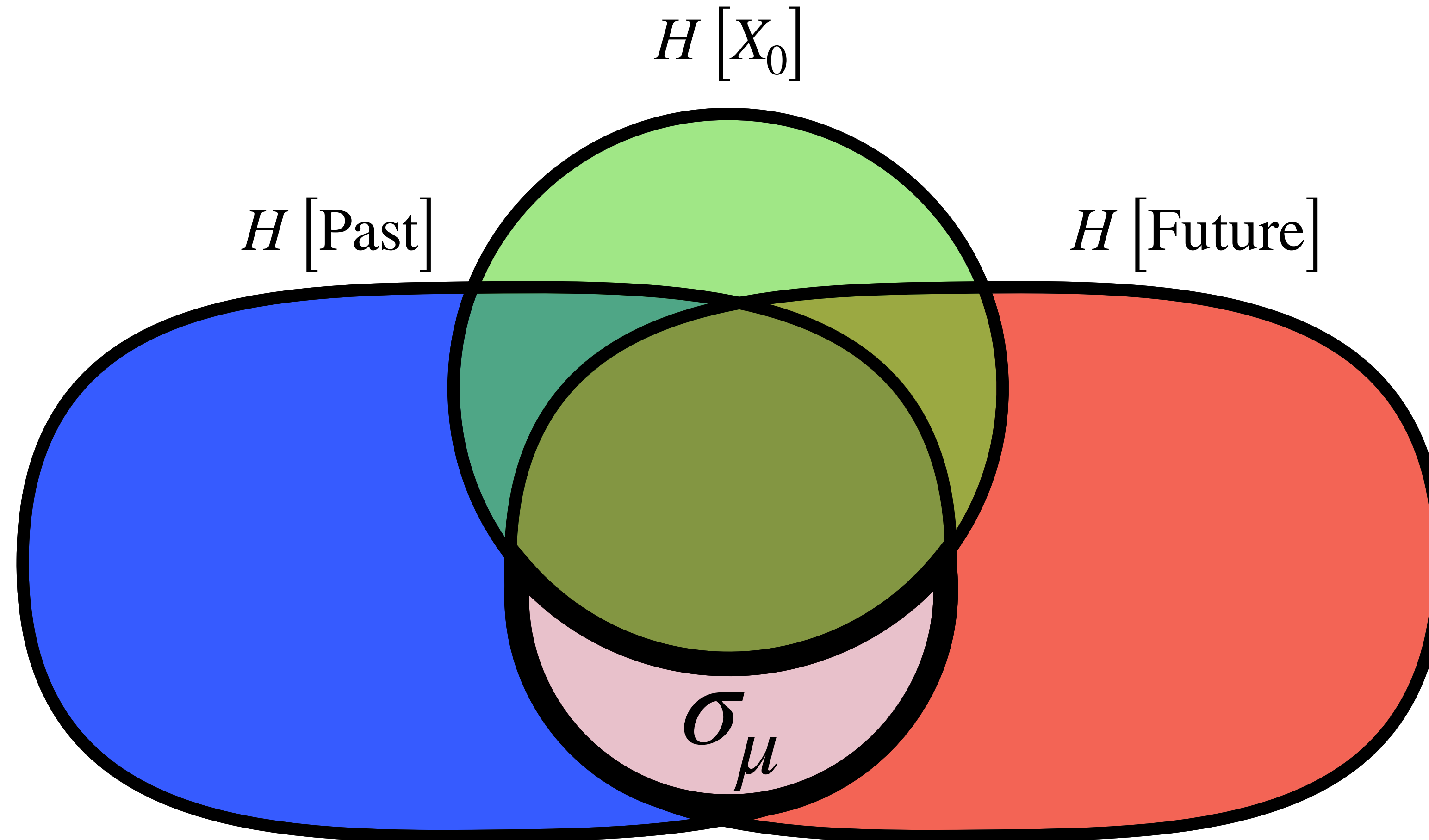
This is the limit of our predictive abilities.

Excess Entropy



Excess entropy: information shared between the past and the present + future.

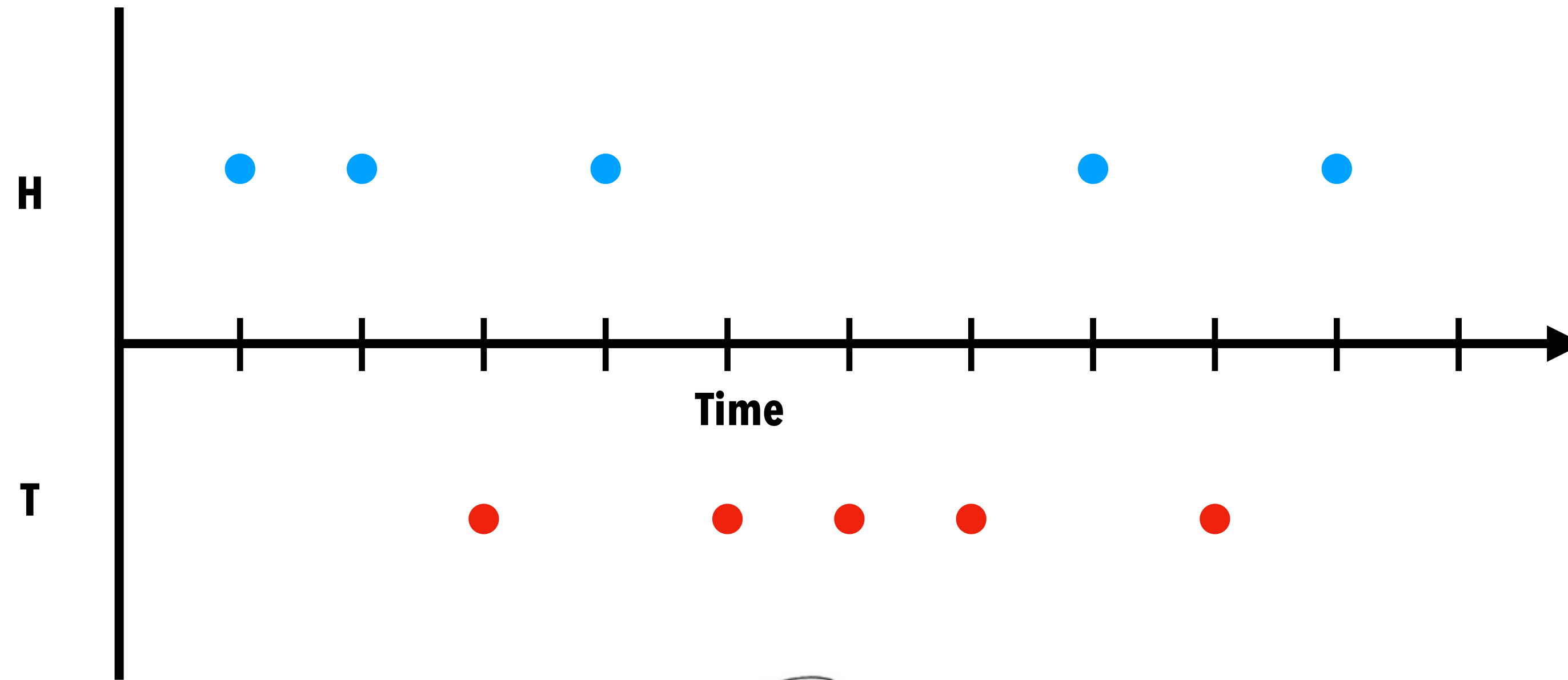
Elusive Information



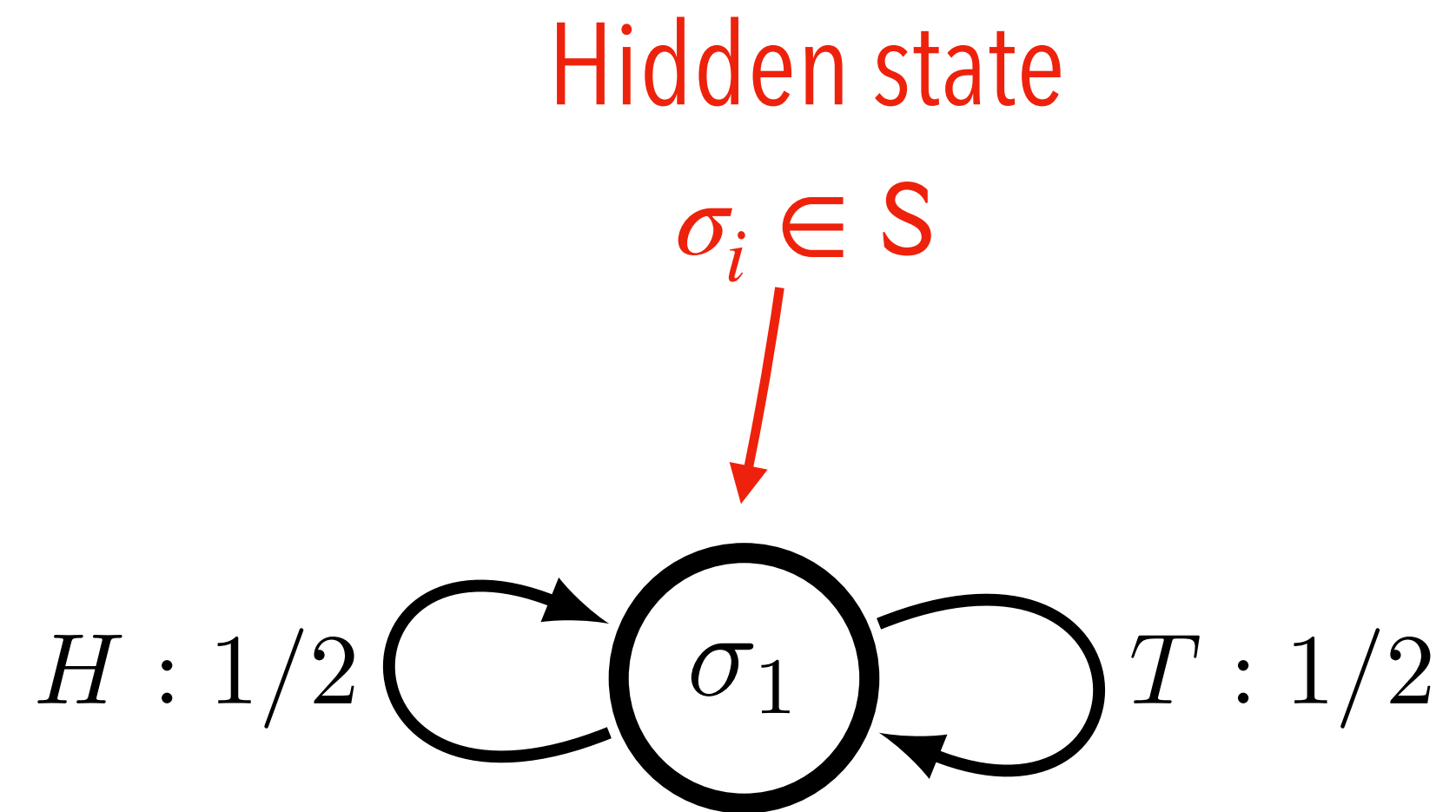
Ryan G. James, Christopher J. Ellison, James P. Crutchfield. *Anatomy of a bit: Information in a time series observation*. Chaos 21, 037109 (2011)

James P. Crutchfield, David P. Feldman. *Regularities unseen, randomness observed: Levels of entropy convergence*. Chaos 1 March 2003; 13 (1): 25-54.

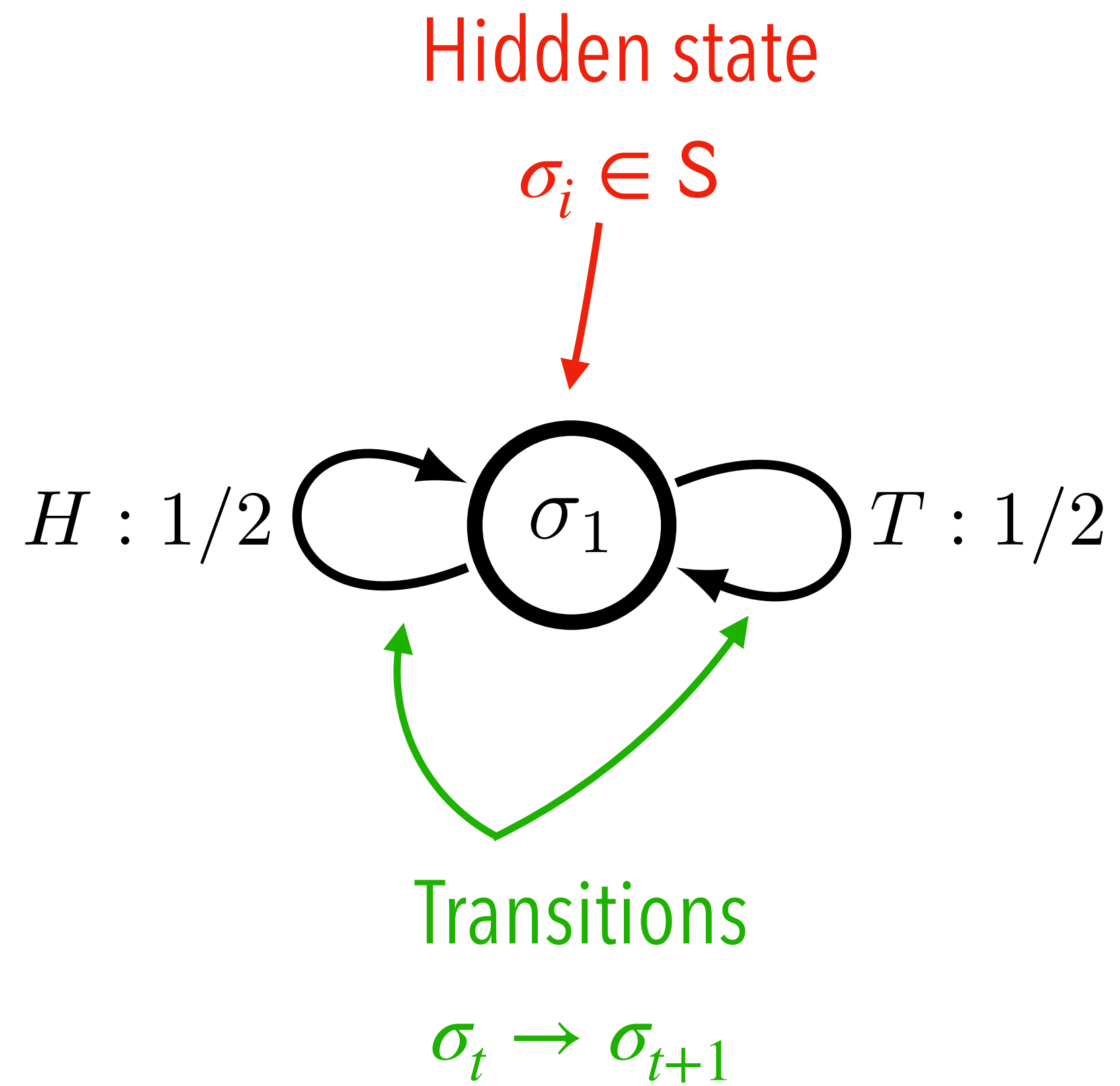
Models of Time Series



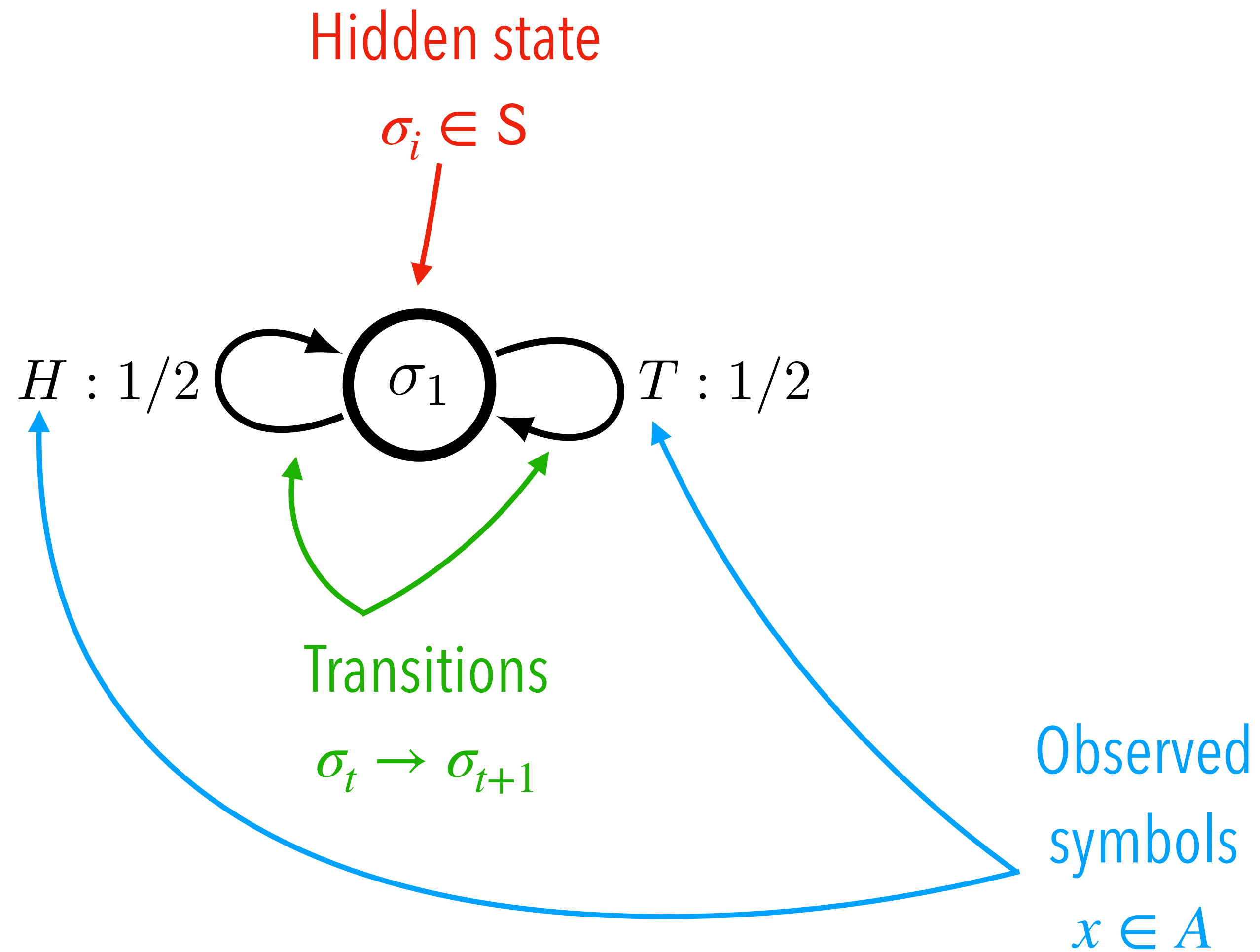
Models of Time Series



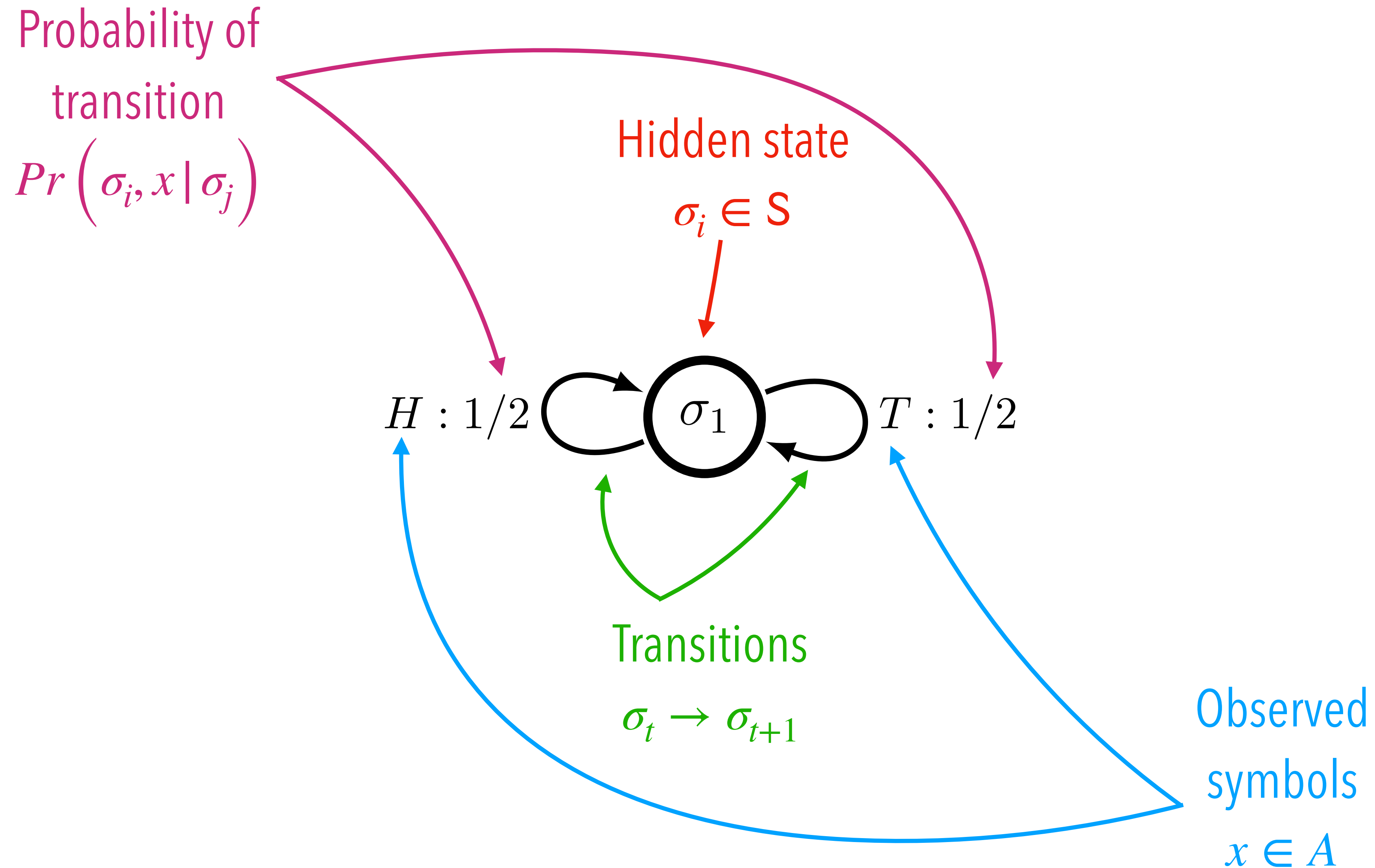
Models of Time Series



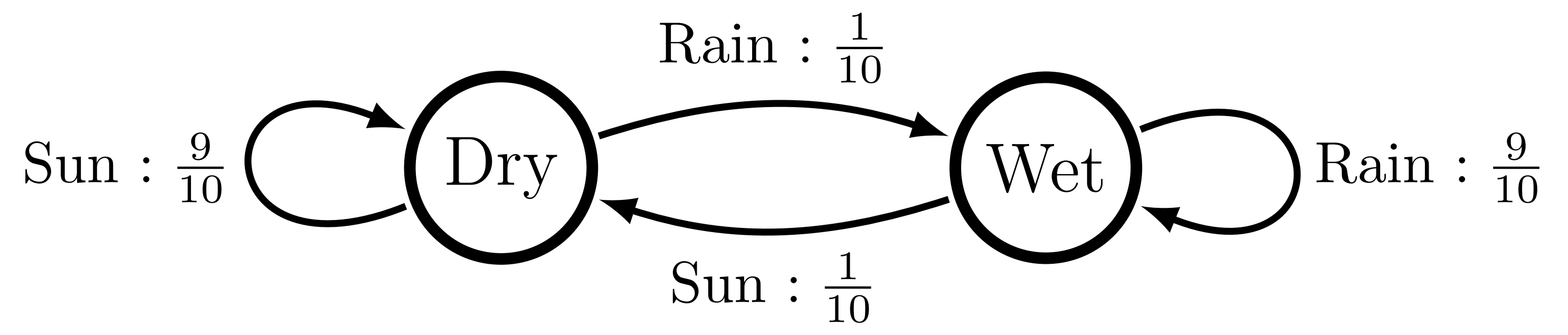
Models of Time Series



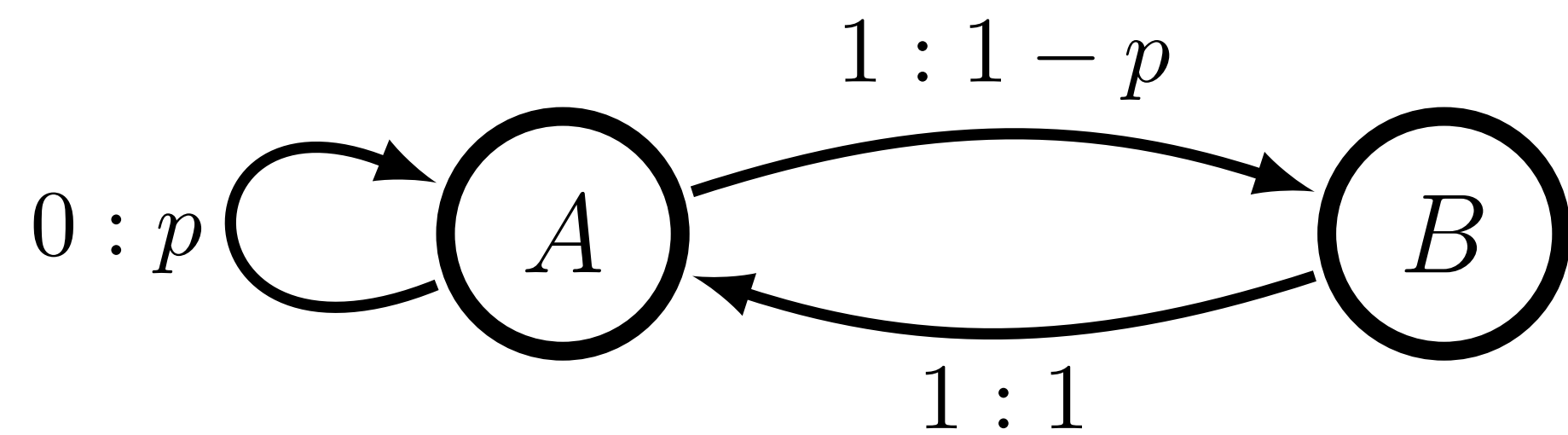
Models of Time Series



Hidden Markov Models

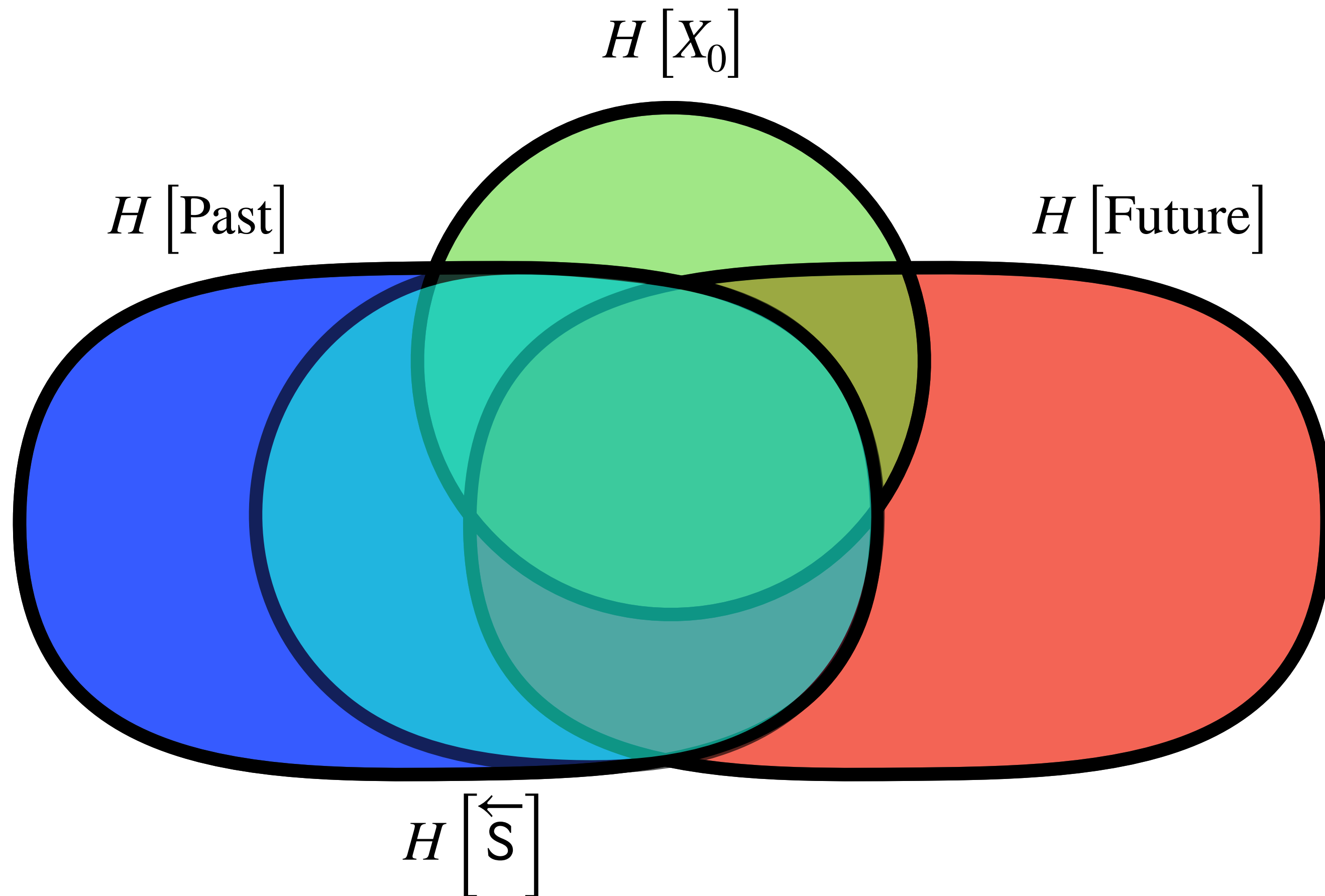


The Utility of (Good) Models



$$h_{\mu} = \sum_{\sigma \in \mathcal{S}} \Pr(\sigma) H[X | \sigma]$$

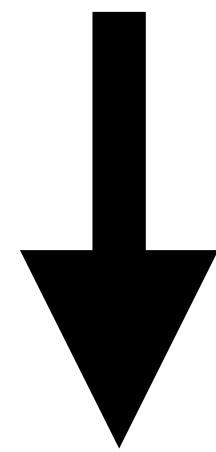
The Utility of (Good) Models



- States are a function of the past
- Optimal predictor
- Minimal in size
- Unique

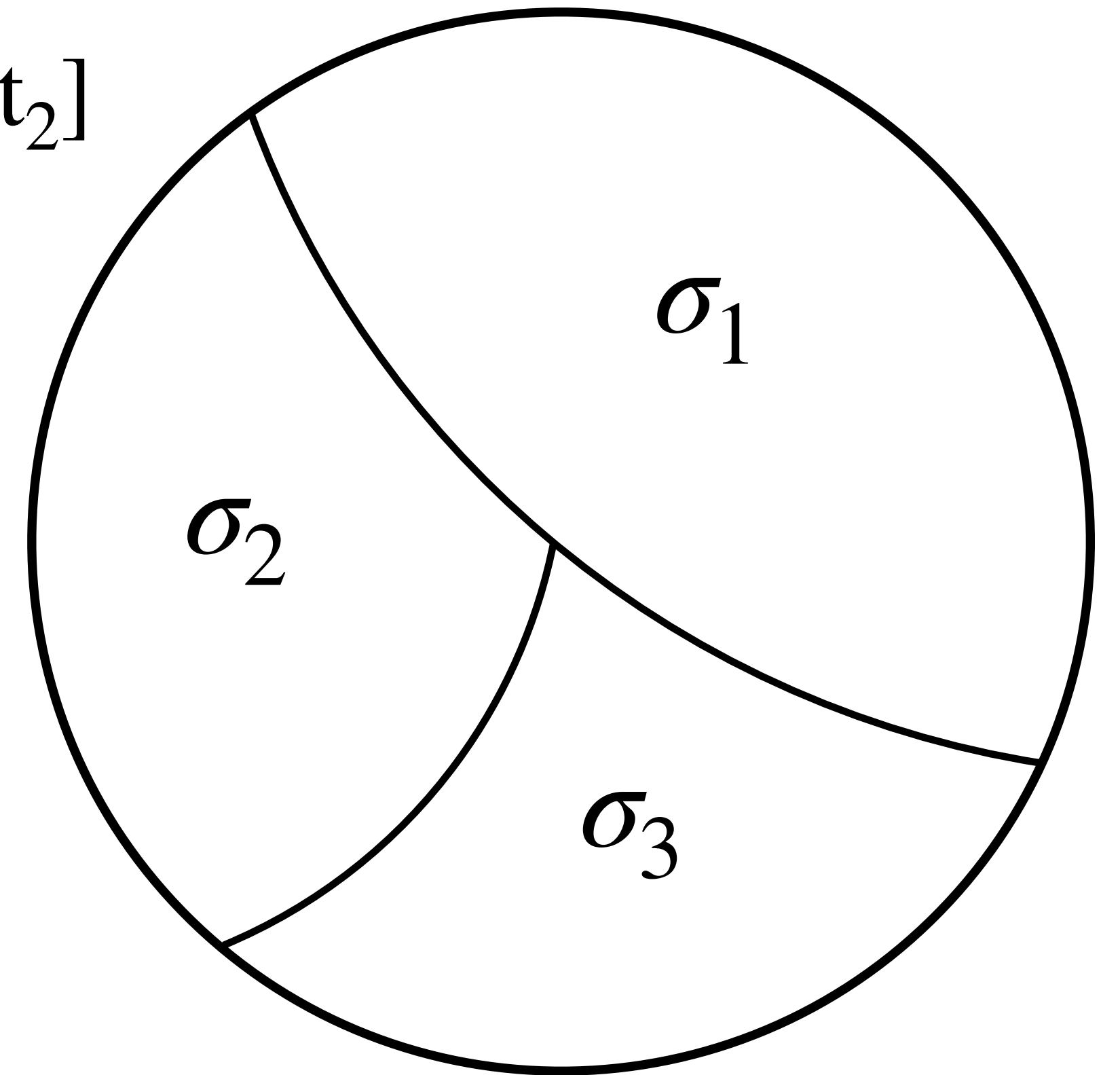
What's "Good"?

$$\text{Past}_1 \sim \text{Past}_2 \iff \Pr[\text{Future} \mid \text{Past}_1] = \Pr[\text{Future} \mid \text{Past}_2]$$



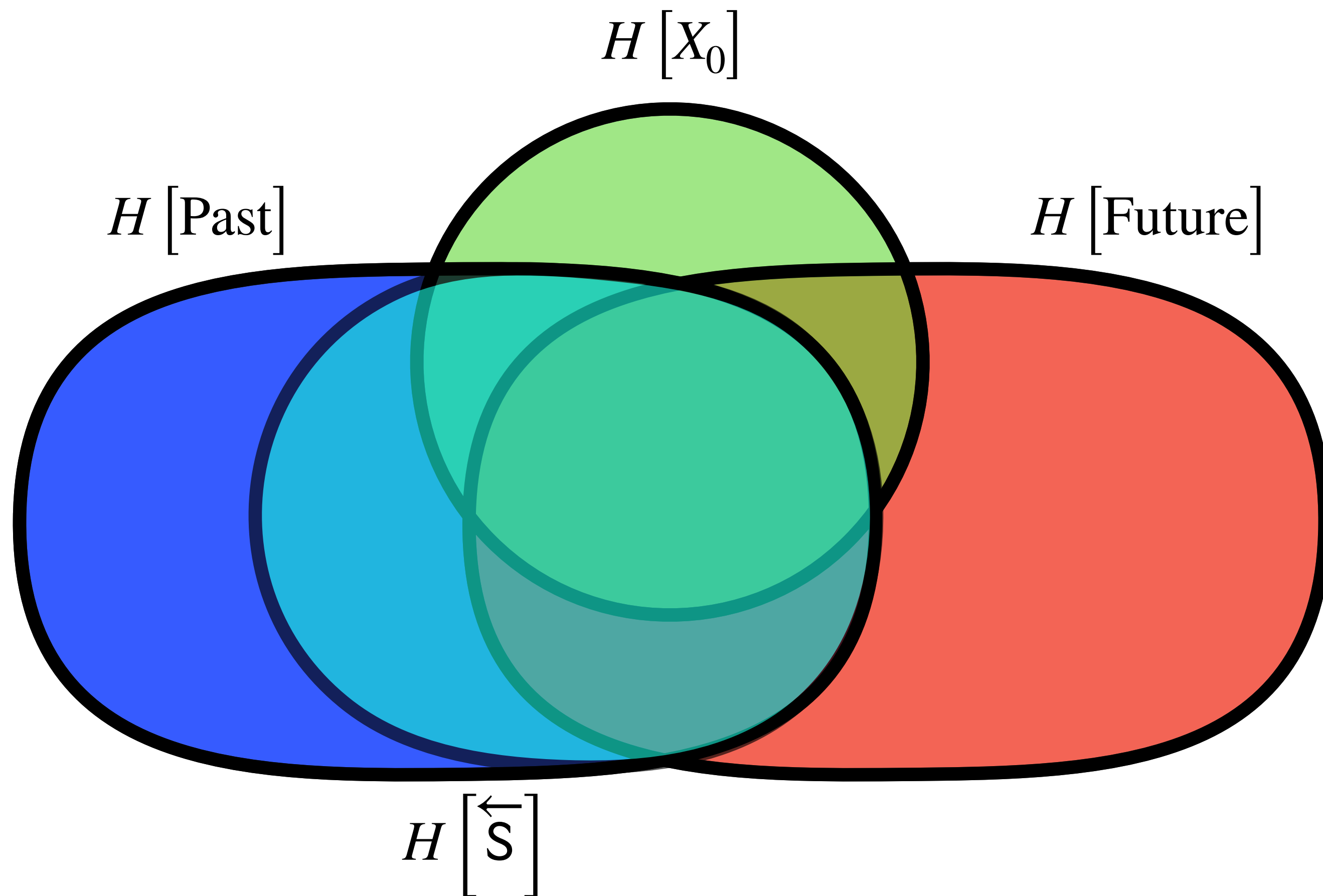
Equivalence class partitions pasts into the *causal states*

$$\mathbf{S} = \mathbf{H} / \sim .$$



Space of all
Pasts \mathbf{H}

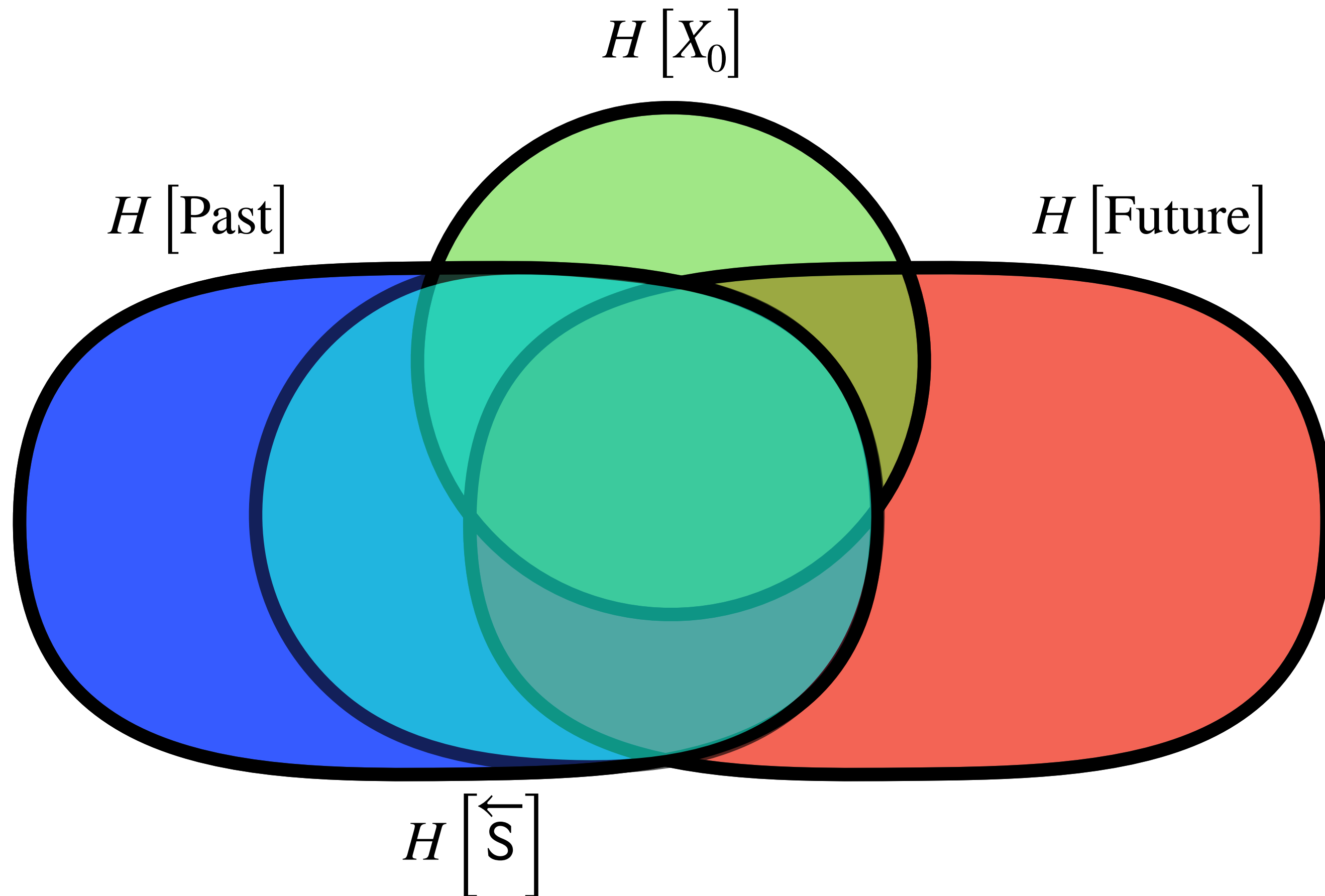
ϵ -Machines are "Good"



- States are a function of the past
- Optimal predictor
- Minimal in size
- Unique

➡ The ϵ machine

ϵ -Machines are "Good"

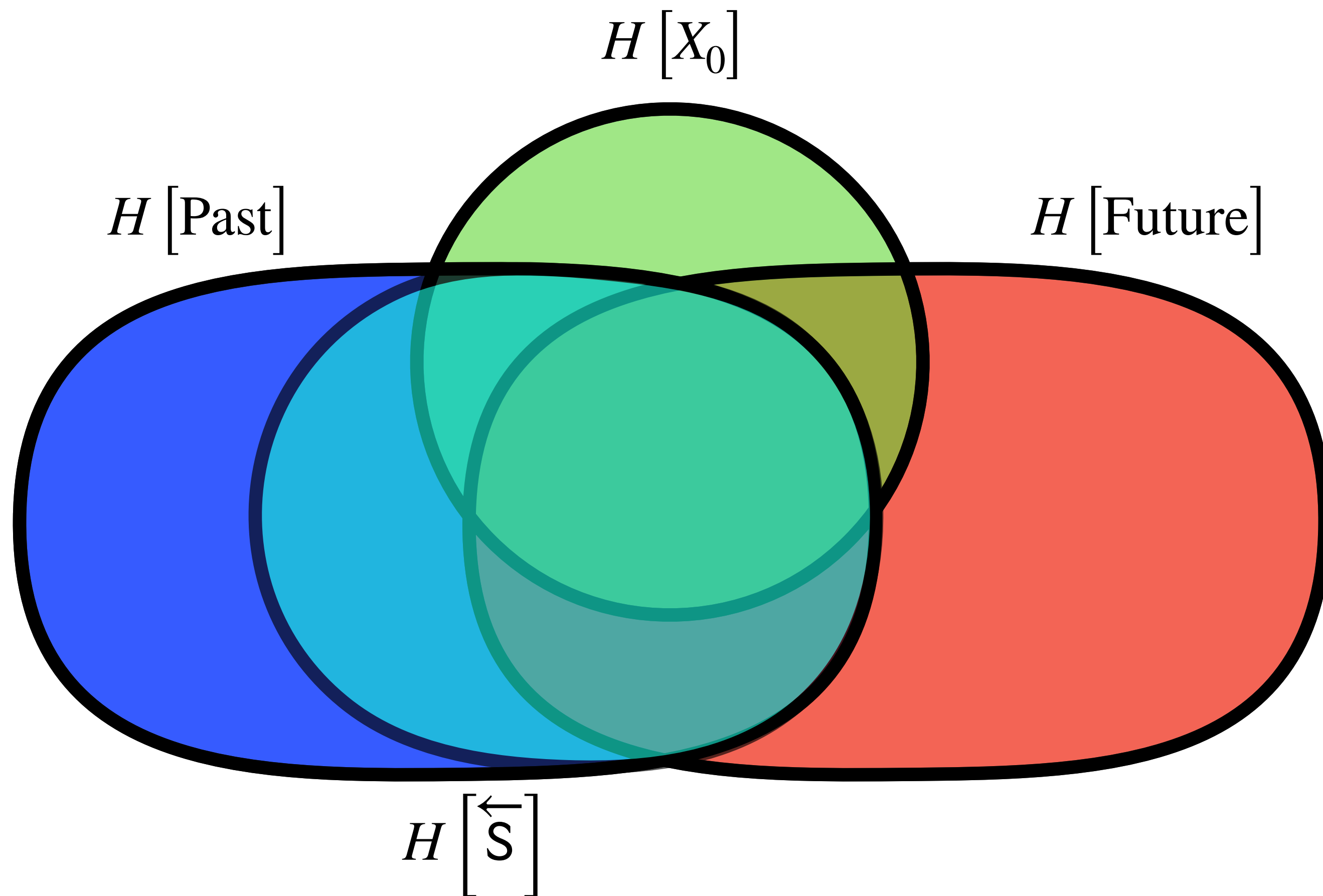


- States are a function of the past
- Optimal predictor
- Minimal in size
- Unique

➔ The ϵ machine

Mechanism?

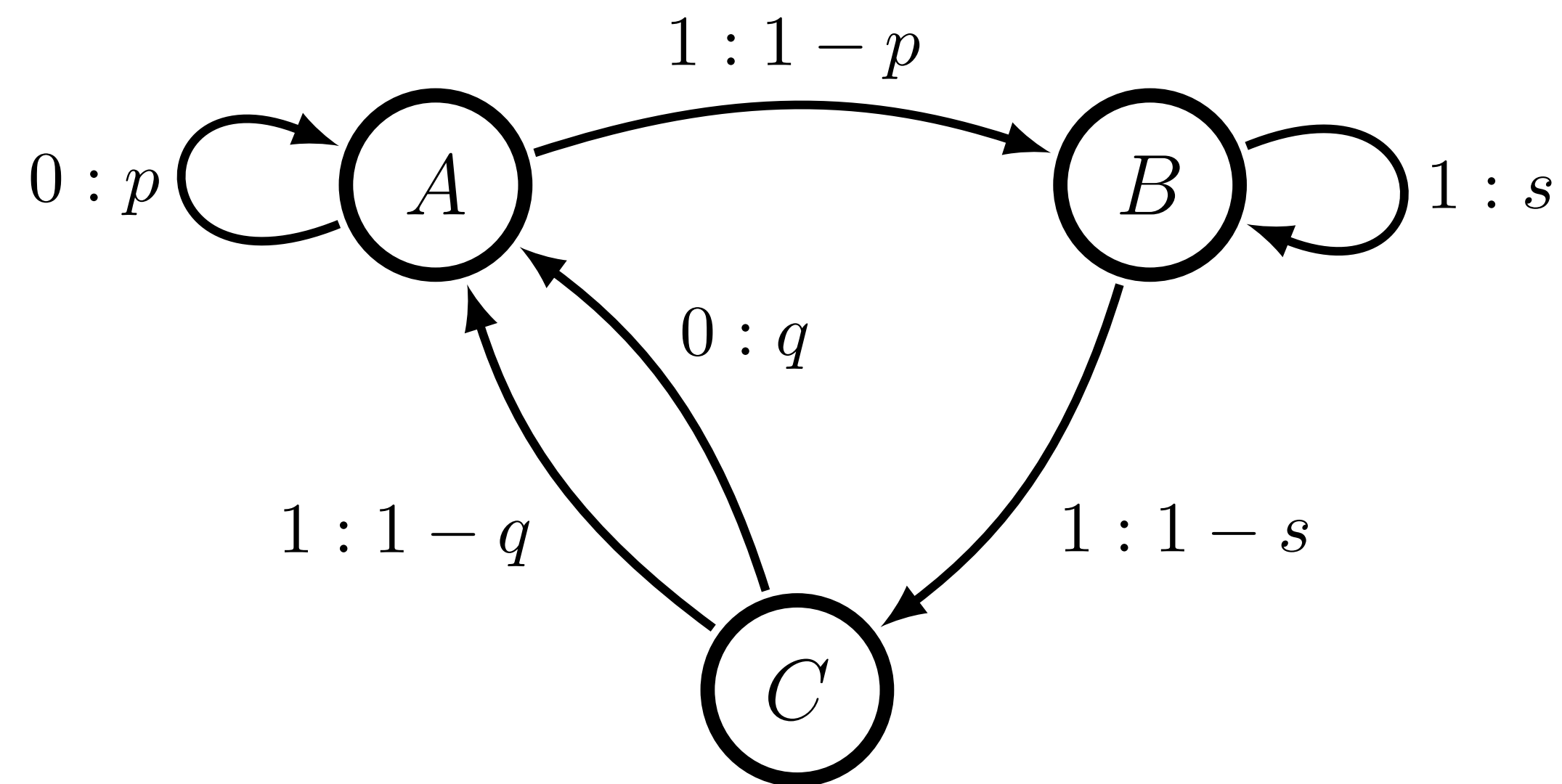
ϵ -Machines are "Good"



The number of bits required to store the ϵ -machine is called the *statistical complexity*:

$$C_\mu = H[\vec{S}]$$

Most HMMs are *Not* Optimal Predictors!



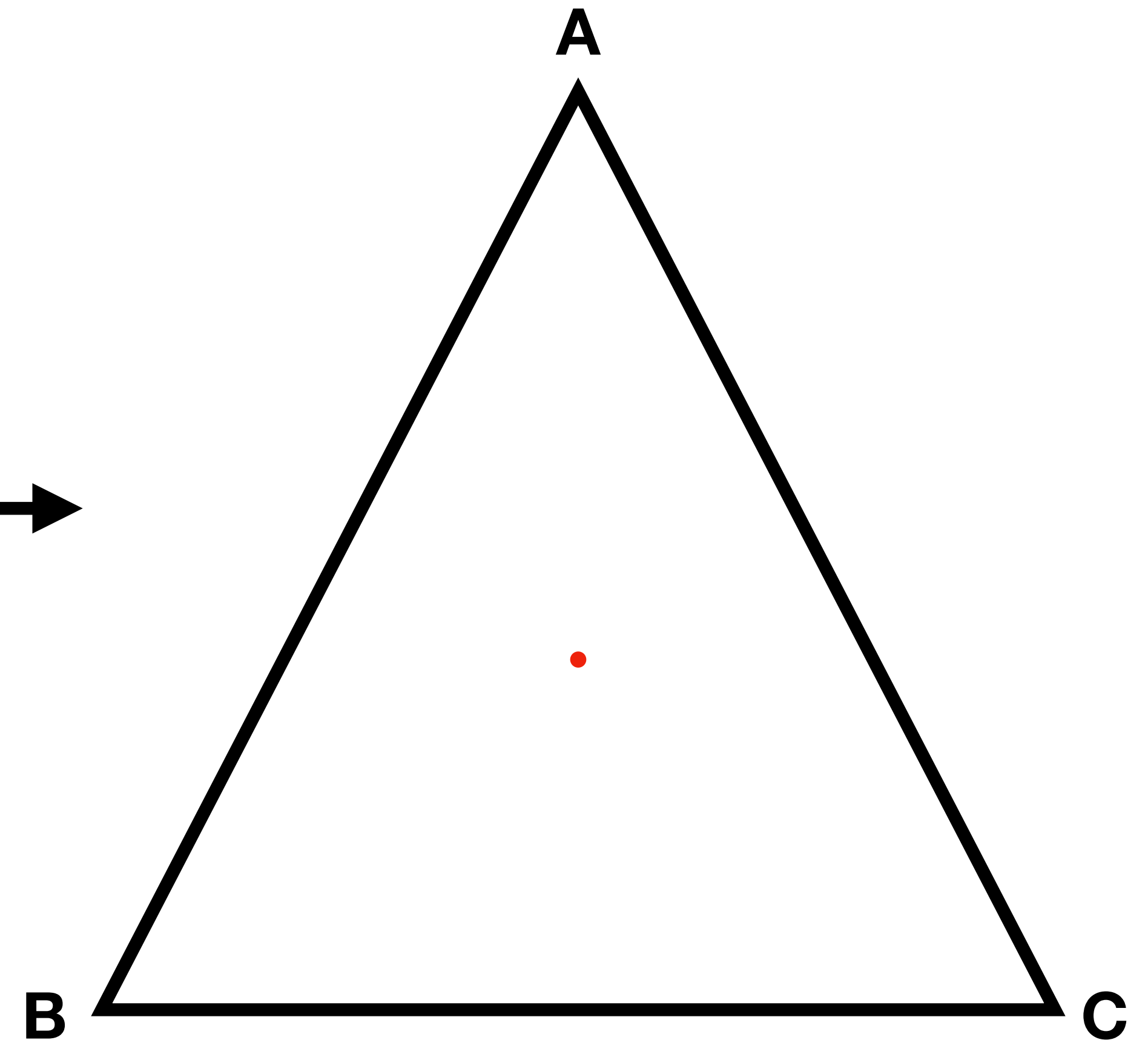
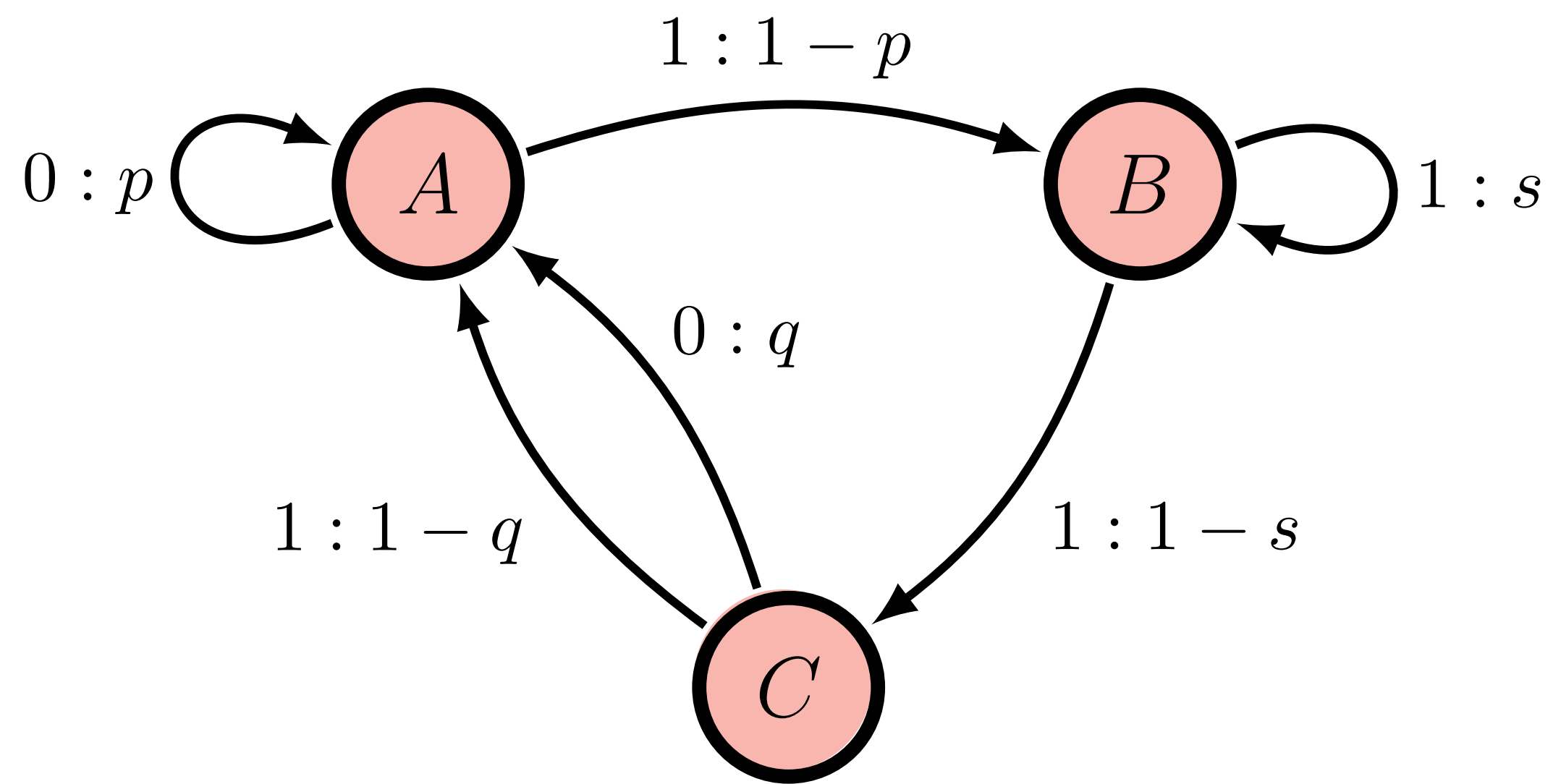
$$h_{\mu} \neq \sum_{\sigma \in \mathcal{S}} \Pr(\sigma) H[X | \sigma]$$

Problem Statement

Given an arbitrary finite-state hidden Markov model,
how to get ϵ -machine?

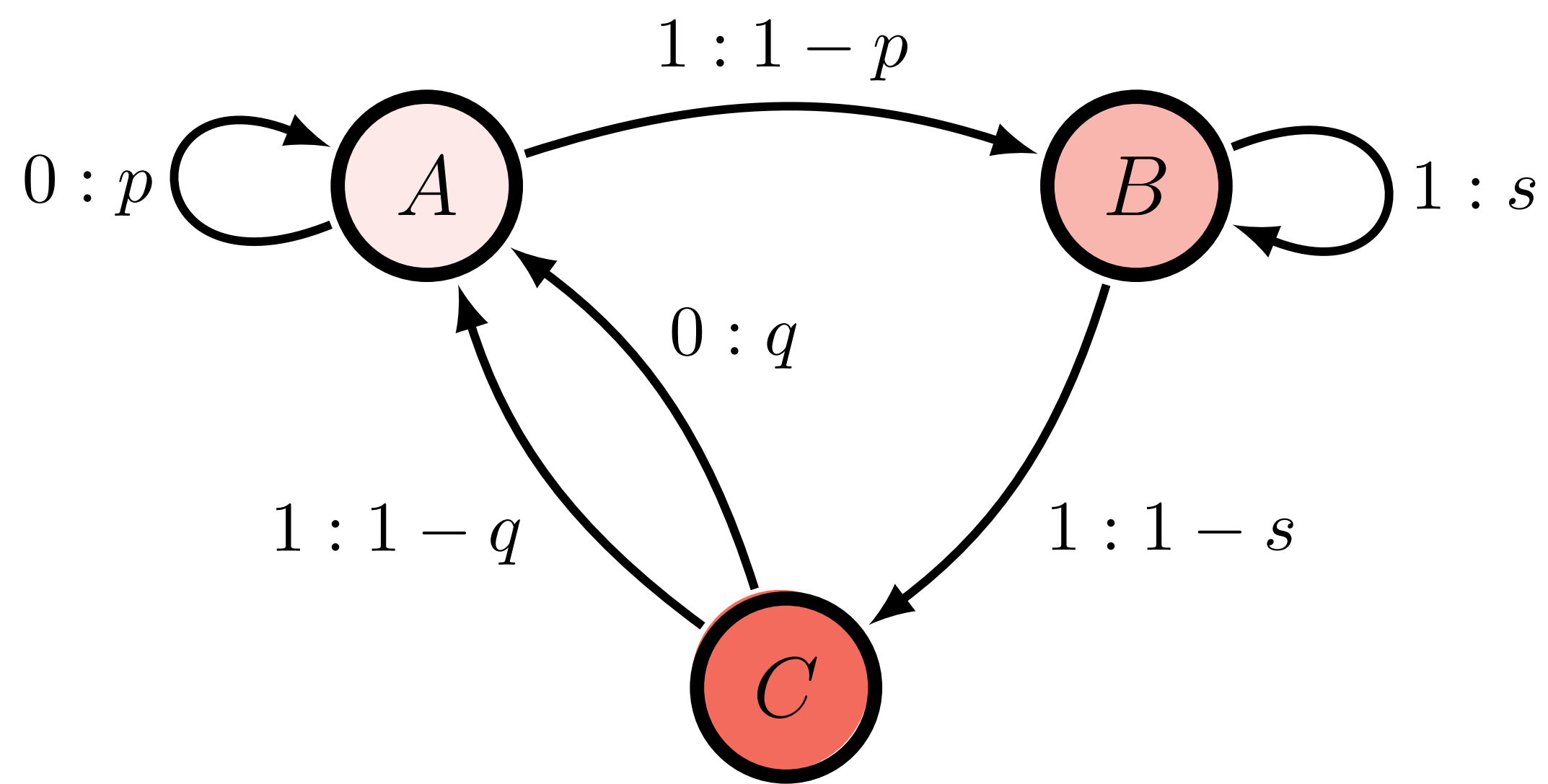
Observe and Update Game

Take our "bad" HMM and track our belief over what state we're in:

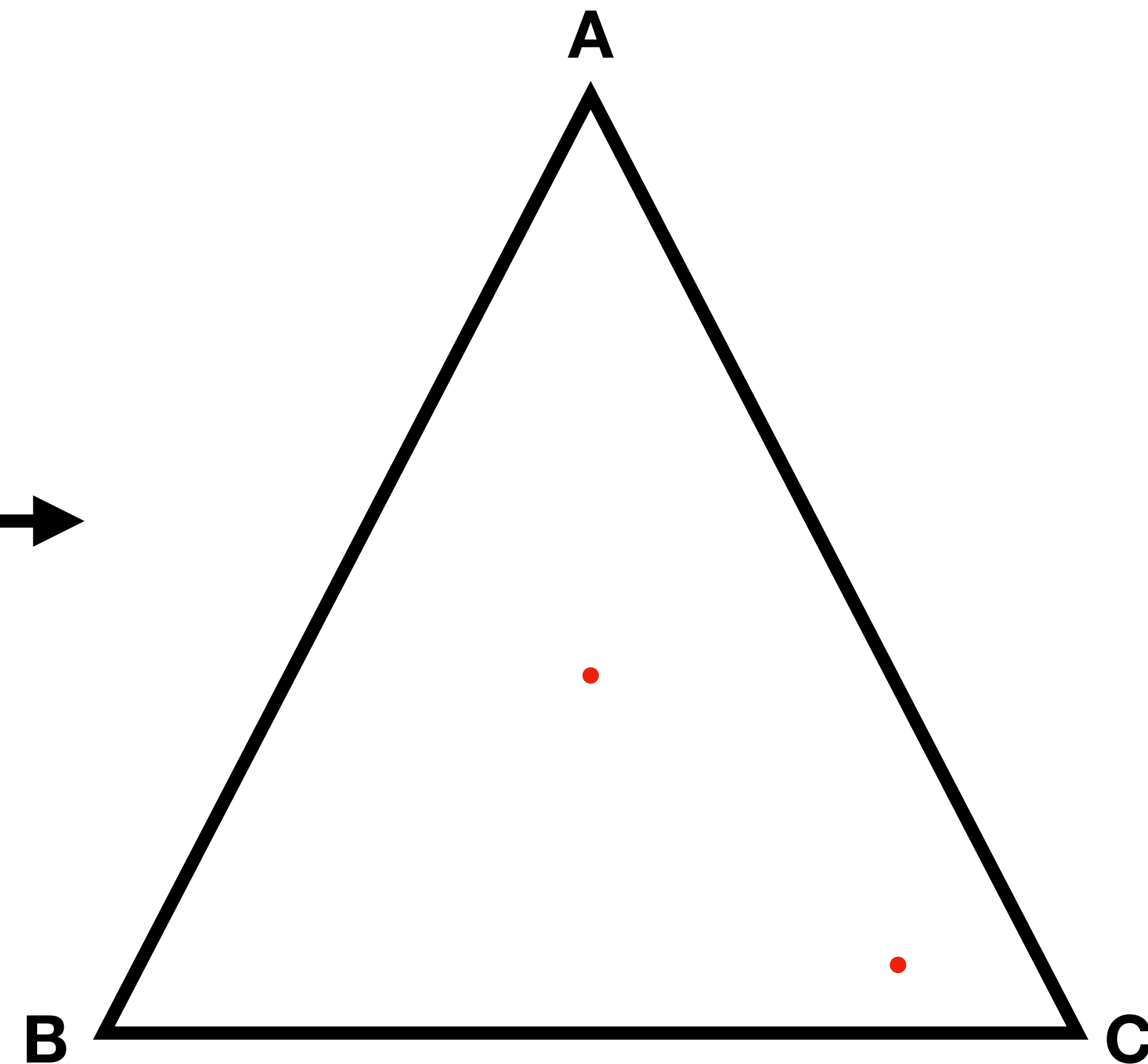


Observed sequence: λ

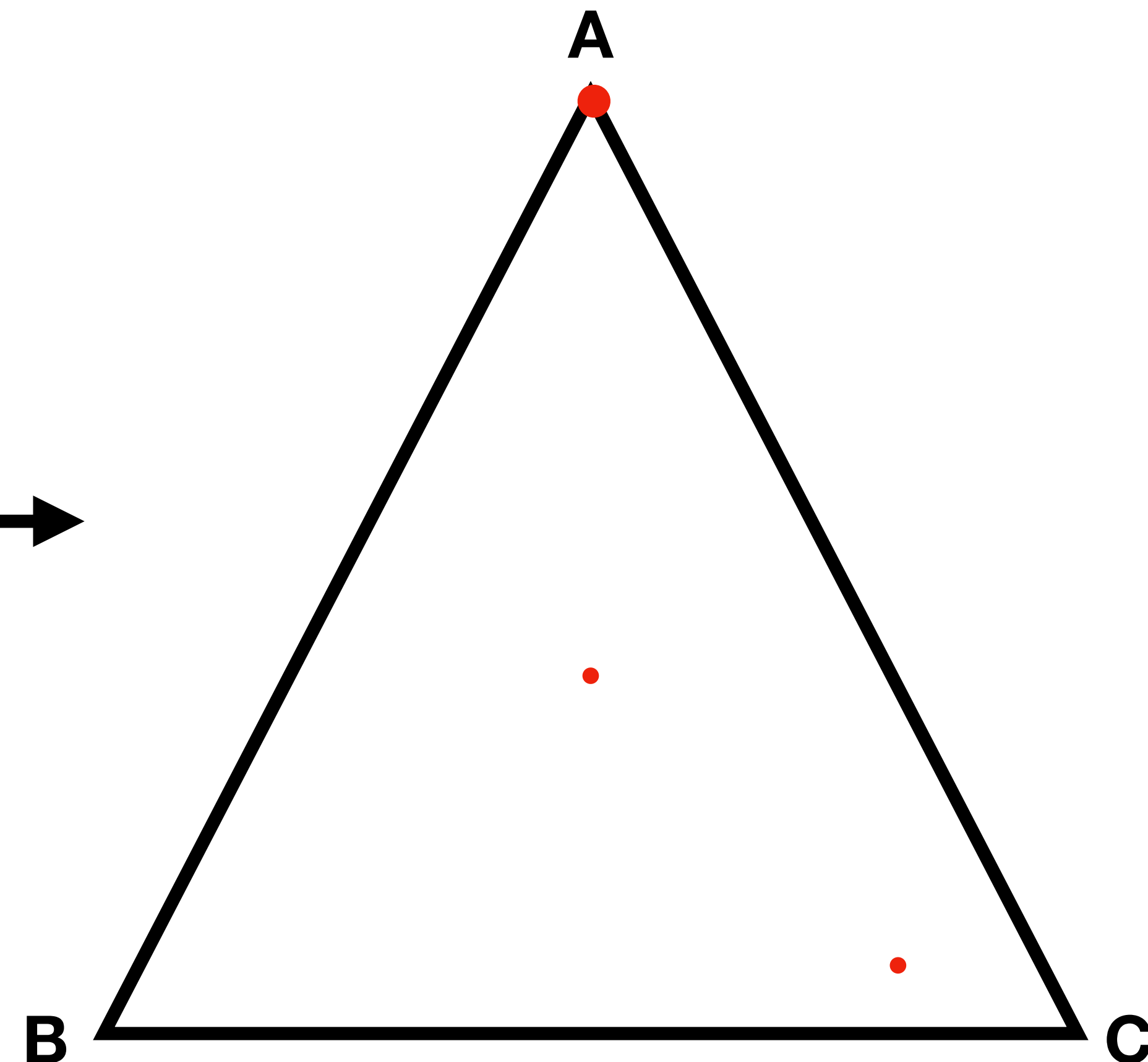
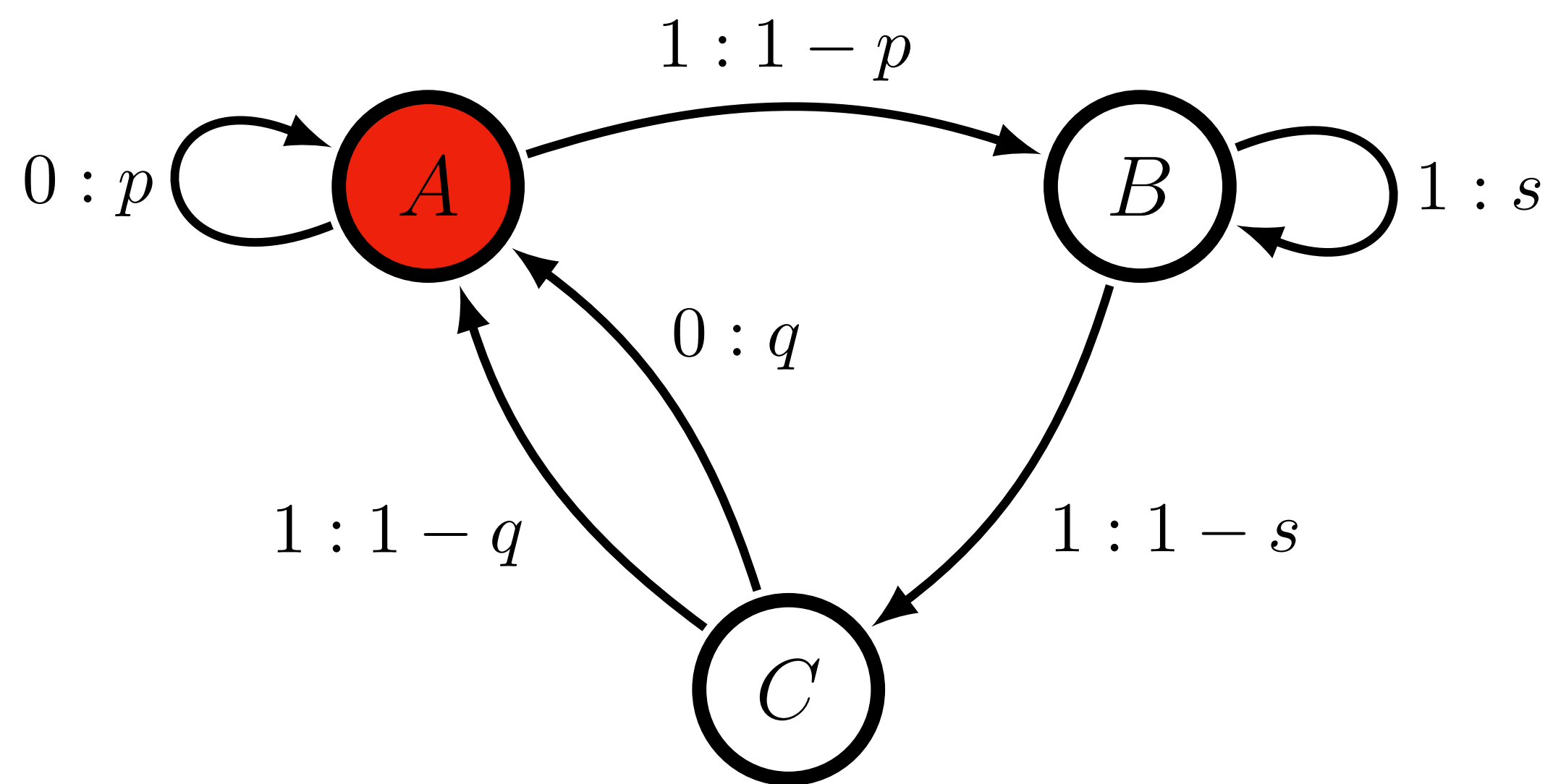
Observe and Update Game



Observed sequence: $\lambda 1$

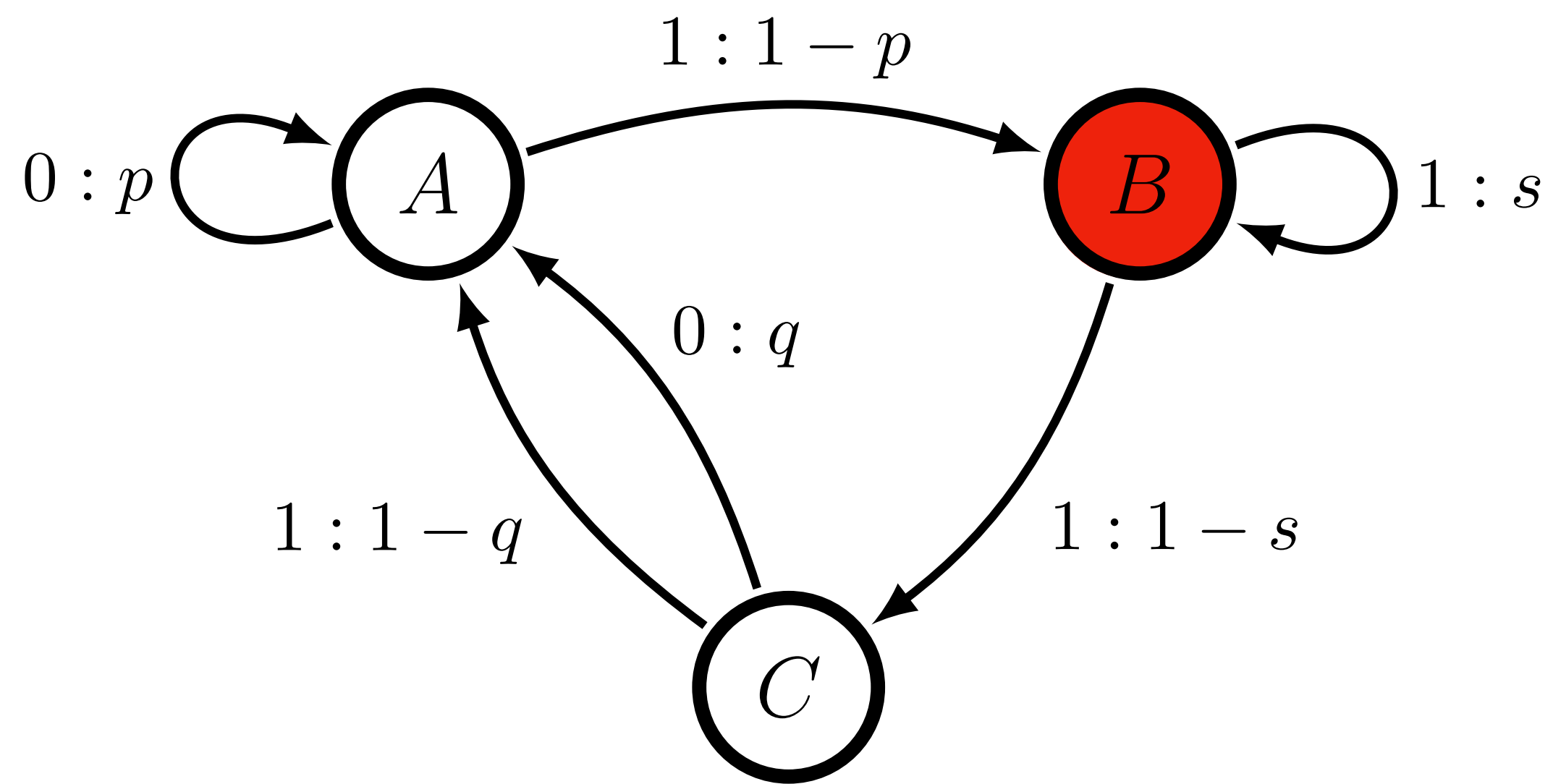


Observe and Update Game

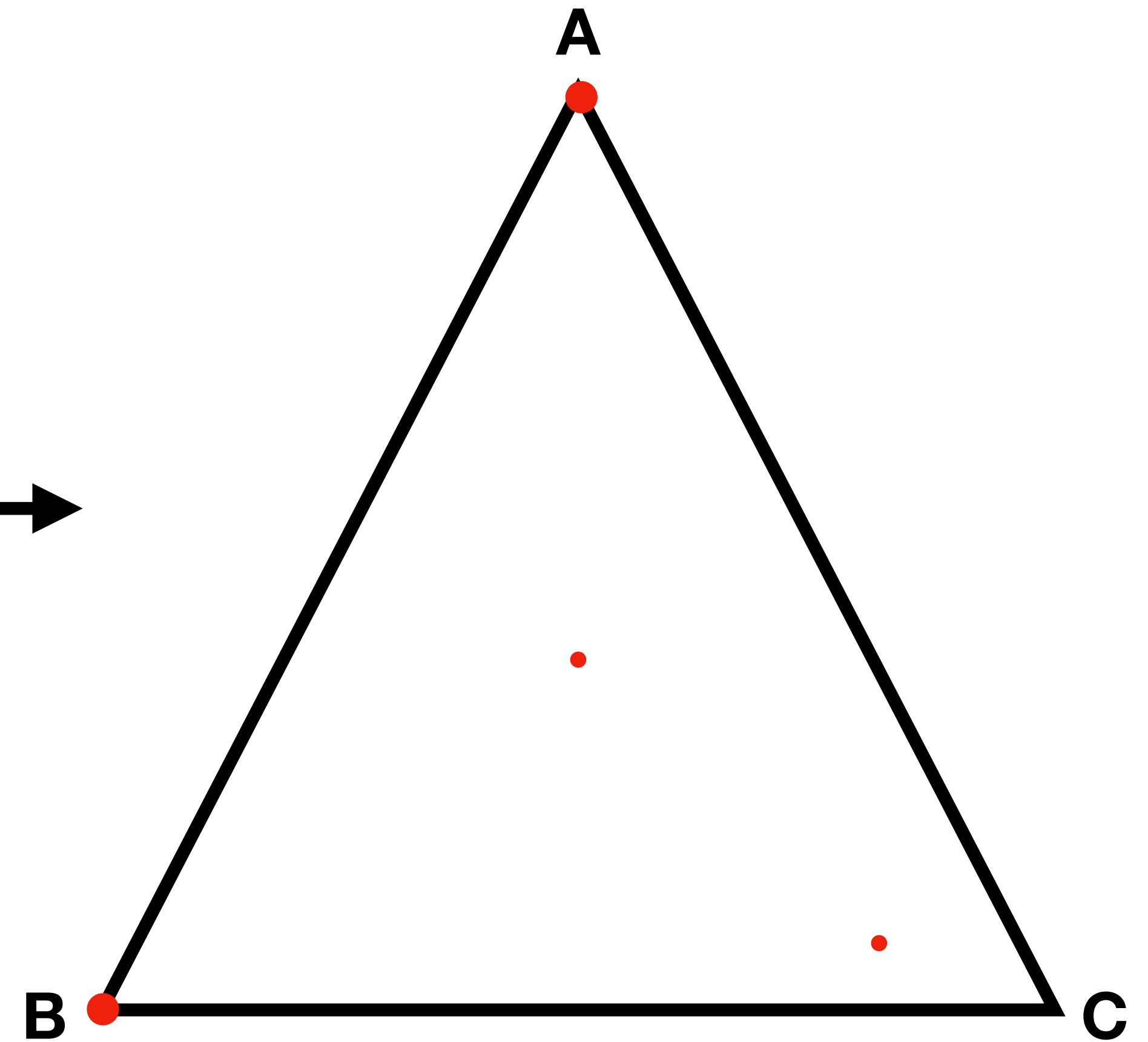


Observed sequence: $\lambda 10$

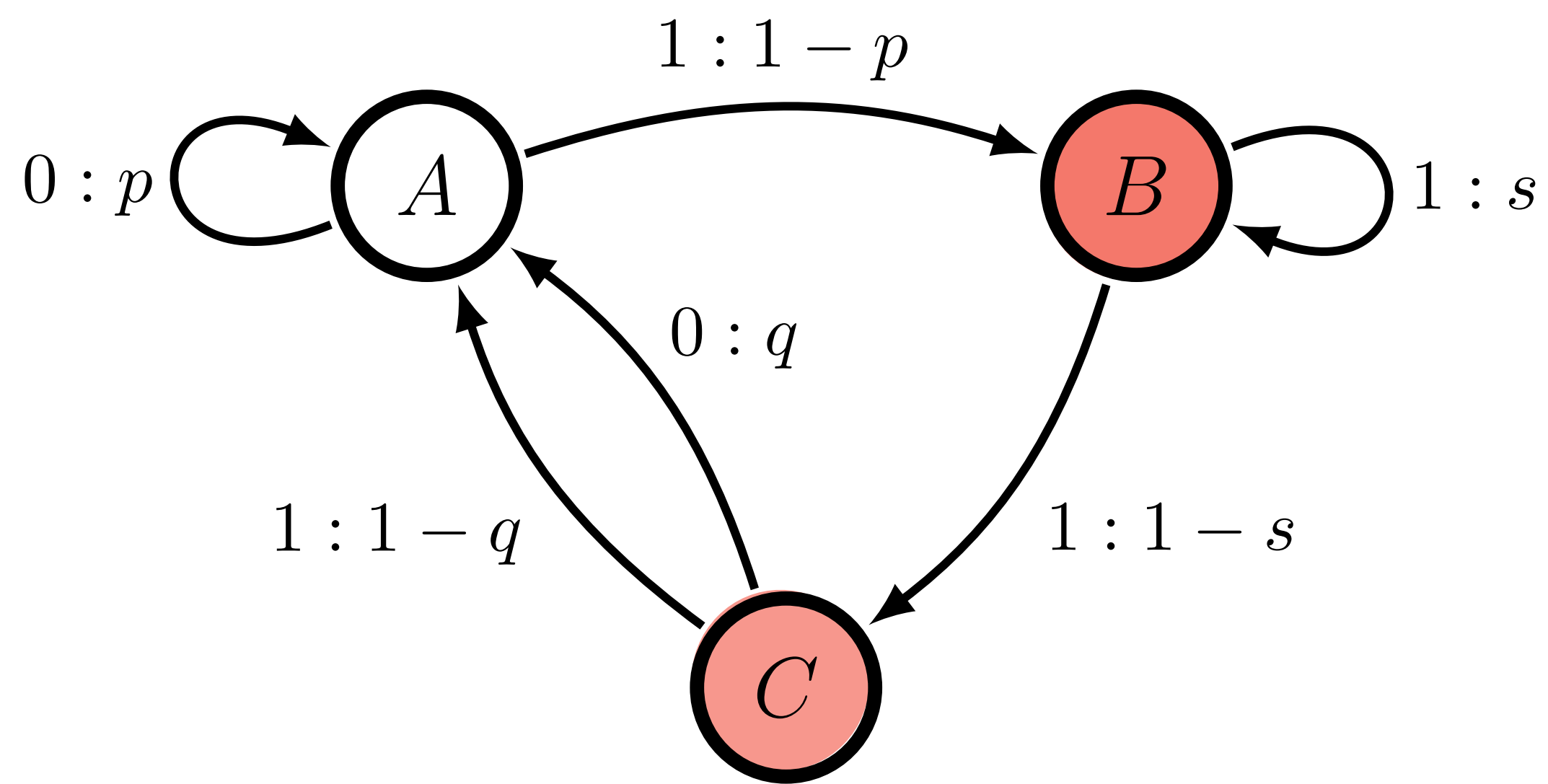
Observe and Update Game



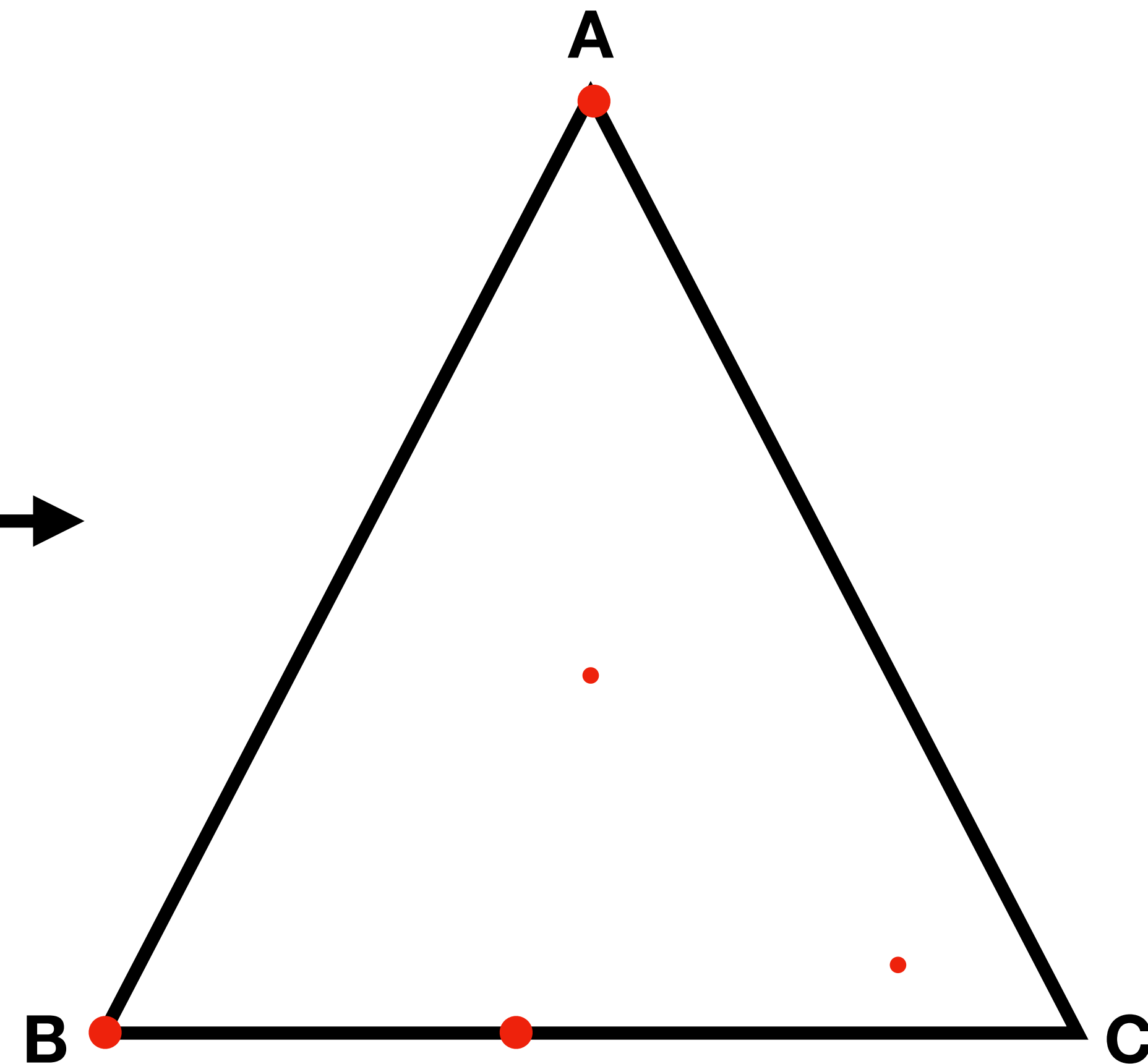
Observed sequence: $\lambda 101$



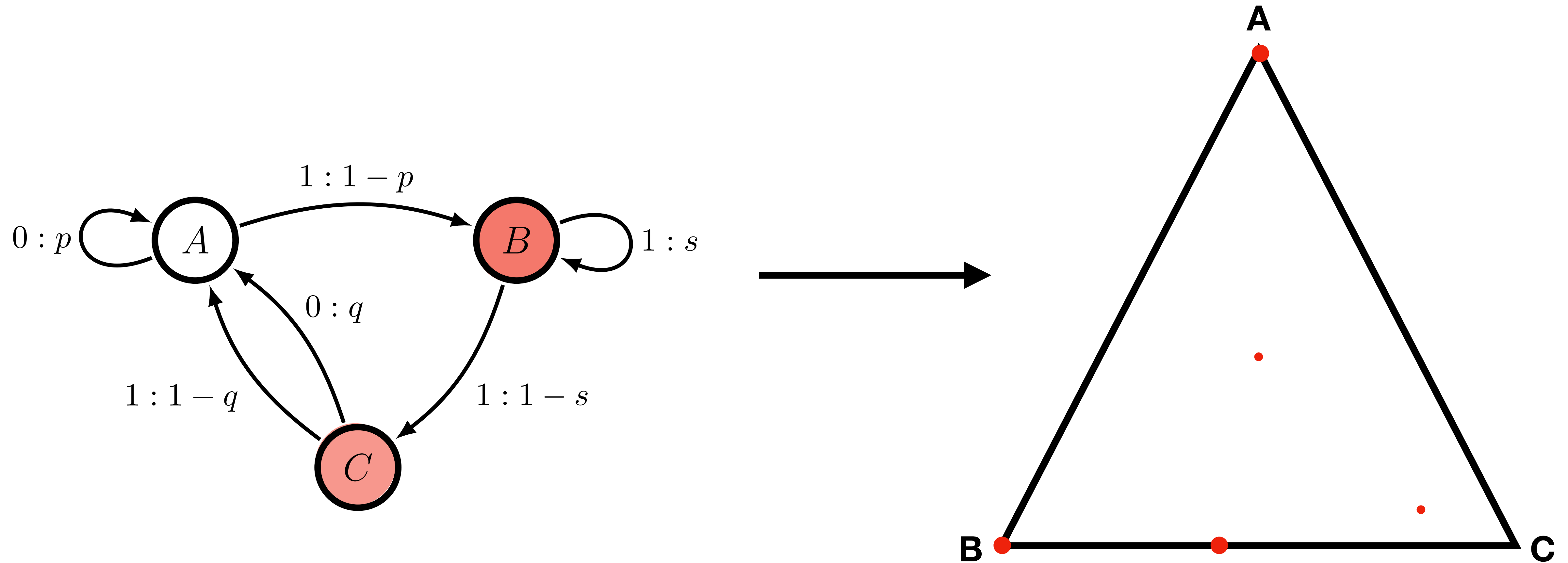
Observe and Update Game



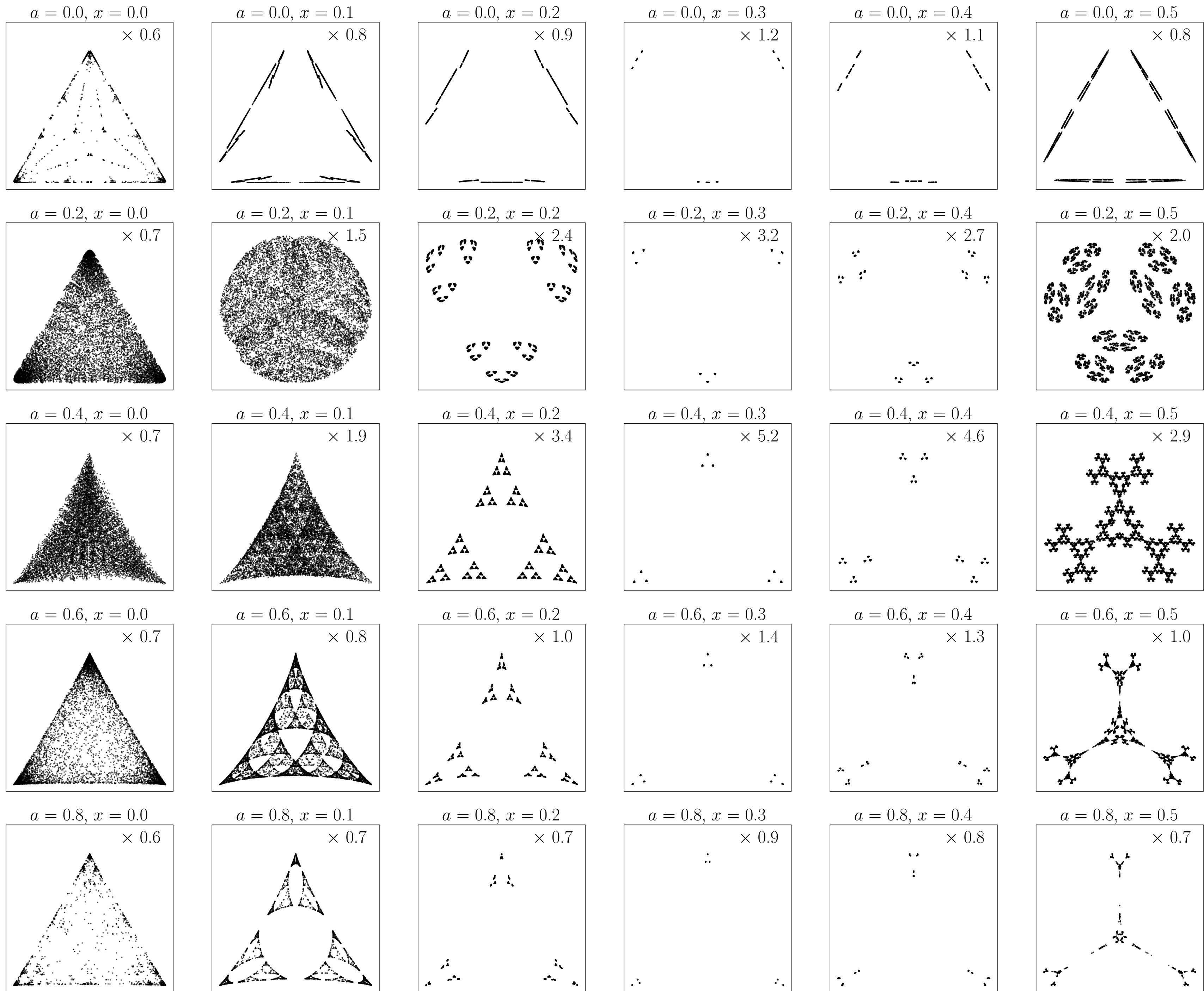
Observed sequence: $\lambda 1011$



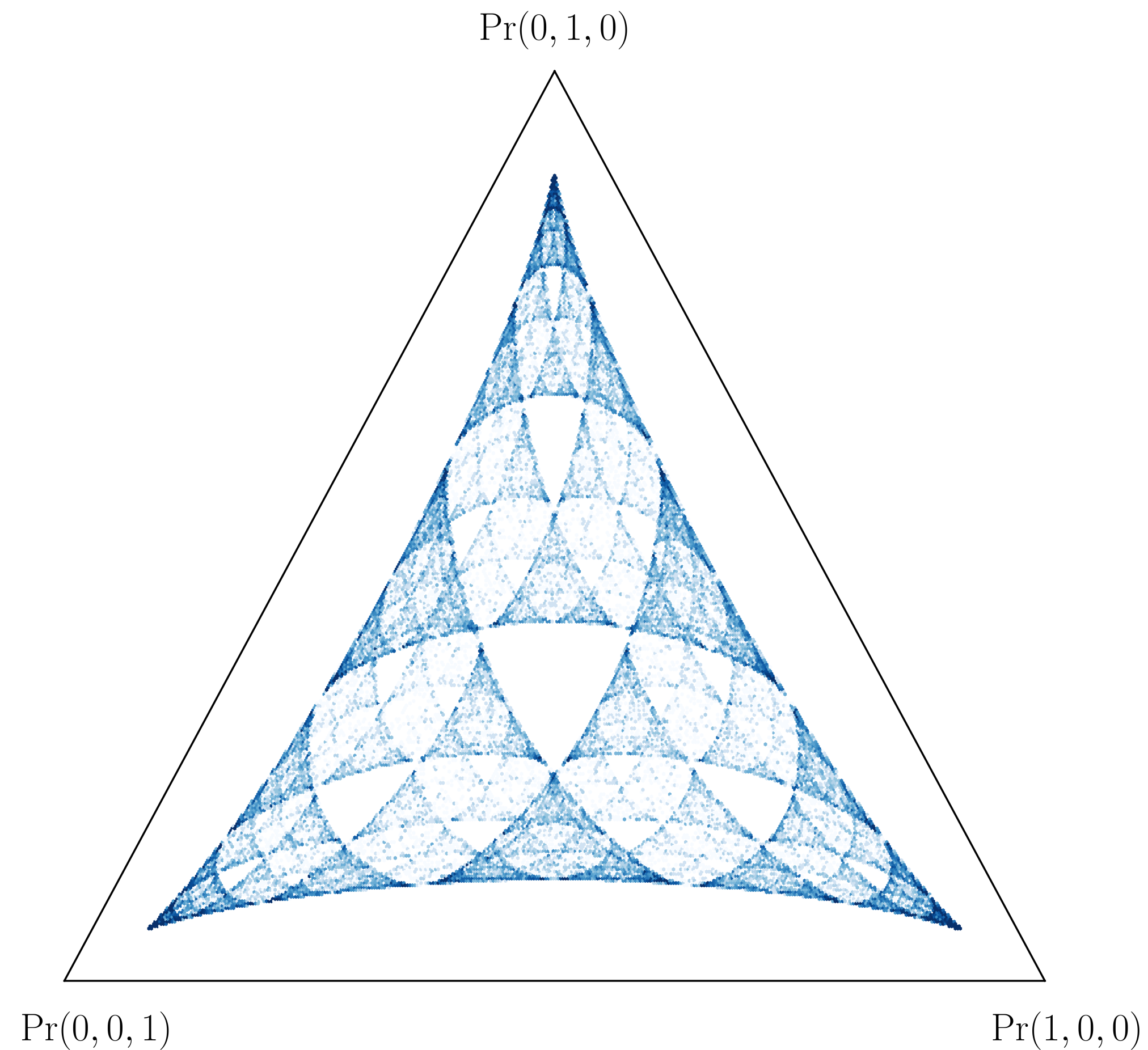
Observe and Update Game



→ Build up a set $R = \left\{ \eta : \eta(w) = \Pr(S_l | X_{0:l} = w, S_0 = \pi) \right\}$ of belief states.

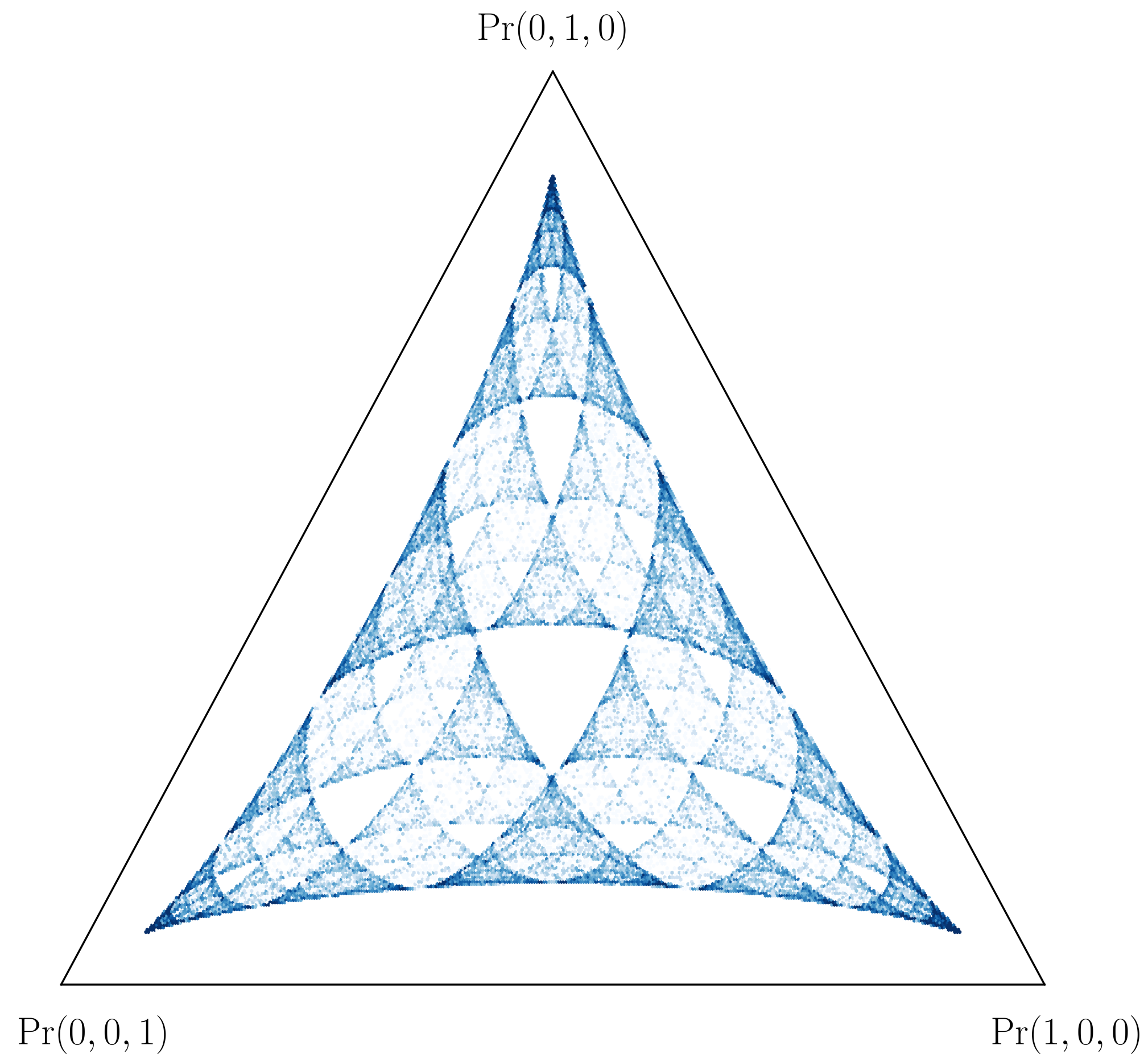


The Infinite State ϵ -Machine



The set of belief states R is the attractor of the "observe and update" stochastic dynamical system (known as an *iterated function system*).

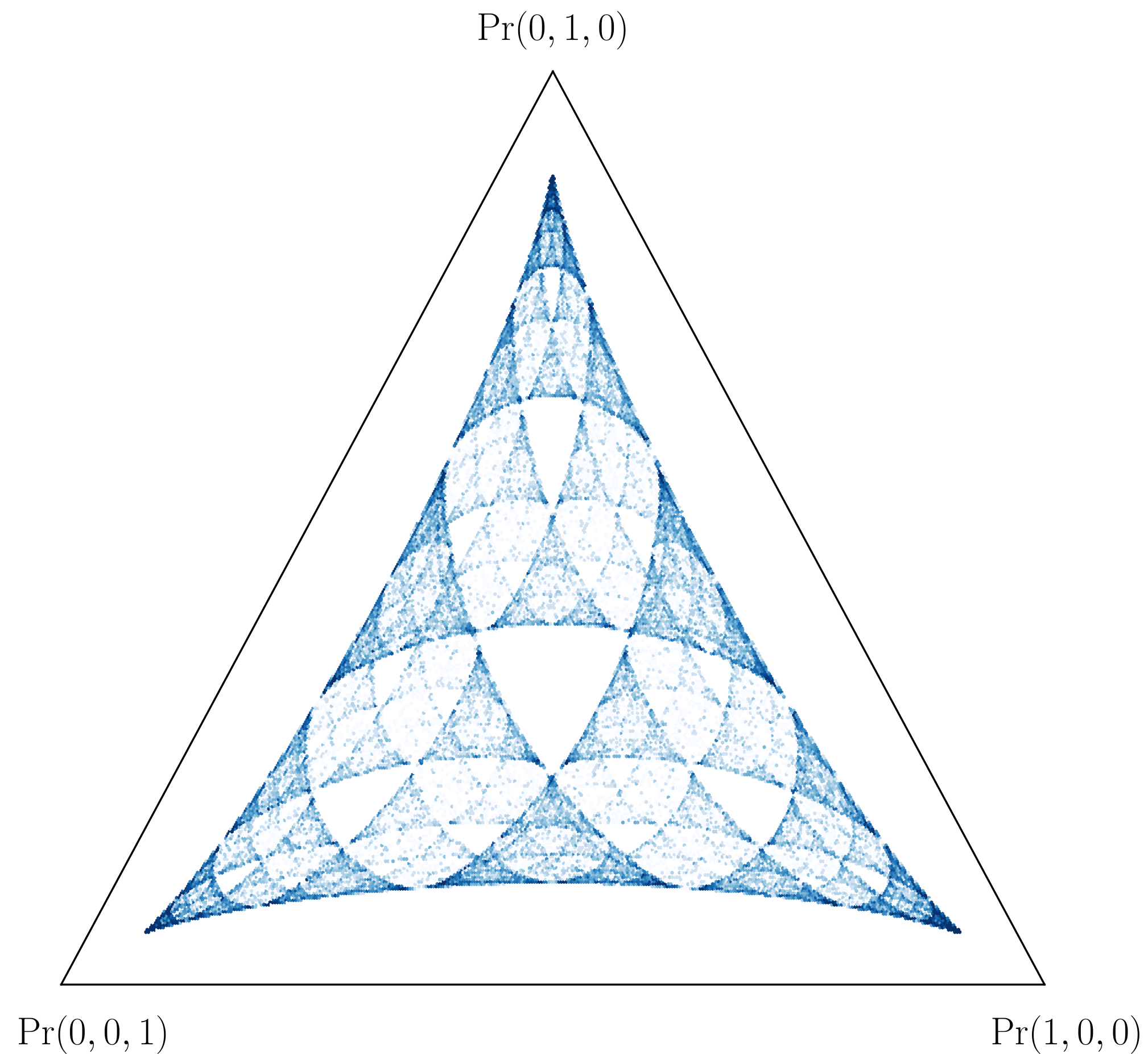
The Infinite State ϵ -Machine



The set of belief states R is the attractor of the “observe and update” stochastic dynamical system (known as an *iterated function system*).

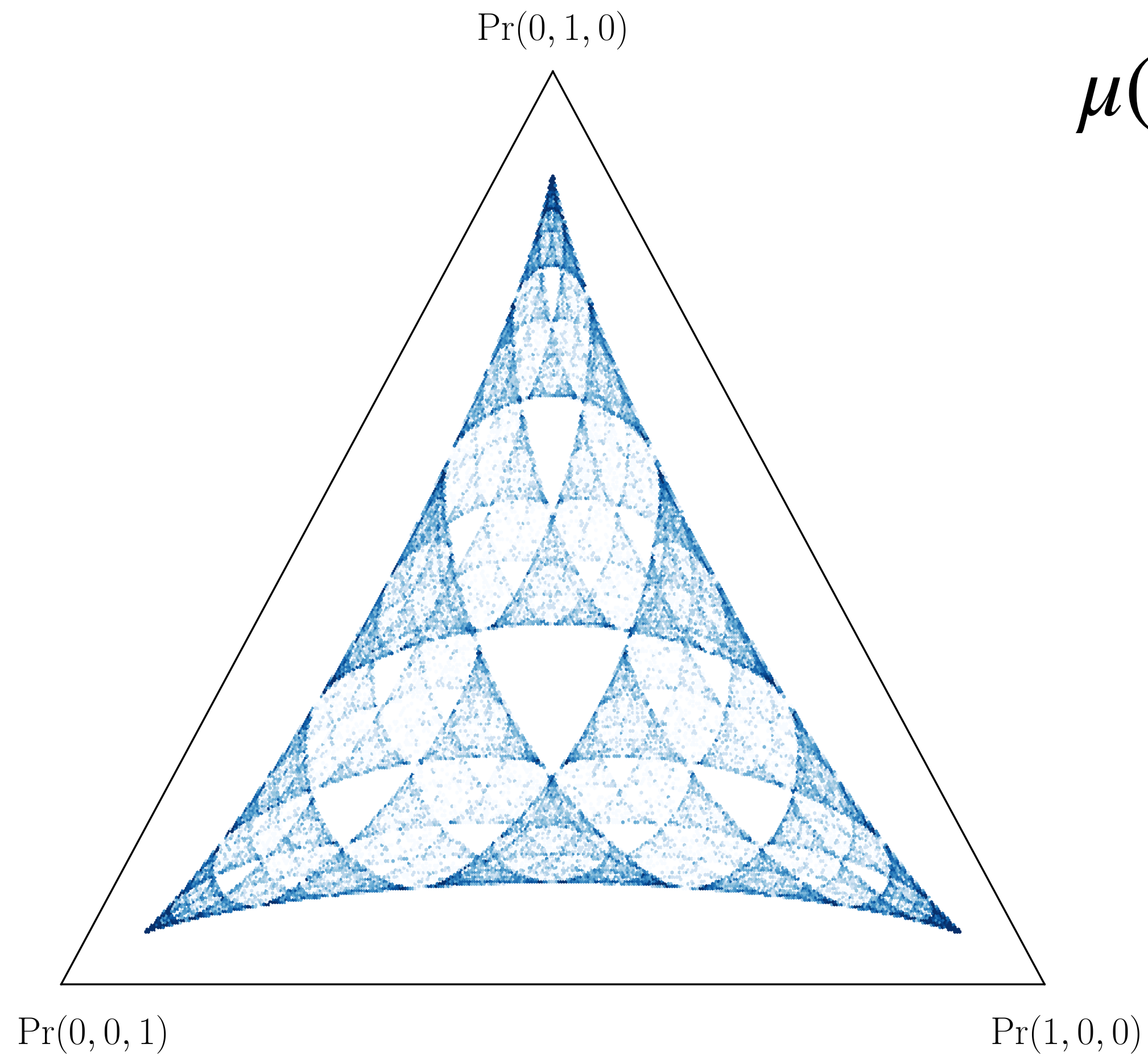
This attractor exists and is unique due to contractivity, but is generically fractal-like.

The Infinite State ϵ -Machine



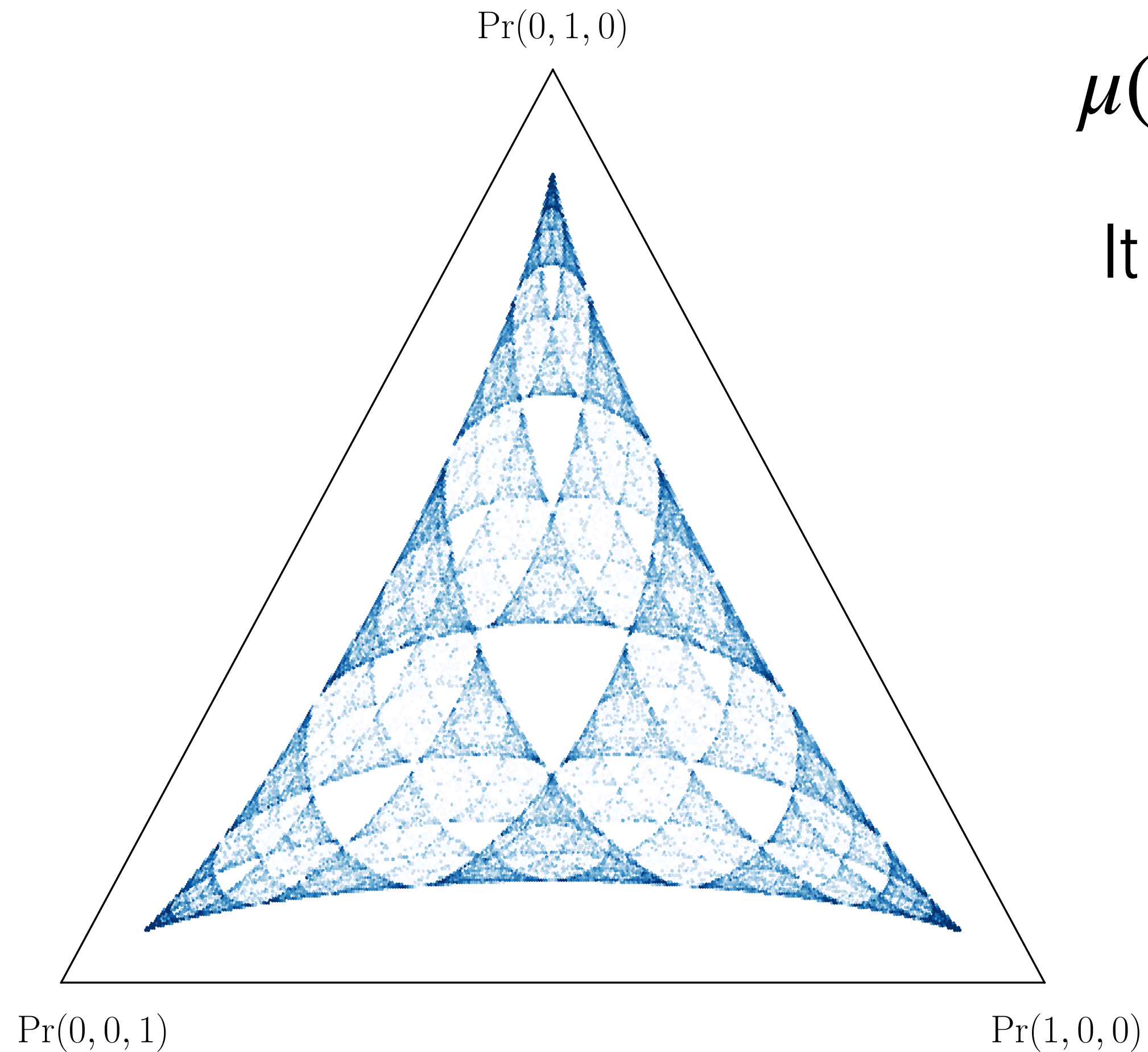
For most *finitely generated* hidden Markov processes, the ϵ machine is an uncountably infinite state set + transitions between these states.

Calculating Shannon Entropy Rate



$\mu(R)$ is called the Blackwell measure.

Calculating Shannon Entropy Rate

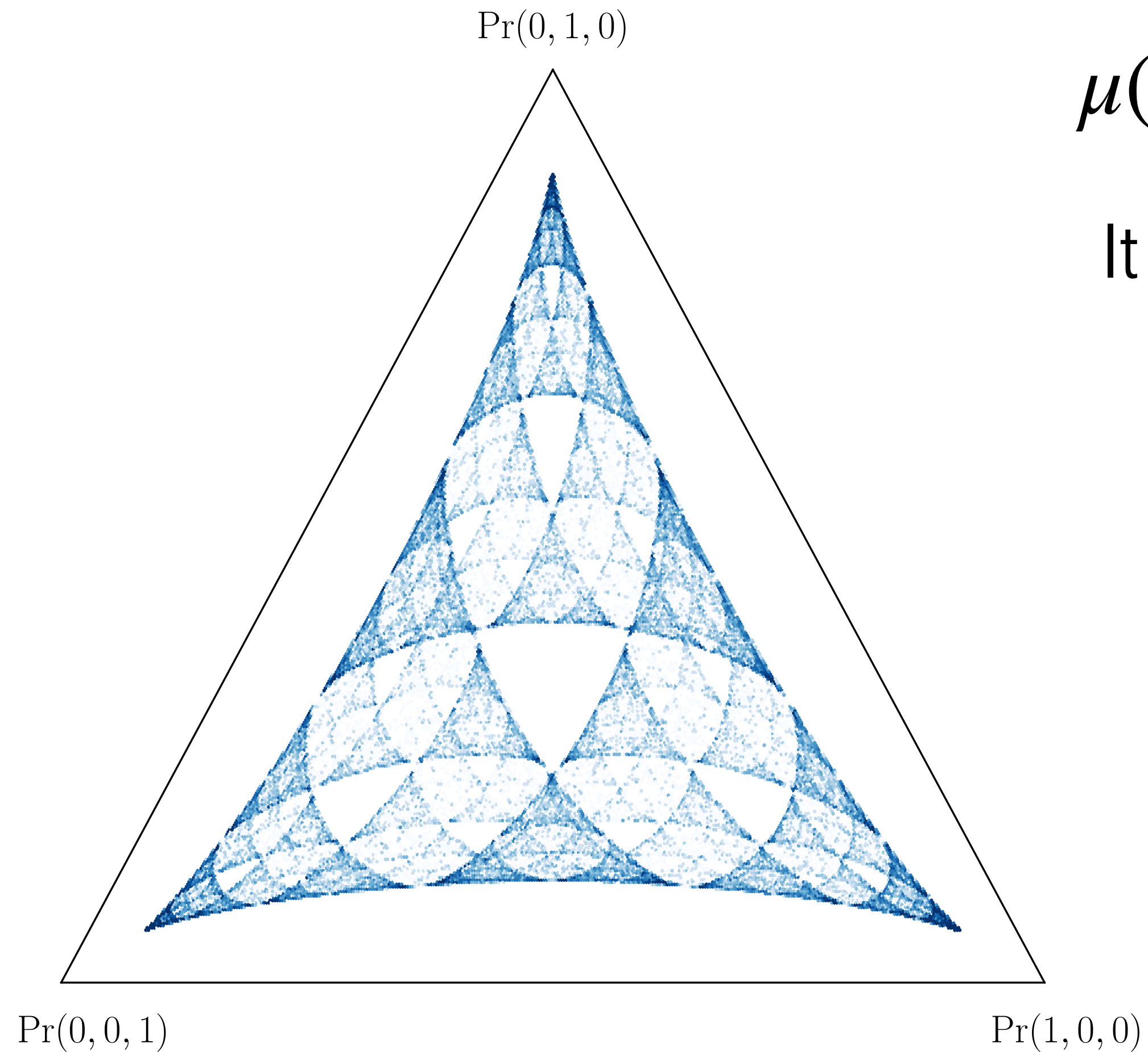


$\mu(R)$ is called the Blackwell measure.

It can be used to calculate the entropy rate:

$$h_{\mu}^B = \int_R d\mu(\eta) H[X | \eta]$$

Calculating Shannon Entropy Rate



$\mu(R)$ is called the Blackwell measure.

It can be used to calculate the entropy rate:

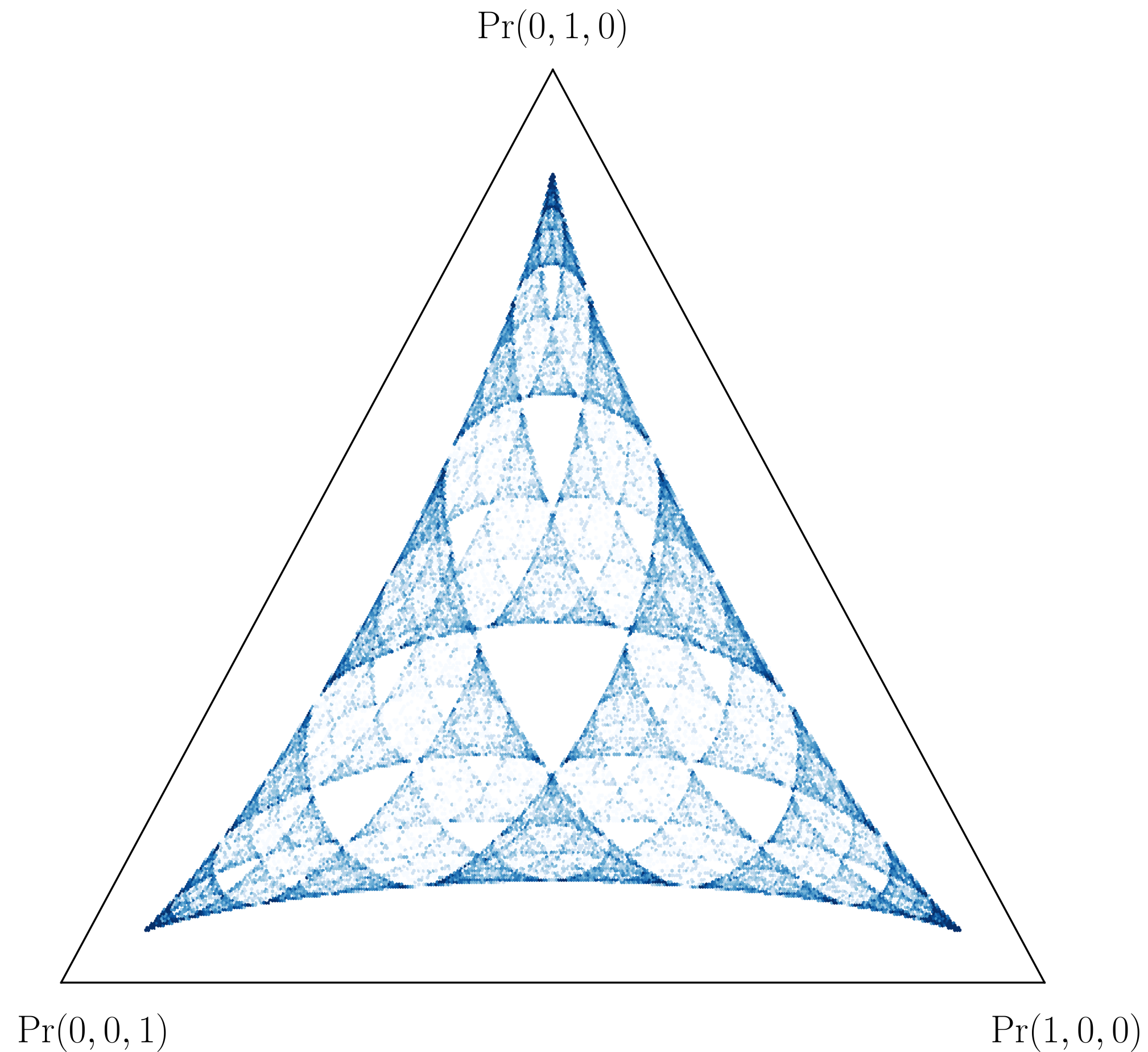
$$h_{\mu}^B = \int_R d\mu(\eta) H[X | \eta]$$

$$\widehat{h}_{\mu}^B = \lim_{L \rightarrow \infty} \frac{1}{L} \sum_{t=0}^L H[X_t | \eta_t]$$

Alex M. Jurgens, James P. Crutchfield. *Shannon Entropy Rate of Hidden Markov Processes*. J. Stat. Phys., 183 (2), 1-18, 2020.

D. Blackwell. *The entropy of functions of finite-state Markov chains*. 1957.

Structure: Statistical Complexity?

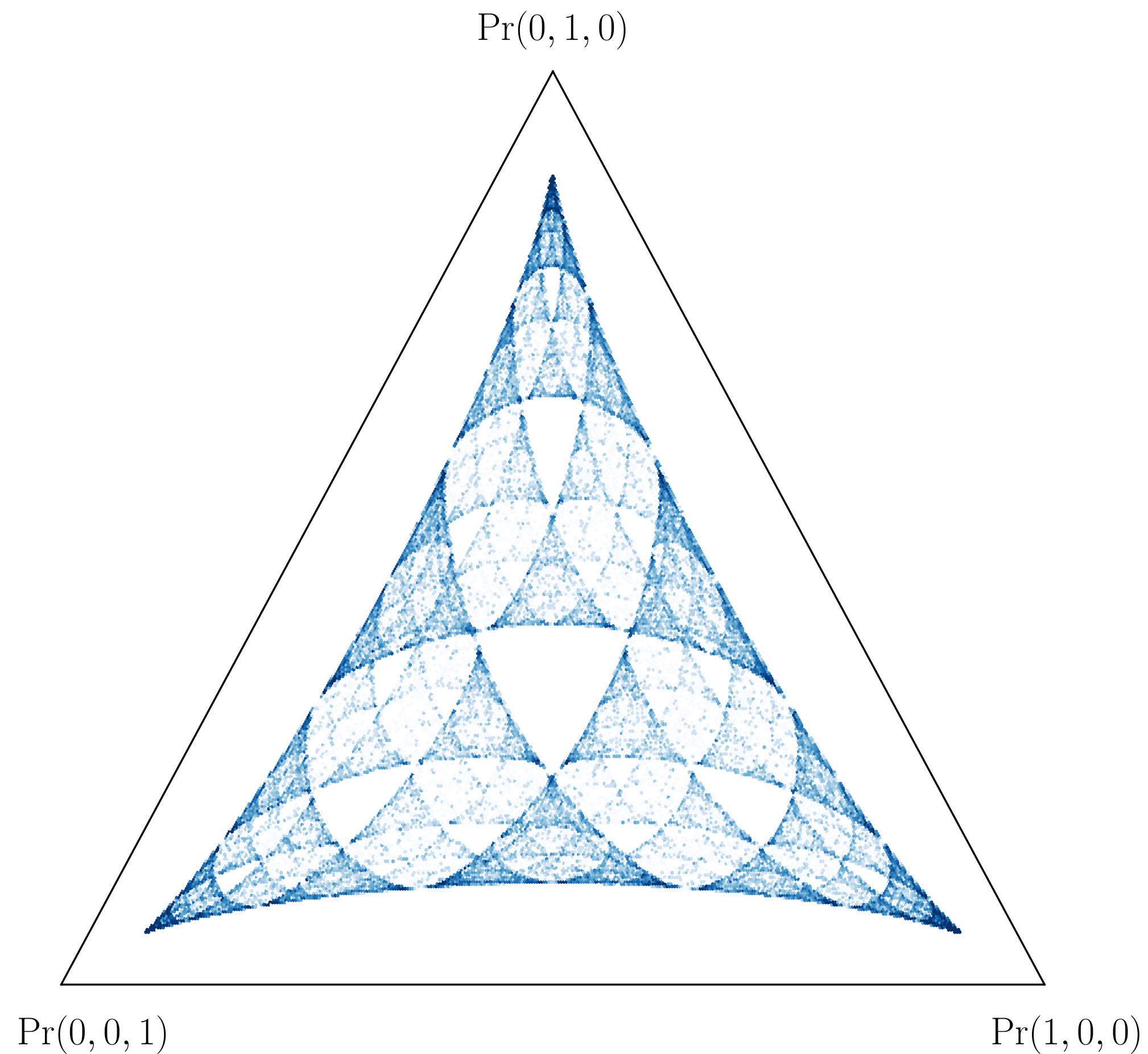


$$C_{\mu} \rightarrow \infty$$

Alexandra M. Jurgens, James P. Crutchfield. *Divergent Predictive Memory: The Statistical Complexity Dimension of Stationary, Ergodic Finite-State Hidden Markov Processes*. *Chaos* 31, 083114, 2021.

Alexandra M. Jurgens, James P. Crutchfield. *Ambiguity rate of hidden Markov processes*. *Phys. Rev. E*, 104 (2021)

Structure: Information Dimension



$$\dim_{\mu}(\mathbb{R}) \sim \frac{\Delta H [R_{\epsilon}]}{\Delta \ln \epsilon}$$

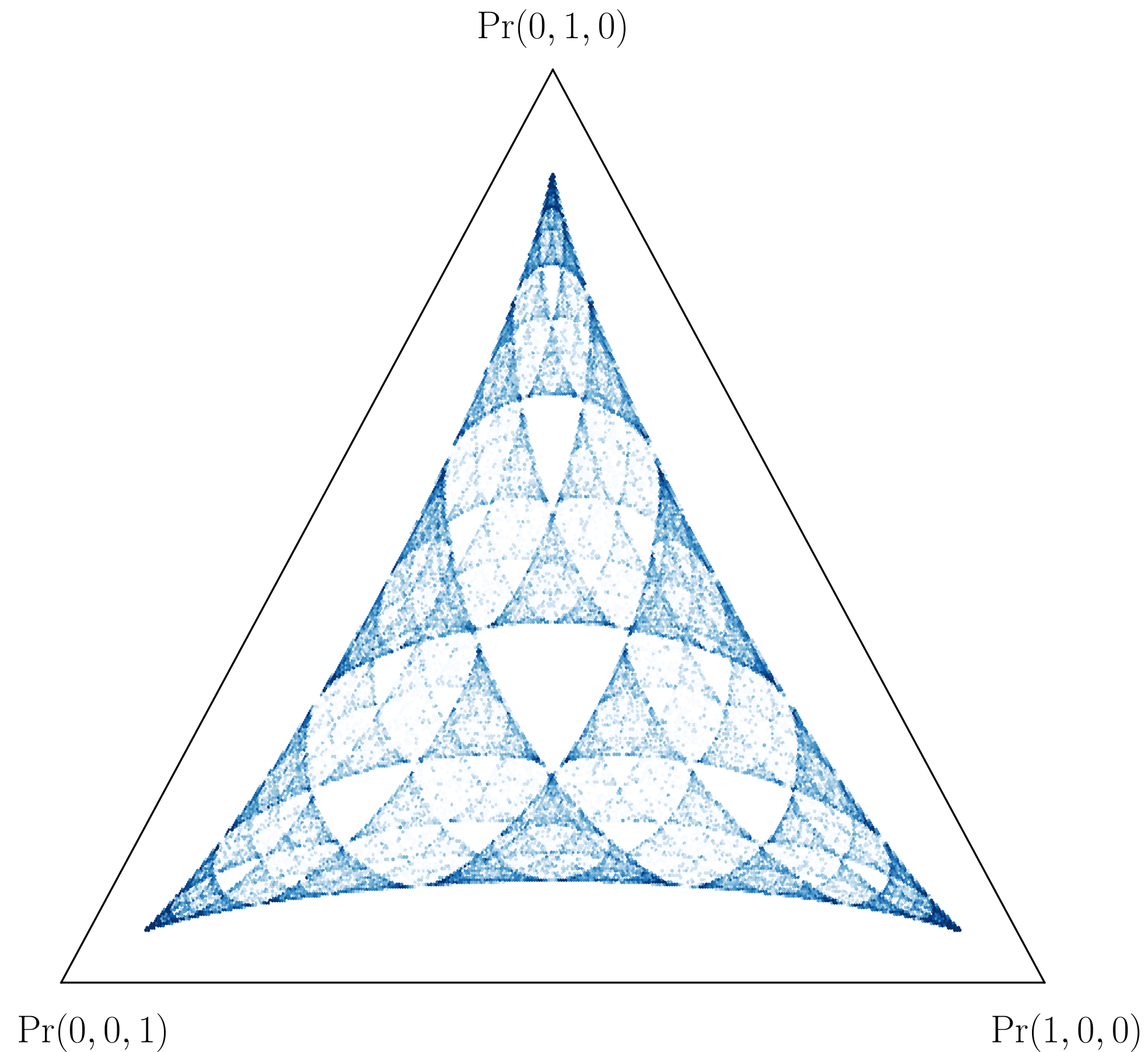
Information in a set

Scale of measurement

Alexandra M. Jurgens, James P. Crutchfield. *Divergent Predictive Memory: The Statistical Complexity Dimension of Stationary, Ergodic Finite-State Hidden Markov Processes*. Chaos 31, 083114, 2021.

Alexandra M. Jurgens, James P. Crutchfield. *Ambiguity rate of hidden Markov processes*. Phys. Rev. E, 104 (2021)

Statistical Complexity Dimension?



$$C_{\mu} \rightarrow \infty$$

$$\dim_{\mu}(R) \sim \frac{\Delta H[R_{\epsilon}]}{\Delta \ln \epsilon} = \frac{\Delta C_{\mu, \epsilon}}{\Delta \ln \epsilon}$$

Alexandra M. Jurgens, James P. Crutchfield. *Divergent Predictive Memory: The Statistical Complexity Dimension of Stationary, Ergodic Finite-State Hidden Markov Processes*. *Chaos* 31, 083114, 2021.

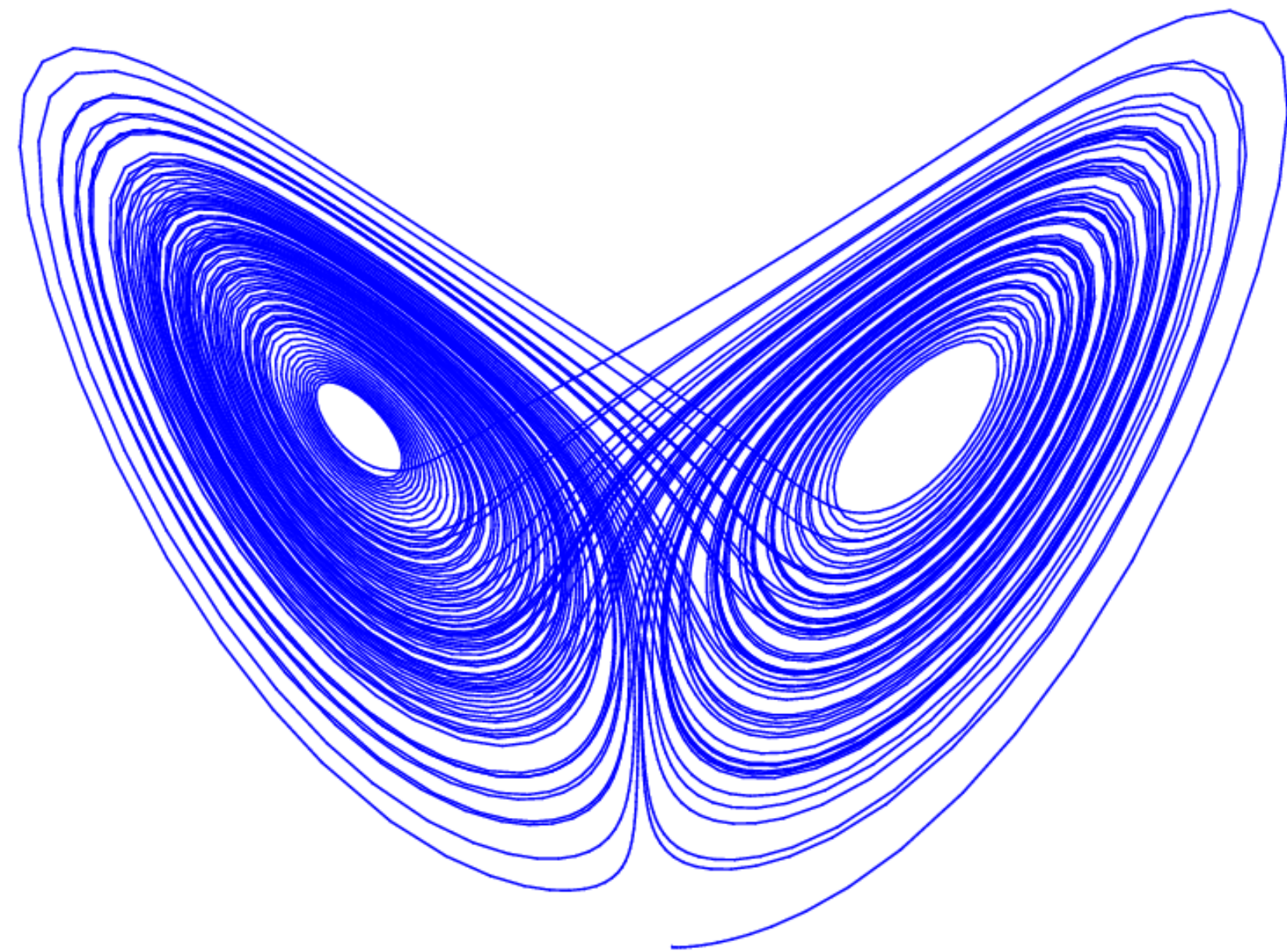
Alexandra M. Jurgens, James P. Crutchfield. *Ambiguity rate of hidden Markov processes*. *Phys. Rev. E*, 104 (2021)

Calculating Information Dimension

Calculate the Lyapunov spectrum:

$$\Gamma = \{ \lambda_1, \lambda_2, \dots, \lambda_N \}$$

$$\text{s.t. } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$$



Kaplan–Yorke conjecture:

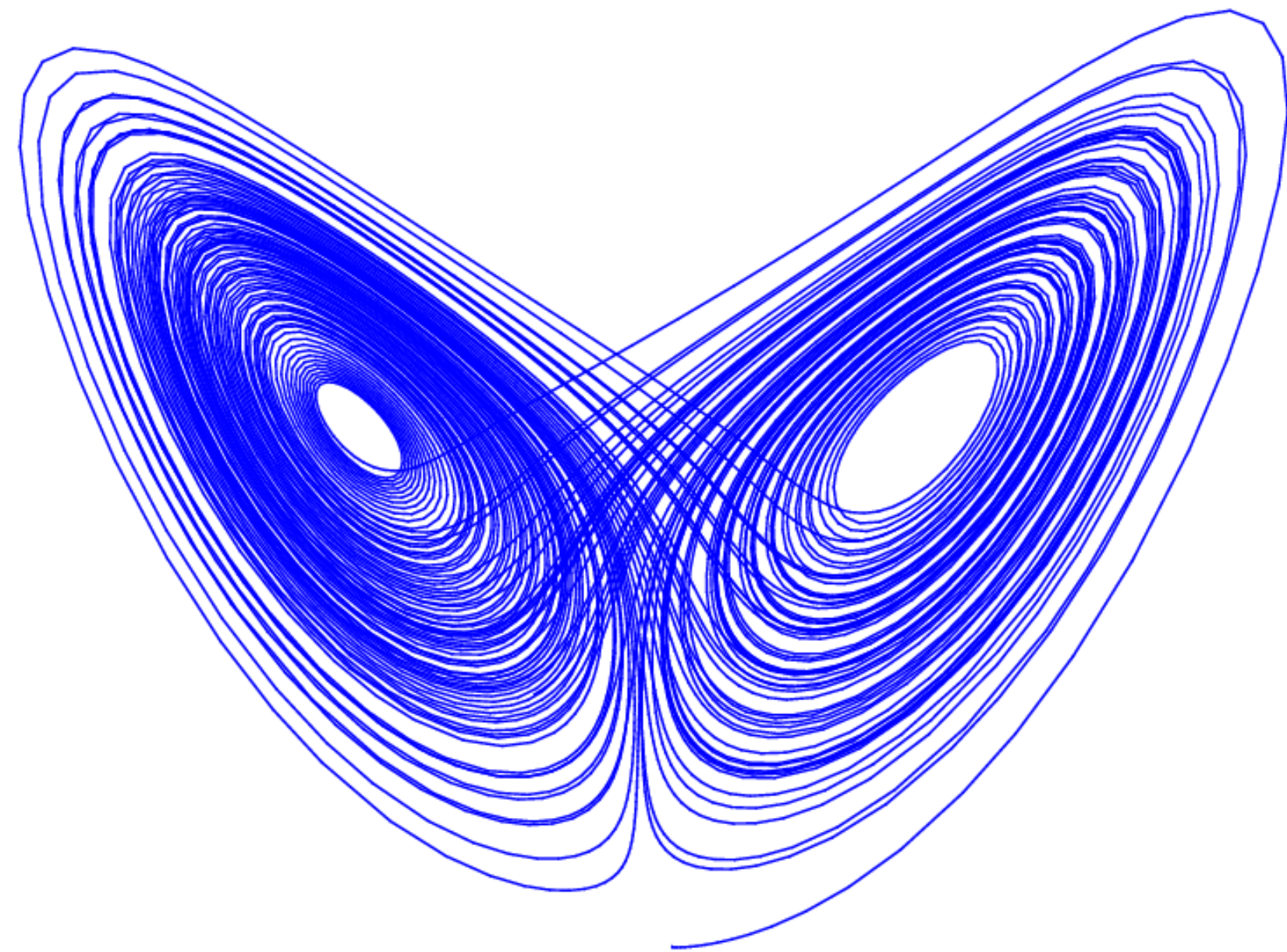
$$\dim_I = k + \frac{\sum_i^k \lambda_i}{|\lambda_{k+1}|}$$

Where k is the largest index for which the sum

$$\sum_i^k \lambda_i \text{ is positive.}$$

Calculating Information Dimension

For the Lorenz attractor with $\sigma = 10$, $r = 28$,
 $b = 8/3$:



$$\Gamma = \{0.90563, 0, -14.57219\}$$

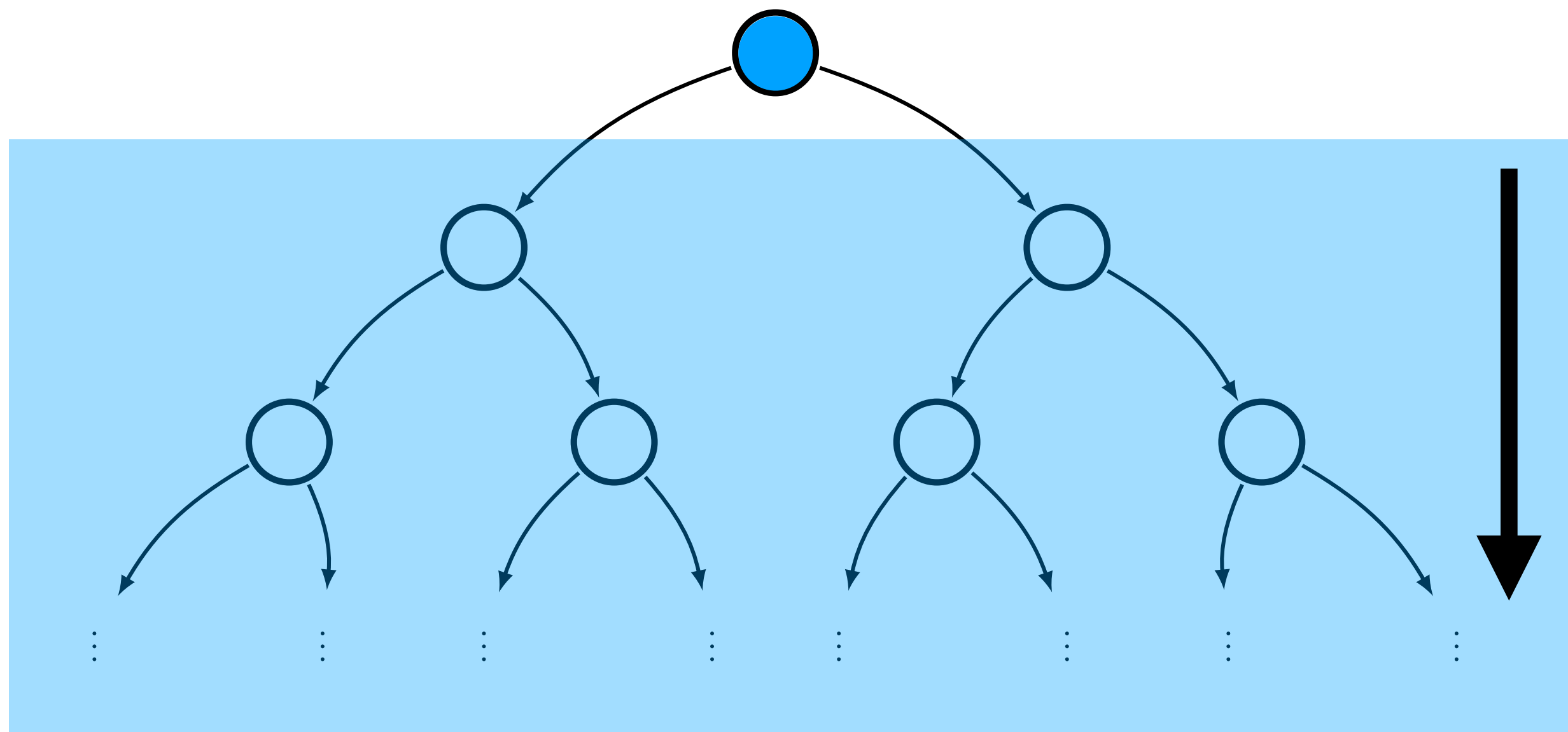
So the Lyapunov dimension is:

$$\dim_I = 2.06215$$

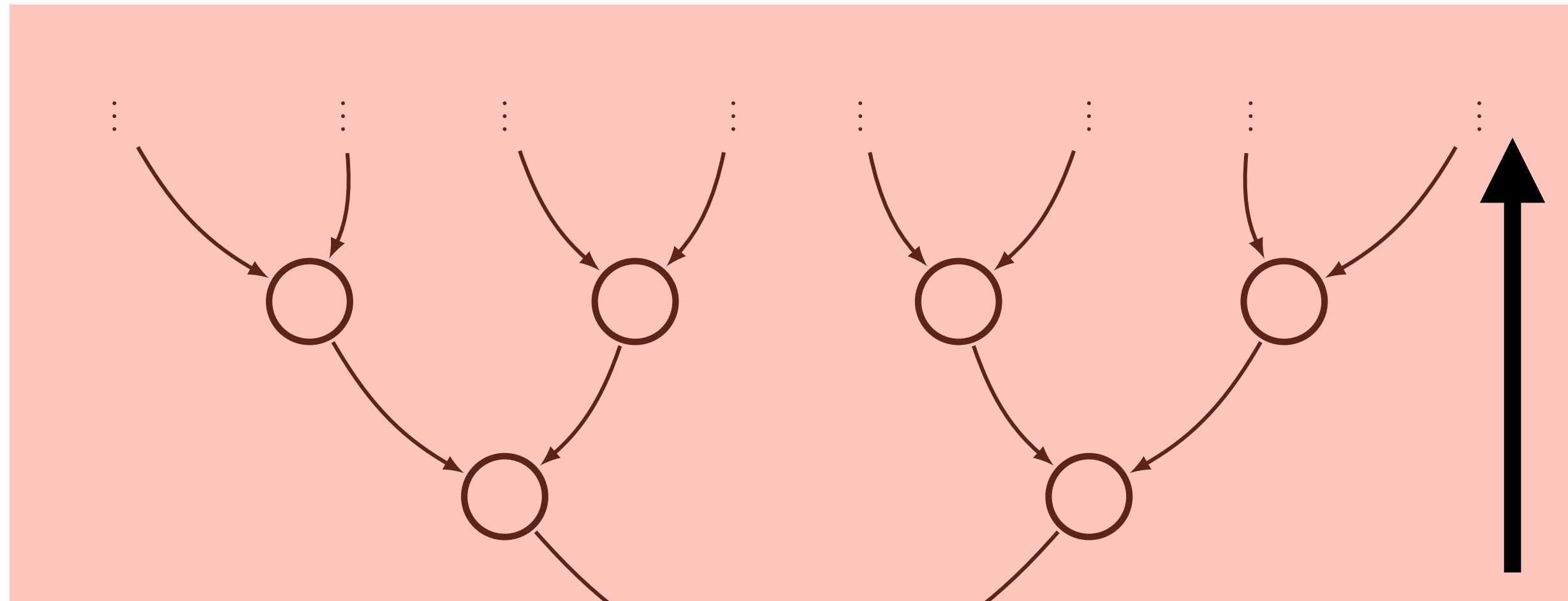
New Quantity: The Ambiguity Rate

Entropy rate measures uncertainty in the next symbol given the present.

$$h_{\mu} = H [X_0, S_1 | S_0]$$



New Quantity: The Ambiguity Rate

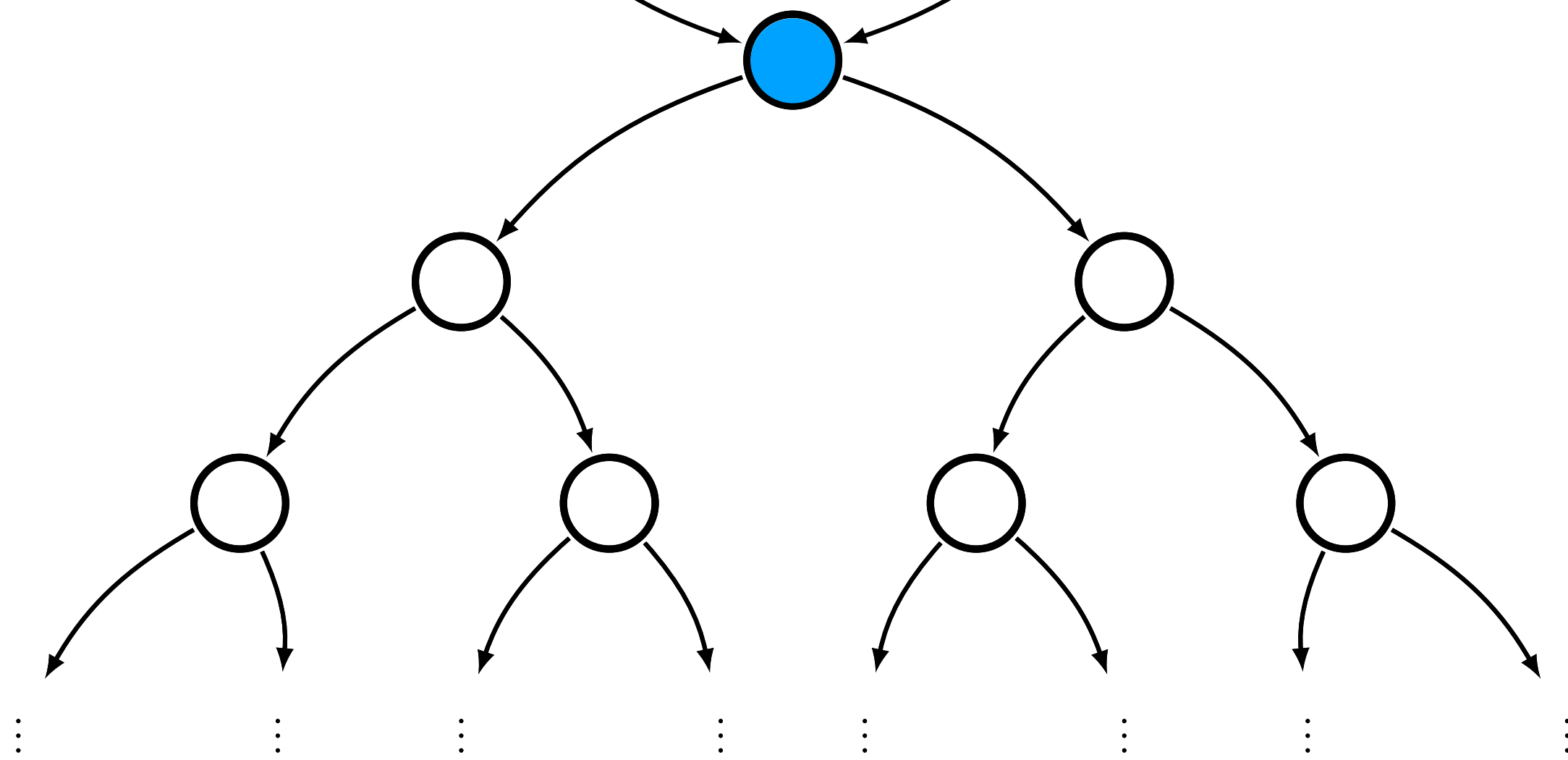


Entropy rate measures uncertainty in the next symbol given the present.

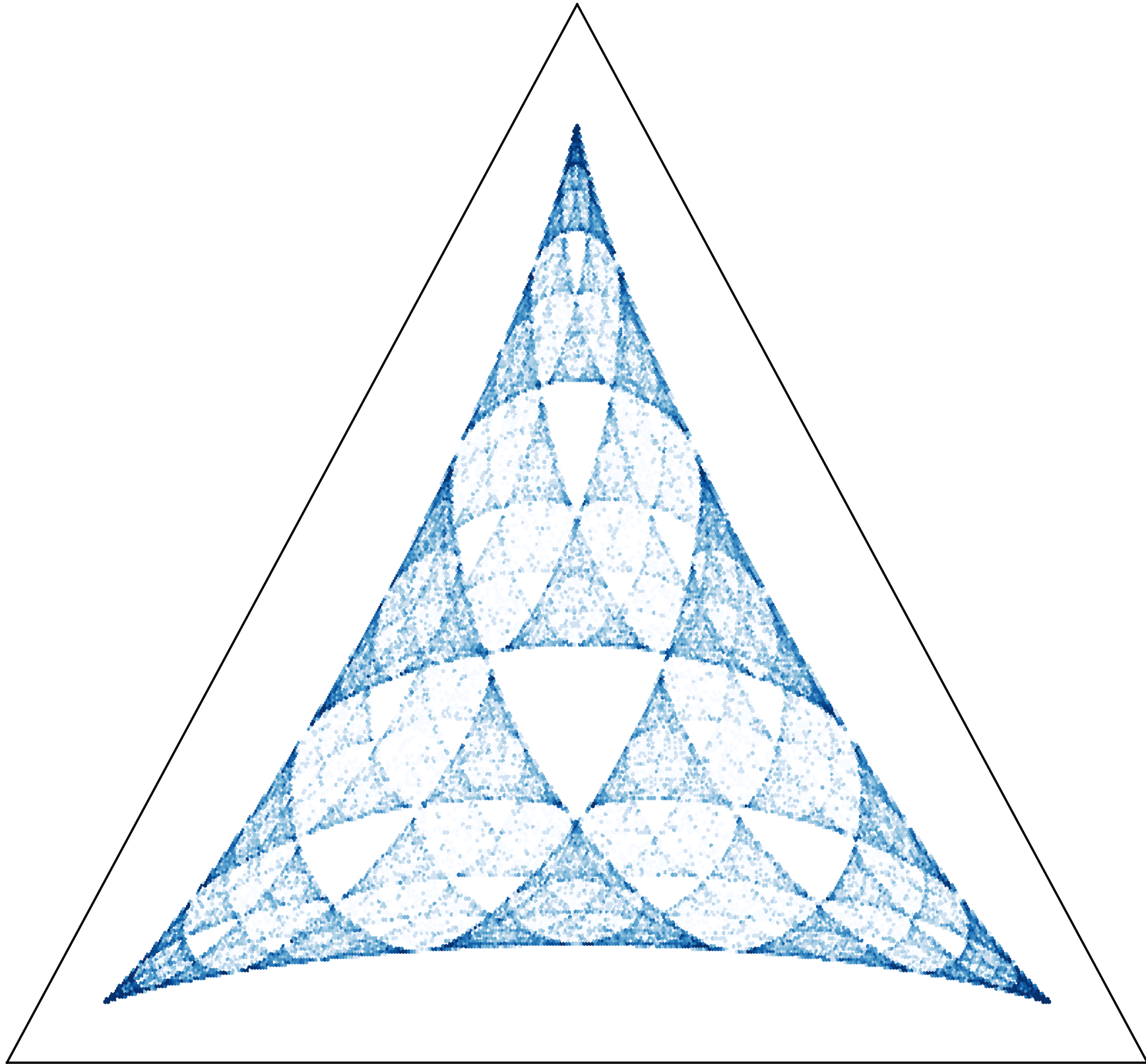
$$h_{\mu} = H [X_0, S_1 | S_0]$$

Ambiguity rate measures uncertainty *in prior* symbol given the present.

$$h_a = H [X_{-1}, S_{-1} | S_0]$$



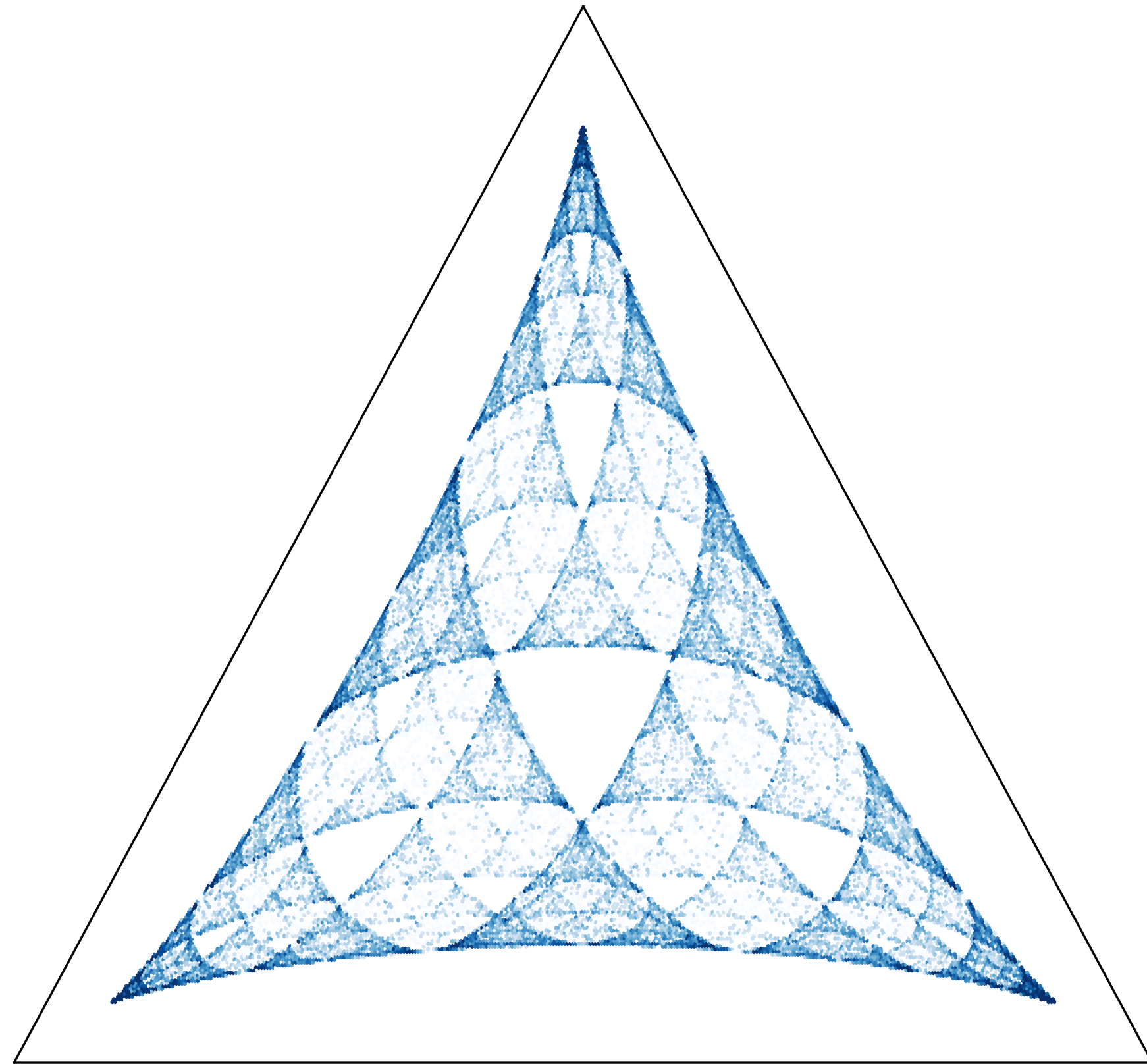
Statistical Complexity Dimension



$$\dim_{\mu}(R) = k + \frac{h_{\mu} - h_a + \sum_i^k \lambda_i}{|\lambda_{k+1}|}$$

Alexandra M. Jurgens, James P. Crutchfield. *Divergent Predictive Memory: The Statistical Complexity Dimension of Stationary, Ergodic Finite-State Hidden Markov Processes*. *Chaos* 31, 083114, 2021.
Alexandra M. Jurgens, James P. Crutchfield. *Ambiguity rate of hidden Markov processes*. *Phys. Rev. E*, 104 (2021)

Statistical Complexity Dimension

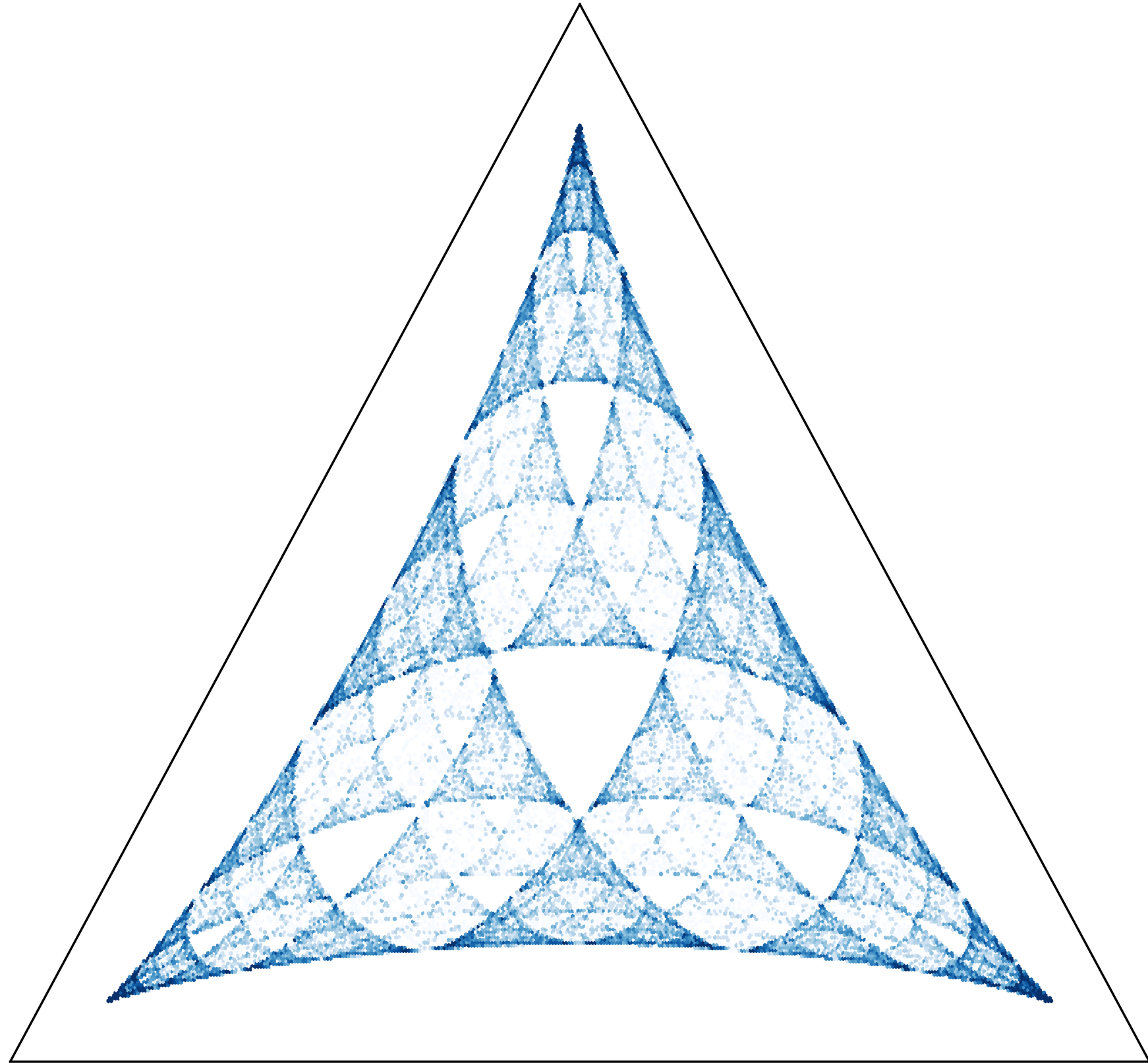


Entropy rate →

Ambiguity rate →

$$\dim_{\mu}(R) = k + \frac{h_{\mu} - h_a + \sum_i^k \lambda_i}{|\lambda_{k+1}|}$$

Statistical Complexity Dimension



$$C_{\mu} \rightarrow \infty$$

$$\dim_{\mu}(R) = \frac{\Delta C_{\mu, \epsilon}}{\Delta \ln \epsilon} = k + \frac{h_{\mu} - h_a + \sum_i^k \lambda_i}{|\lambda_{k+1}|}$$

→ Can calculate randomness and structure,
now for infinite states!

Alexandra M. Jurgens, James P. Crutchfield. *Divergent Predictive Memory: The Statistical Complexity Dimension of Stationary, Ergodic Finite-State Hidden Markov Processes*. *Chaos* 31, 083114, 2021.

Alexandra M. Jurgens, James P. Crutchfield. *Ambiguity rate of hidden Markov processes*. *Phys. Rev. E*, 104 (2021)

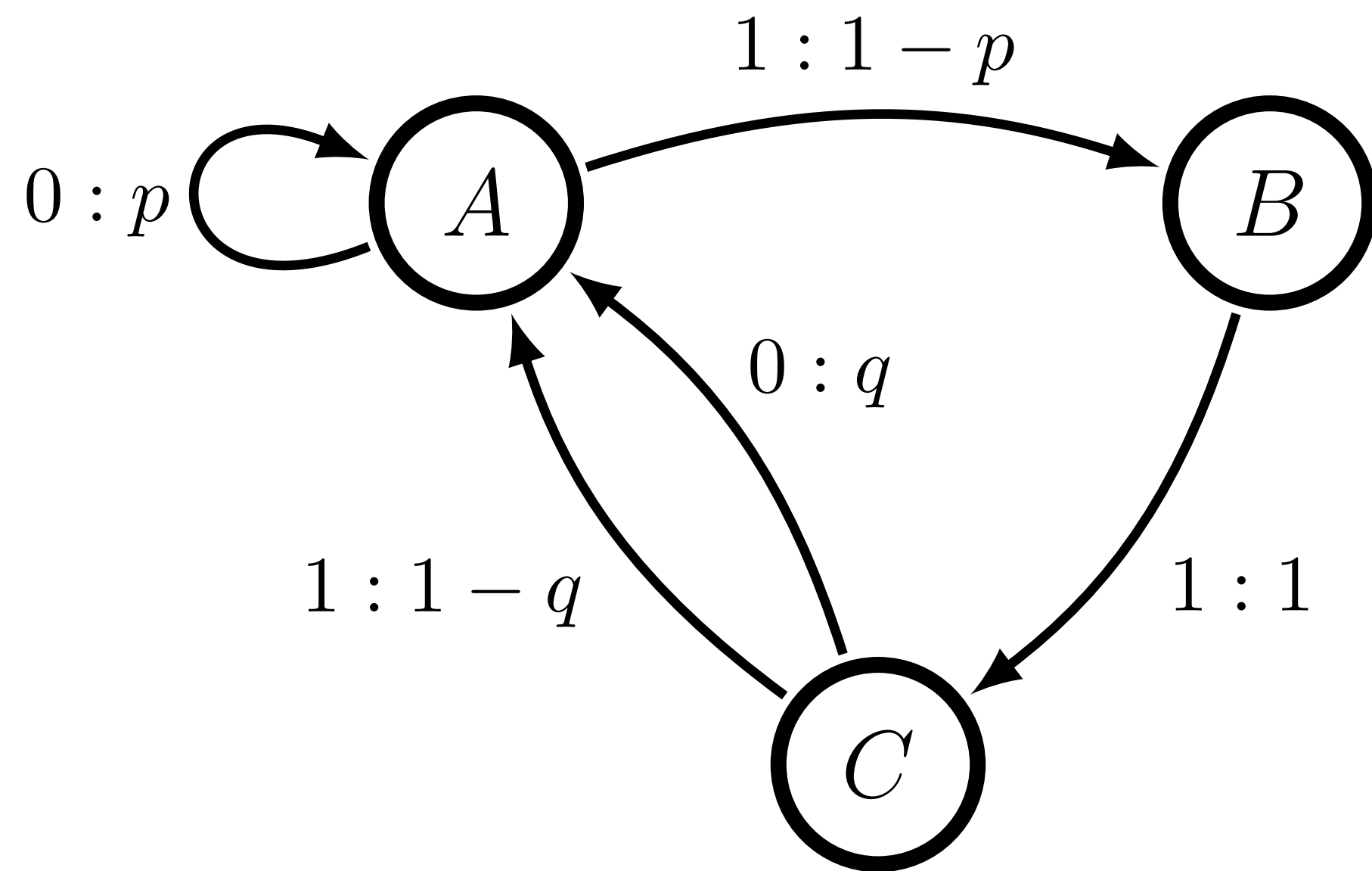
Growth Rate of Predictive Models

$$\Delta H [\text{Predictive states}] = h_{\mu} - h_a$$

Entropy rate

Ambiguity rate

Growth Rate of Predictive Models



For a finite state model:

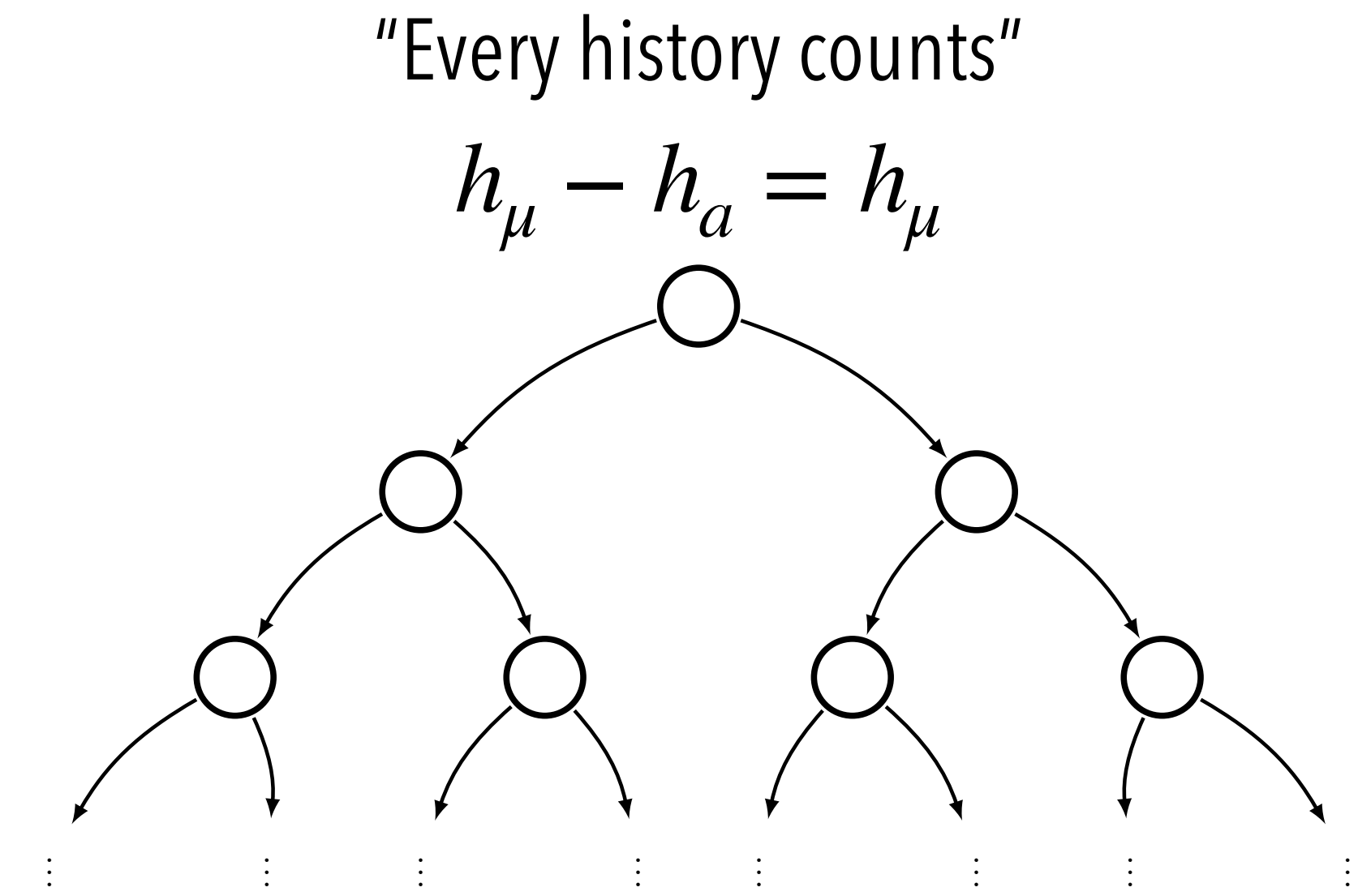
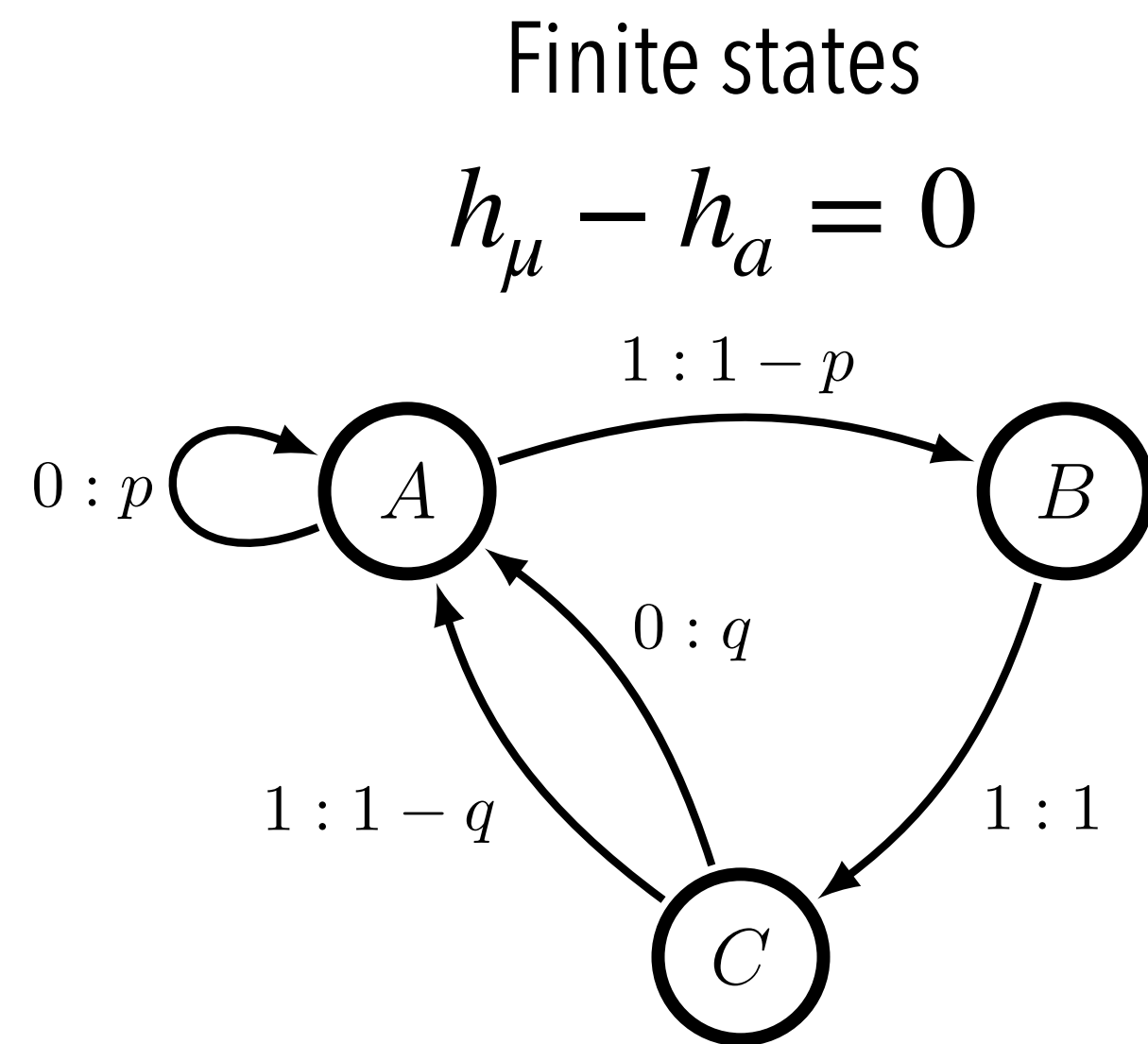
$$h_{\mu} - h_a = 0$$

Alexandra M. Jurgens, James P. Crutchfield. *Divergent Predictive Memory: The Statistical Complexity Dimension of Stationary, Ergodic Finite-State Hidden Markov Processes*. Chaos 31, 083114, 2021.

Alexandra M. Jurgens, James P. Crutchfield. *Ambiguity rate of hidden Markov processes*. Phys. Rev. E, 104 (2021)

Growth Rate of Predictive Models

$$\Delta H [\text{Predictive states}] = h_\mu - h_a$$



Alexandra M. Jurgens, James P. Crutchfield. *Divergent Predictive Memory: The Statistical Complexity Dimension of Stationary, Ergodic Finite-State Hidden Markov Processes*. Chaos 31, 083114, 2021.

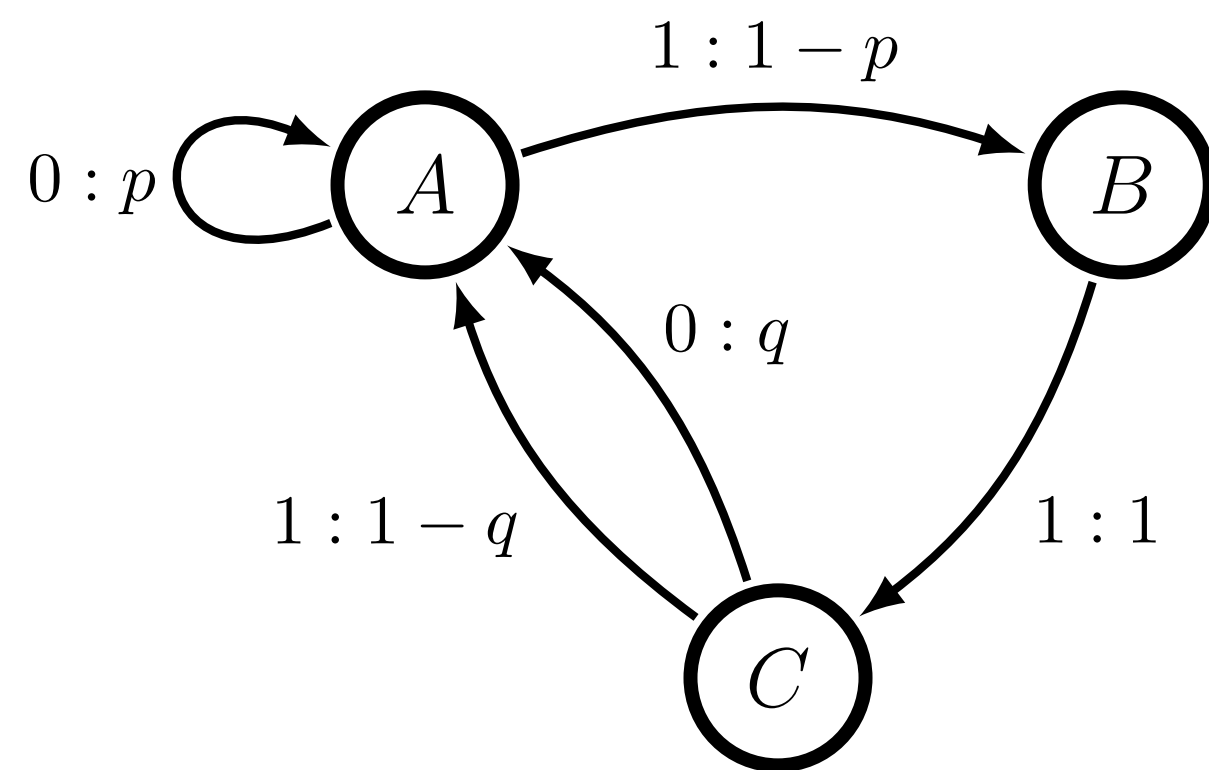
Alexandra M. Jurgens, James P. Crutchfield. *Ambiguity rate of hidden Markov processes*. Phys. Rev. E, 104 (2021)

Growth Rate of Predictive Models

$$\Delta H [\text{Predictive states}] = h_\mu - h_a$$

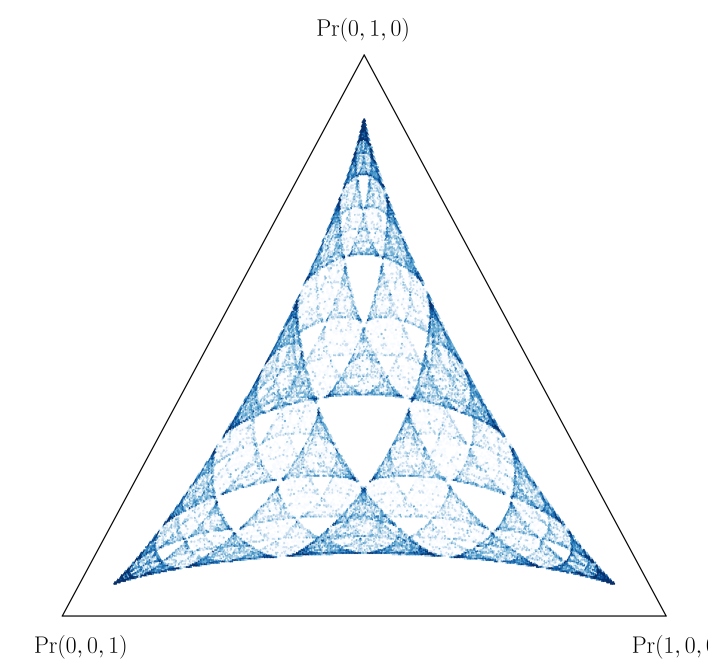
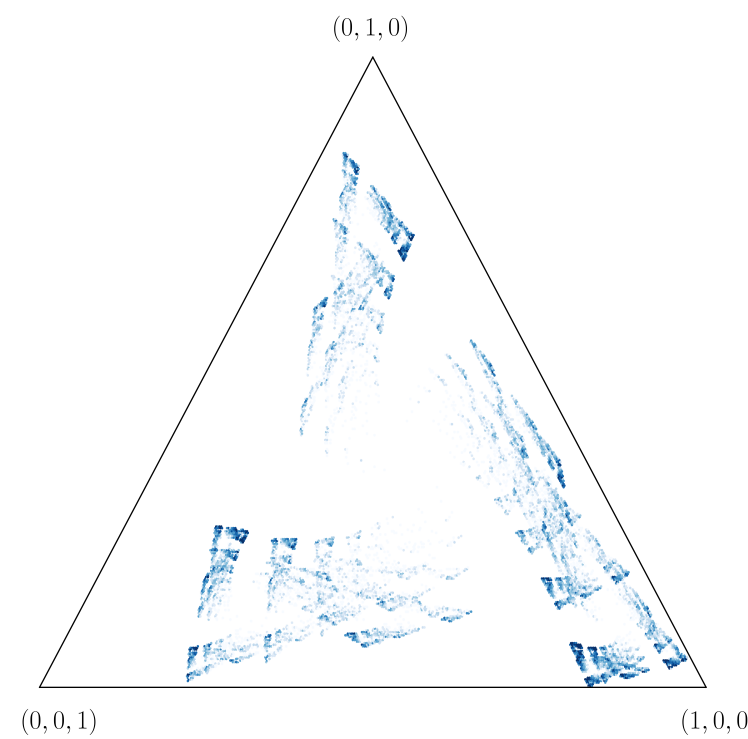
Finite states

$$h_\mu - h_a = 0$$



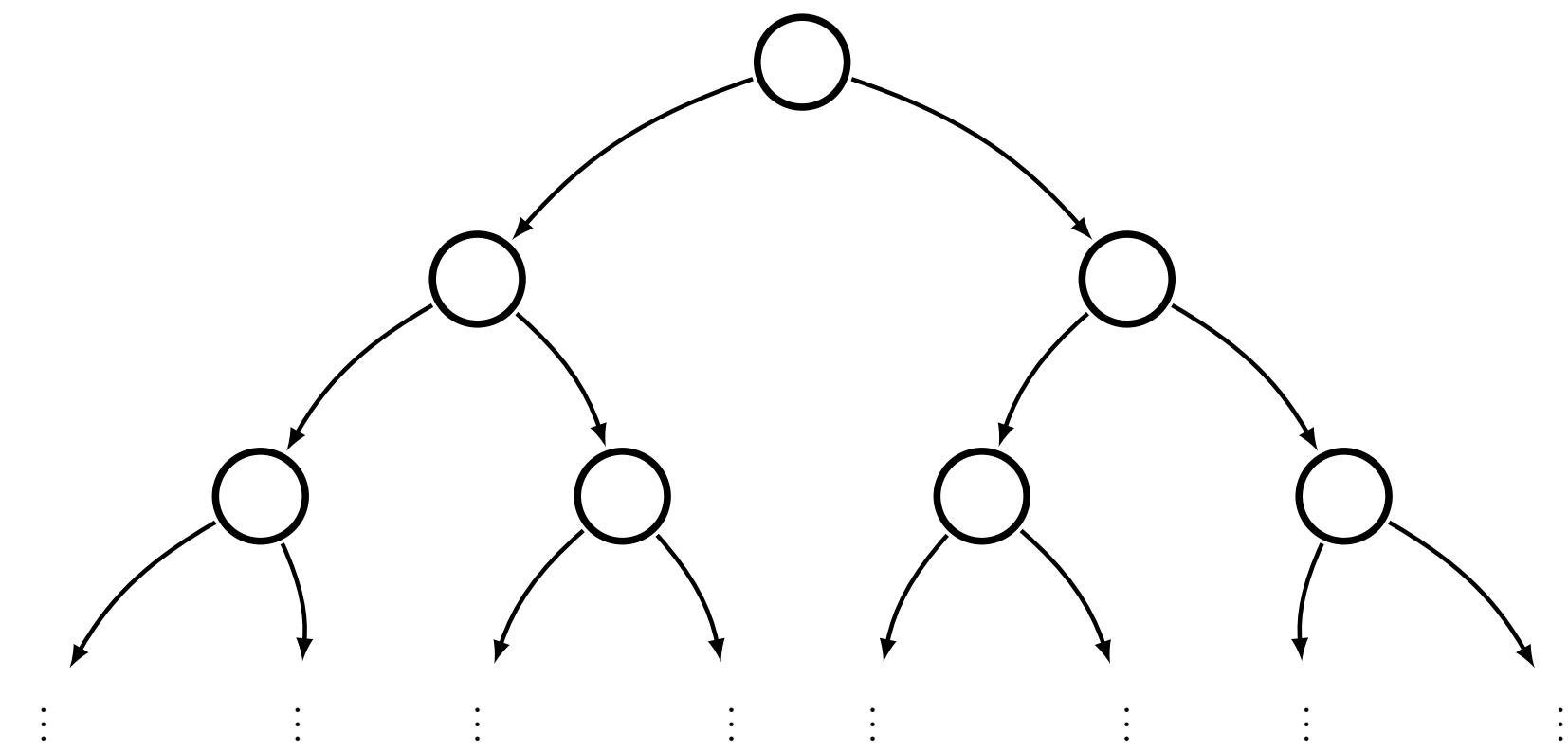
In general:

$$h_\mu > h_\mu - h_a > 0$$

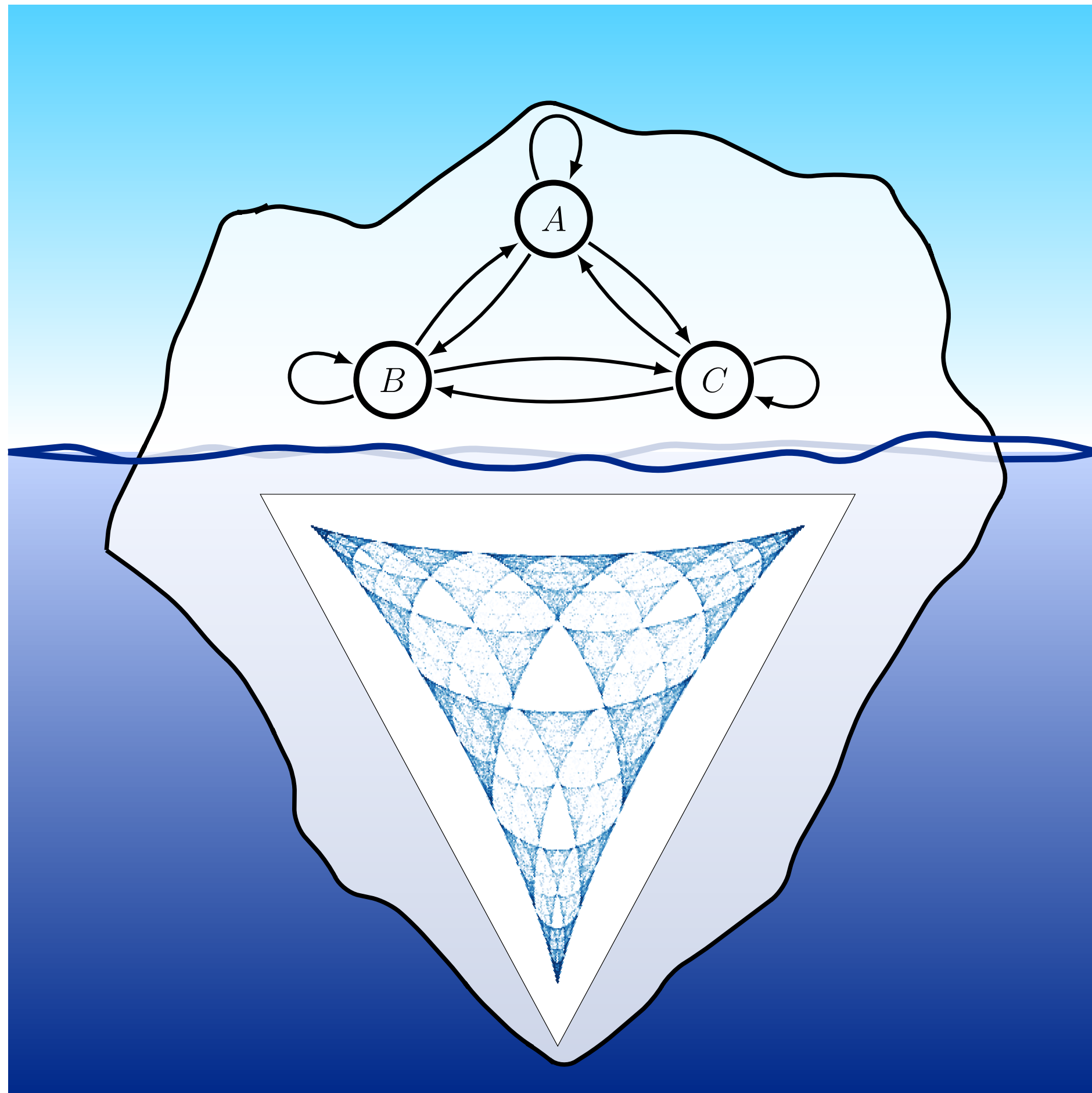


"Every history counts"

$$h_\mu - h_a = h_\mu$$



Acknowledgements + Questions



Thank you to my PhD advisor, Jim Crutchfield, and all those at the Complexity Sciences Center, especially Alec Boyd, Adam Rupe, Sam Loomis and for helpful discussions.

Thank you as well to my current group leader at INRIA Bordeaux, Nicolas Brodu.

Contact:

Website: <https://csc.ucdavis.edu/~ajurgens/>

Email: alexandra.jurgens@inria.fr

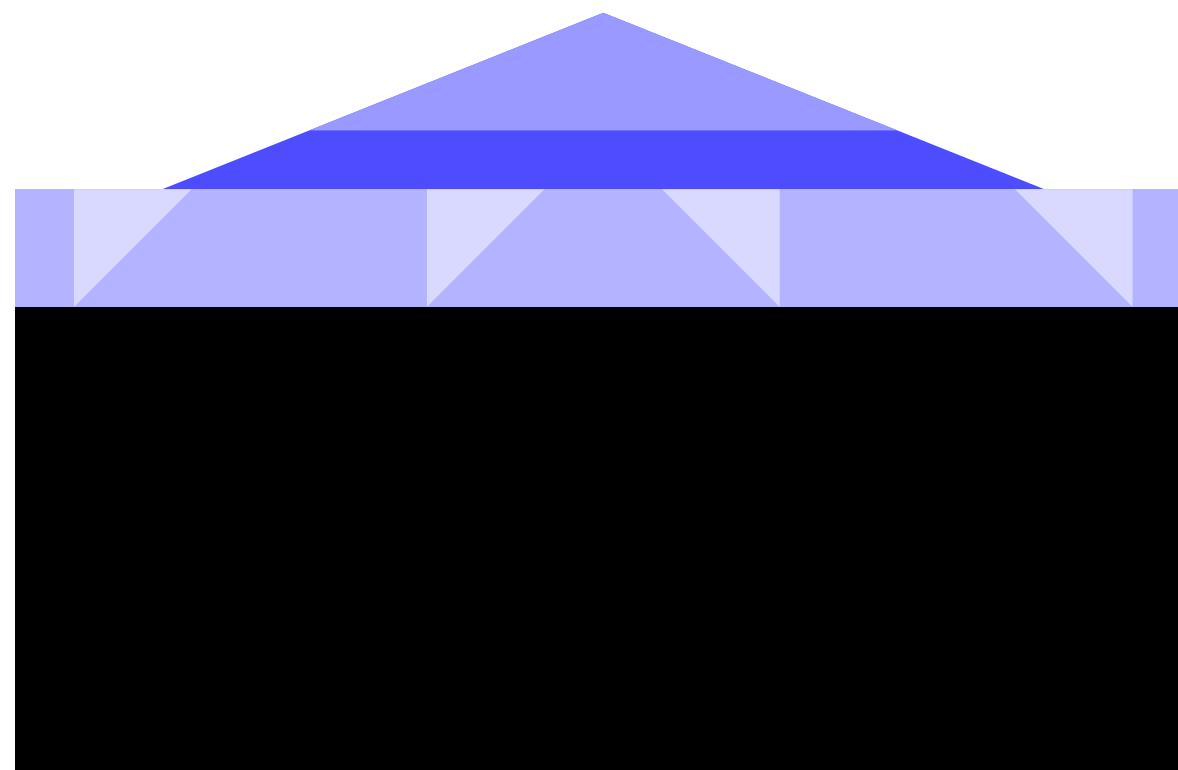
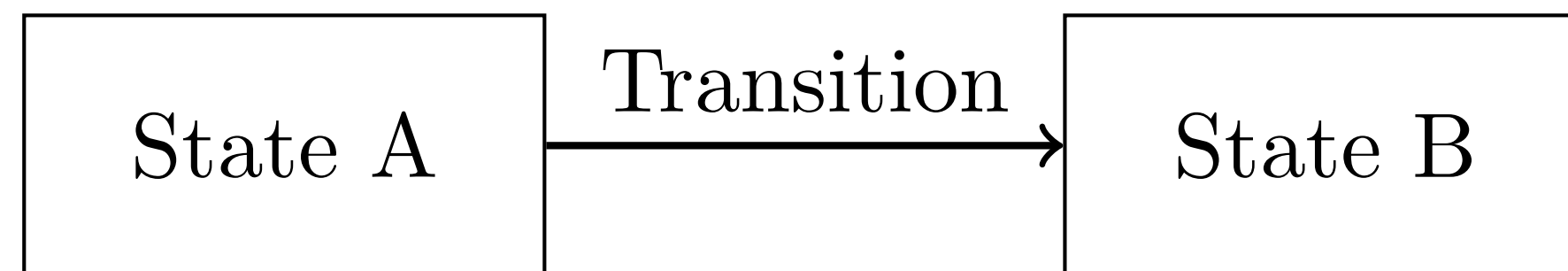
Twitter: @amjurgens



Acknowledgements + Questions

Create tikz code depicting an iceberg, but overlay a finite state machine on the above water portion and a black rectangle on the underwater portion. Include code to replace the rectangle by an inserted image.

Certainly! Below is the TikZ code to depict an iceberg with a finite state machine (FSM) overlay on the above water portion and a black rectangle on the underwater portion. Additionally, it includes the code to replace the black rectangle with an inserted image.



Thank you to my PhD advisor, Jim Crutchfield, and all those at the Complexity Sciences Center, especially Alec Boyd, Adam Rupe, Sam Loomis and for helpful discussions.

Thank you as well to my current group leader at INRIA Bordeaux, Nicolas Brodu.

Contact:

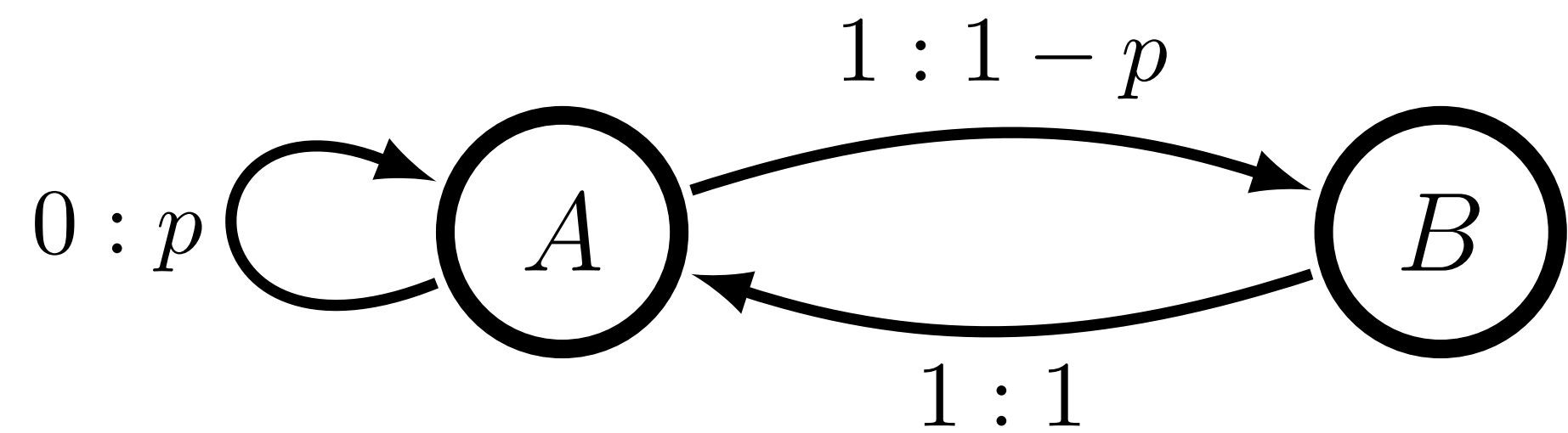
Website: <https://csc.ucdavis.edu/~ajurgens/>

Email: alexandra.jurgens@inria.fr

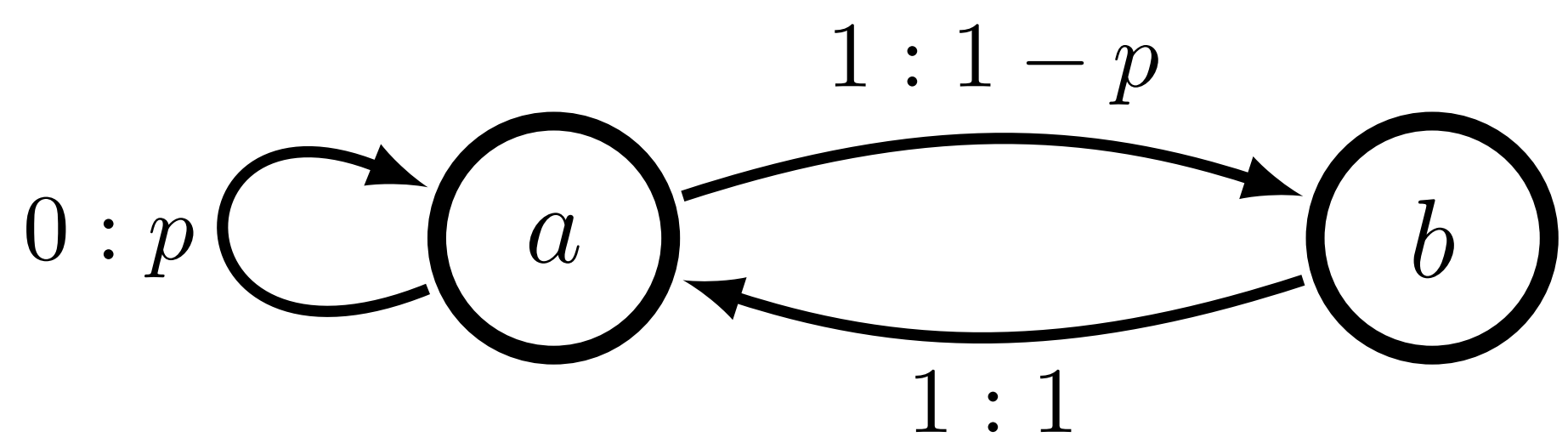
Twitter: @amjurgens



The Utility of (Good) Models



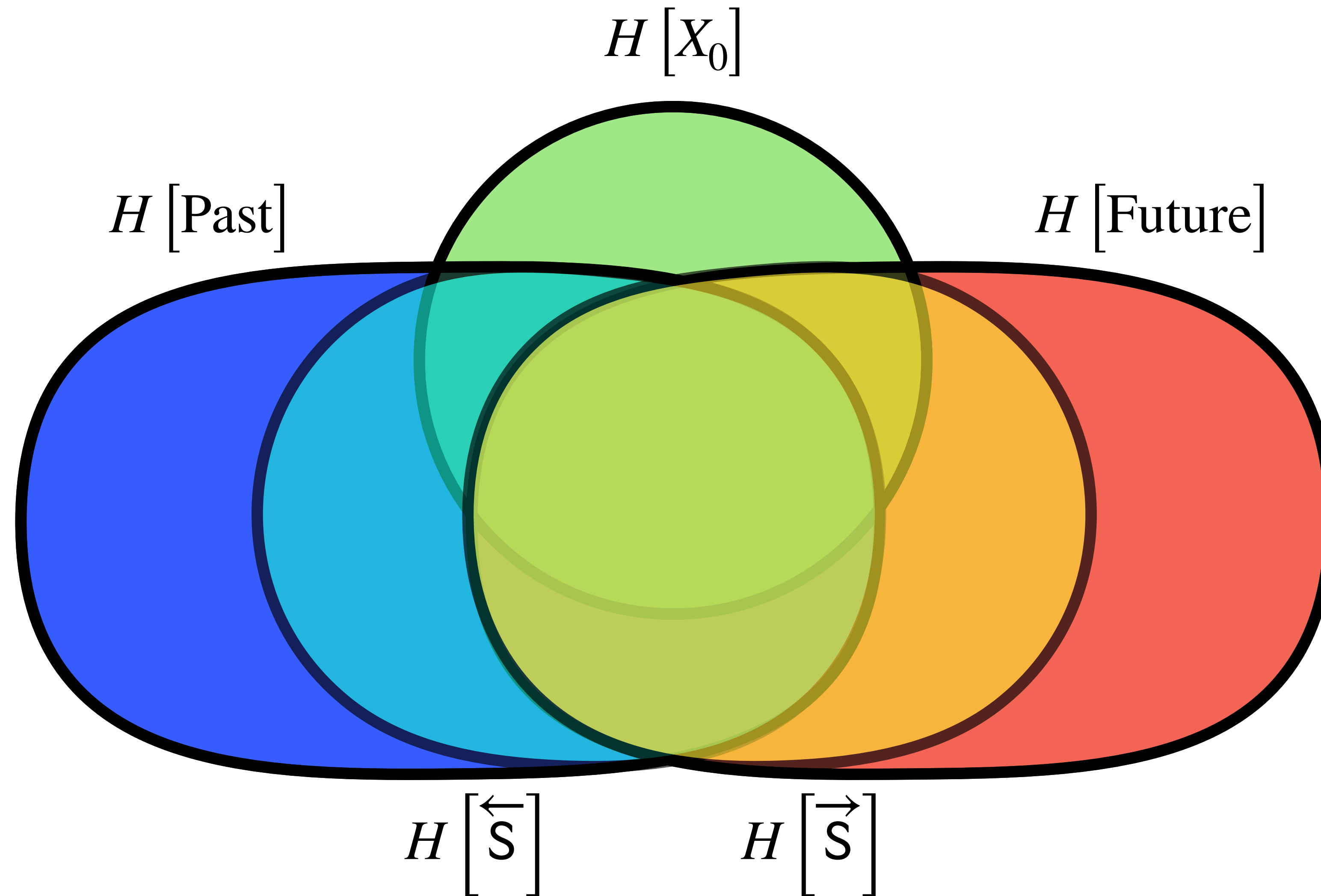
$$h_\mu = \sum_{\sigma \in \mathcal{S}} \Pr(\sigma) H[X | \sigma]$$



$$E = I \left[\begin{array}{c} \vec{\mathcal{S}} \\ \overleftarrow{\mathcal{S}} \end{array} \right]$$

$$\sigma_\mu = I \left[\begin{array}{c} \vec{\mathcal{S}} \\ \overleftarrow{\mathcal{S}} \end{array} \middle| X_0 \right]$$

The Utility of (Good) Models



Ryan G. James, Christopher J. Ellison, James P. Crutchfield. *Anatomy of a bit: Information in a time series observation*. *Chaos* 21, 037109 (2011)

James P. Crutchfield, David P. Feldman. *Regularities unseen, randomness observed: Levels of entropy convergence*. *Chaos* 1 March 2003; 13 (1): 25-54.