

# PHYS 256: POCI

## Physics of Information & Computation



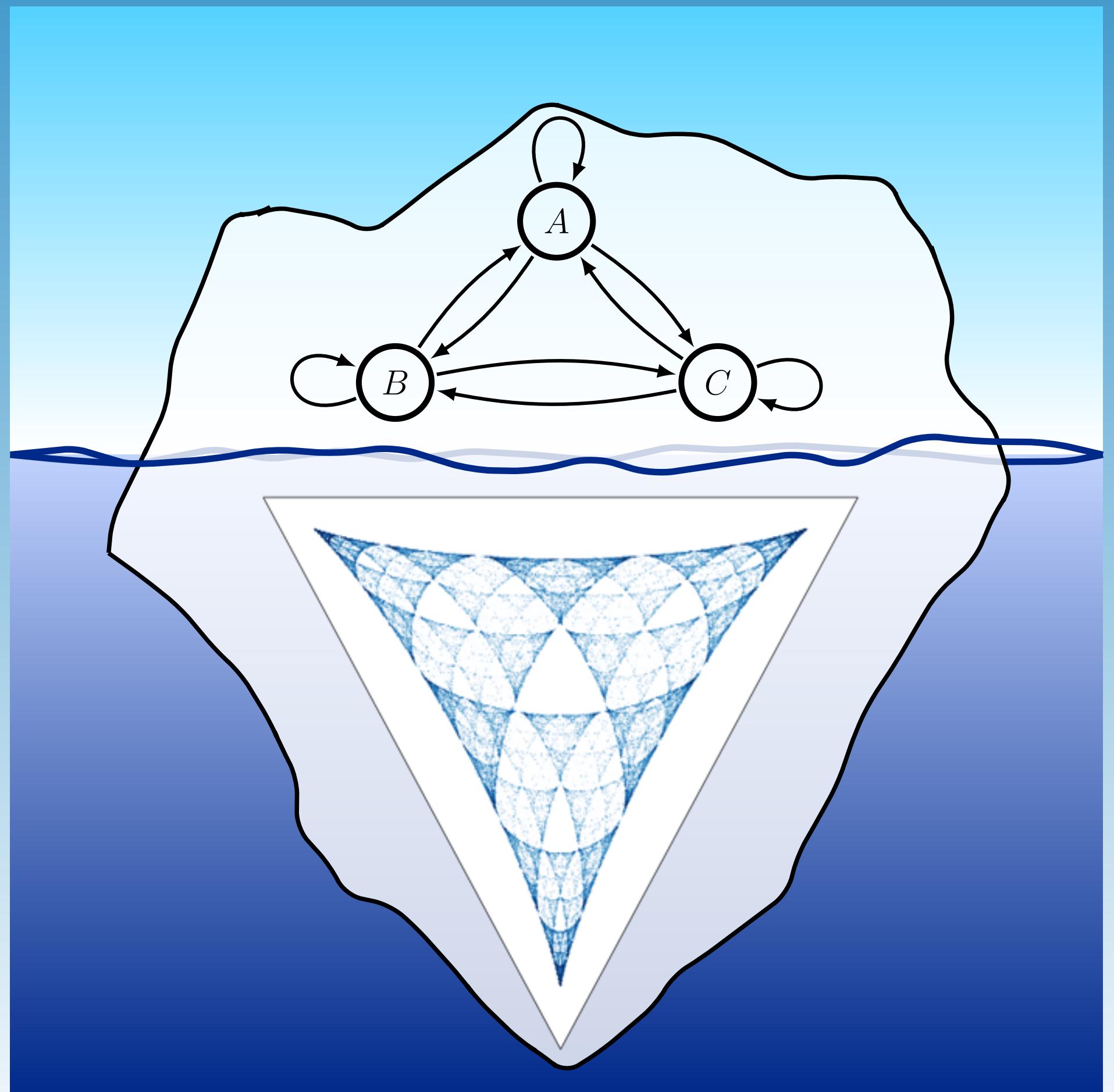
Alexandra Jurgens  
Inria Centre at the University of Bordeaux  
05/20/2025

# Infinite State Processes II

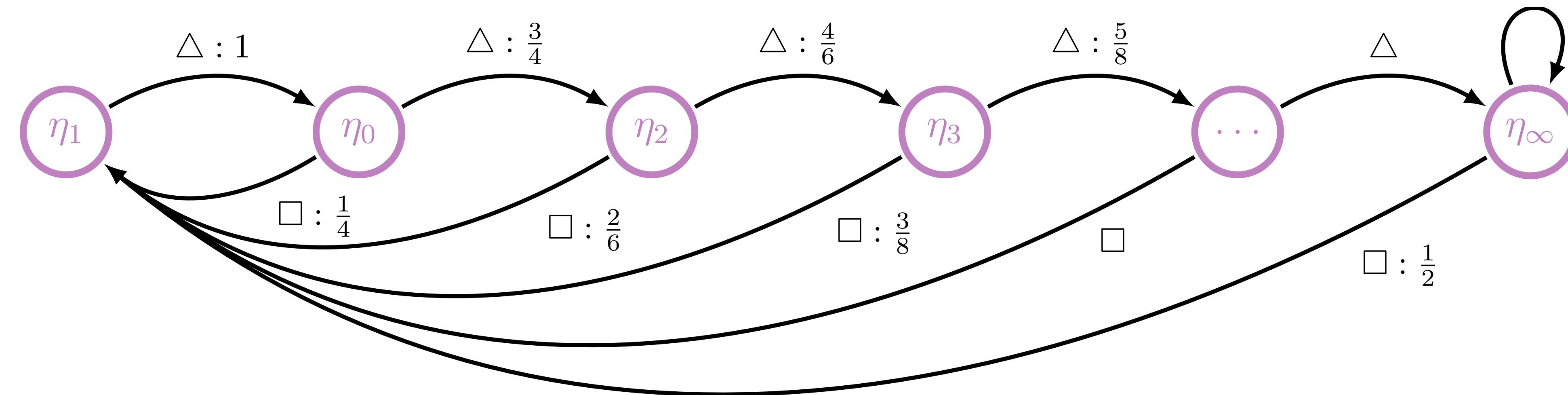
Alexandra Jurgens  
Inria Centre at the University of Bordeaux  
05/20/2025



*Inria*



# Blackwell Entropy Formula



$$h_\mu = \int_R d\mu(\eta) H[x|\eta] = \lim_{N \rightarrow \infty} - \sum_n \Pr(\eta_n) \sum_{x \in A} \Pr(x|\eta_n) \log_2 \Pr(x|\eta_n)$$

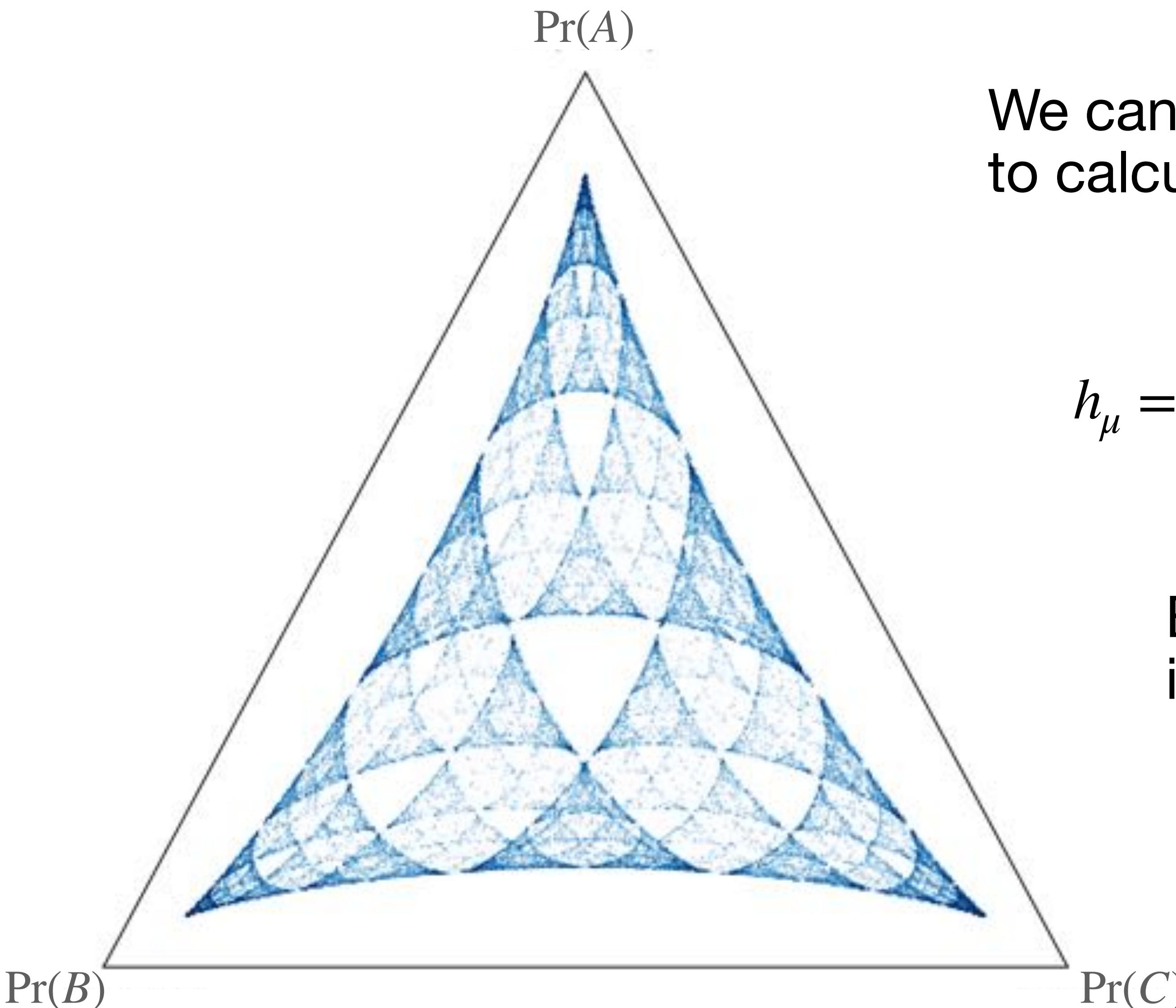
$$\rightarrow h_\mu = 0.67$$

Blackwell, D. (1957) The entropy of functions of finite-state Markov chains.

Jurgens, A. M., & Crutchfield, J. P. (2021) Shannon entropy rate of hidden Markov processes. *Journal of Statistical Physics*, 183(2), 32.

# Mixed State Sets are Fractals

---

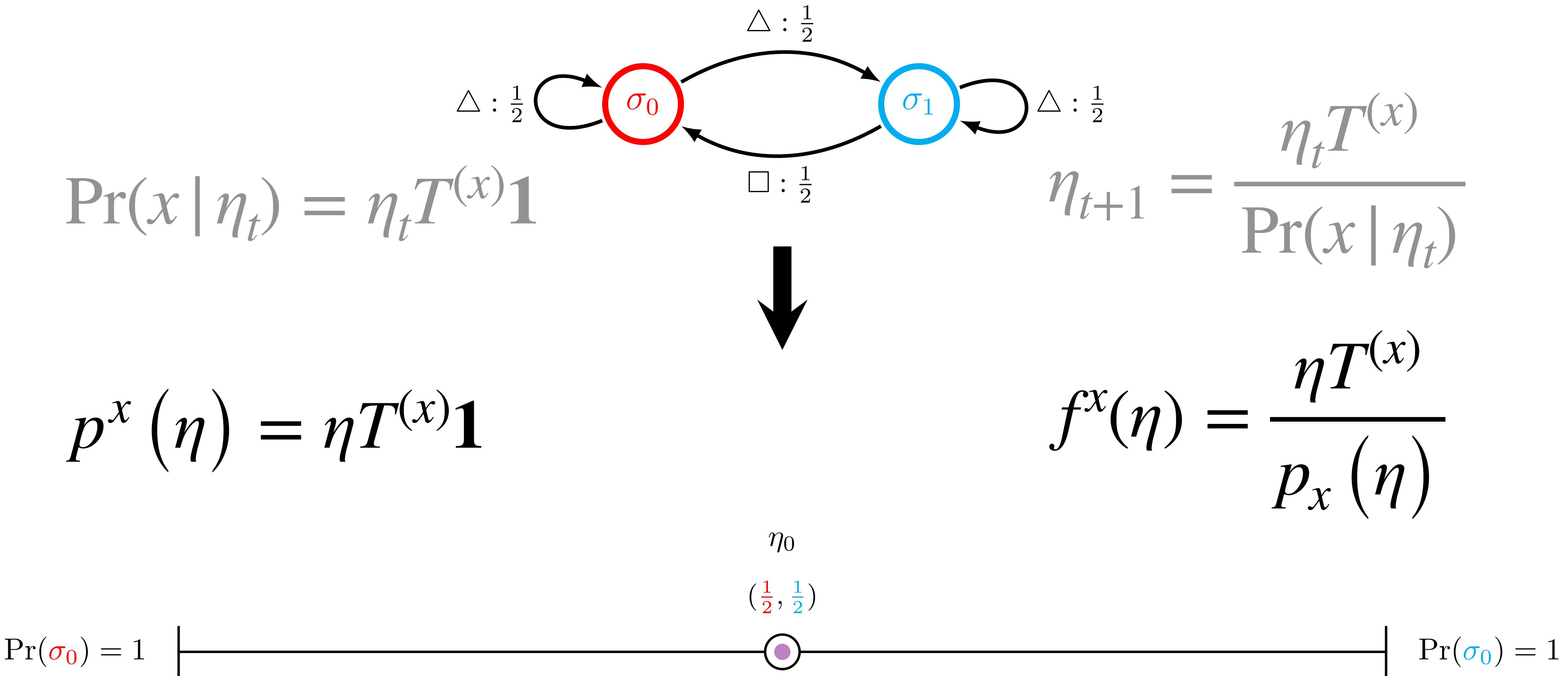


We can still use the Blackwell formula  
to calculate entropy:

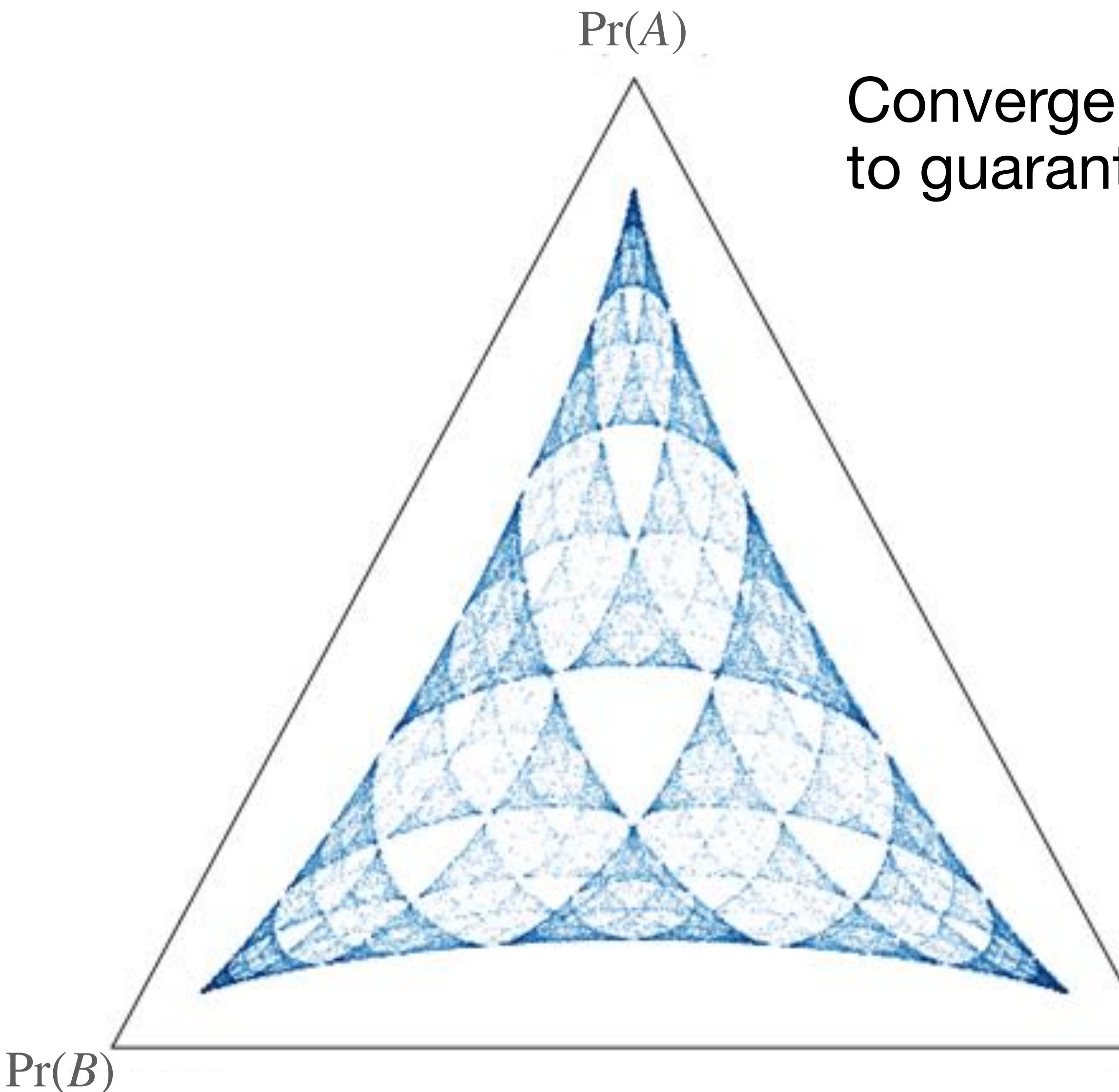
$$h_\mu = \int_R d\mu(\eta) H[x | \eta]$$

Except finding the Blackwell measure  
is much more difficult.

# Mixed State Generation → IFS



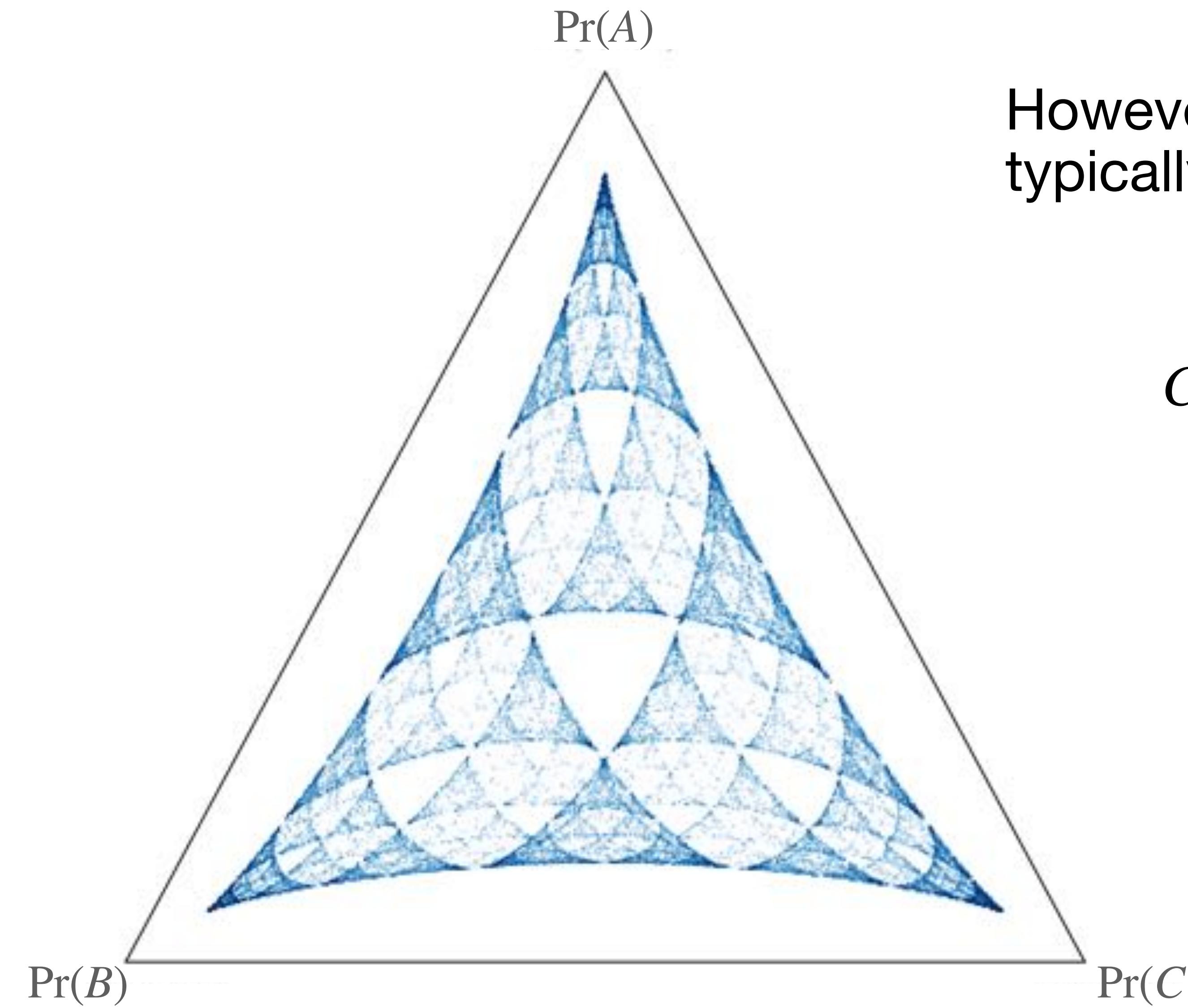
# Entropy Rate of Fractal HMMs



Convergence theorems about IFSs can now be applied to guarantee we can apply the ergodic theorem:

$$\begin{aligned} h_\mu &= \int_R d\mu(\eta) H[x | \eta] \\ &= \lim_{T \rightarrow \infty} -\frac{1}{T} \sum_t^T H[x | \eta_t] \\ &= \lim_{T \rightarrow \infty} -\frac{1}{T} \sum_t^T \sum_x p^x(\eta_t) \log_2 p^x(\eta_t) \end{aligned}$$

# Complexity of Fractal HMMs



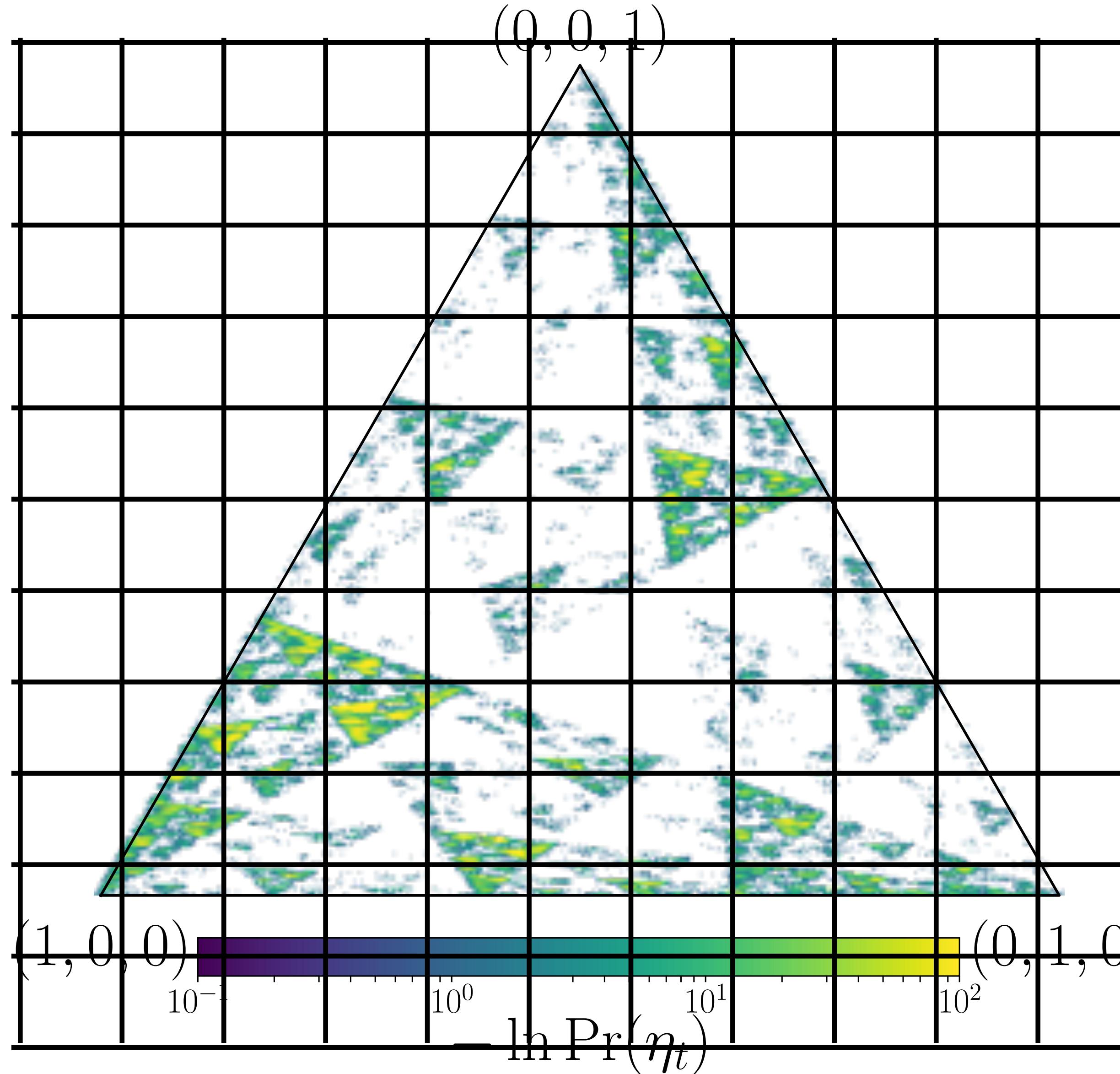
However, the statistical complexity still typically diverges.

$$C_\mu = - \int_R \Pr(\eta) \log_2 \Pr(\eta) d\mu(\eta)$$

Infinite complexity?

Jurgens, A. M., & Crutchfield, J. P. (2021). Divergent predictive states: The statistical complexity dimension of stationary, ergodic hidden Markov processes. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 31(8).

# Statistical Complexity Dimension



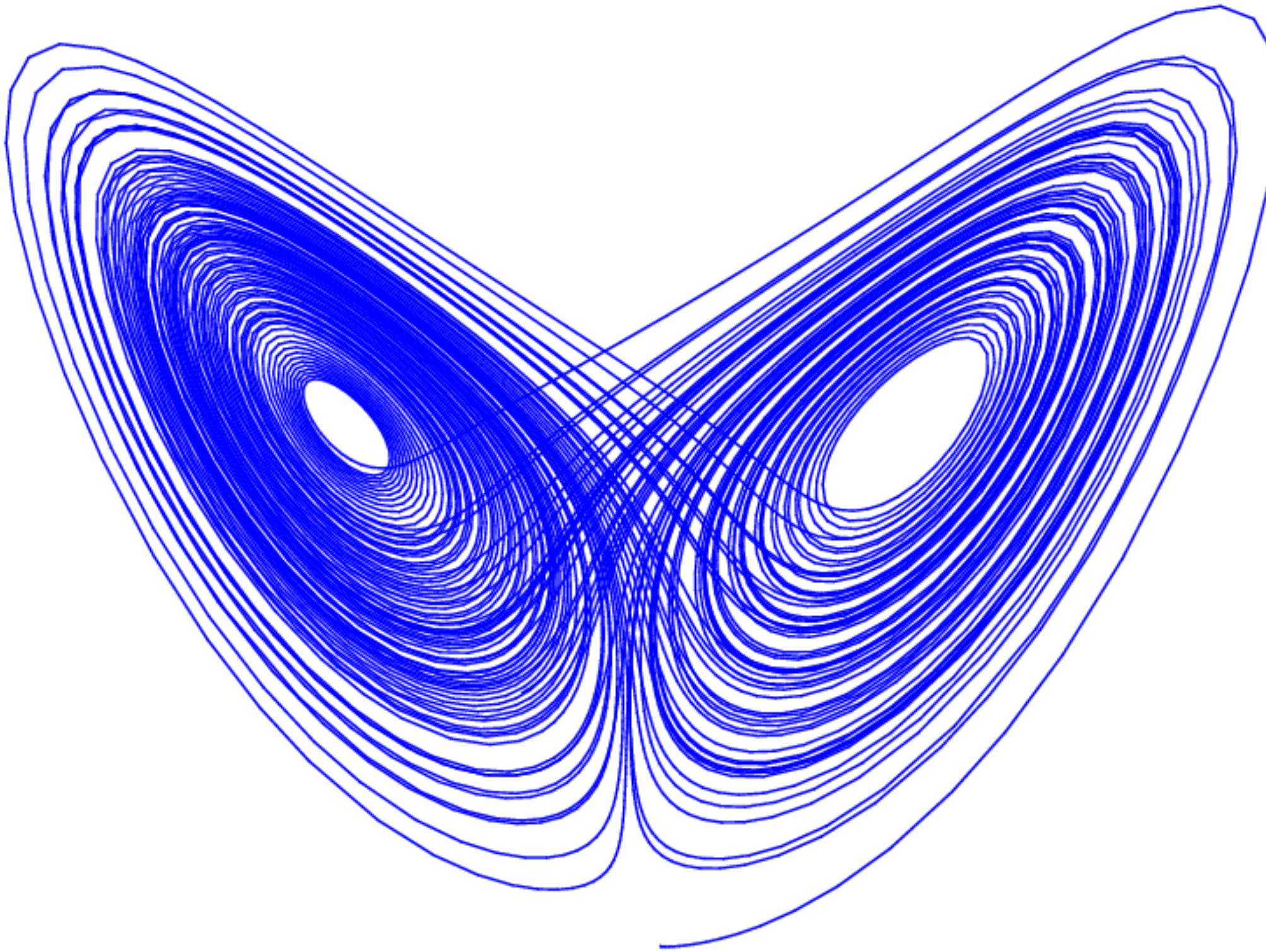
So, we define the *statistical complexity dimension* as the information dimension of the mixed state set:

$$d_\mu = \lim_{\epsilon \rightarrow 0} -\frac{H[R]}{\log(\epsilon)}$$

It gives the rate at which memory resources diverge as we increase predictive power of a finitized  $\epsilon$ -machine.

# Kaplan-Yorke Conjecture

---

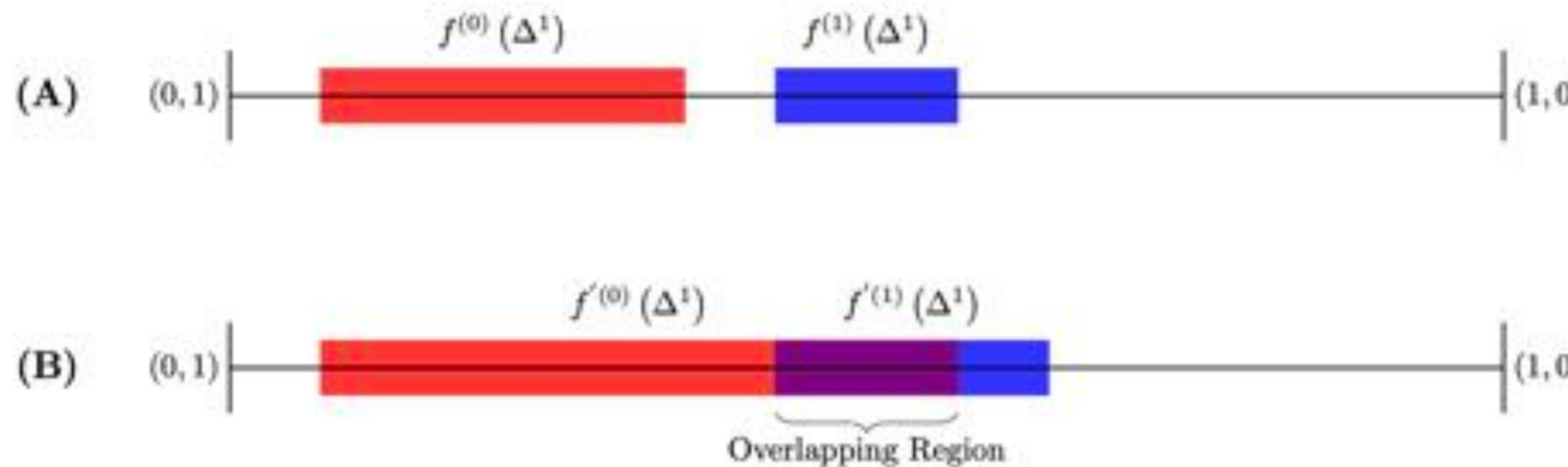


Kaplan—Yorke conjecture:

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

Where  $k$  is the largest index for which  
the sum  $\sum_i^k \lambda_i$  is positive.

# The Overlapping Problem



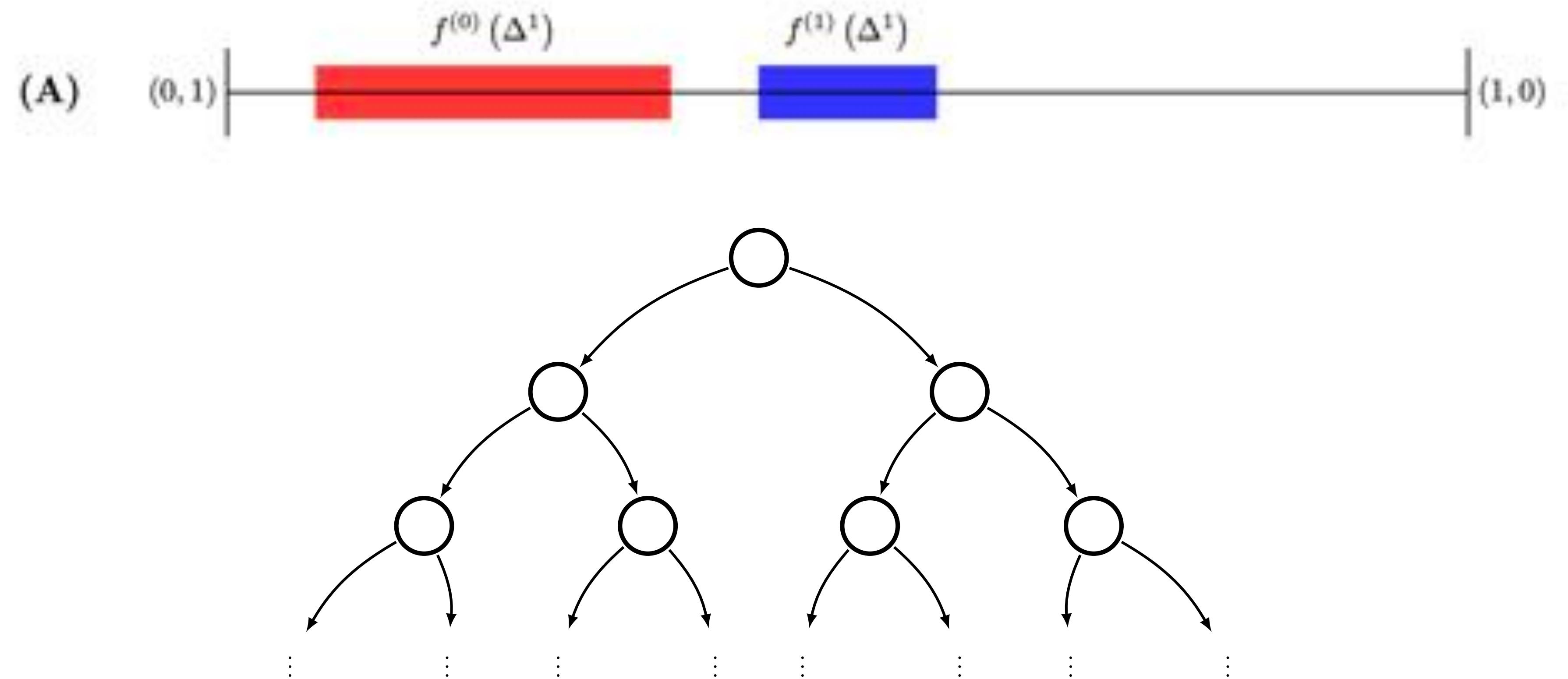
Problem: this solution does not work for “overlapping” IFSs.

We end up “double counting” some of the states.

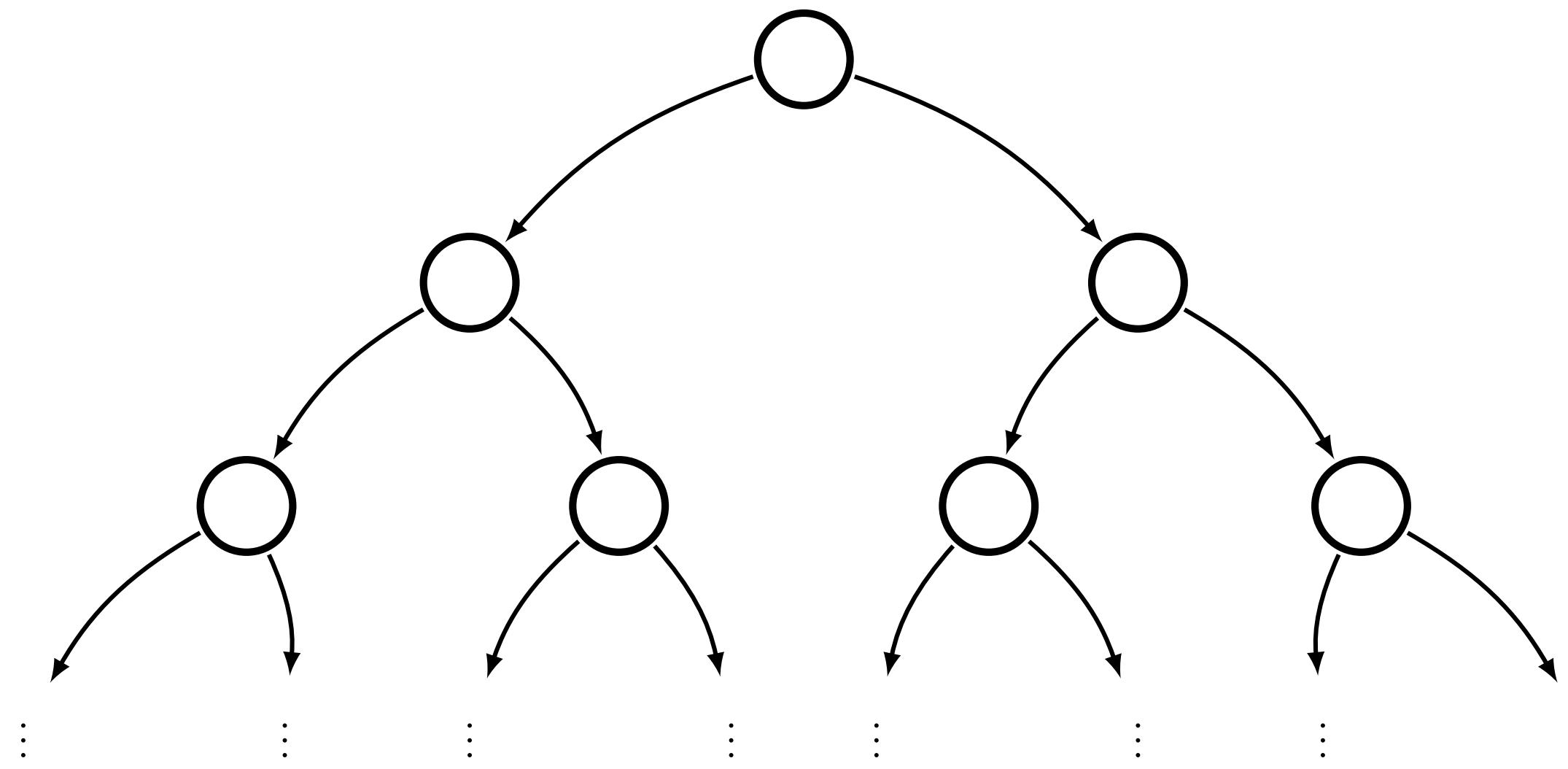
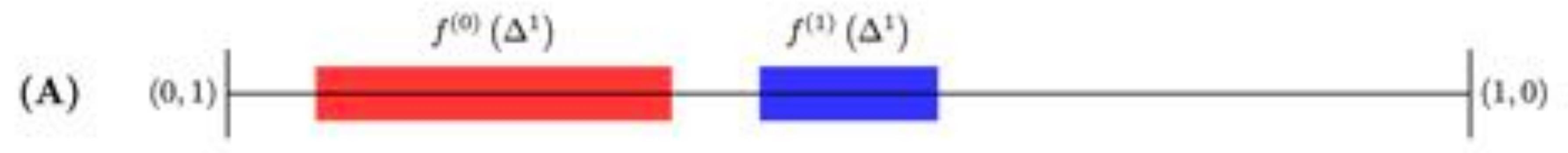
$$\dim_{\mu}(R) \leq k + \frac{h_{\mu} + \sum_i^k \lambda_i}{|\lambda_{k+1}|}$$

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)

# When KY Conjecture Works

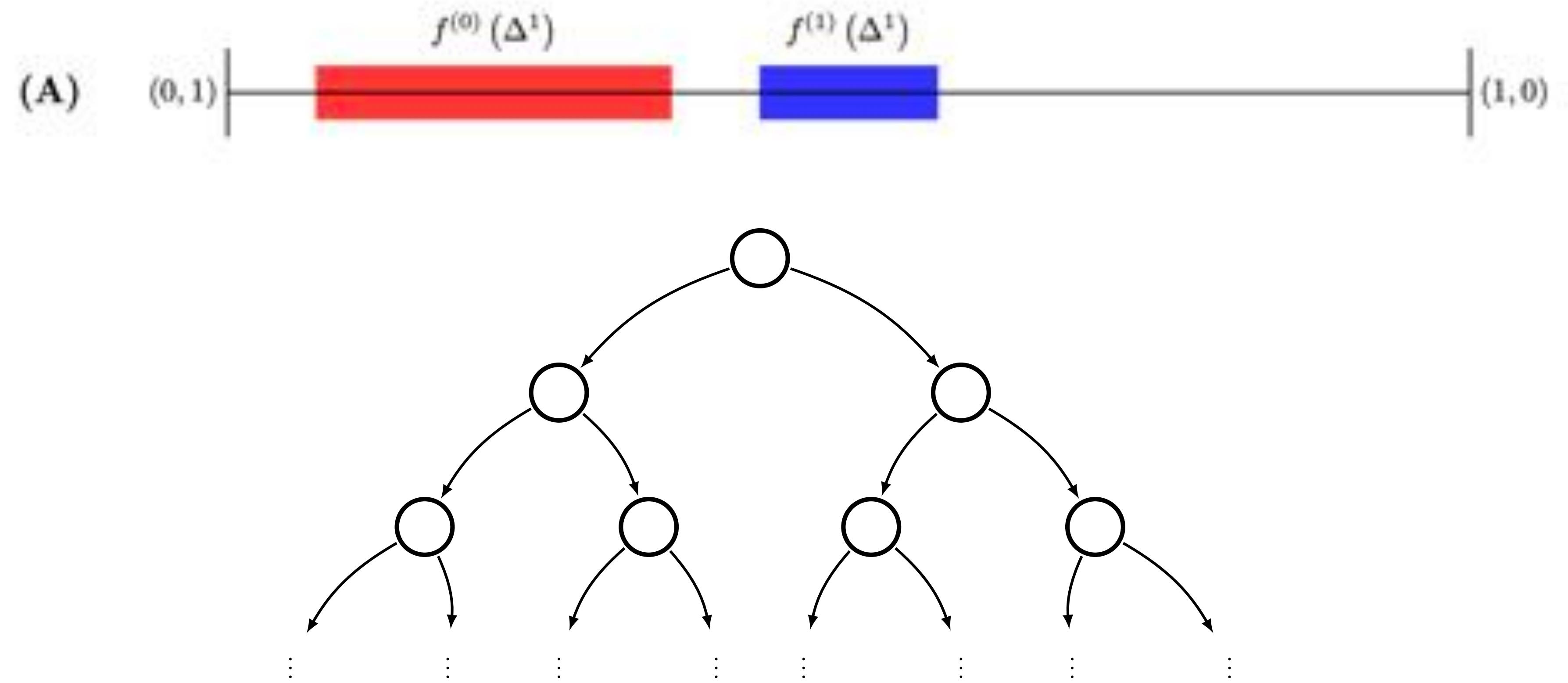


# When KY Conjecture Works



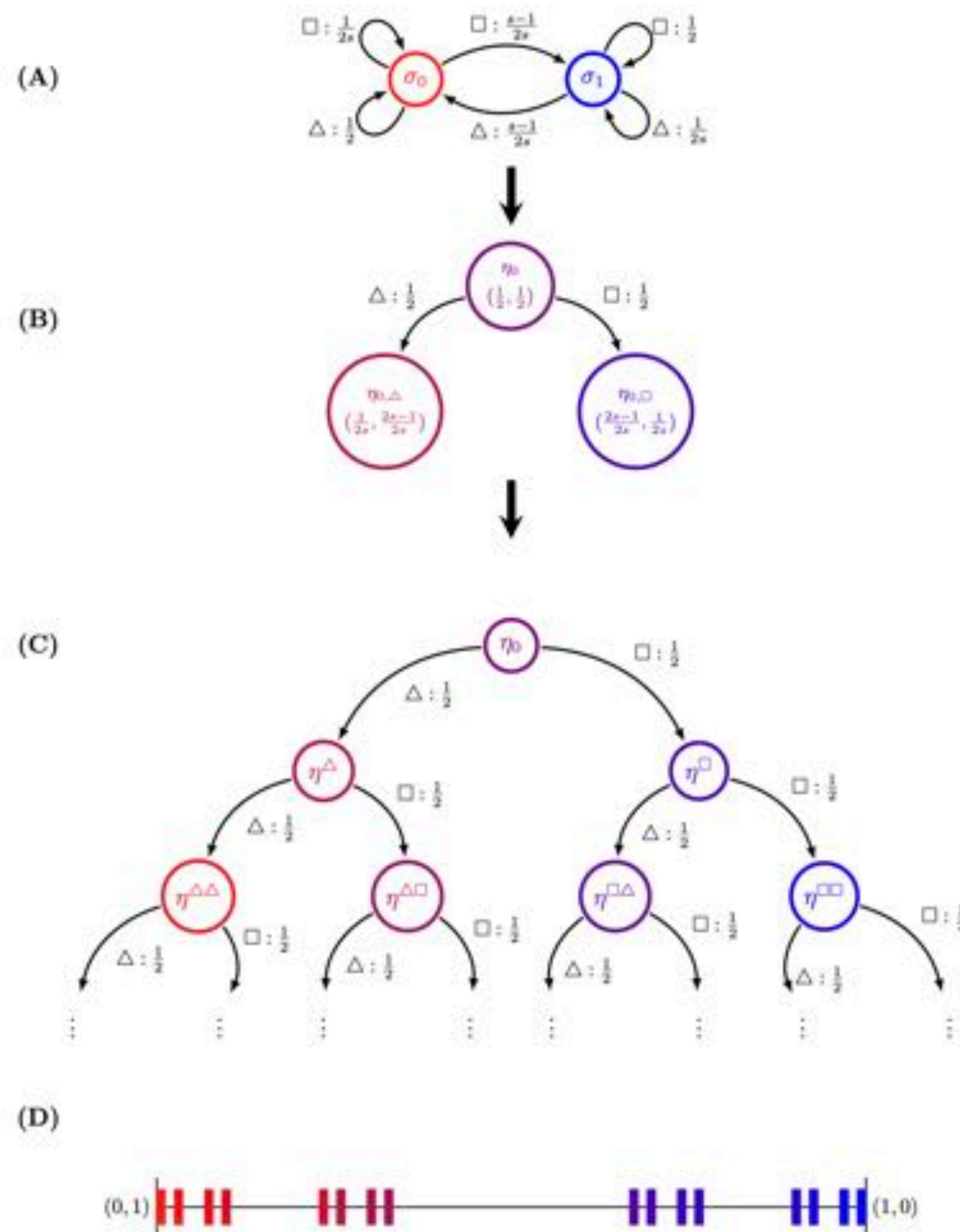
$$\dim_{\mu} (R) = k + \frac{h_{\mu} + \sum_i^k \lambda_i}{|\lambda_{k+1}|}$$

# When KY Conjecture Works



$$\dim_{\mu}(R) = \frac{h_{\mu}}{|\lambda_1|}$$

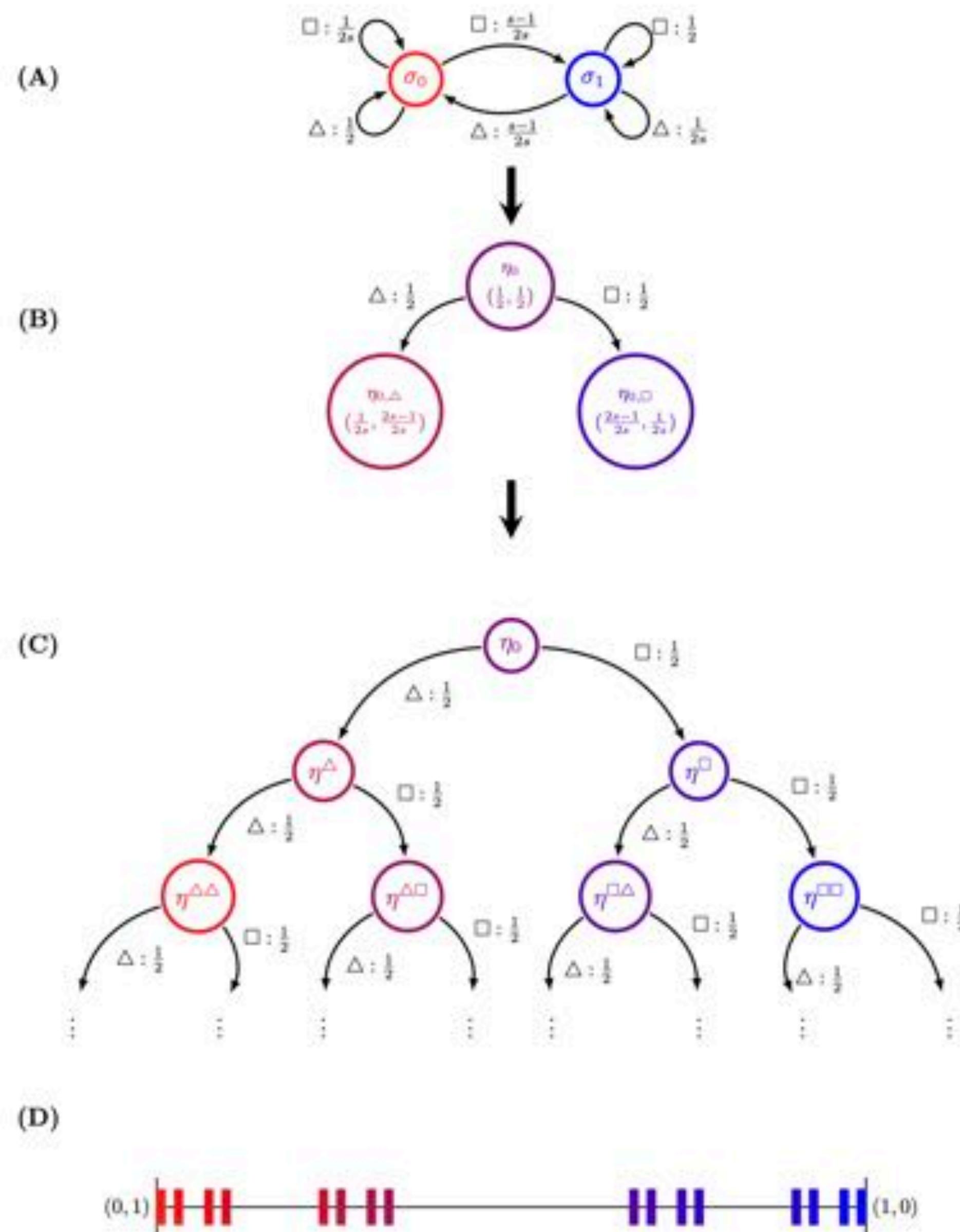
# When KY Conjecture Works



$$\dim_{\mu} (R) = \frac{h_{\mu}}{|\lambda_1|}$$

Where  $\lambda_1$  is the time-averaged Lyapunov exponent of each map.

# When KY Conjecture Works



$$\dim_{\mu} (R) = \frac{h_{\mu}}{|\lambda_1|}$$

Let  $w$  be an  $L$  length word.

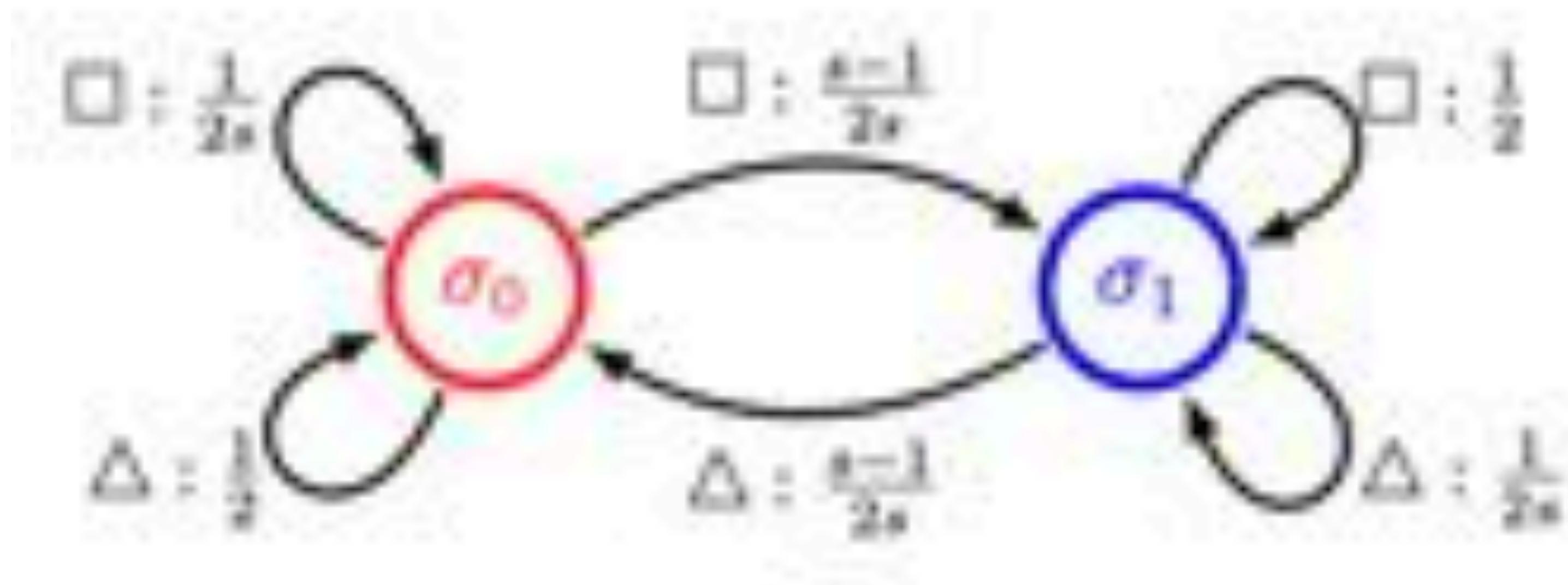
Let  $\eta_t$  be the mixed state at time  $t$ , which emits observation  $w_t$ .

Let  $\eta_t = (x_t, y_t)$ .

$$\lambda(w) = \lim_{L \rightarrow \infty} \frac{1}{L} \sum_{i=0}^{n-1} \ln \left| \frac{df^{(w_i)}(x)}{dx} \right|_{x=x_t}$$

# Cantor Set Machine

---

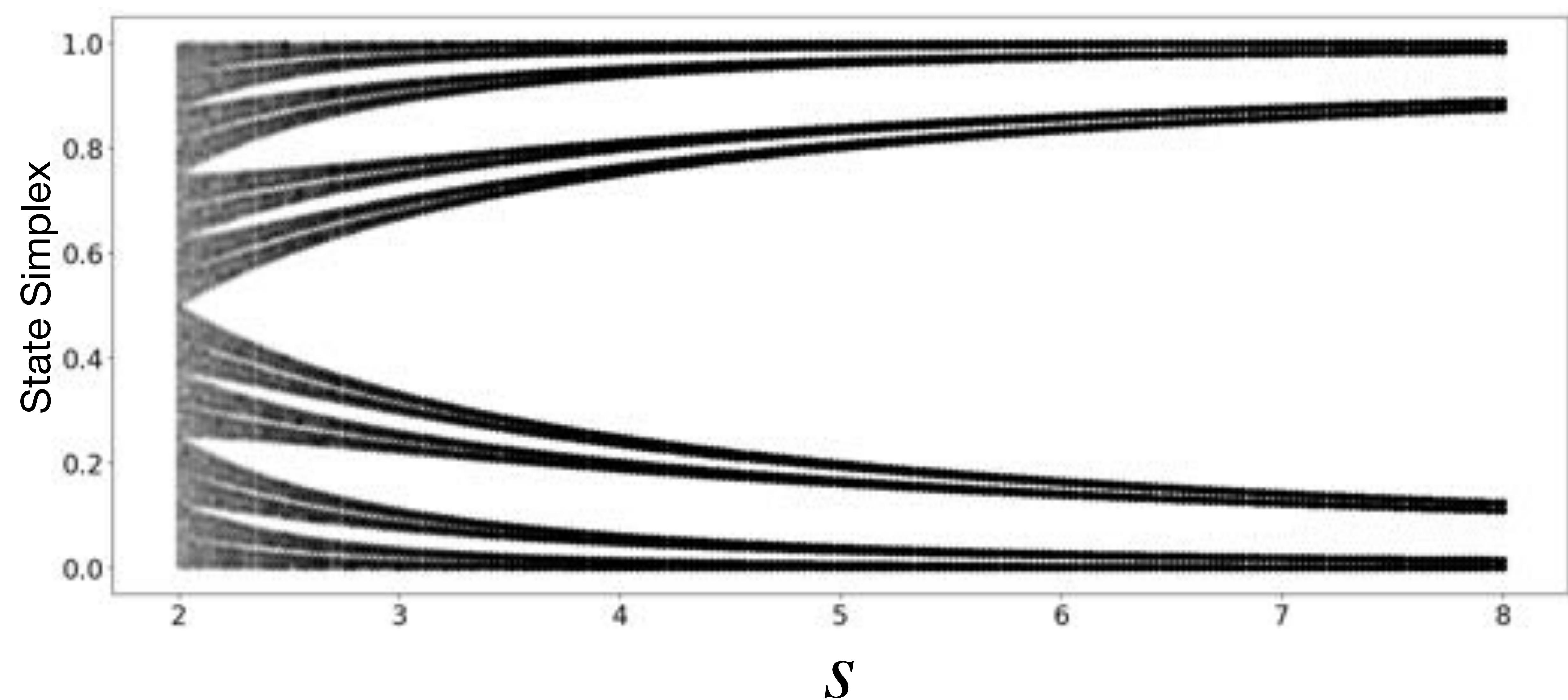


This machine produces a missing  $s$  Cantor set.



# Cantor Set Machine

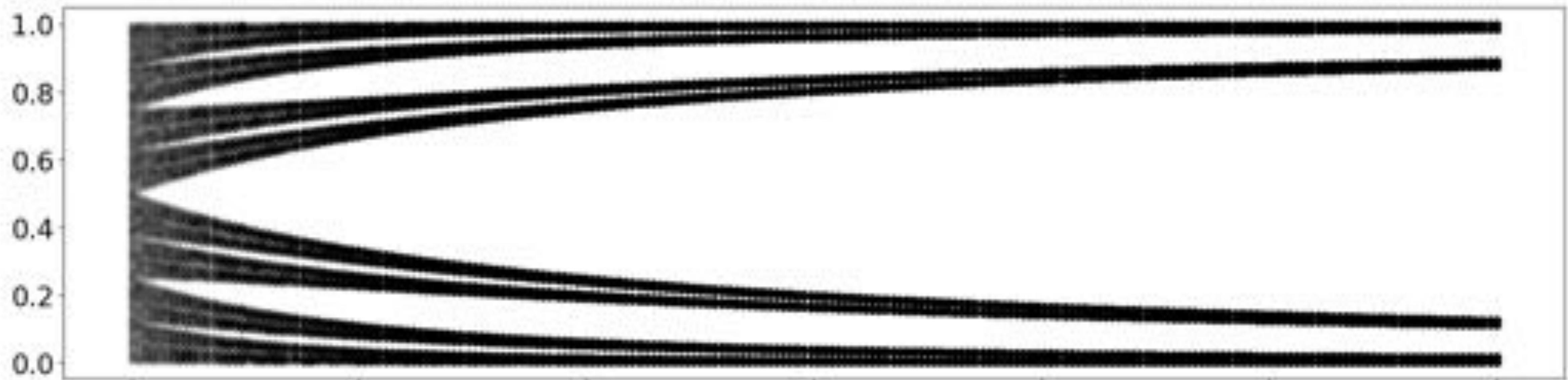
---



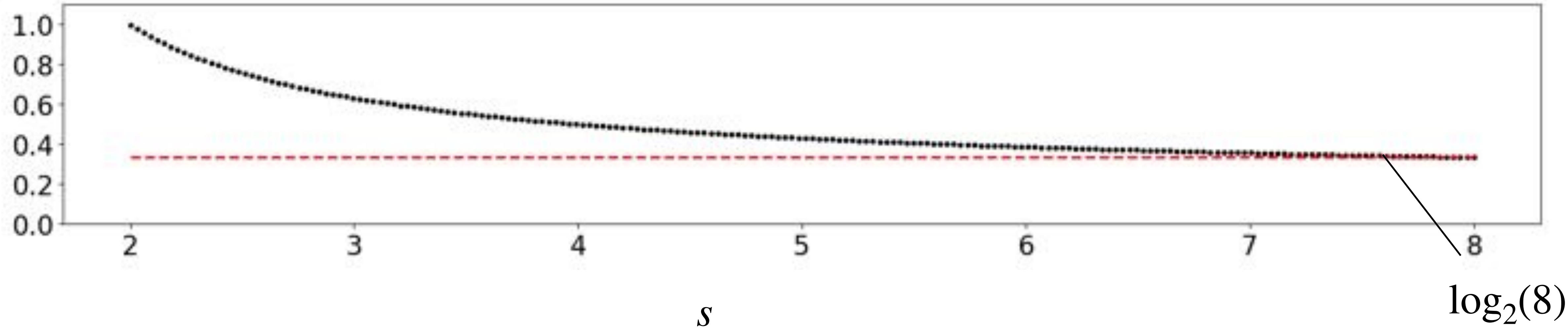
# Cantor Set Machine

---

State Simplex



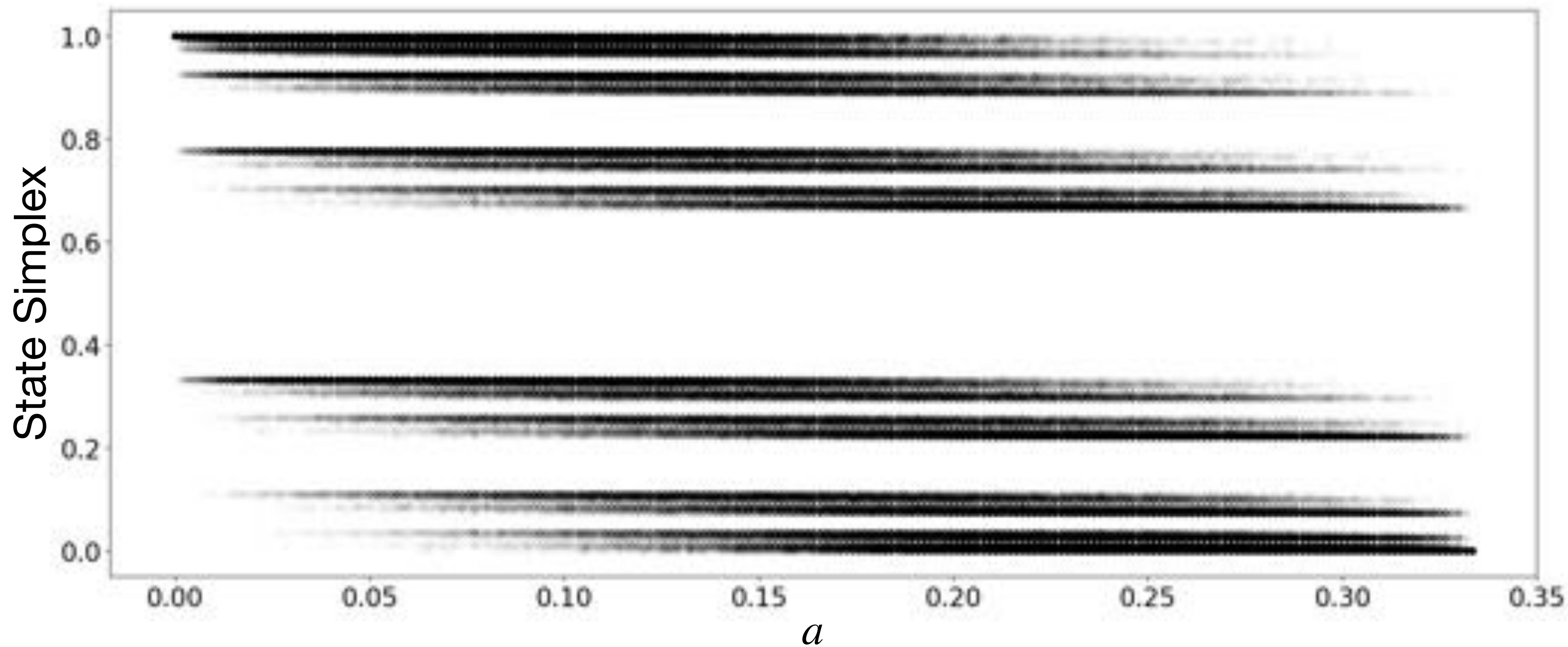
$\dim_{\mu}$



# Cantor Set Machine

---

$$s = 3$$

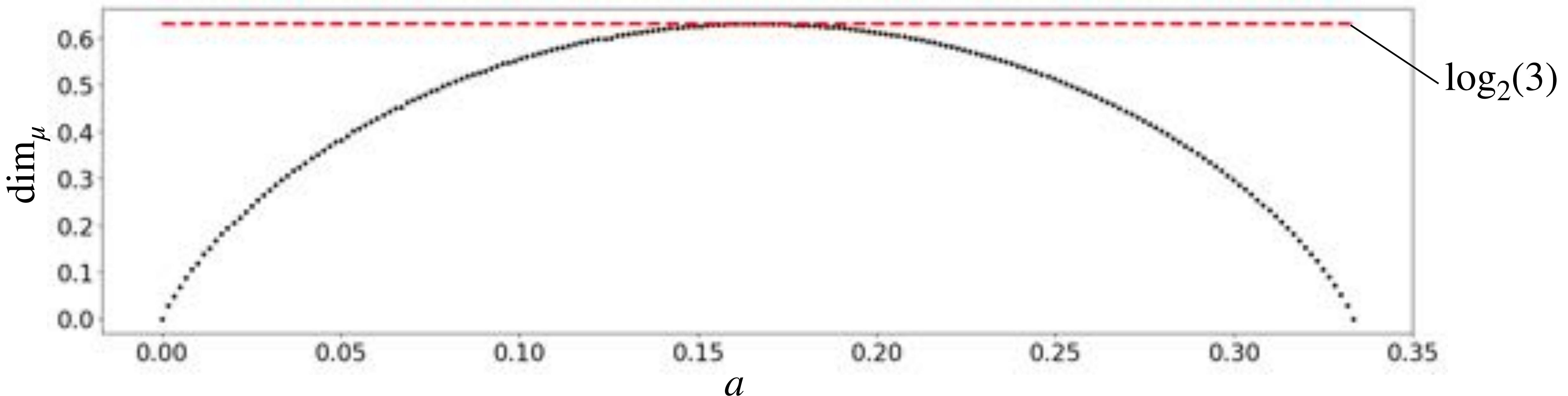
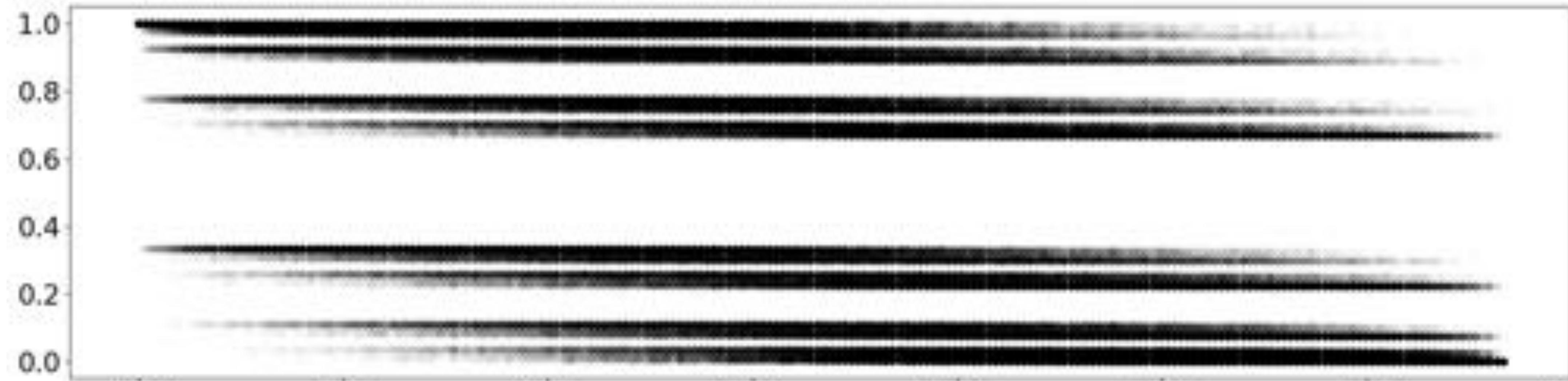


# Cantor Set Machine

---

$s = 3$

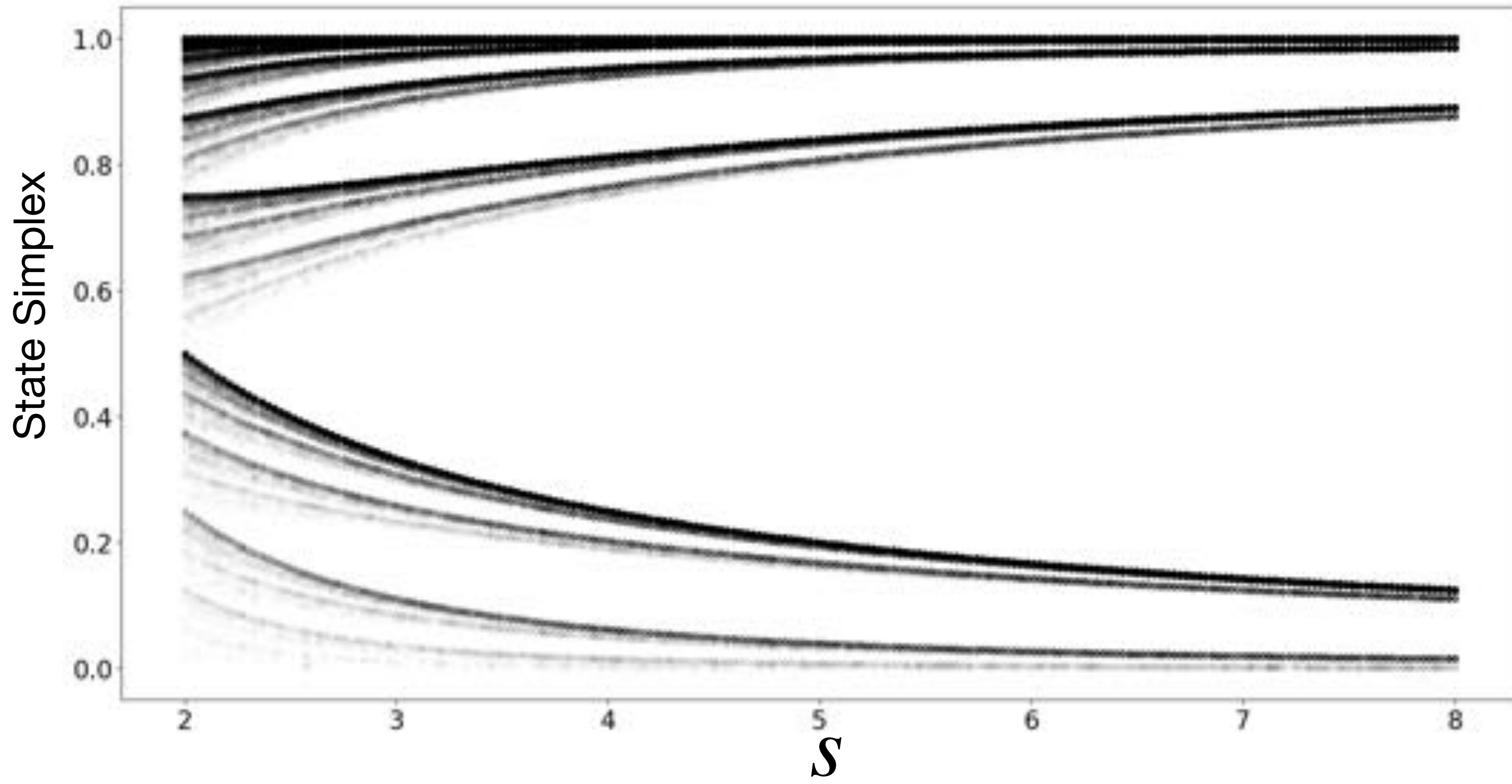
State Simplex



# Cantor Set Machine

---

$$a = 1/6s$$

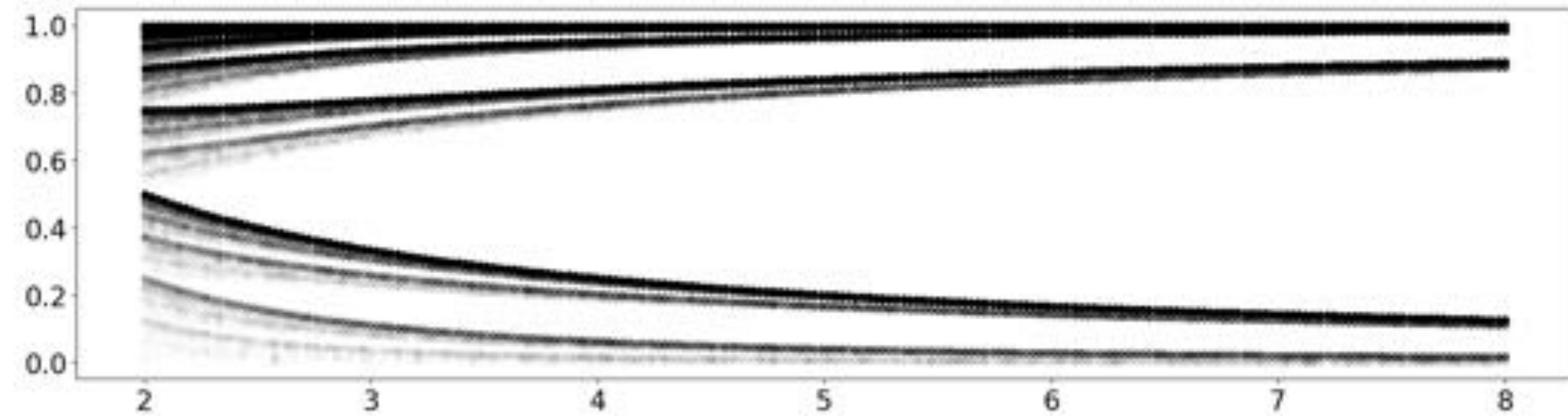


# Cantor Set Machine

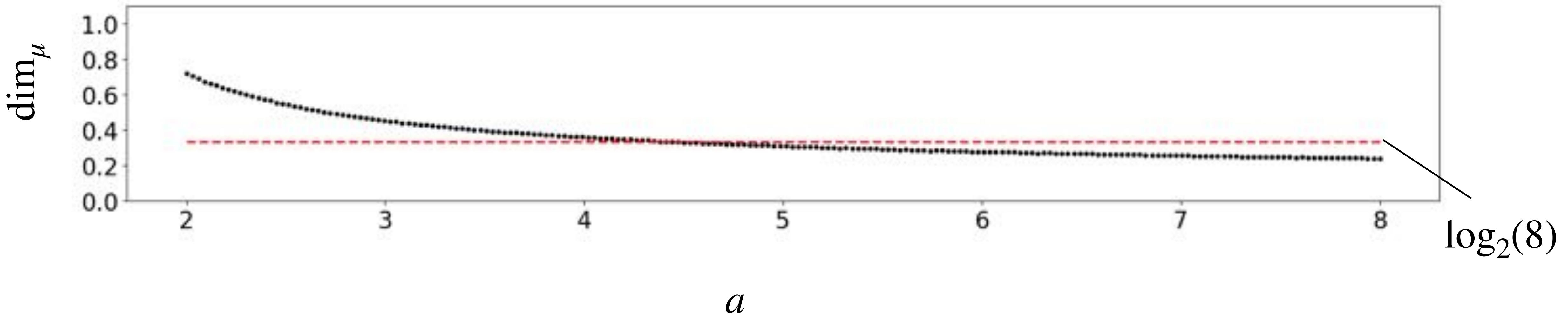
---

$s = 3$

State Simplex

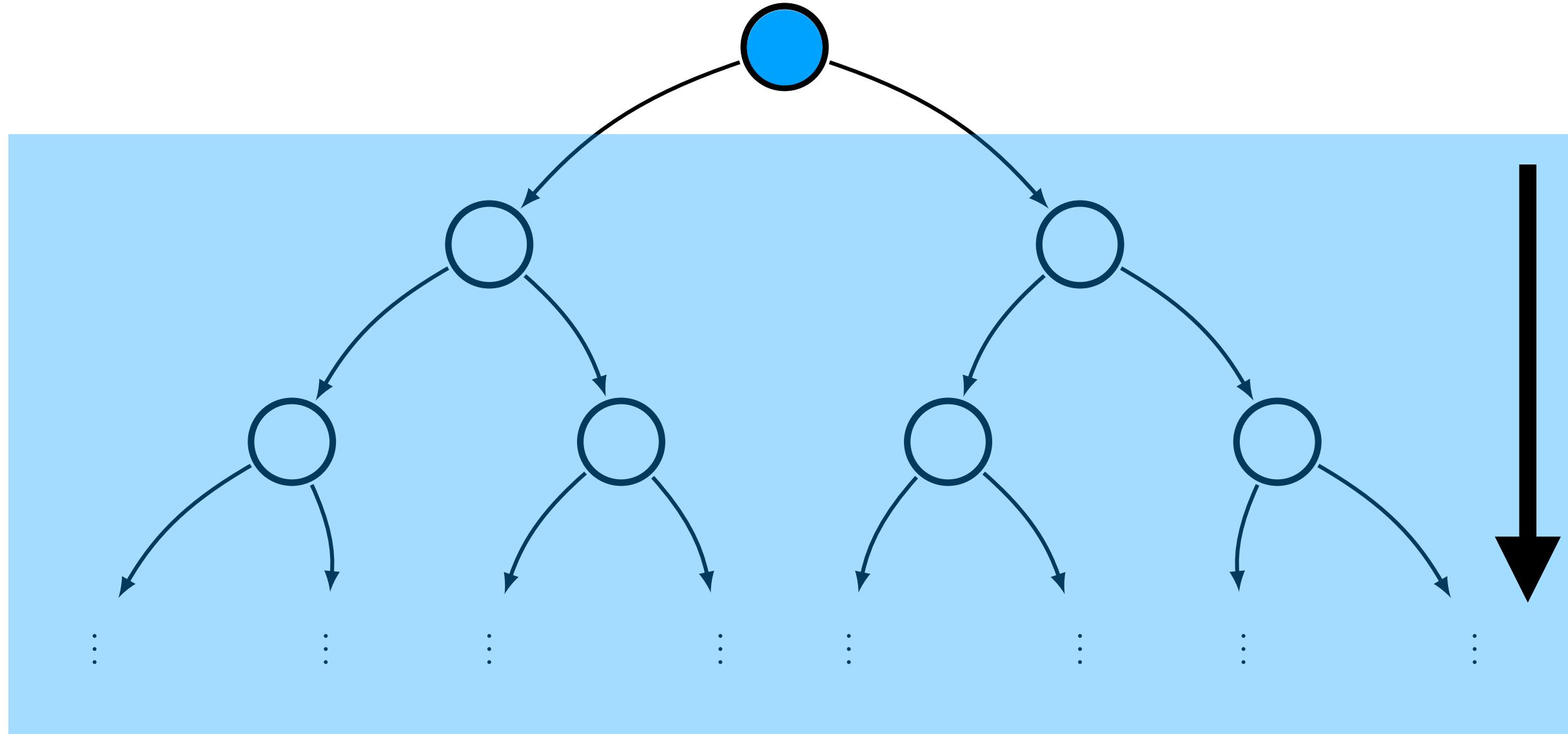


$\dim_{\mu}$



# Overlap Problem

---



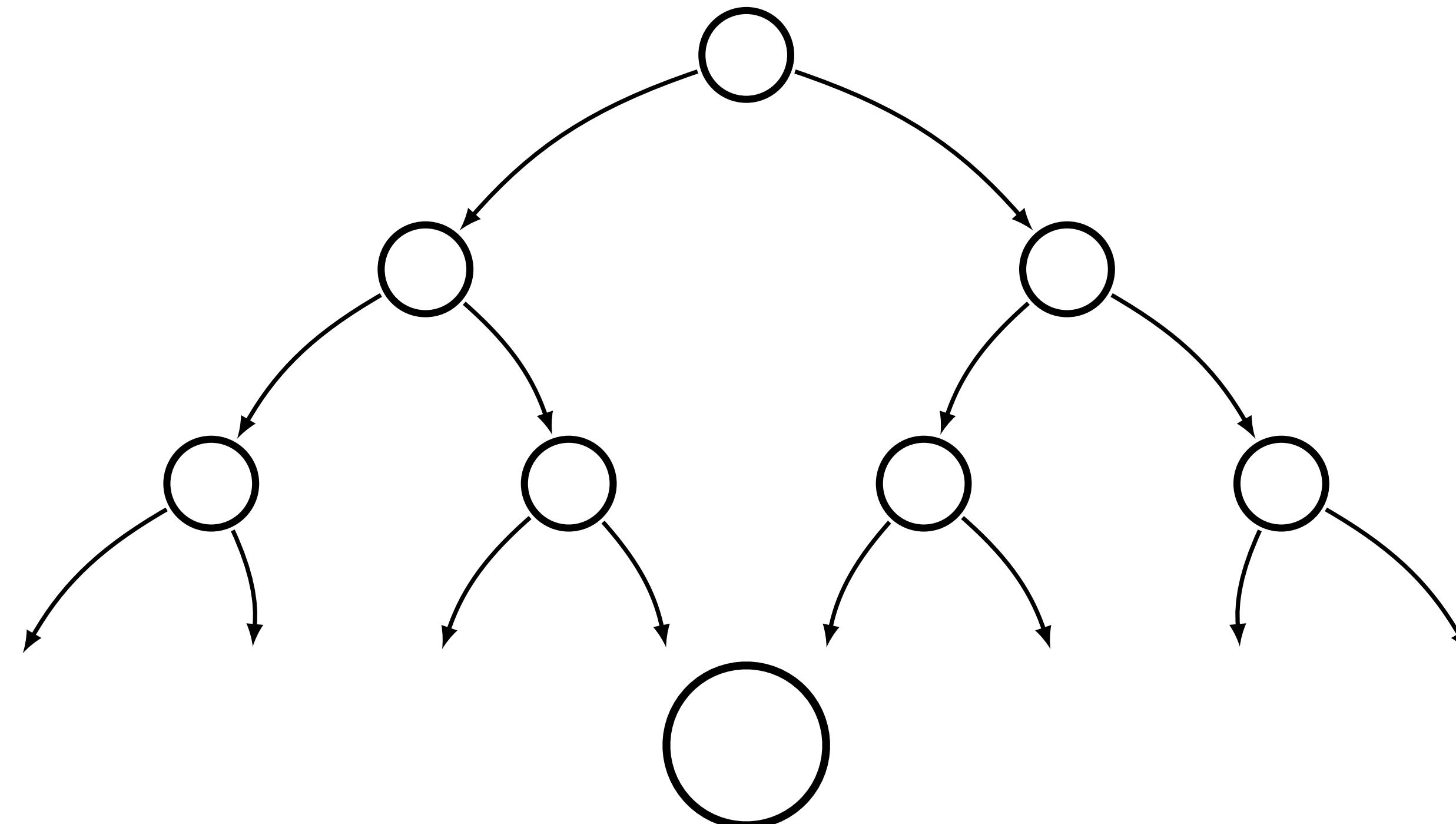
Entropy rate measures uncertainty in the next symbol given the present.

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

Rate of new symbols, on average...

# Overlap Problem

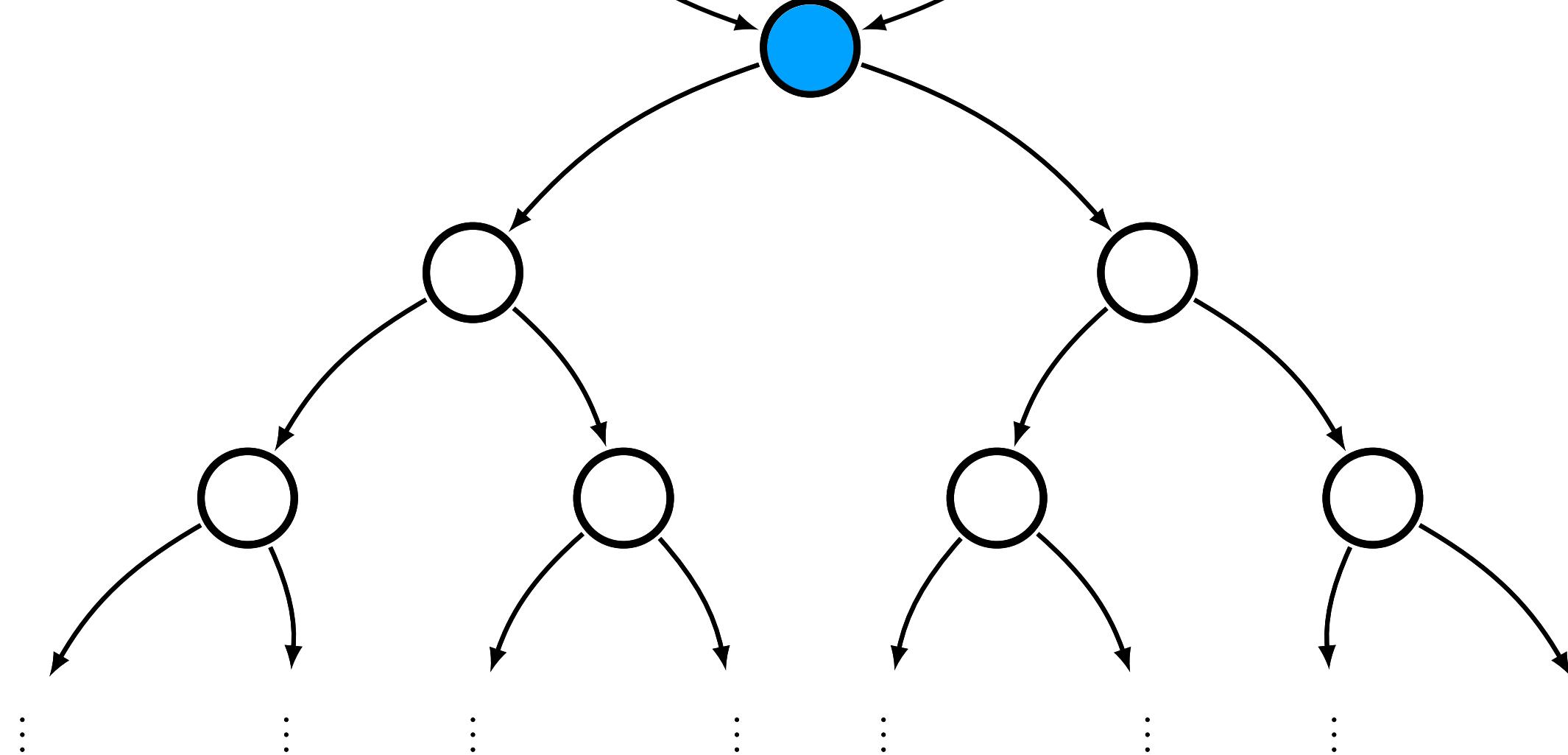
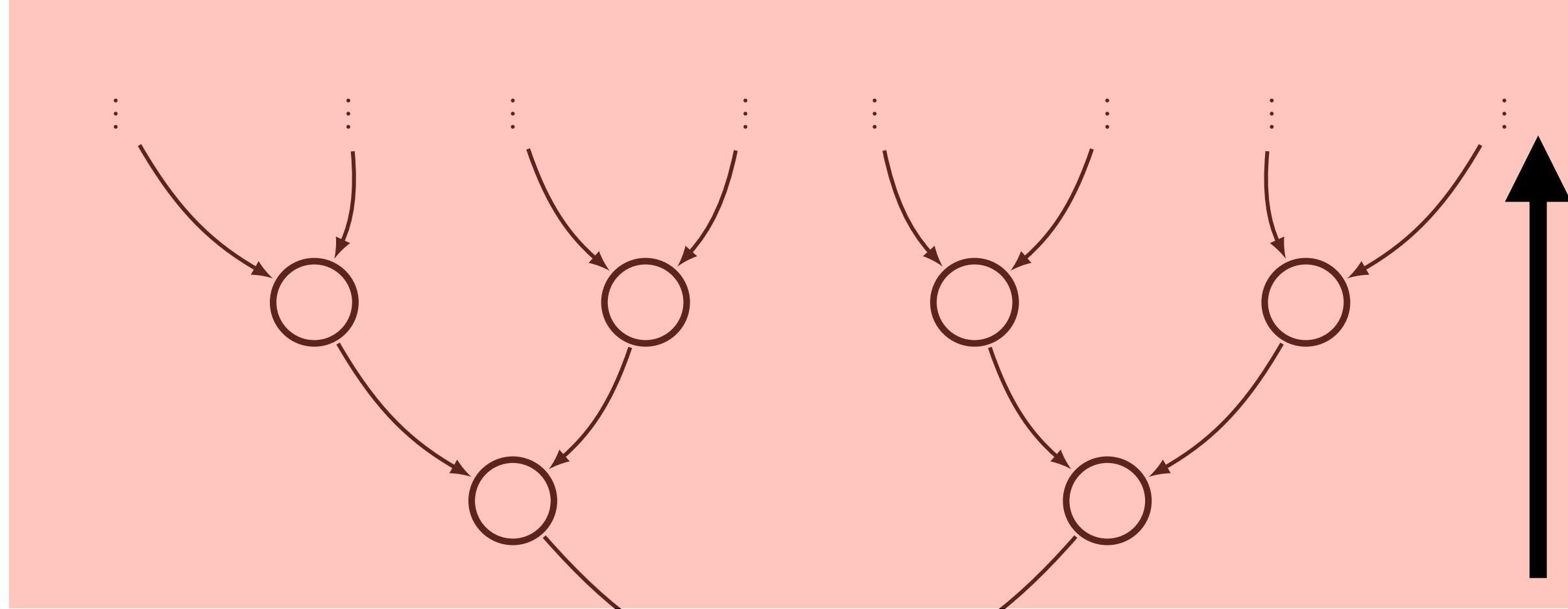
---



But, the branching entropy *overcounts* the average number of states generated when it is possible for two states to map into the same state.

We need a correction to count the average number of new states properly.

# New Quantity: 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*.

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

# Proposed Correction to Kaplan Yorke Conjecture

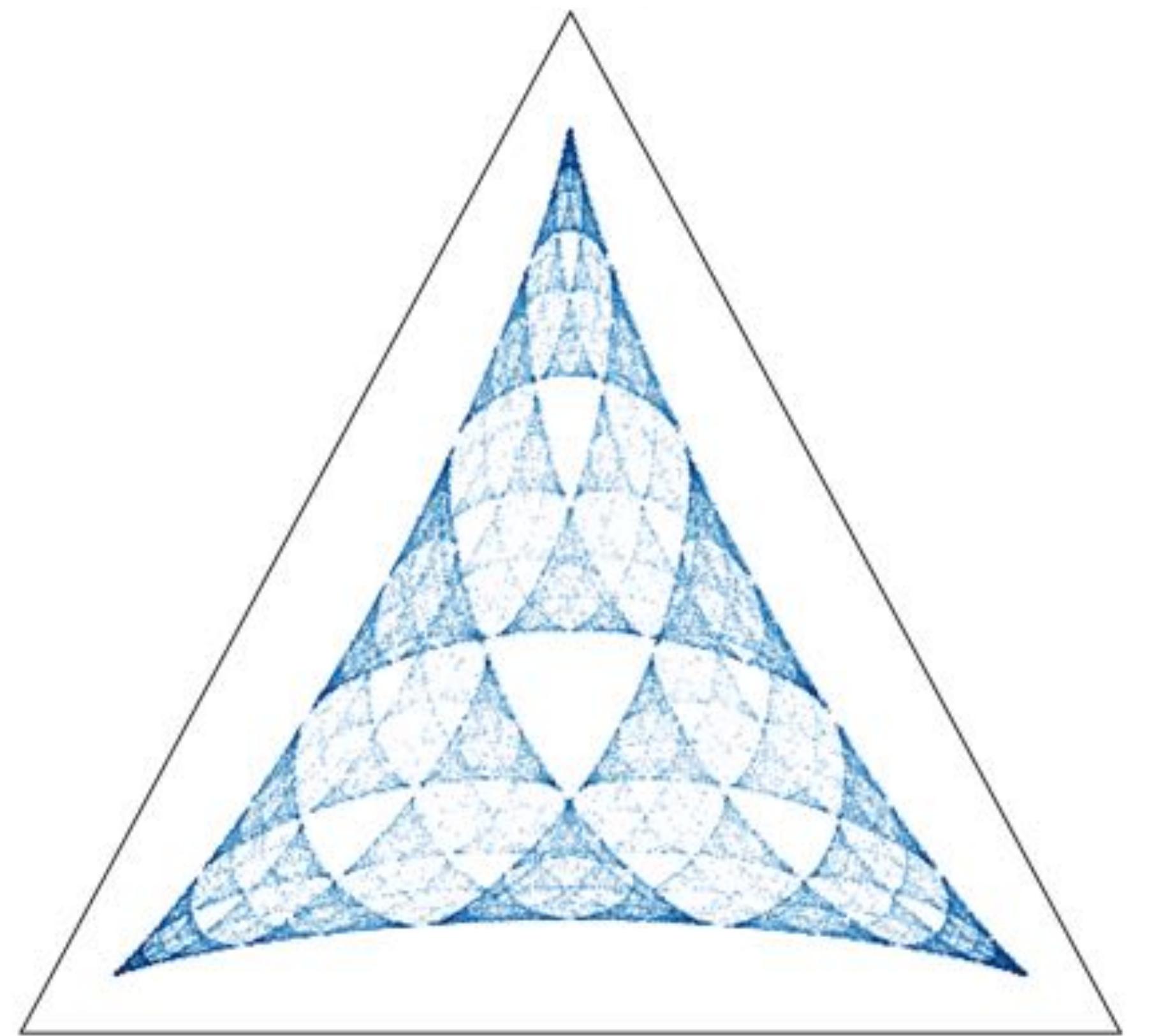
---

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

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)

# Proposed Correction to Kaplan-Yorke Conjecture



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

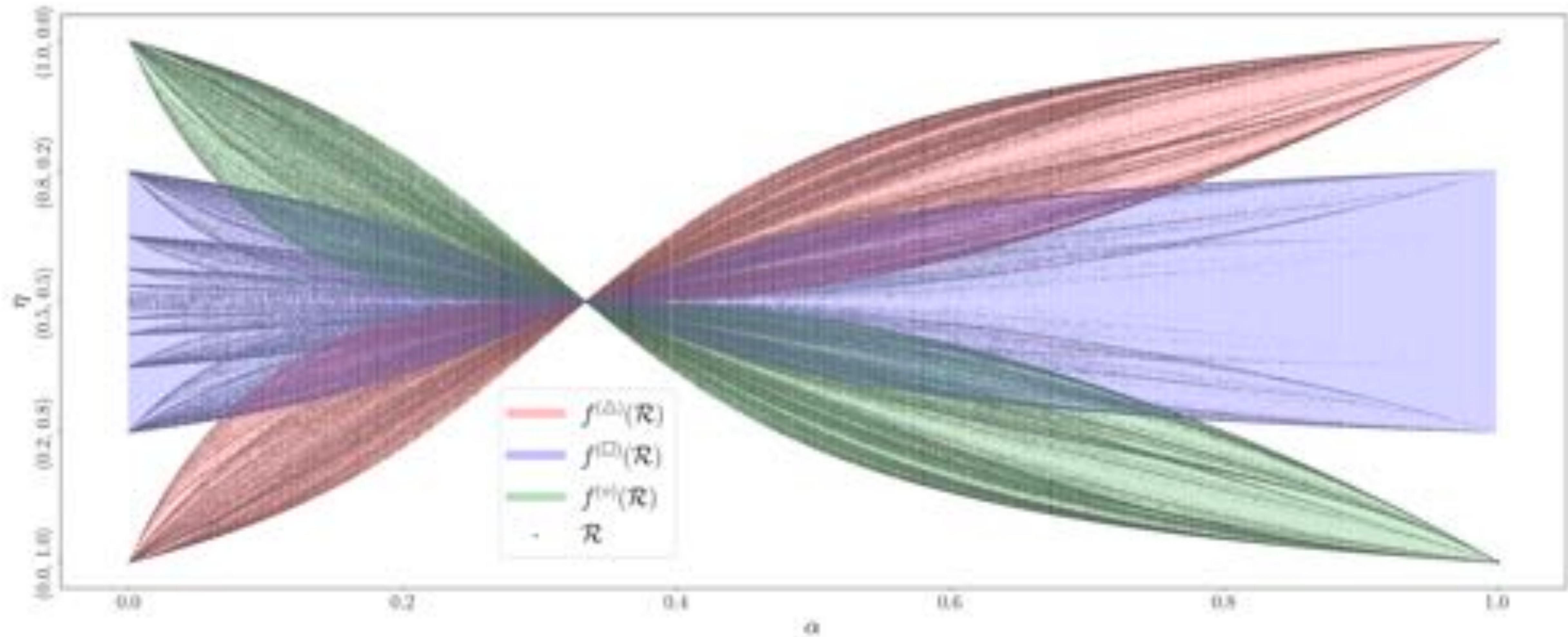
Entropy rate

Ambiguity rate

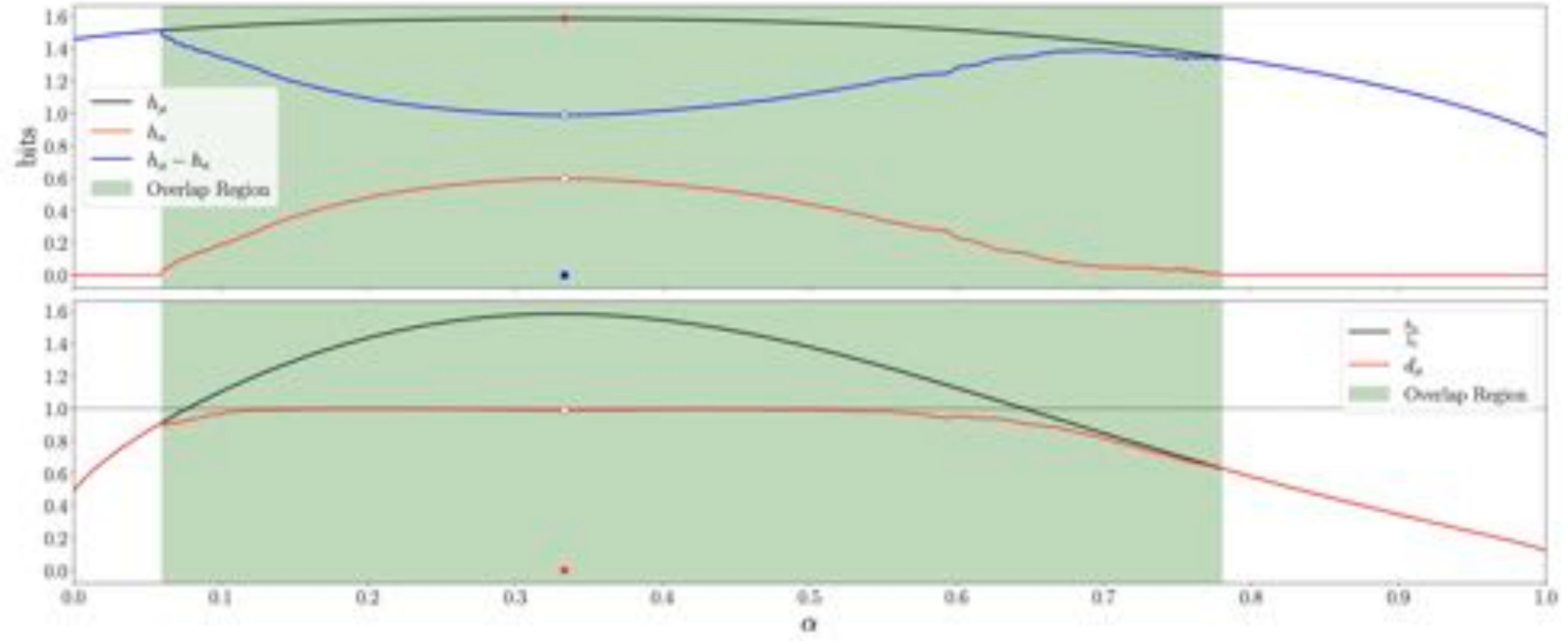
Alexandra M. Jurgens, James P. Crutchfield. *Divergent Predictive Memory: The Statistical Complexity Dimension of Stationary, Ergodic Finite-State Hidden Markov Processes*. Chaos 31, 083114, 2021

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

# Proposed Correction to Kaplan Yorke Conjecture



# Proposed Correction to Kaplan Yorke Conjecture

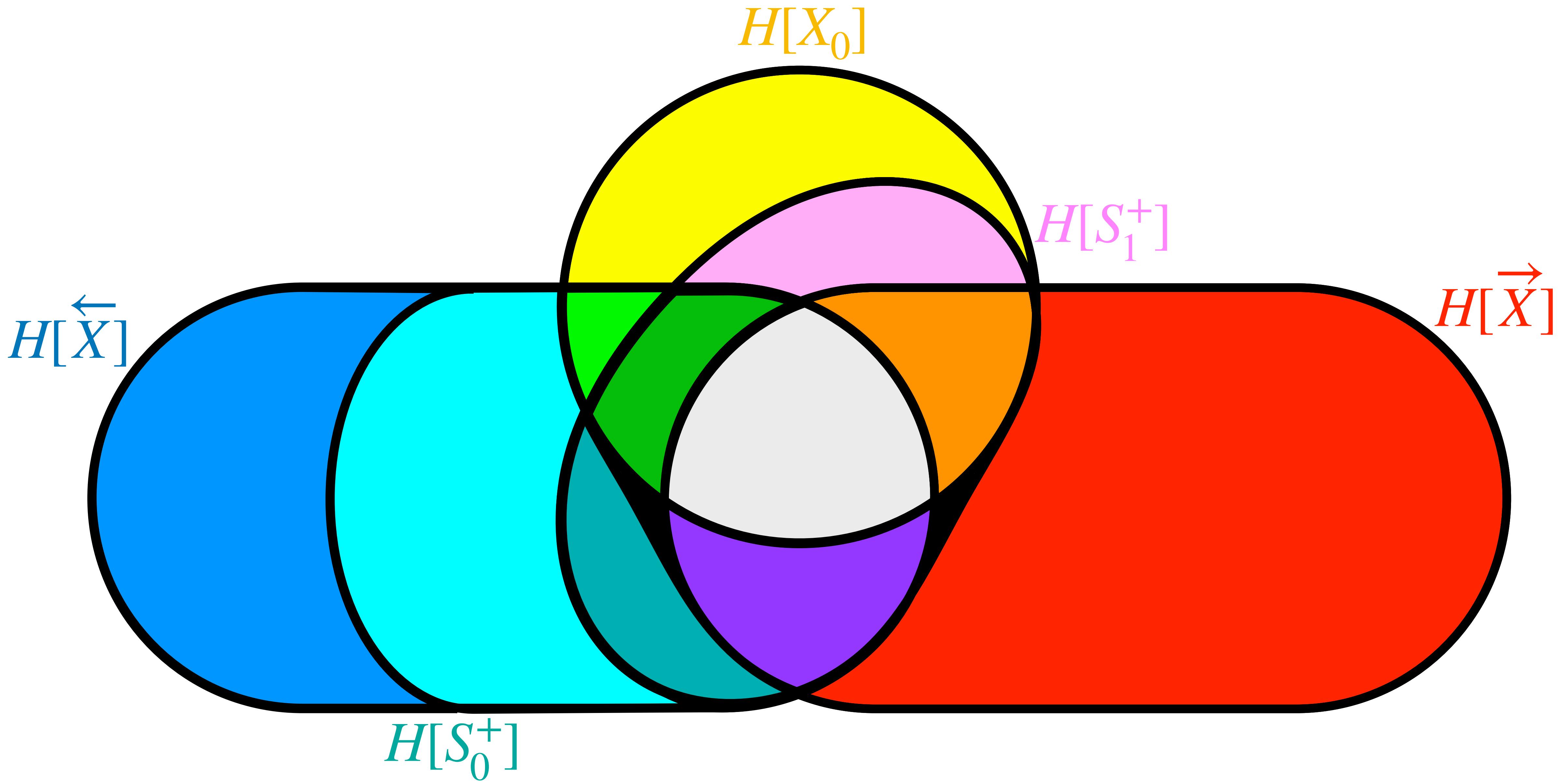


---

# What is $h_a$ ?

# What is $h_a$ ?

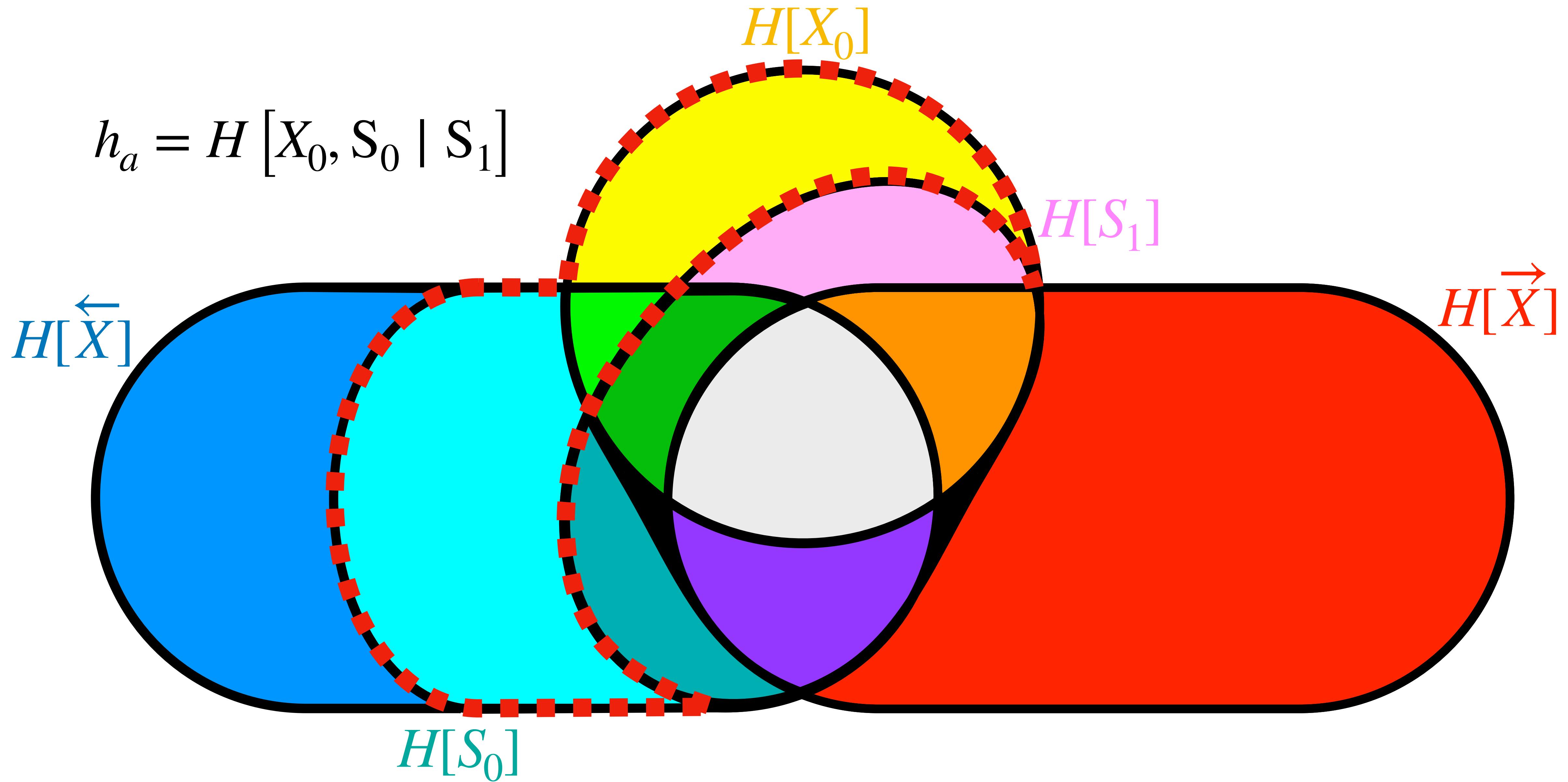
---



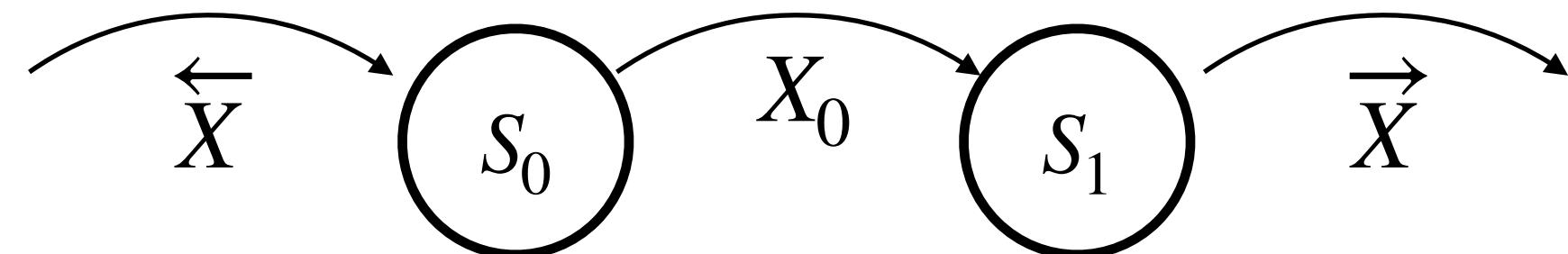
# What is $h_a$ ?

---

$$h_a = H [X_0, S_0 \mid S_1]$$



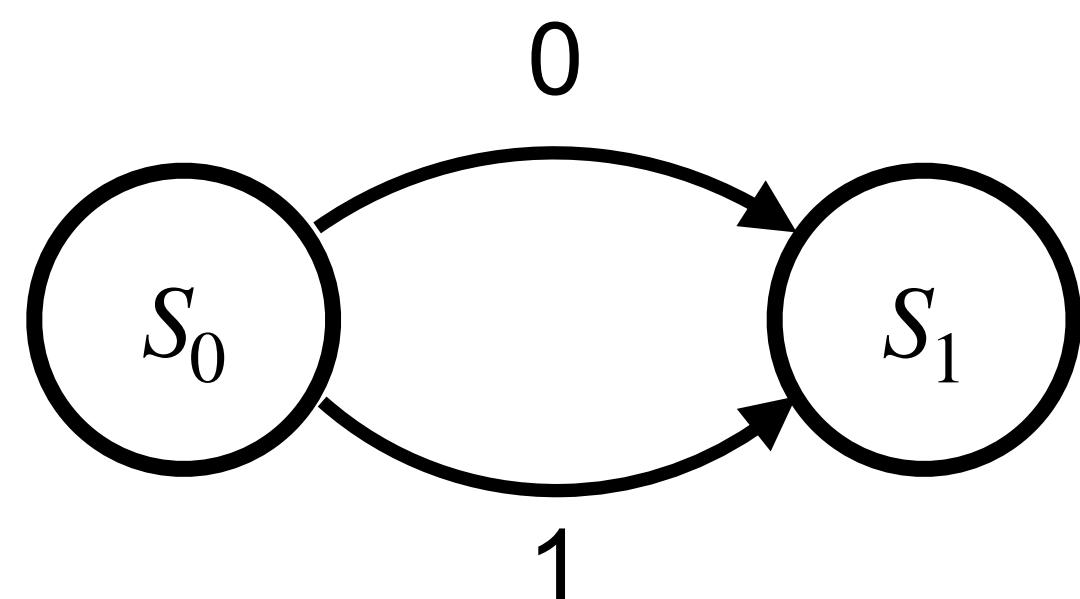
# Components of Ambiguity Rate



*Ephemeral information*

$$H[X_0 | S_0, S_1]$$

Information generated in  
present, not used for  
prediction.



# Components of Ambiguity Rate

---

$$h_a = H [X_0, S_0 \mid S_1]$$

$$H[\vec{X}]$$

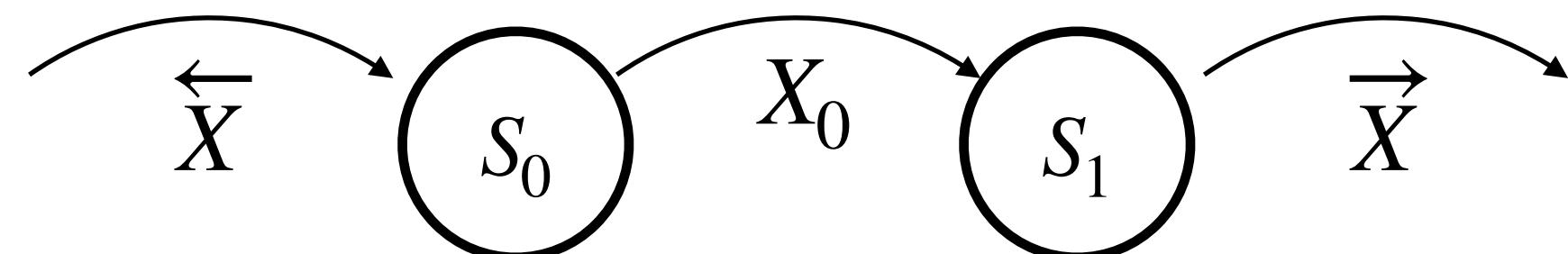
$$H[X_0]$$

$$H[S_1]$$

$$H[\vec{X}]$$

$$H[S_0]$$

# Components of Ambiguity Rate

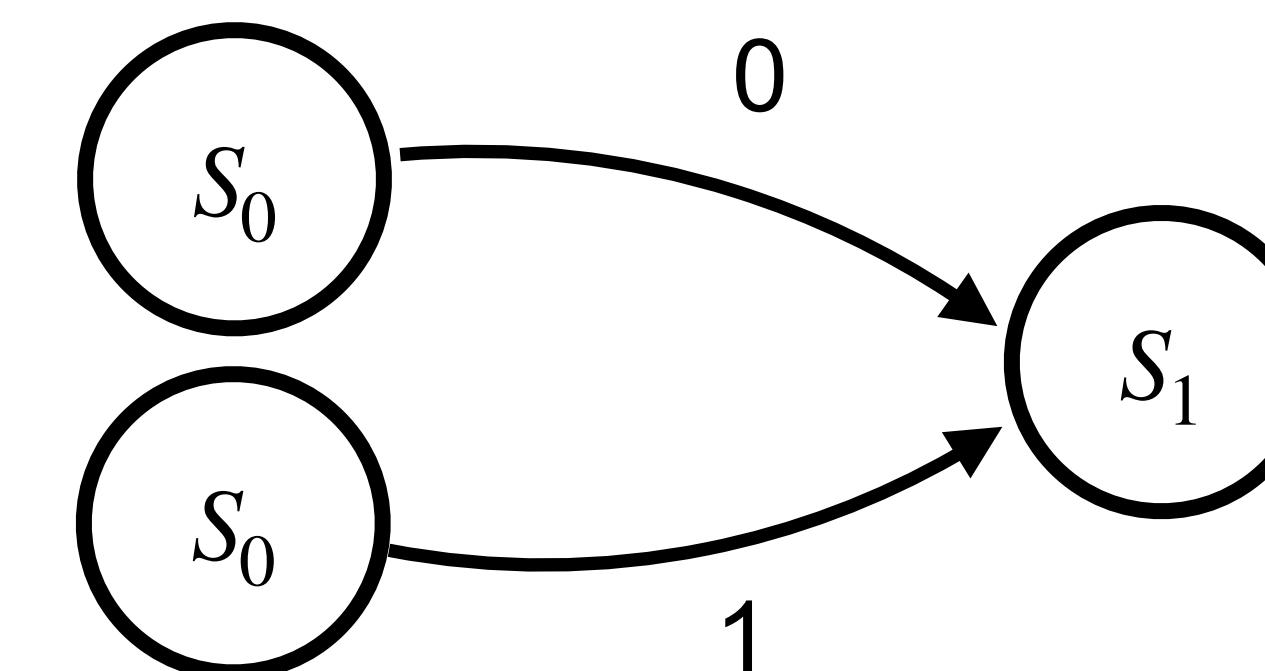


*Binding information*

$$I[X_0; S_0 | S_1]$$



Information in the past and present that doesn't get used for future prediction.



# Components of Ambiguity Rate

---

$$h_a = H [X_0, S_0 \mid S_1]$$

$$H[\vec{X}]$$

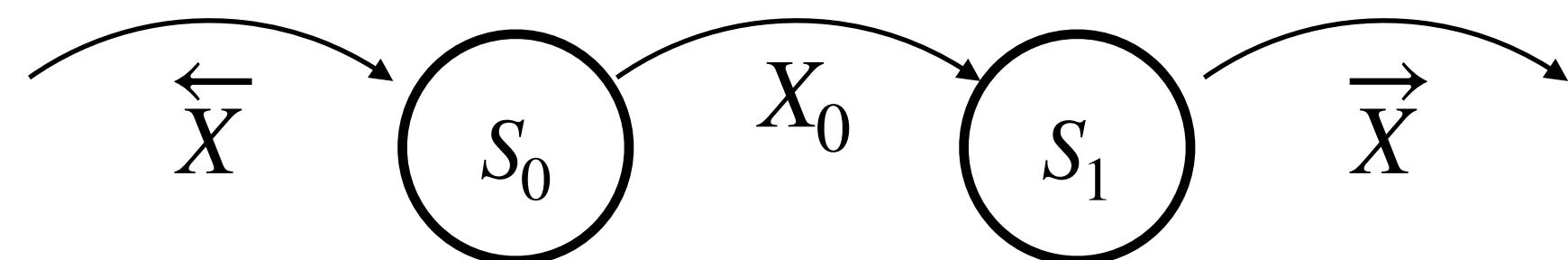
$$H[X_0]$$

$$H[S_1]$$

$$H[\vec{X}]$$

$$H[S_0]$$

# Components of Ambiguity Rate

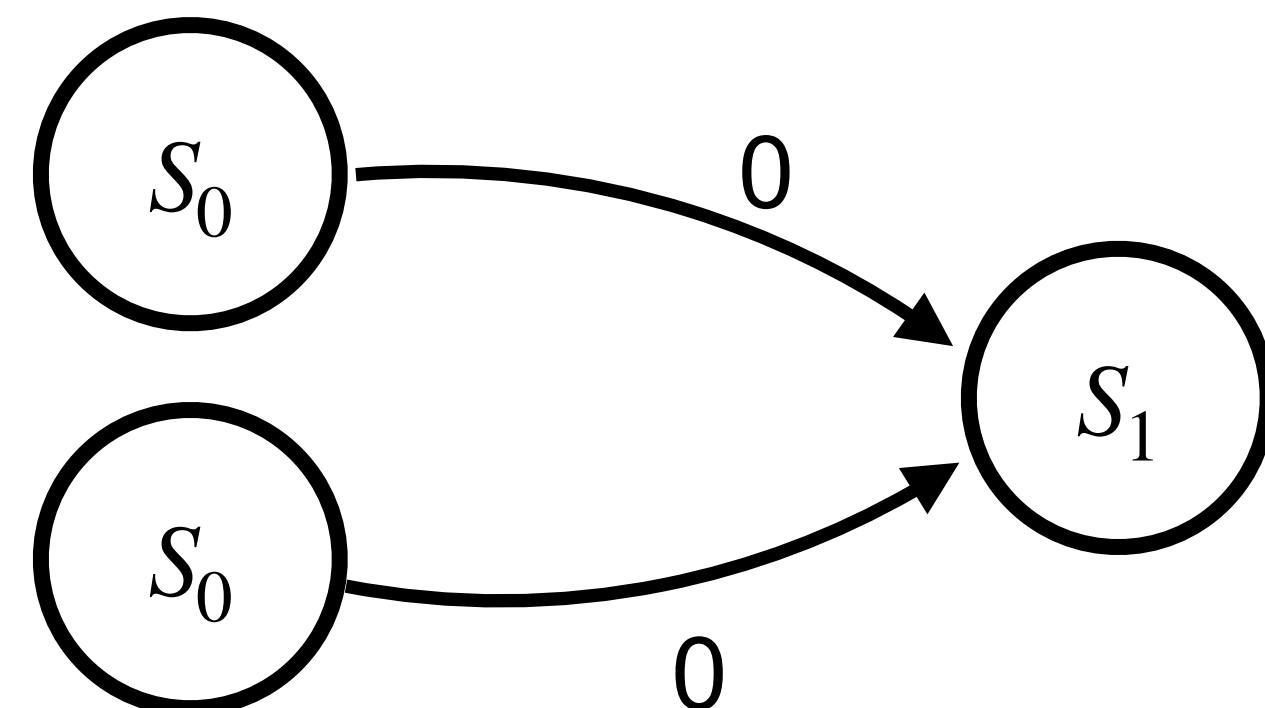


*Crypticity*

$$H[S_0 | X_0, S_1]$$



Information in past used for prediction but then forgotten for future prediction.



A Jurgens, JP Crutchfield. Taxonomy of Prediction. arXiv preprint arXiv:2504.11371, 2025

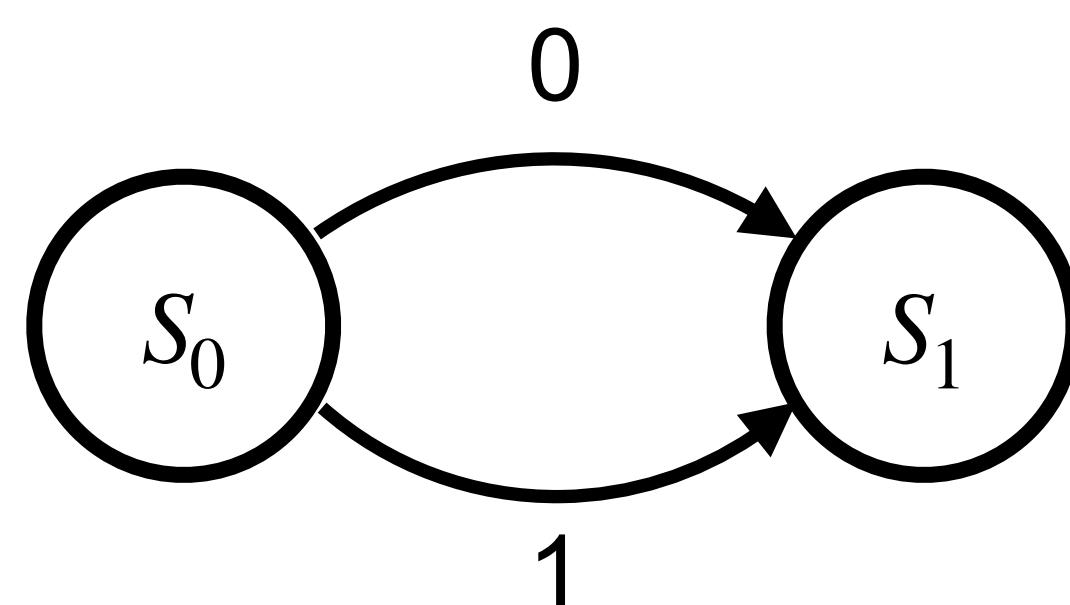
R. G. James, C. J. Ellison, and J. P. Crutchfield. Anatomy of a bit: Information in a time series observation. CHAOS. 17. 21(3):037109, 2011. doi:10.1063/1.3637494

# Components of Ambiguity Rate

*Ephemeral information*

$$H[X_0 | S_0, S_1]$$

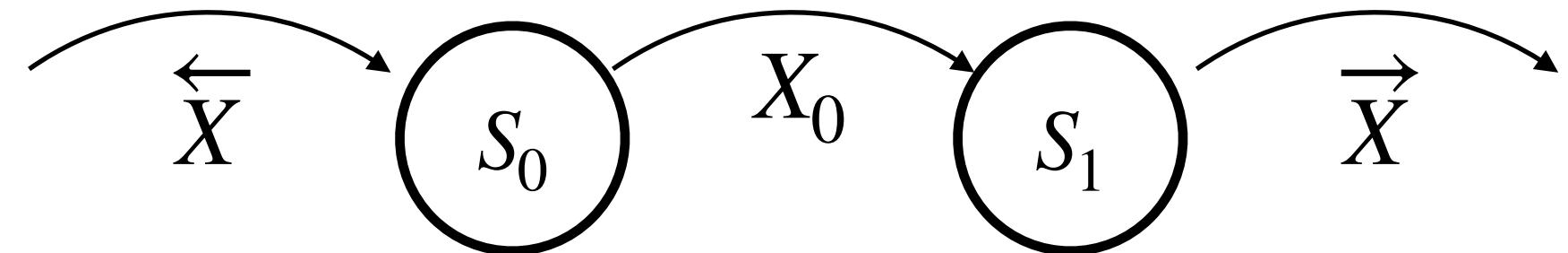
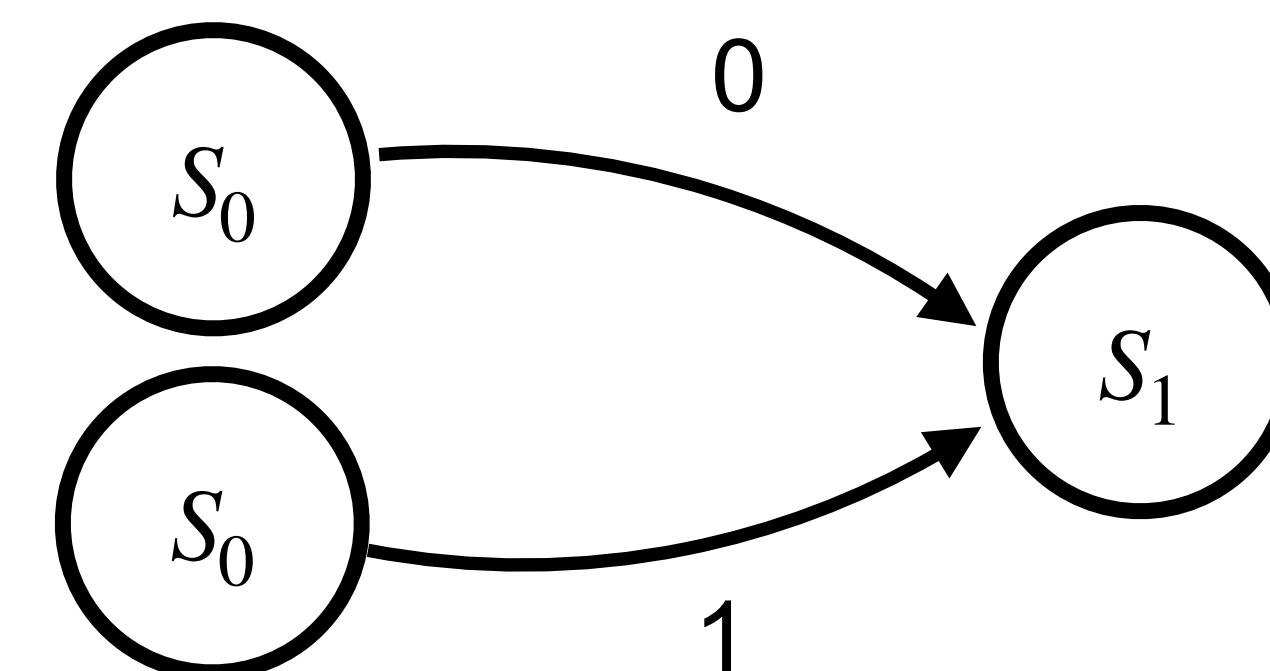
Information generated in present, not used for prediction.



*Binding information*

$$I[X_0; S_0 | S_1]$$

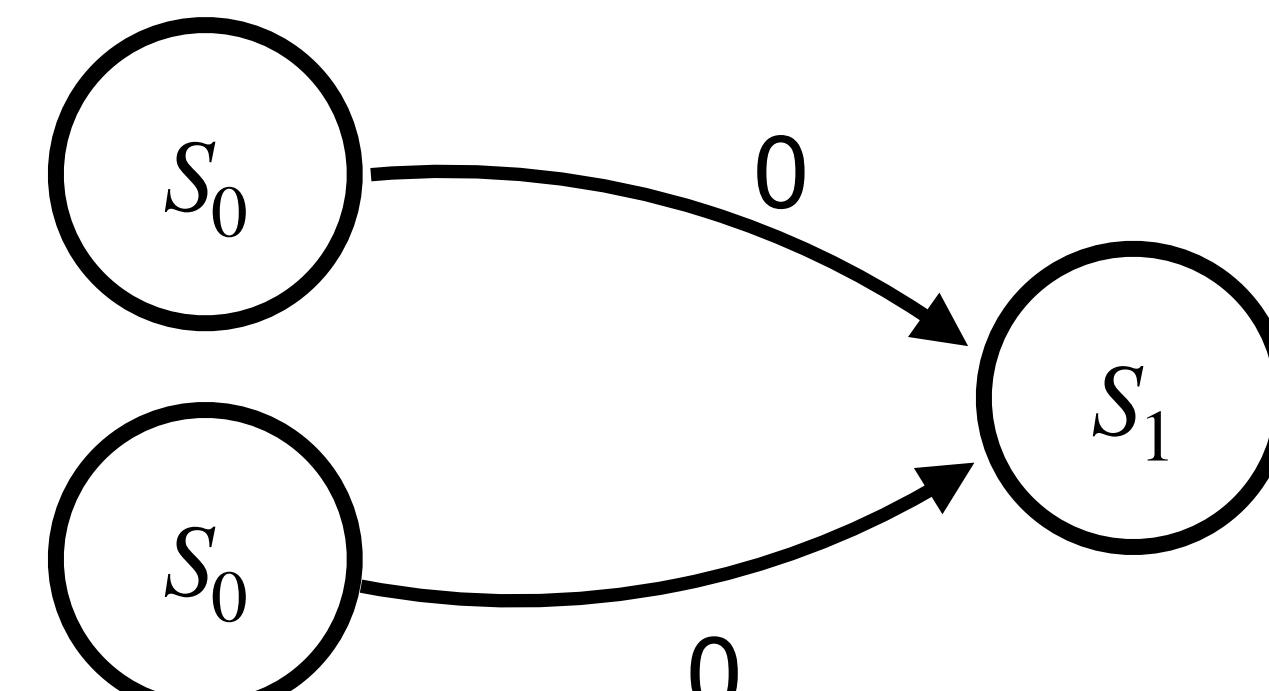
Information in the past and present that doesn't get used for future prediction.



*Crypticity*

$$H[S_0 | X_0, S_1]$$

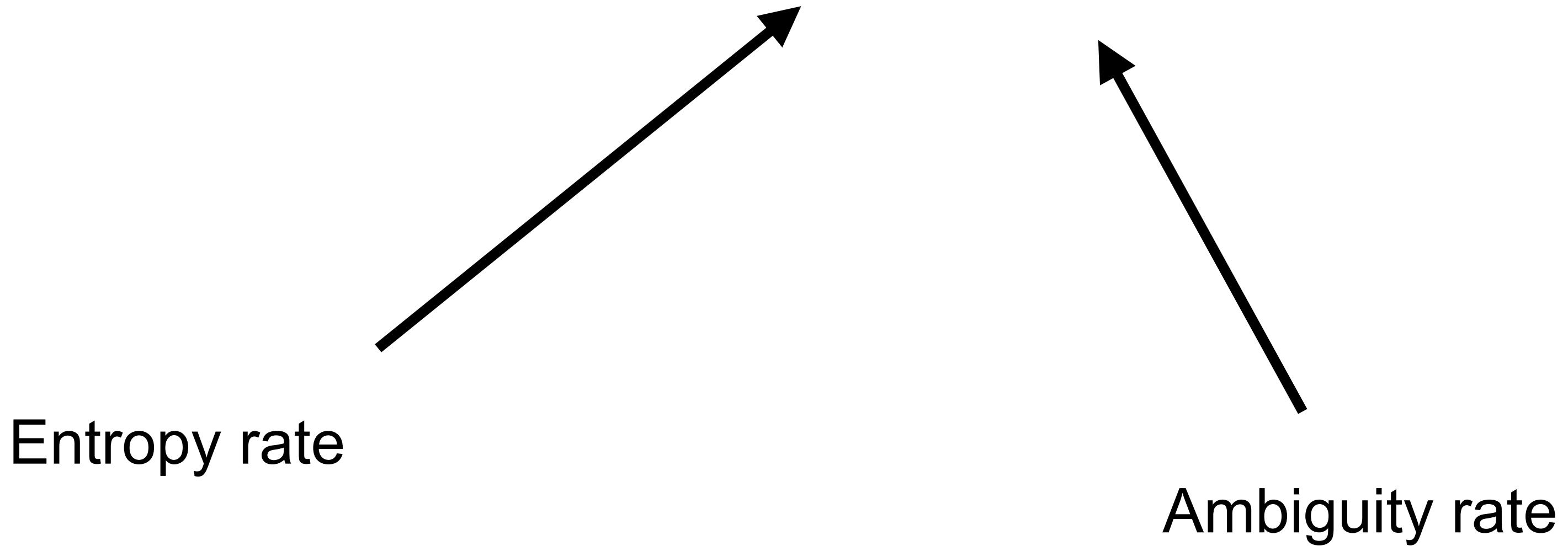
Information in past used for prediction but then forgotten for future prediction.



# $h_a$ as Model Growth Rate

---

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



# $h_a$ as Model Growth Rate

---

$$\begin{aligned} h_\mu - h_a &= H[X_t | S_t] - H[X_t, S_t | S_{t+1}] \\ &= H[S_{t+1} | S_t, X_t] + H[S_{t+1}] - H[S_t] \\ &= \Delta H[S] \end{aligned}$$

# $h_a$ as Model Growth Rate

---

$$\begin{aligned} h_\mu - h_a &= H[X_t | S_t] - H[X_t, S_t | S_{t+1}] \\ &= H[S_{t+1} | S_t, X_t] + H[S_{t+1}] - H[S_t] \\ &= \Delta H[S] \end{aligned}$$

# What is $h_a$ ?

---

$$h_a = H [X_0, S_0 \mid S_1]$$

$$H[X_0]$$

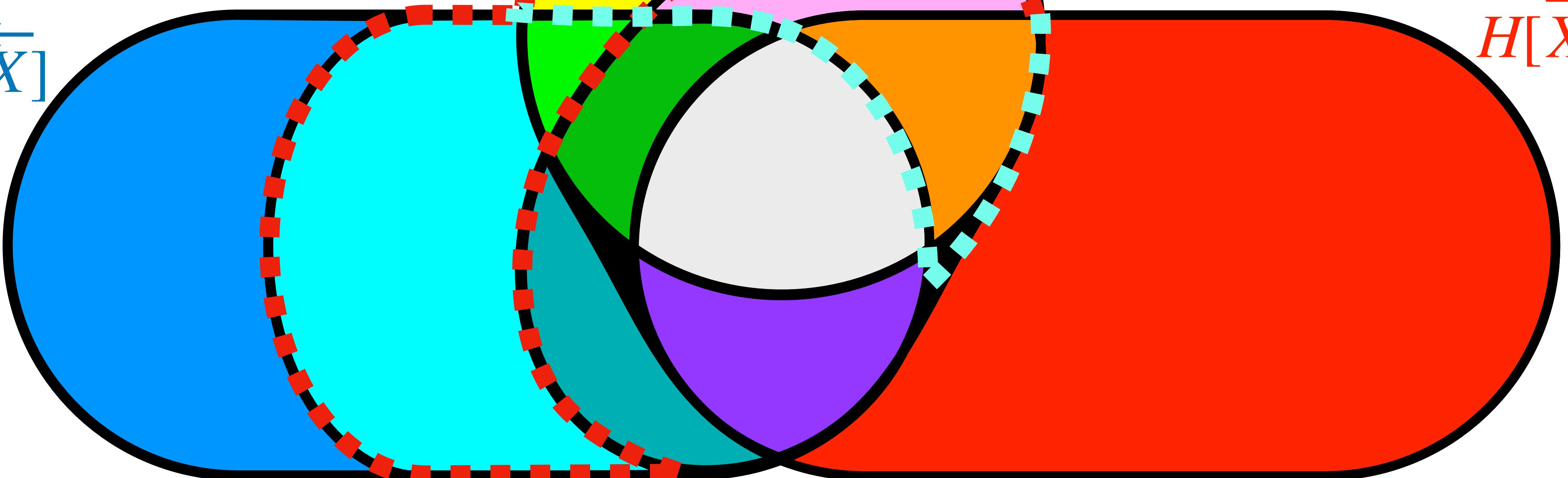
$$h_\mu = H [X_0 \mid S_0]$$

$$H[\vec{X}]$$

$$H[S_1]$$

$$H[\vec{X}]$$

$$H[S_0]$$

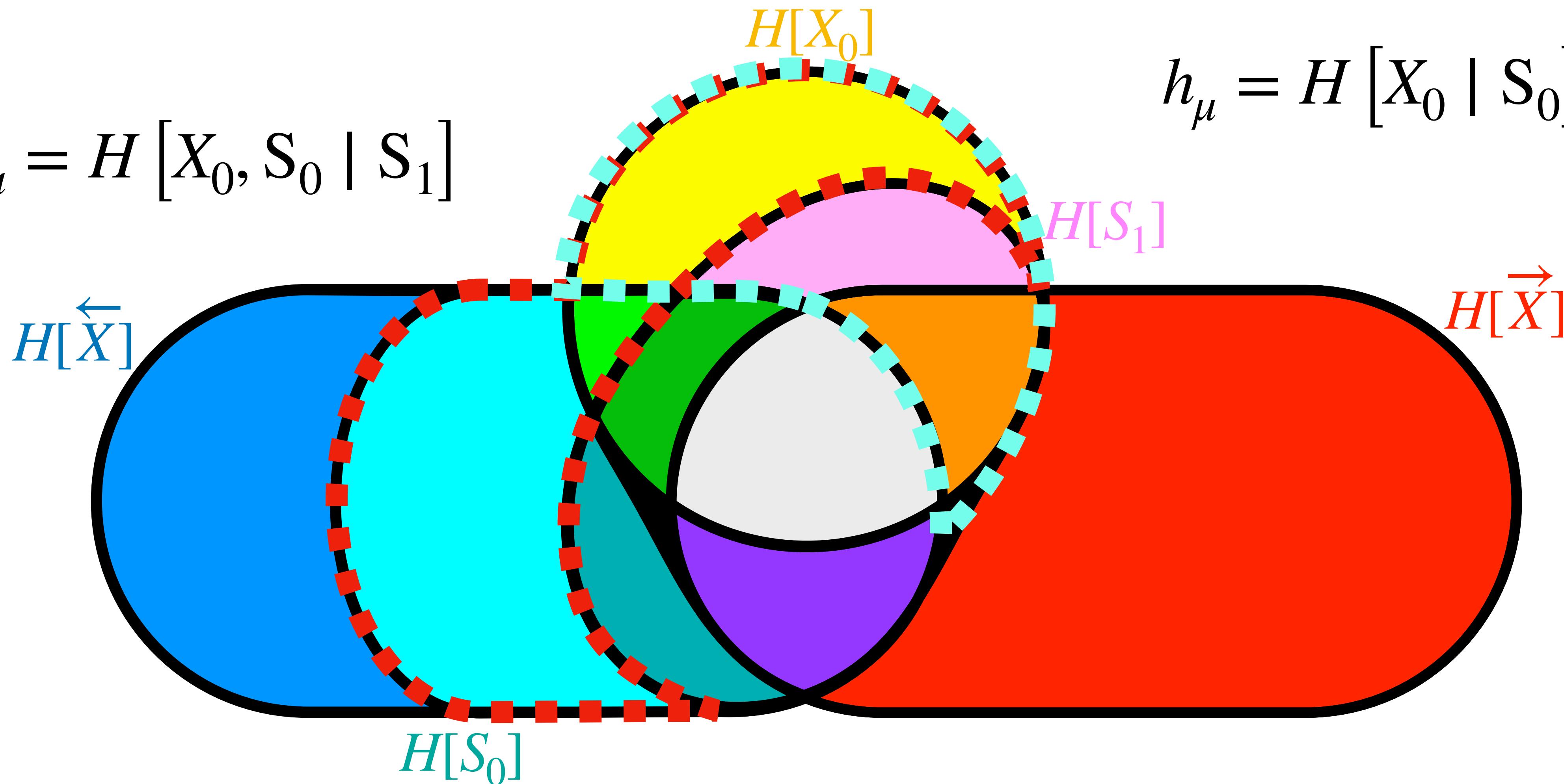


# What is $h_a$ ?

---

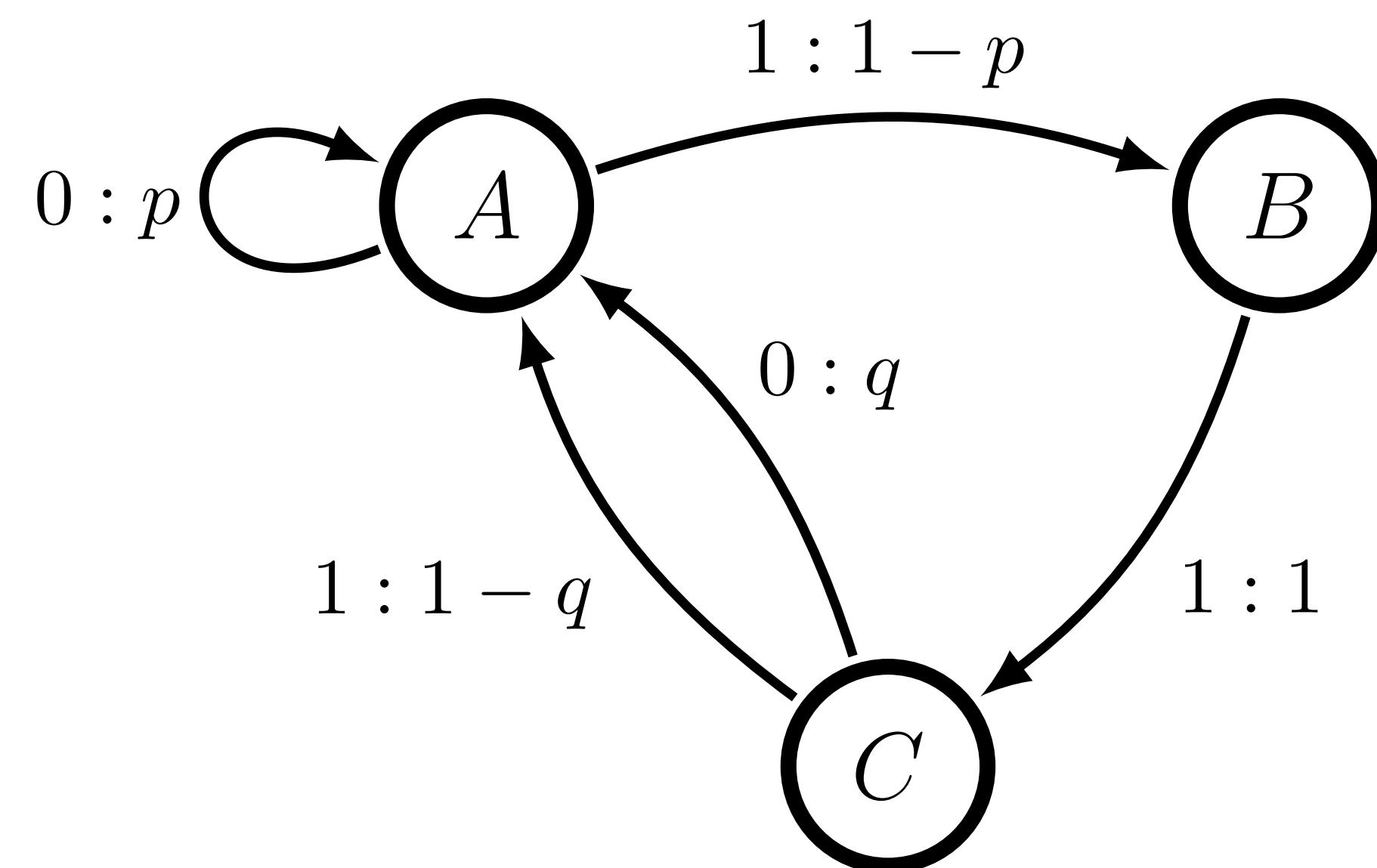
$$h_a = H [X_0, S_0 \mid S_1]$$

$$h_\mu = H [X_0 \mid S_0]$$



If  $\Delta H [S] = 0$ , then  $h_\mu = h_a$

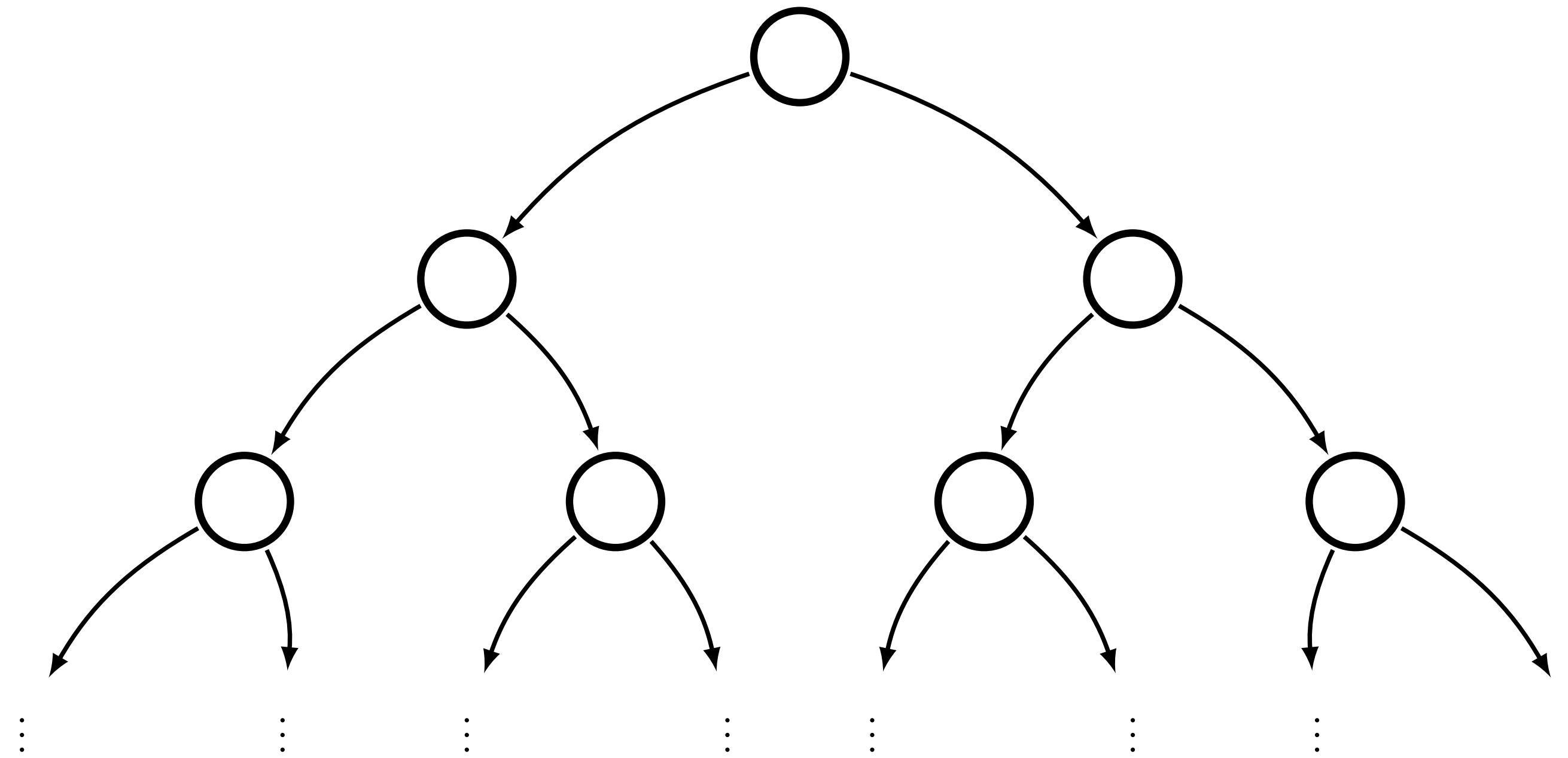
# Growth Rate of Optimally Predictive Models



For a finite  $C_\mu$  model:

$$h_\mu - h_a = 0$$

# Growth Rate of Optimally Predictive Models



# No uncertainty in past:

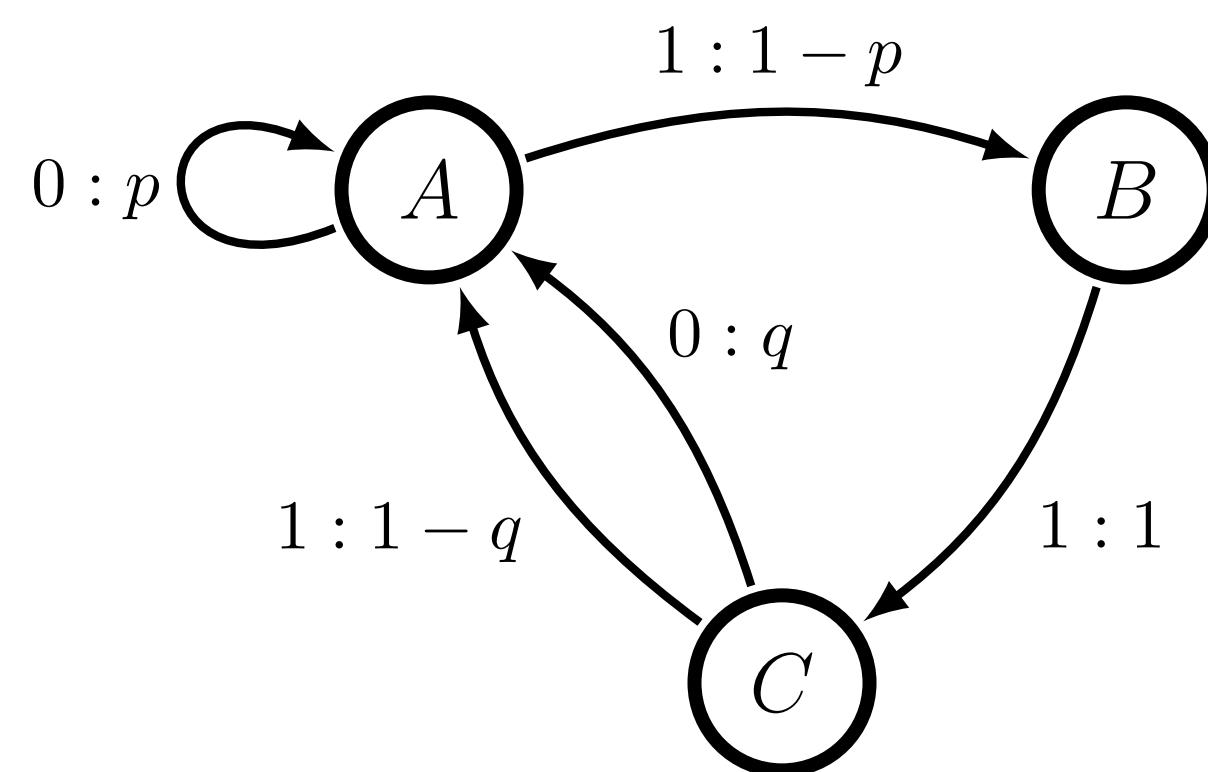
$$h_\mu - h_a = h_\mu$$

# Growth Rate of Optimally Predictive Models

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

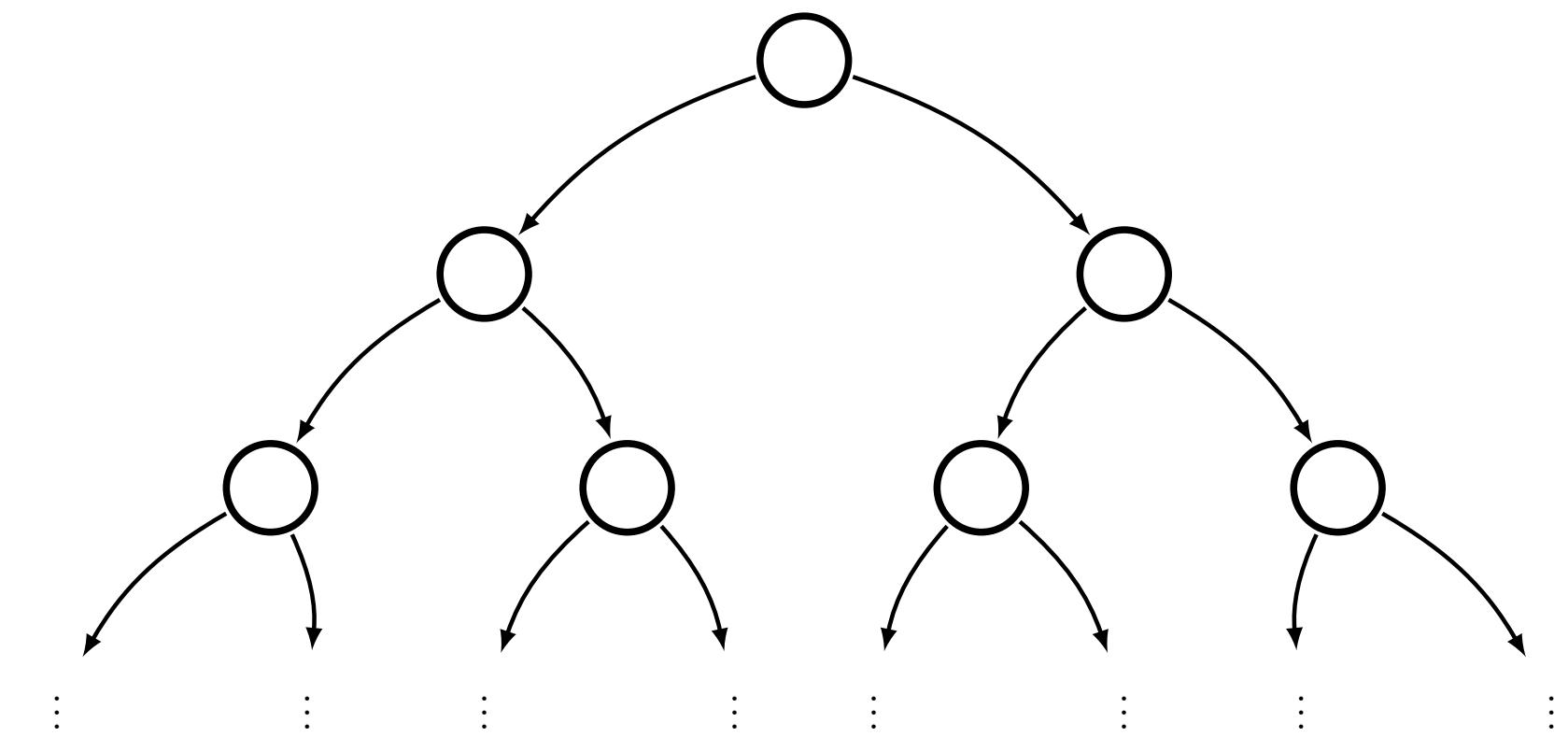
Finite states

$$h_\mu - h_a = 0$$



“Every history counts”

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

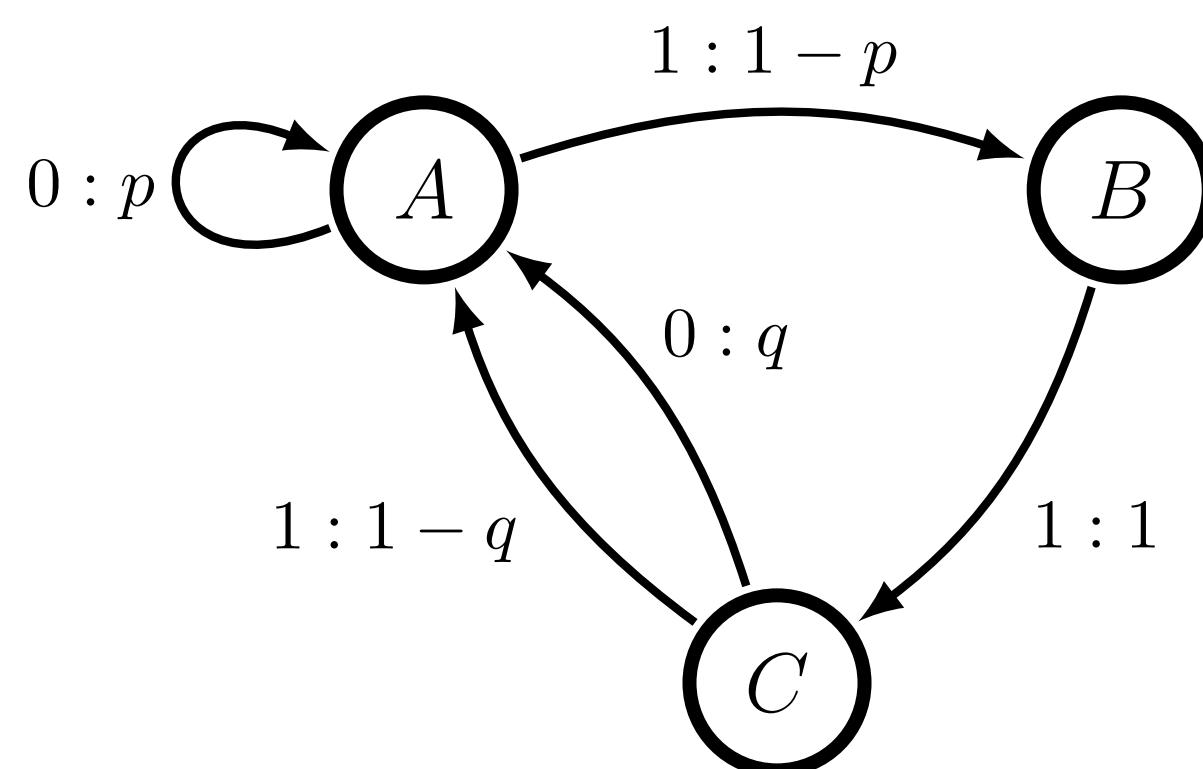


# Growth Rate of Optimally Predictive Models

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

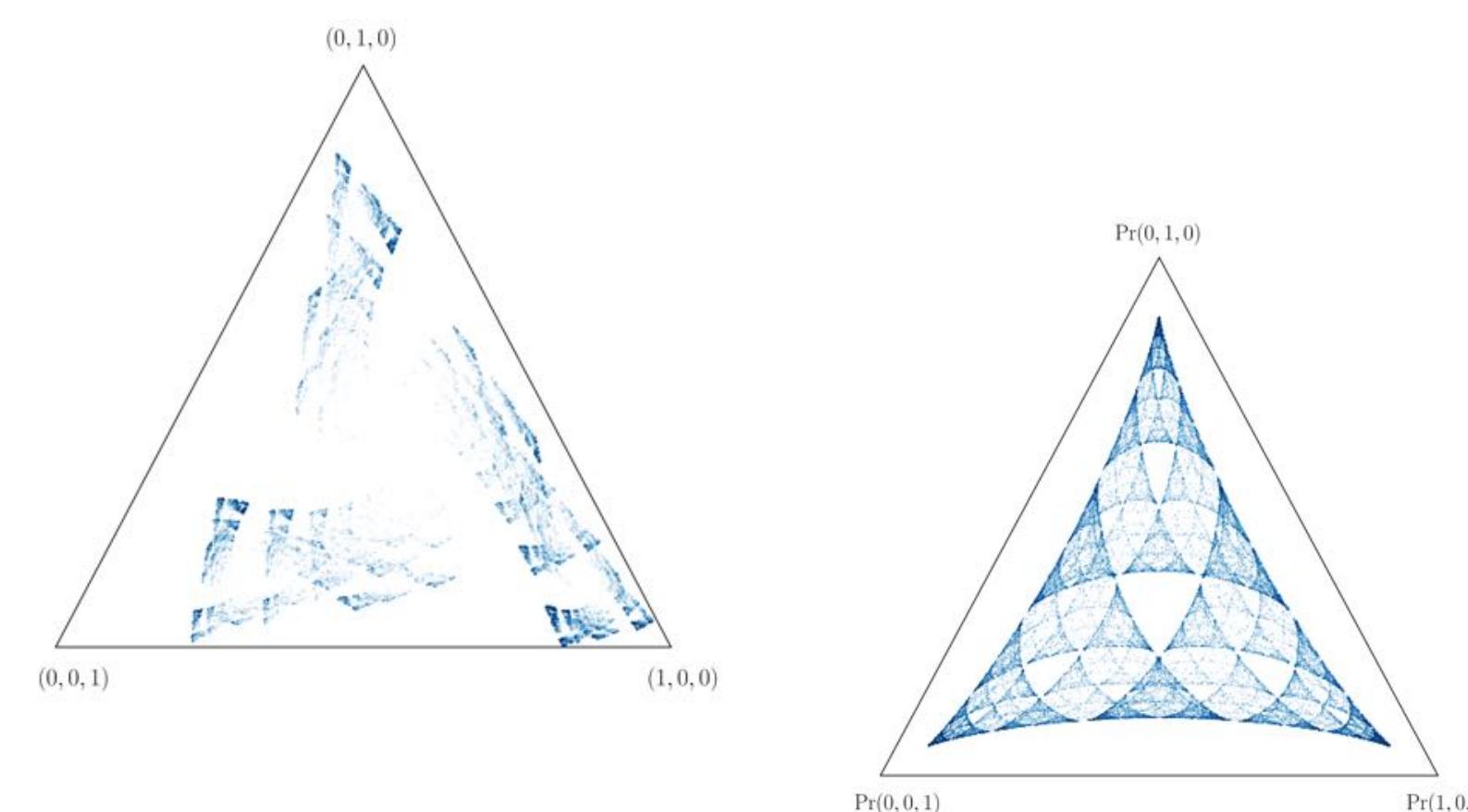
Finite states

$$h_\mu - h_a = 0$$



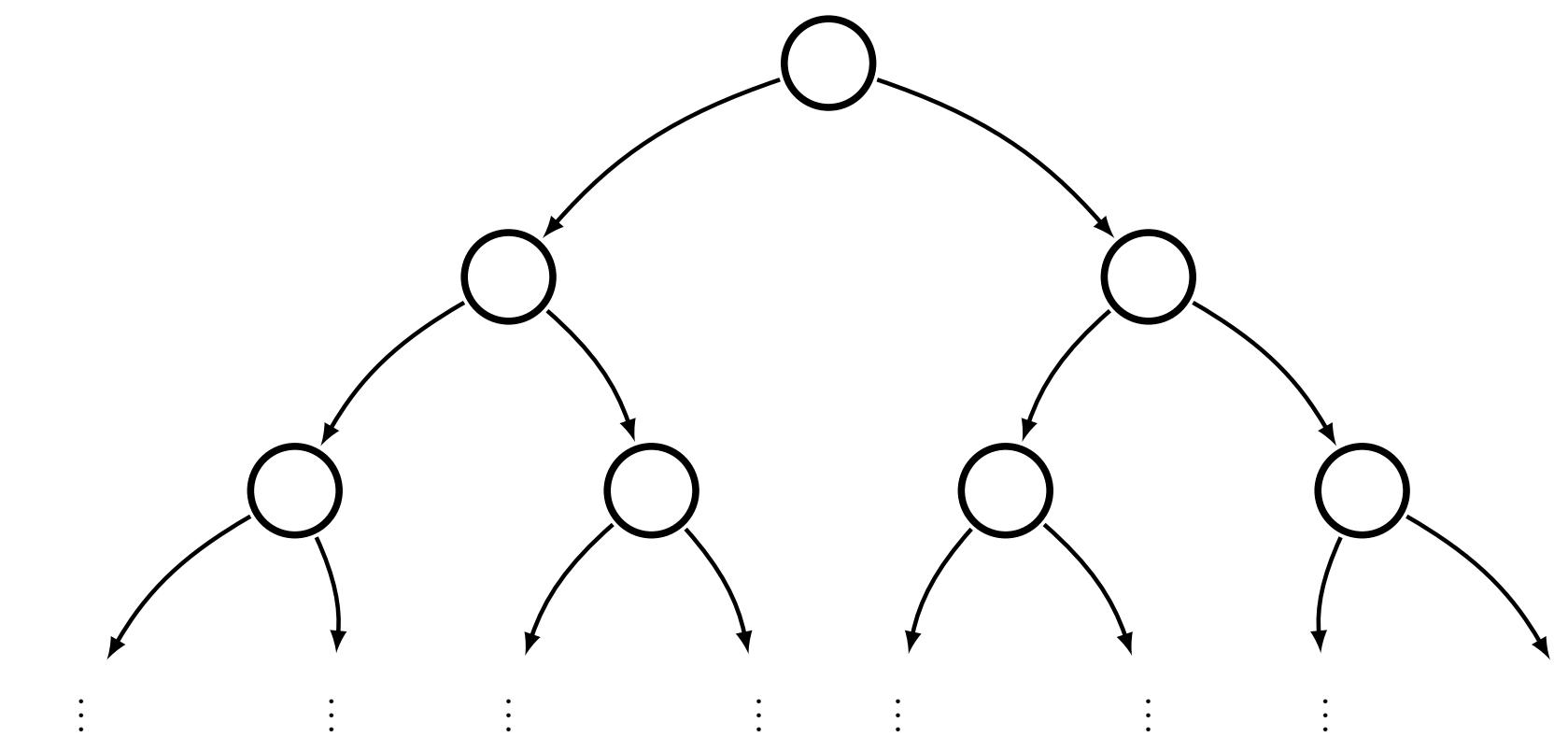
In general:

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

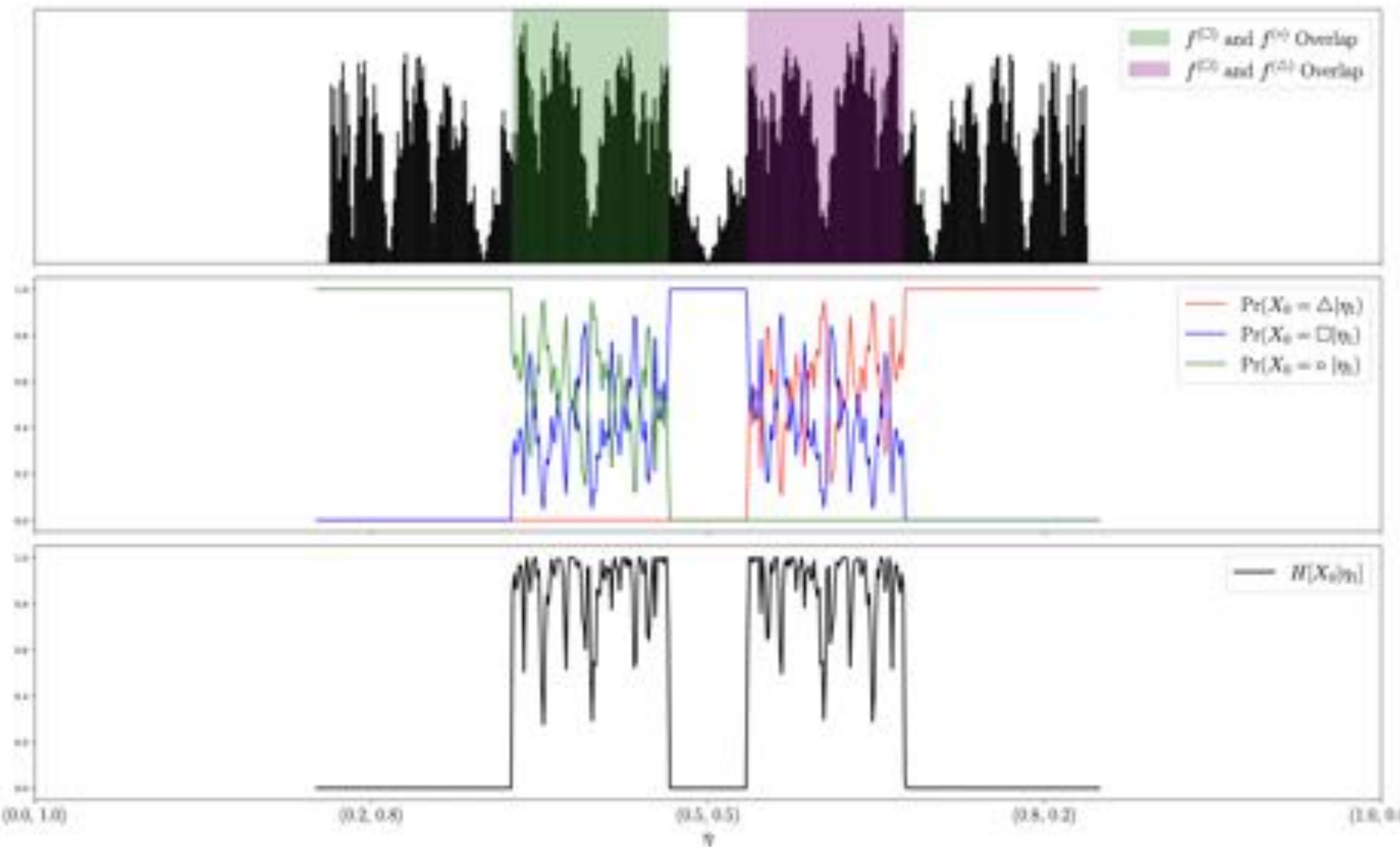


“Every history counts”

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



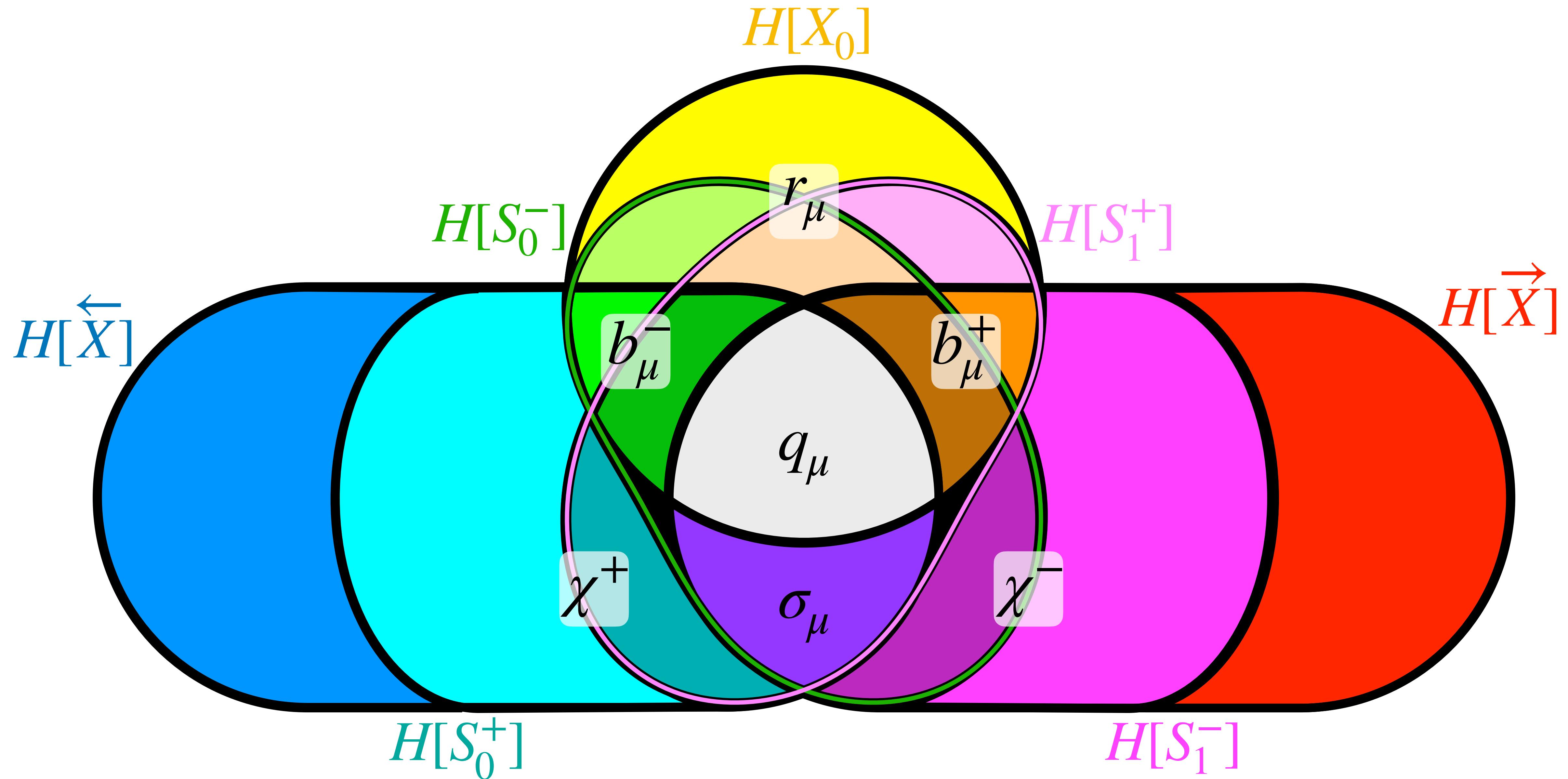
# Calculation of Ambiguity Rate?



Ulam's method....



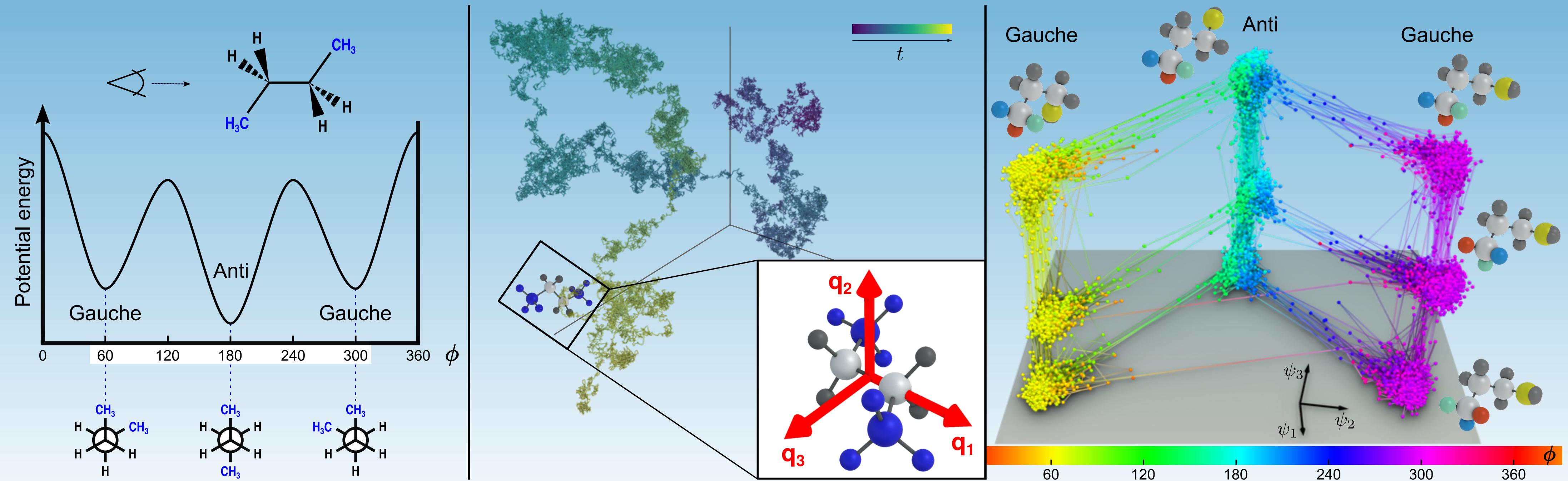
# This work inspired...





*(Secret Bonus Talk)*

# Geometric Perspectives on $\epsilon$ -Machines



# Context for this Talk

## Geometry from a Time Series

N. H. Packard, J. P. Crutchfield, J. D. Farmer, and R. S. Shaw

Dynamical Systems Collective, Physics Department, University of California, Santa Cruz, California 95064

(Received 13 November 1979)

It is shown how the existence of low-dimensional chaotic dynamical systems describing turbulent fluid flow might be determined experimentally. Techniques are outlined for reconstructing phase-space pictures from the observation of a single coordinate of any dissipative dynamical system, and for determining the dimensionality of the system's attractor. These techniques are applied to a well-known simple three-dimensional chaotic dynamical system.

PACS numbers: 47.25.-c



FIG. 1.  $(x,y)$  projection of Rossler (Ref. 7).



FIG. 2.  $(x,\dot{x})$  reconstruction from the time series.

Continuous and deterministic dynamic  $\rightarrow$  attractor  
reconstruction

# Context for this Talk

**Geometry from a Time Series**

N. H. Packard, J. P. Crutchfield, J. D. Farmer, and R. S. Shaw  
*Dynamical Systems Collective, Physics Department, University of California, Santa Cruz, California 95064*  
(Received 13 November 1979)

It is shown how the existence of low-dimensional chaotic dynamical systems describing turbulent fluid flow might be determined experimentally. Techniques are outlined for reconstructing phase-space pictures from the observation of a single coordinate of any dissipative dynamical system, and for determining the dimensionality of the system's attractor. These techniques are applied to a well-known simple three-dimensional chaotic dynamical system.

PACS numbers: 47.25.-c



FIG. 1.  $(x,y)$  projection of Rossler (Ref. 7).

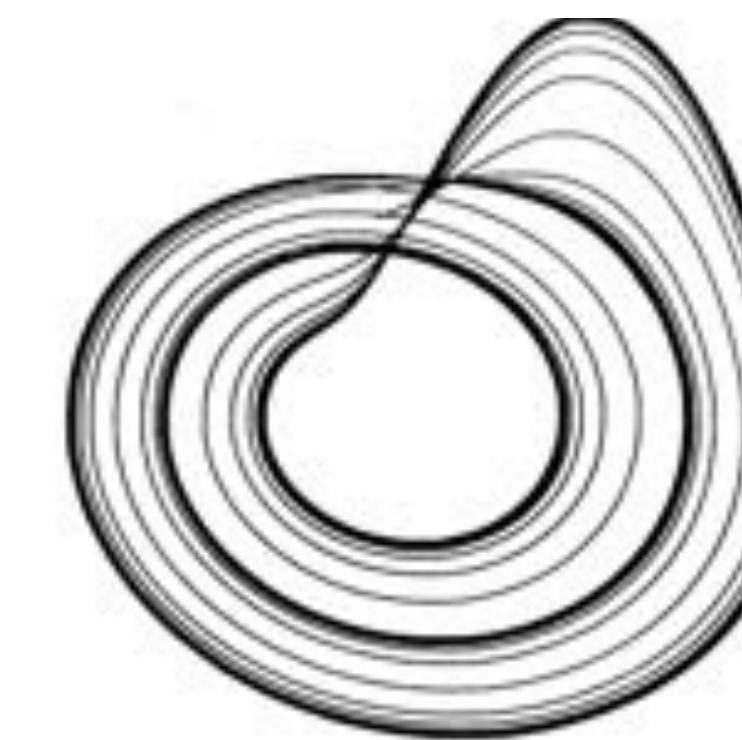


FIG. 2.  $(x,\dot{x})$  reconstruction from the time series.

Continuous and deterministic dynamic  $\rightarrow$  attractor reconstruction

## Transformers Represent Belief State Geometry in their Residual Stream

Adam S. Shai  
Simplex  
PIBBSS\*

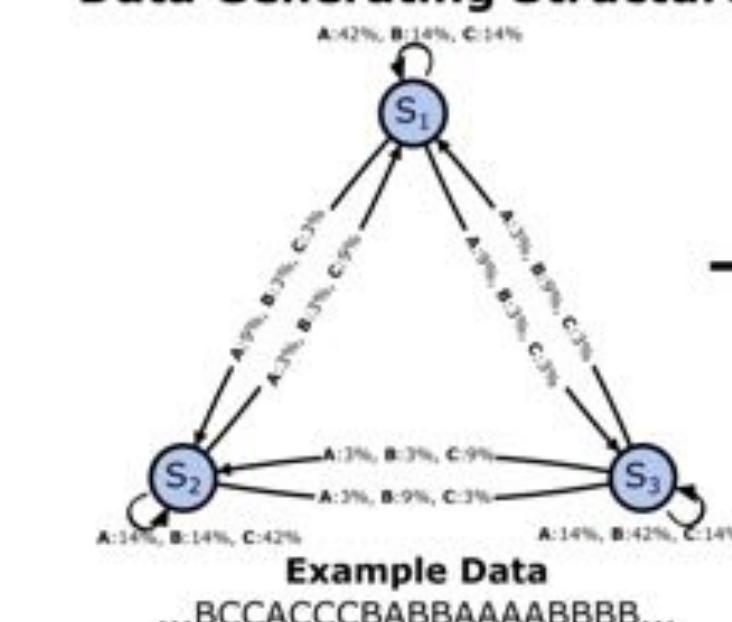
Sarah E. Marzen  
W. M. Keck Science Department  
Pitzer, Scripps, and Claremont McKenna College

Lucas Teixeira  
PIBBSS

Alexander Gietelink Oldenziel  
University College London  
Timaeus

Paul M. Riechers  
Simplex  
BITS

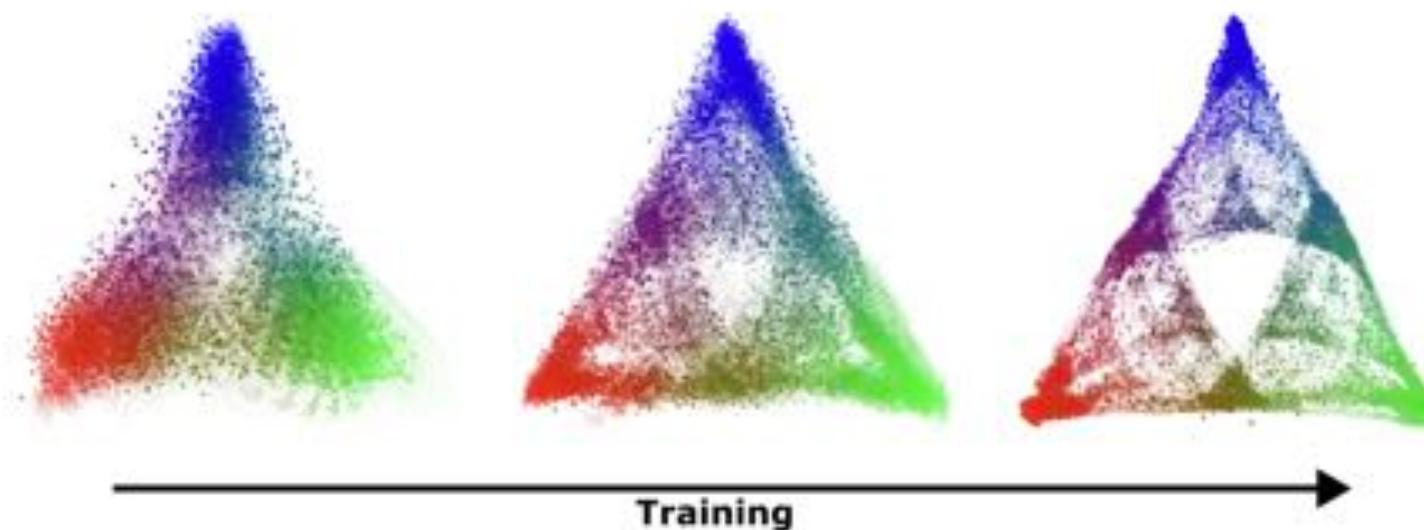
### Data Generating Structure



### Theoretical Prediction

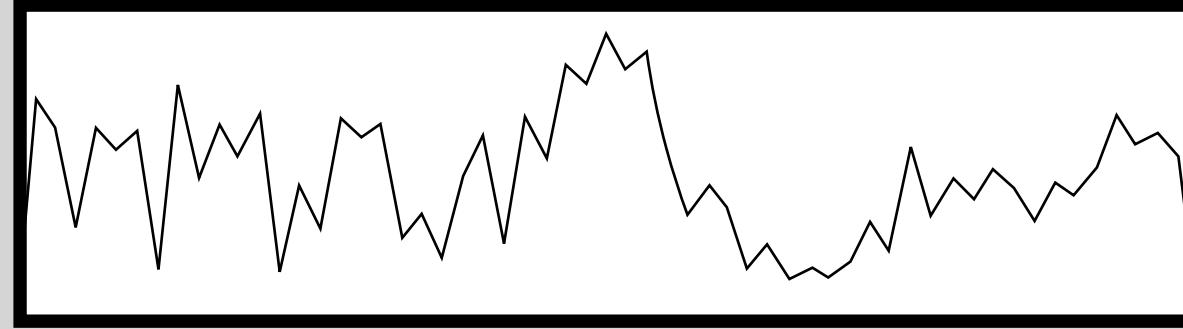


### Residual Stream Geometry



# Writing Down $\epsilon$ -Machines

Data



0100010000101110...

...  
Lorem ipsum dolor sit amen...

⋮

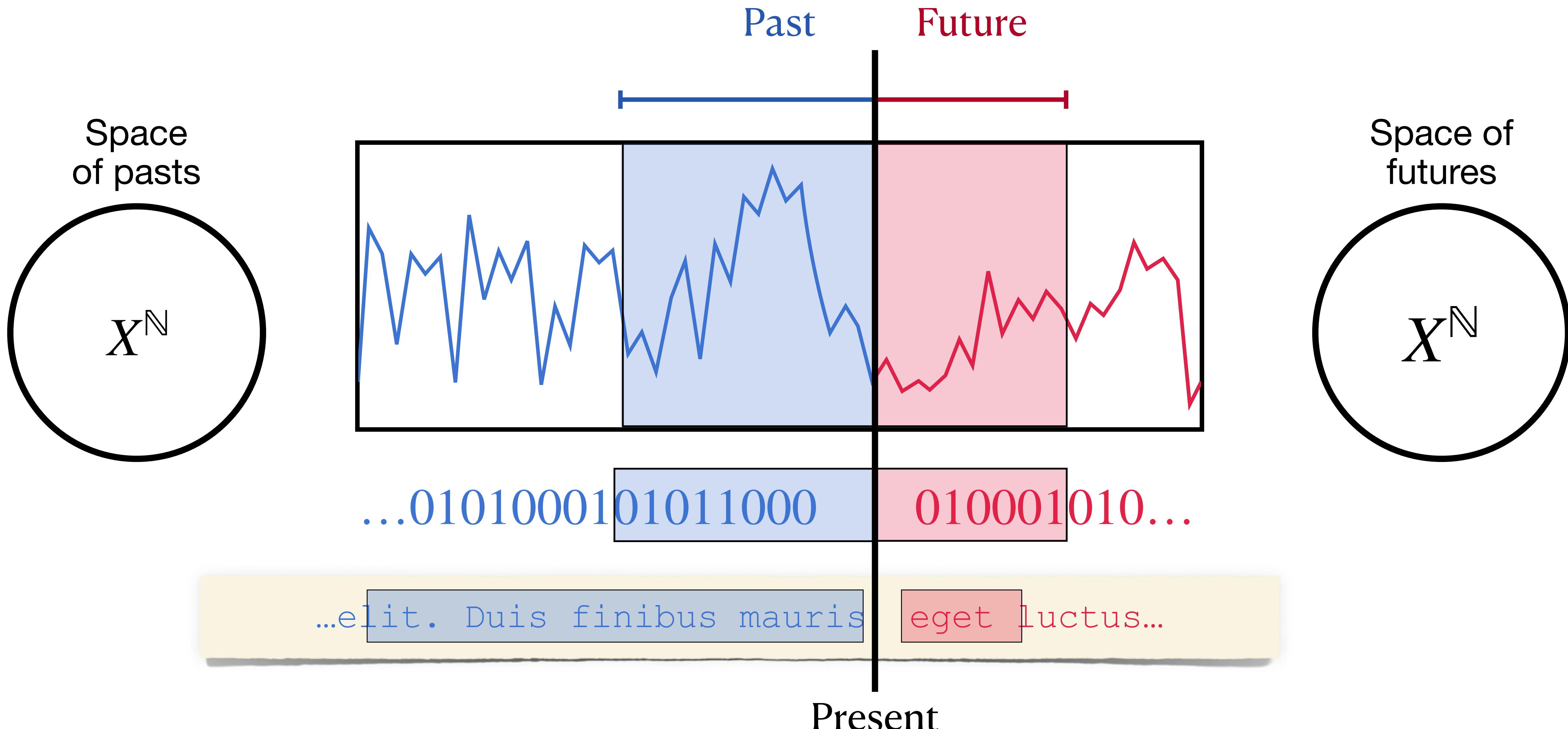
Causal States  $S$

- Via the predictive equivalence relation

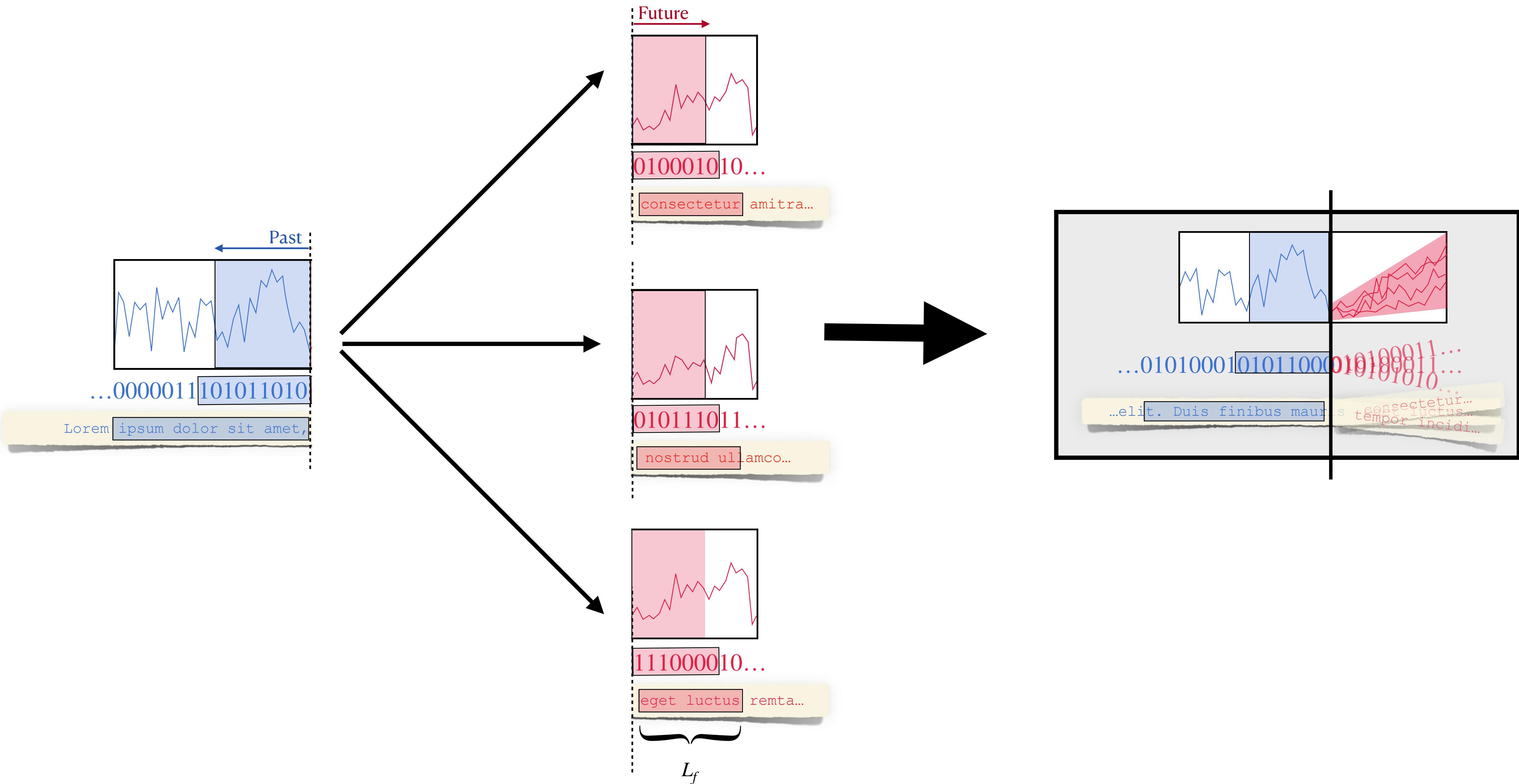
$\epsilon$ -Machine

- Dynamic over the causal states

# Pasts & Futures

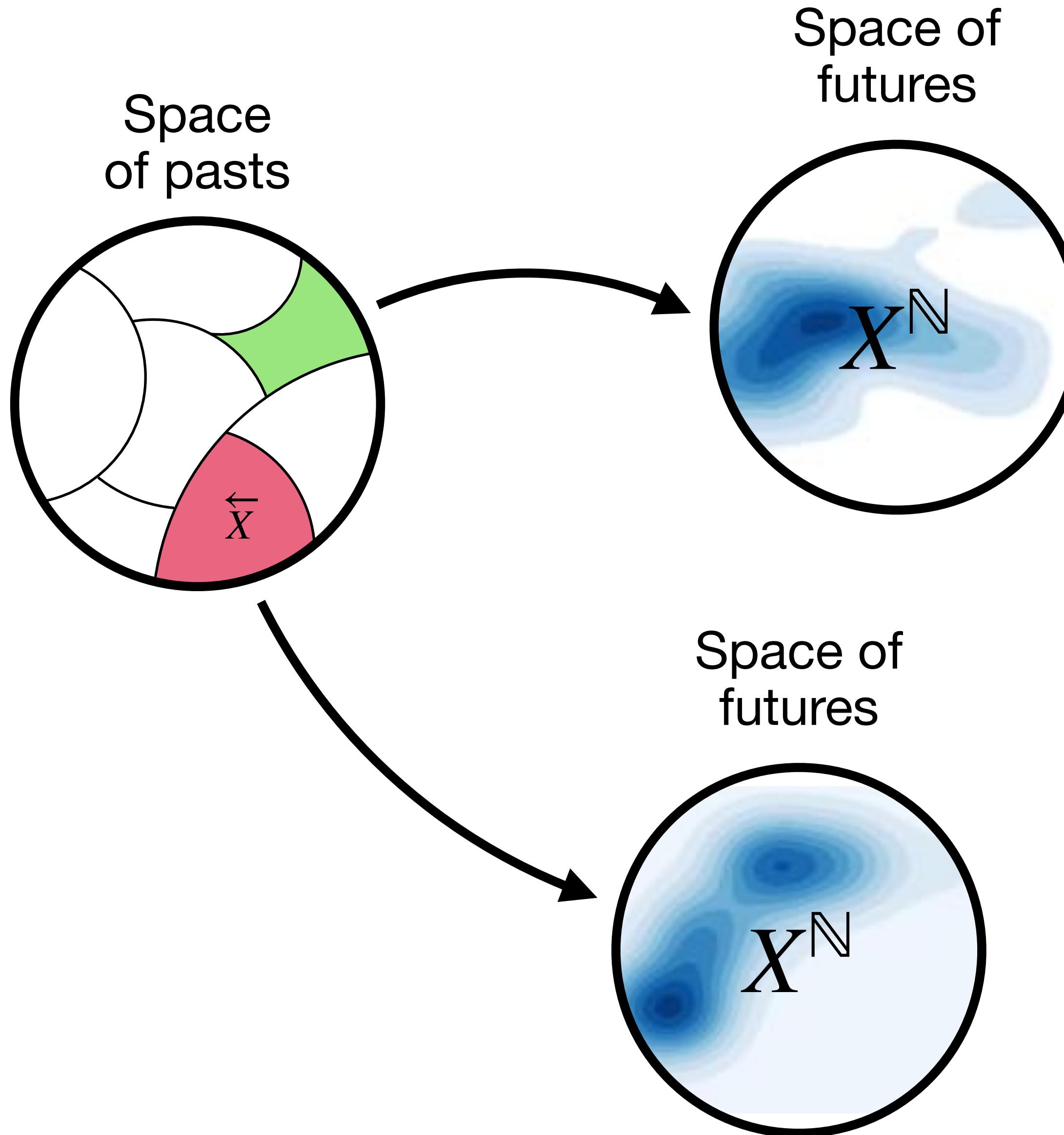


# Pasts Induce Distribution Over Futures



# Past Partitioning

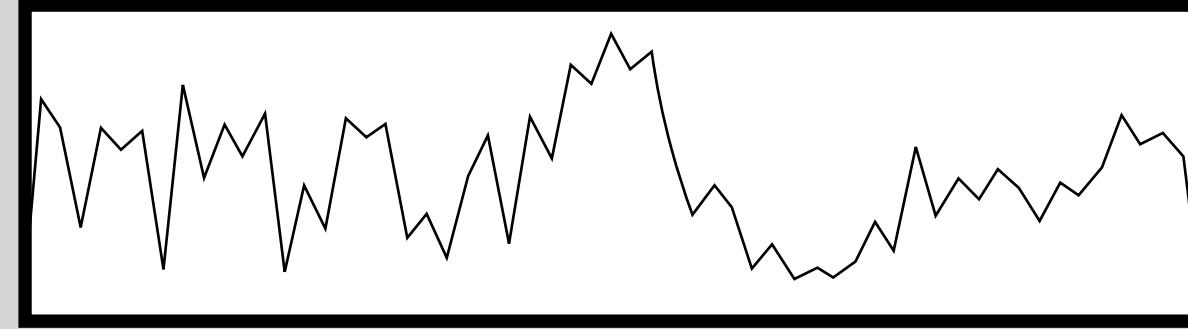
Predictive equivalence relation induces partition over pasts....



...where each element of the partition induces a unique distribution over futures.

# Writing Down $\epsilon$ -Machines

Data



0100010000101110...

...  
Lorem ipsum dolor sit amen...

⋮

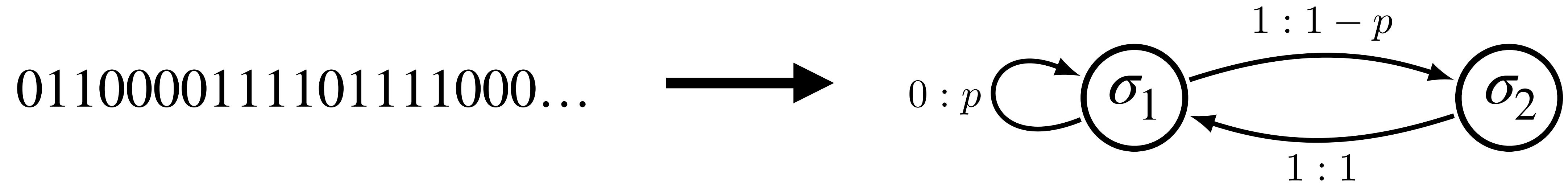
Causal States  $S$

- Via the predictive equivalence relation
- 1. By inspection

$\epsilon$ -Machine

- Dynamic over the causal states

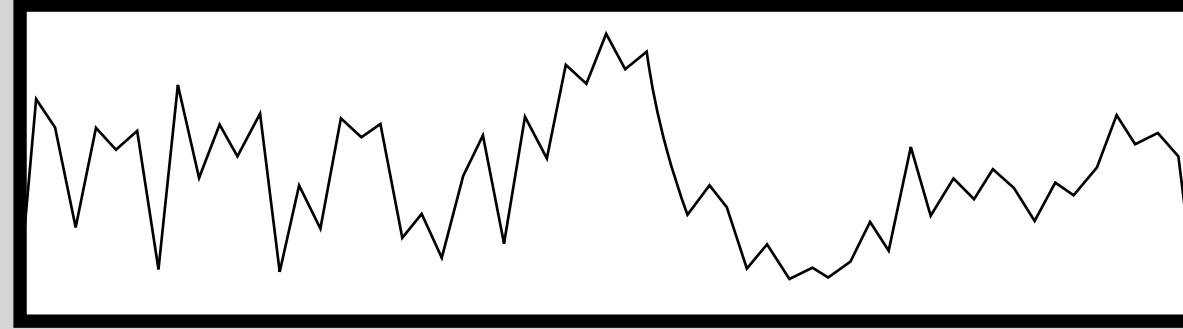
# $\epsilon$ -Machine: Finite States



Causal states + dynamic over states:  $\epsilon$ -machine

# Writing Down $\epsilon$ -Machines

Data



0100010000101110...

...  
Lorem ipsum dolor sit amen...

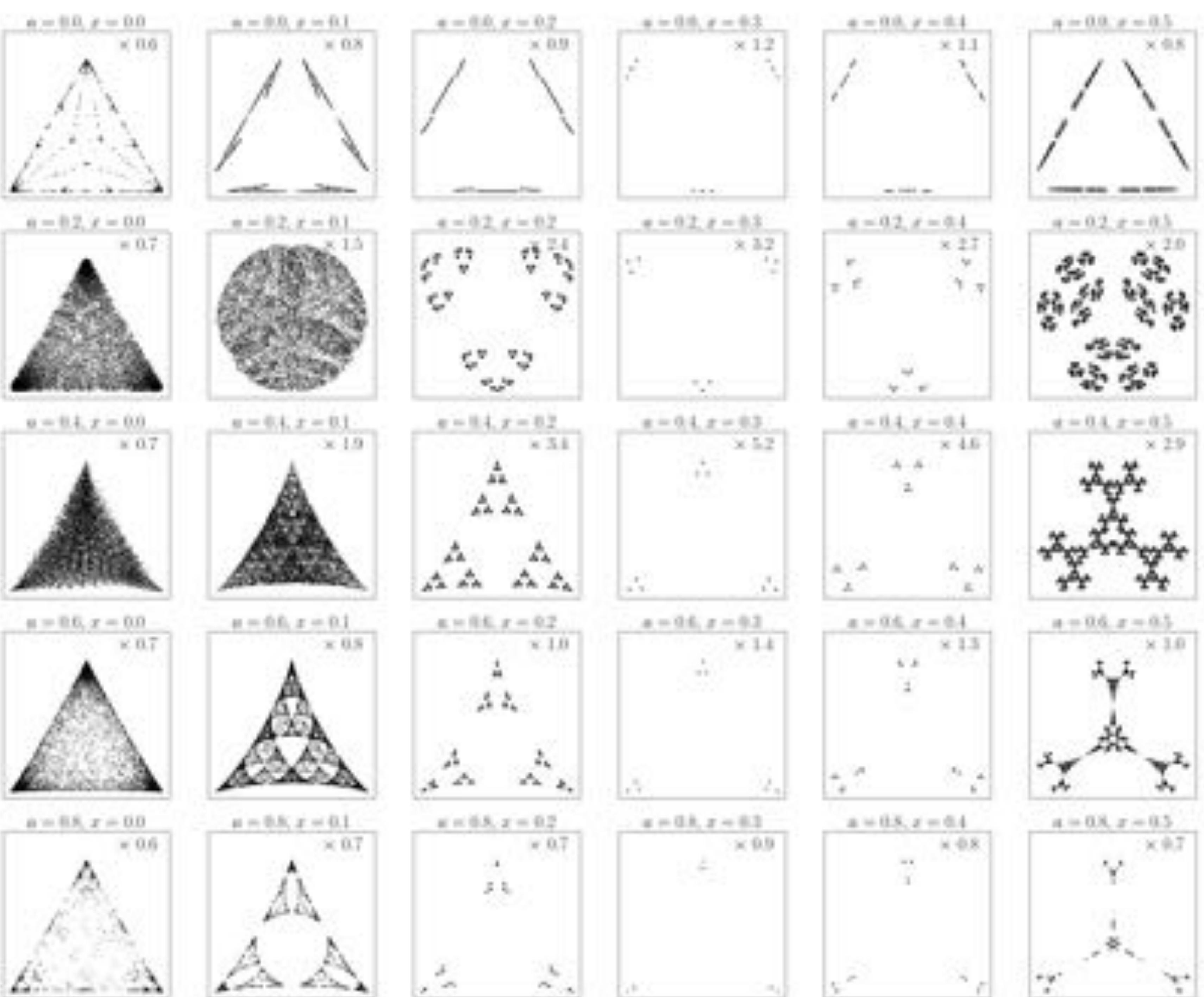
⋮

Causal States  $S$

- Via the predictive equivalence relation
- 1. By inspection
- 2. By generation in the space spanned by a generating model

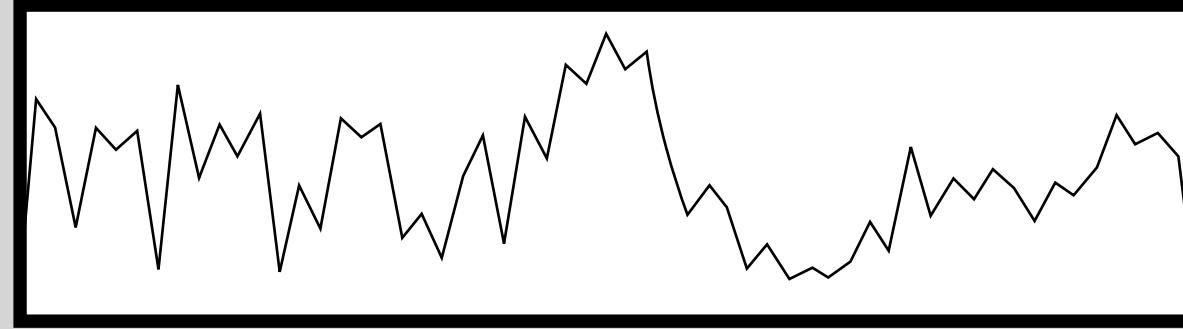
$\epsilon$ -Machine

- Dynamic over the causal states



# Writing Down $\epsilon$ -Machines

Data



0100010000101110...

...  
Lorem ipsum dolor sit amen...

⋮

Causal States  $S$

- Via the predictive equivalence relation
- 1. By inspection
- 2. By generation in the space spanned by a generating model
- 3. By embedding in a space

$\epsilon$ -Machine

- Dynamic over the causal states

# Cantor Embedded $\epsilon$ -Machines

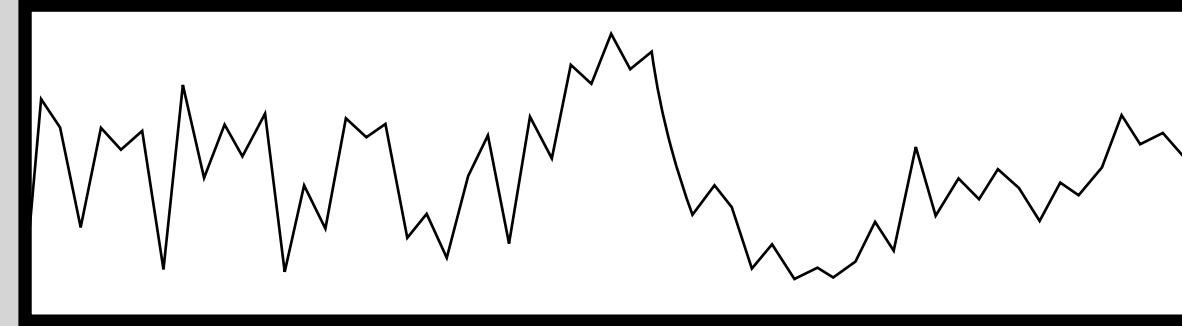
The screenshot shows a journal article page from the Chaos journal. The title of the article is "Exploring predictive states via Cantor embeddings and Wasserstein distance". The authors listed are Samuel P. Loomis<sup>a</sup> and James P. Crutchfield<sup>b</sup>. The article was submitted on June 10, 2022, and accepted on November 2, 2022, published online on December 5, 2022. The DOI is 10.1063/5.0102603. There are links for "View Article", "Export Citation", and "CrossRef". The text "Cite as: Chaos 32, 123115 (2022); doi: 10.1063/5.0102603" is also present.

Given a sequence of symbols  $x_1, x_2, x_3 \dots$  where each  $x_k$  is drawn from the alphabet  $X = \{0, 1, 2, \dots, N\}$ :

$$C(x_1, x_2, \dots) = \sum_{k=1}^{\infty} \frac{2x_k}{(2|X|-1)^k}$$

# Writing Down $\epsilon$ -Machines

Data



0100010000101110...

...  
Lorem ipsum dolor sit amen...

⋮

Causal States  $S$

- Via the predictive equivalence relation
- 1. By inspection
- 2. By generation in the space spanned by a generating model
- 3. By embedding in a space
  - A. Which space? What metric properties do we need to recreate the predictive equivalence relation?
  - B. Dimension reduction?

$\epsilon$ -Machine

- Dynamic over the causal states

# $\epsilon$ -Machines in a Hilbert space

Chaos

ARTICLE

[scitation.org/journal/cha](https://scitation.org/journal/cha)

## Discovering causal structure with reproducing-kernel Hilbert space $\epsilon$ -machines

Cite as: Chaos 32, 023103 (2022); doi: 10.1063/5.0062829  
 Submitted: 8 July 2021 · Accepted: 10 January 2022 ·  
 Published Online: 1 February 2022 · Publisher error corrected: 28 March 2022

Nicolas Brodu<sup>1,✉</sup>  and James P. Crutchfield<sup>2,✉</sup> 

AFFILIATIONS

<sup>1</sup>Geostat Team—Geometry and Statistics in Acquisition Data, INRIA 33405 Talence Cedex, France  
<sup>2</sup>Complexity Sciences Center and Department of Physics and Astronomy, University of California at Davis, California 95616, USA

<sup>✉</sup>Author to whom correspondence should be addressed: [nicolas.brodu@inria.fr](mailto:nicolas.brodu@inria.fr)  
<sup>✉</sup>Electronic mail: [chaos@ucdavis.edu](mailto:chaos@ucdavis.edu)

 ELSEVIER

Content

Physica D

journal homepage: [www.elsevier.com/locate/physd](http://www.elsevier.com/locate/physd)

Chaos

ARTICLE

[pubs.aip.org/aip/cha](https://pubs.aip.org/aip/cha)

## Inferring kernel $\epsilon$ -machines: Discovering structure in complex systems

Cite as: Chaos 35, 033162 (2025); doi: 10.1063/5.0242981  
 Submitted: 8 October 2024 · Accepted: 10 March 2025 ·  
 Published Online: 31 March 2025

Alexandra M. Jurgens<sup>1,✉</sup>  and Nicolas Brodu<sup>1,✉</sup> 

AFFILIATIONS

INRIA Bordeaux Sud Ouest, 33405 Talence Cedex, France

<sup>✉</sup>[alexandra.jurgens@inria.fr](mailto:alexandra.jurgens@inria.fr)  
<sup>✉</sup>Author to whom correspondence should be addressed: [nicolas.brodu@inria.fr](mailto:nicolas.brodu@inria.fr)



## Topology, convergence, and reconstruction of predictive states

Samuel P. Loomis, James P. Crutchfield<sup>\*</sup>

Complexity Sciences Center and Department of Physics and Astronomy, University of California at Davis, One Shields Avenue, Davis, CA 95616, United States of America

ARTICLE INFO

Article history:  
 Received 20 September 2021  
 Received in revised form 10 December 2022  
 Accepted 17 December 2022  
 Available online 3 January 2023  
 Communicated by G. Froyland

Keywords:  
 Stochastic process  
 Symbolic dynamics  
 Dynamical systems

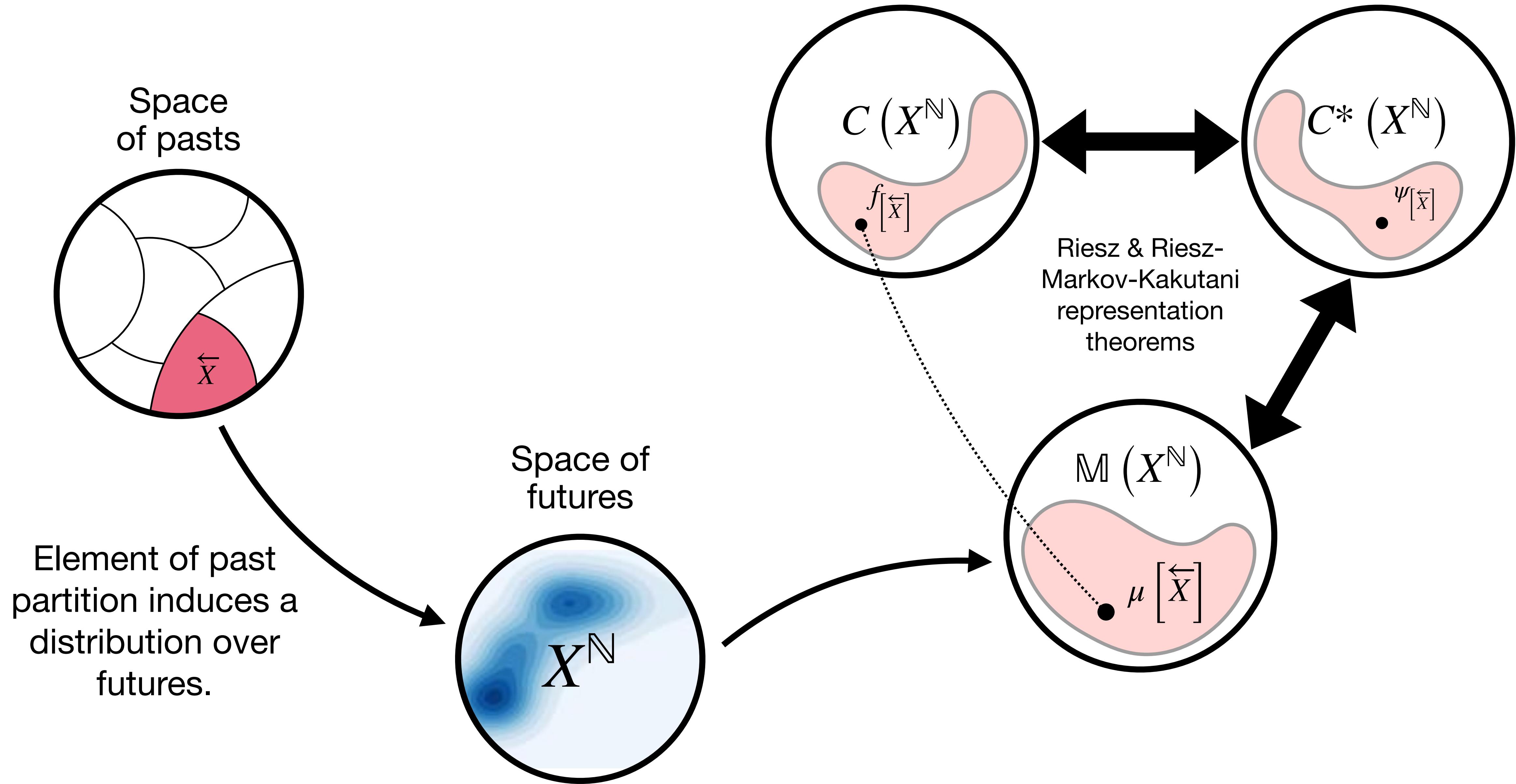
ABSTRACT

Predictive equivalence in discrete stochastic processes has been applied with great success to identify randomness and structure in statistical physics and chaotic dynamical systems and to inferring hidden Markov models. We examine the conditions under which predictive states can be reliably reconstructed from time-series data, showing that convergence of predictive states can be achieved from empirical samples in the weak topology of measures. Moreover, predictive states may be represented in Hilbert spaces that replicate the weak topology. We mathematically explain how these representations are particularly beneficial when reconstructing high-memory processes and connect them to reproducing kernel Hilbert spaces.

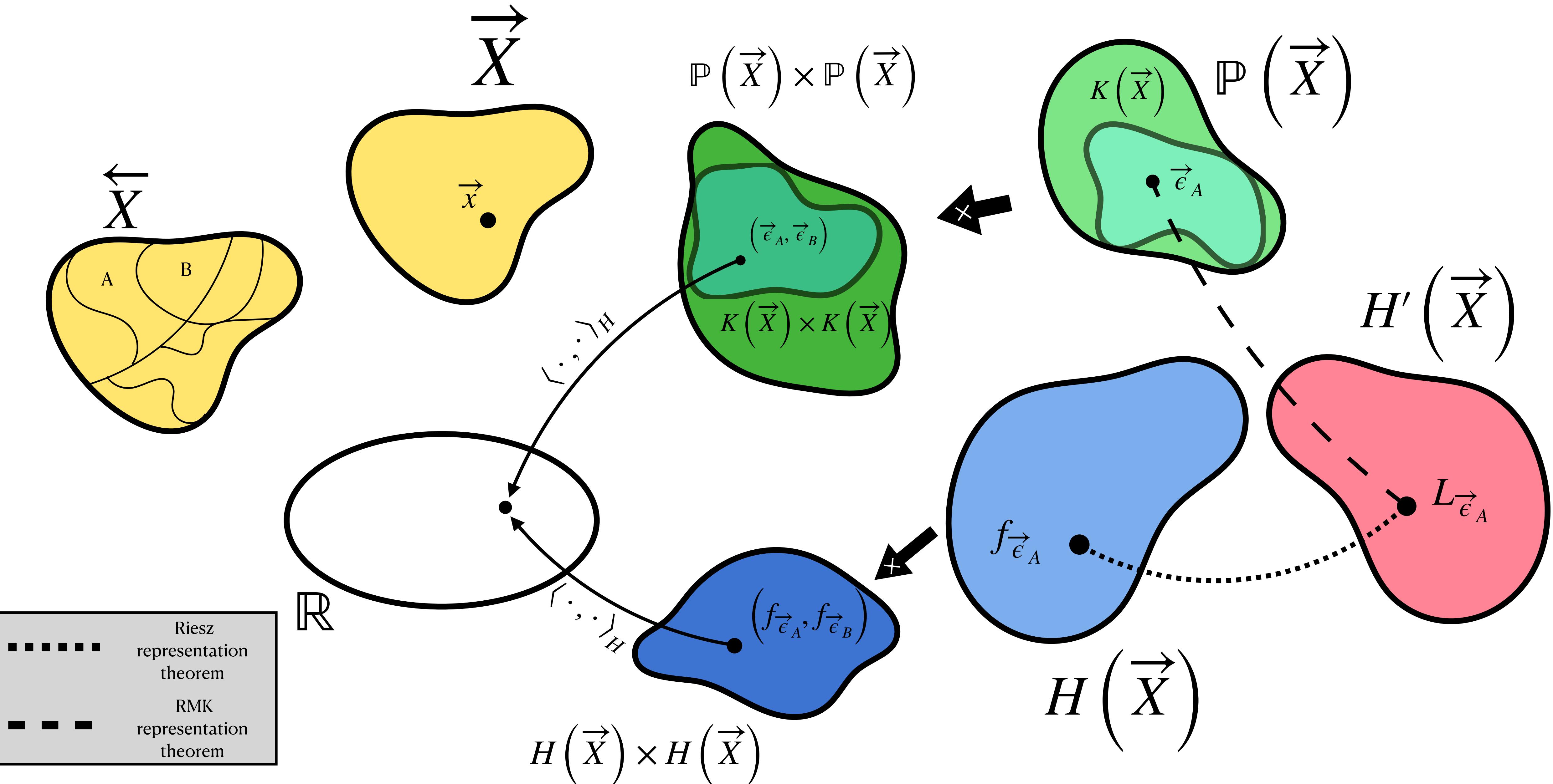
© 2022 The Author(s). Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<https://creativecommons.org/licenses/by-nd/4.0/>)

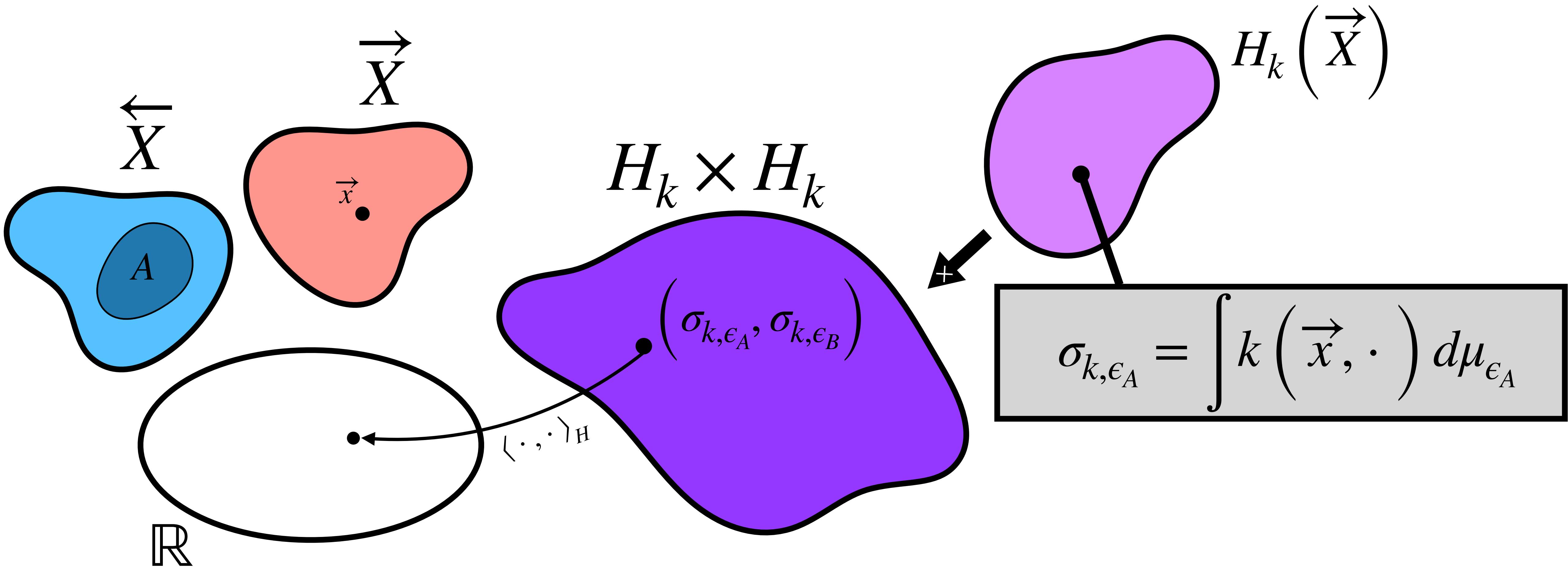
# Causal States in a Hilbert space



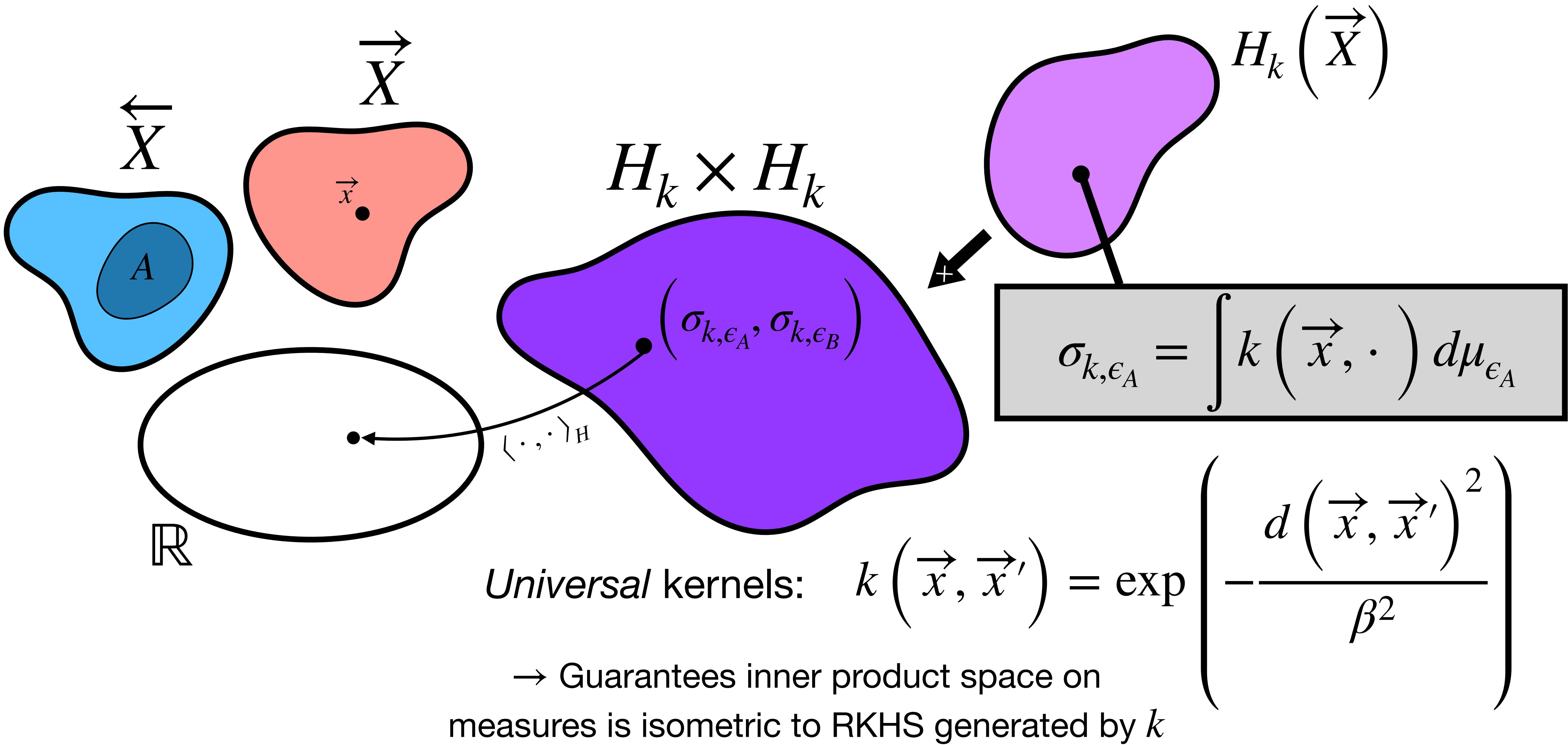
# Kernel Causal States



# Kernel Causal States



# Kernel Causal States

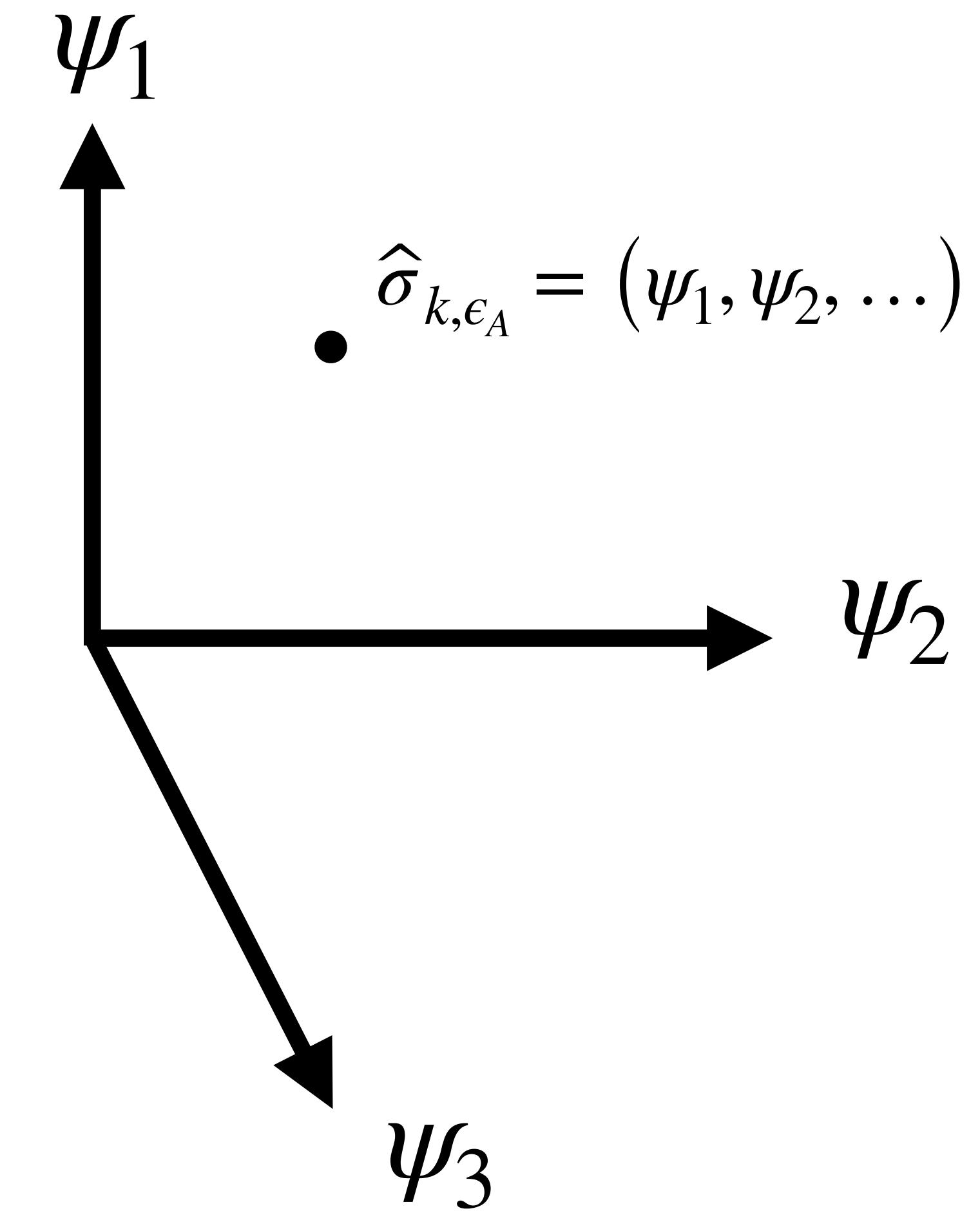


# Dimension Reduction...?

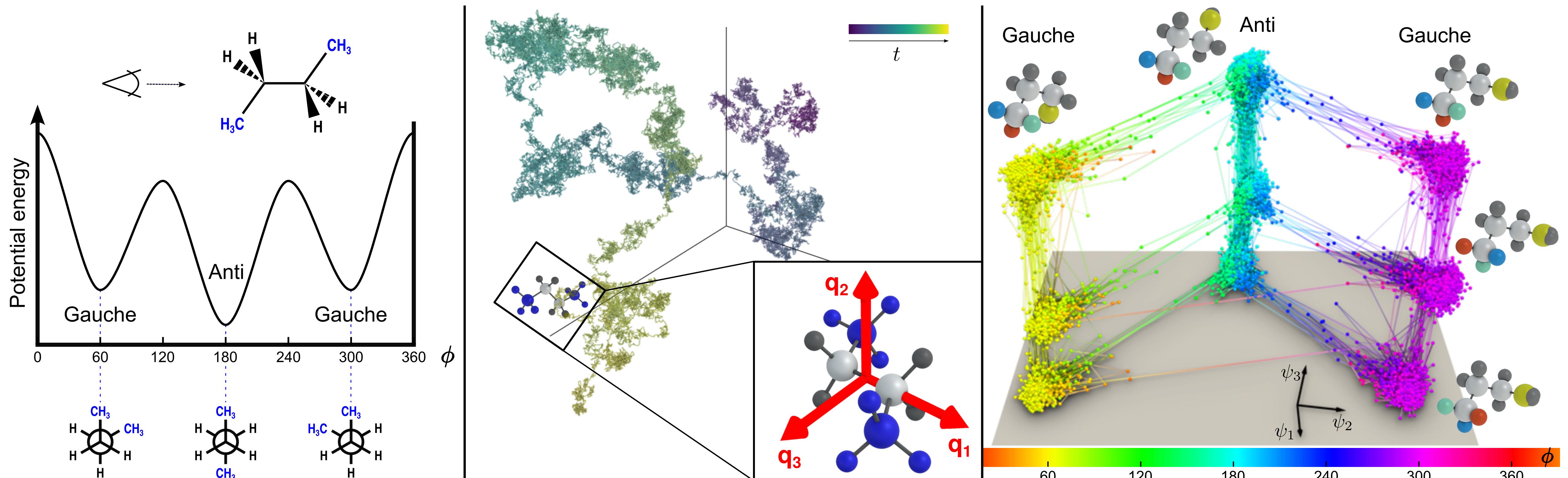
A diagram showing a purple-shaded region  $S$  with a black boundary. Above it, the expression  $H_k(\vec{X})$  is written. Below the region  $S$ , a box contains the formula:

$$\sigma_{k,\epsilon_A} = \int k\left(\vec{x}, \cdot\right) d\mu_{\epsilon_A}$$

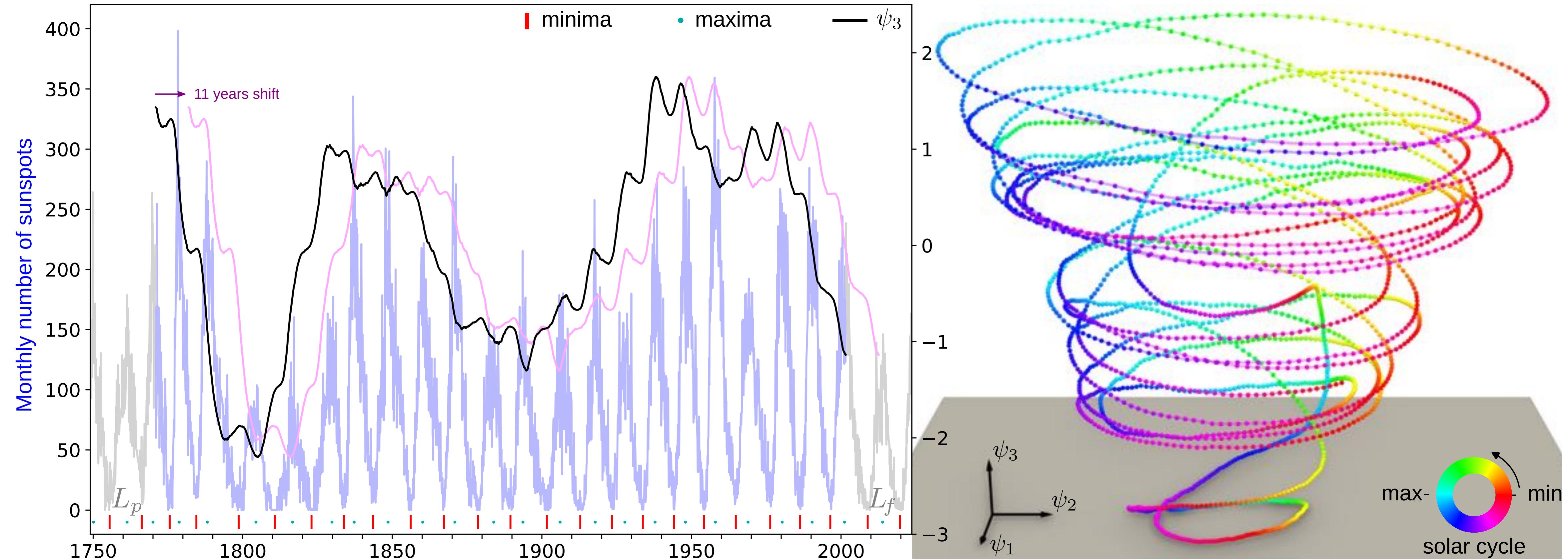
Current technique:  
- diffusion mapping



# Kernel Causal States: *n*-Butane

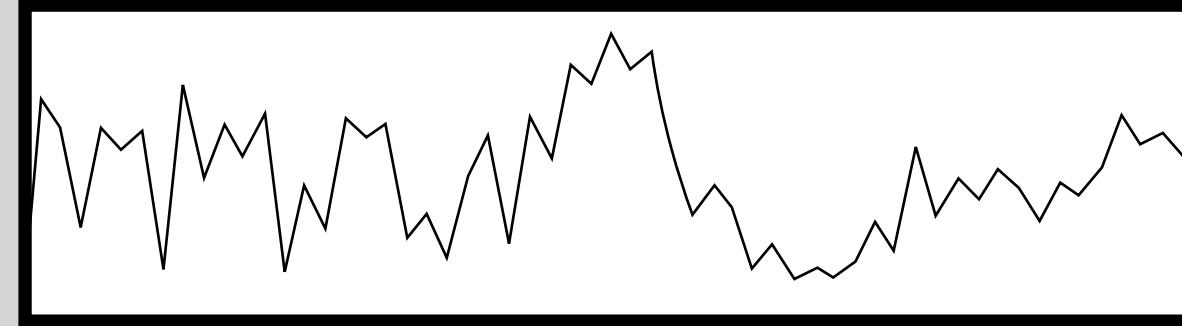


# Kernel Causal States: Sunspot Sequence



# Writing Down $\epsilon$ -Machines

Data



0100010000101110...

...  
Lorem ipsum dolor sit amen...

⋮

Causal States  $S$

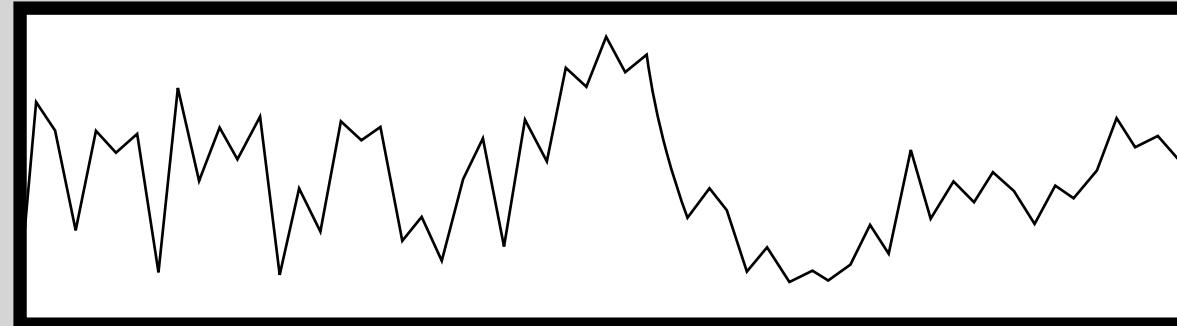
- Via the predictive equivalence relation
- 1. By inspection
- 2. By generation in the space spanned by a generating model
- 3. By embedding in a space
  - A. Which space? What metric properties do we need to recreate the predictive equivalence relation?
  - B. Dimension reduction?

$\epsilon$ -Machine

- Dynamic over the causal states?

# Writing Down $\epsilon$ -Machines

Data



0100010000101110...

...  
Lorem ipsum dolor sit amen...

⋮

Causal States  $S$

- Via the predictive equivalence relation
- 1. By inspection
- 2. By generation in the space spanned by a generating model
- 3. By embedding in a space
  - A. Which space? What metric properties do we need to recreate the predictive equivalence relation?
  - B. Dimension reduction?

$\epsilon$ -Machine

- Dynamic over the causal states?

Complexity & Information Measures

- Information measure zoo?

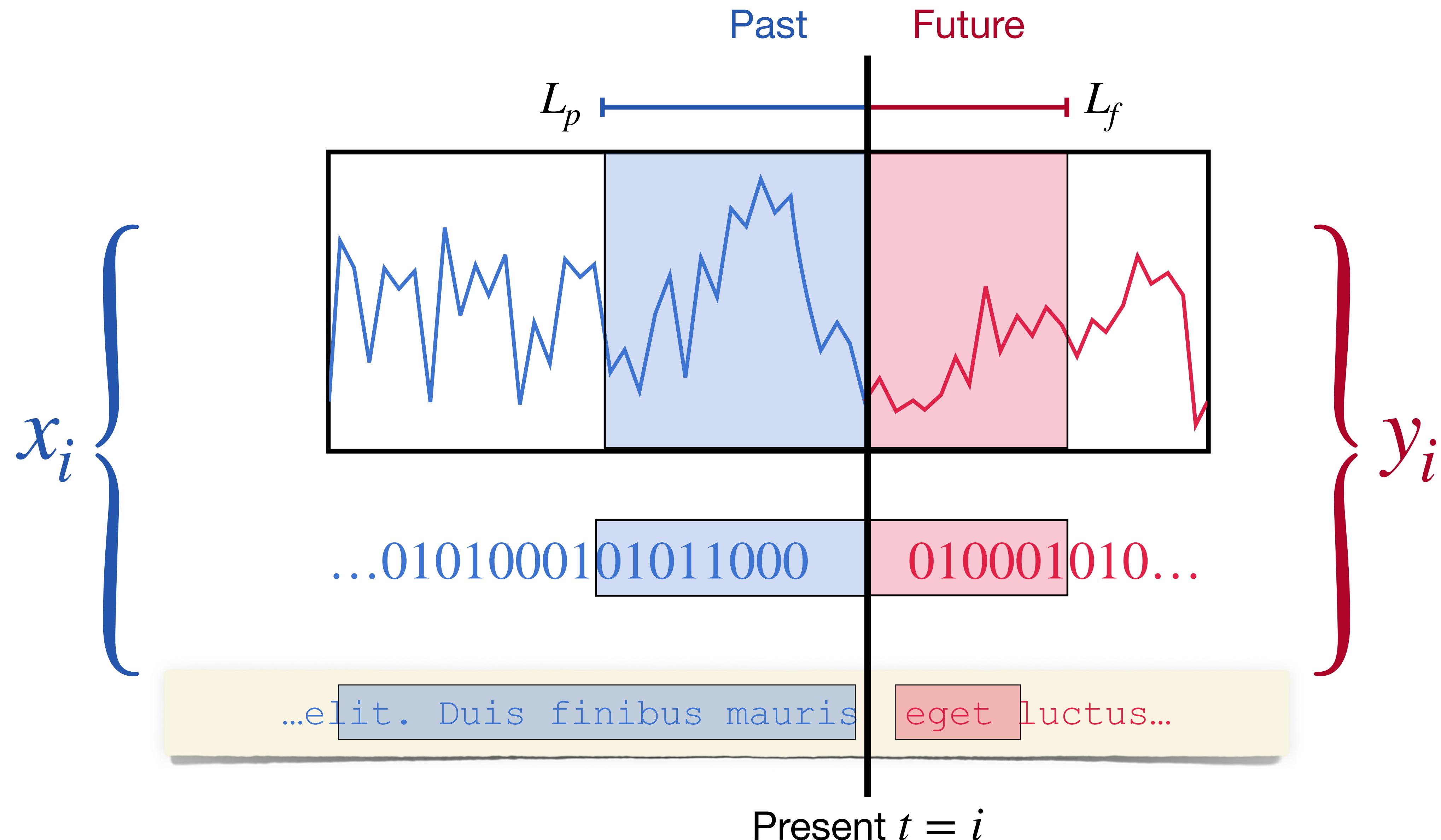
# Thank you and Questions

---

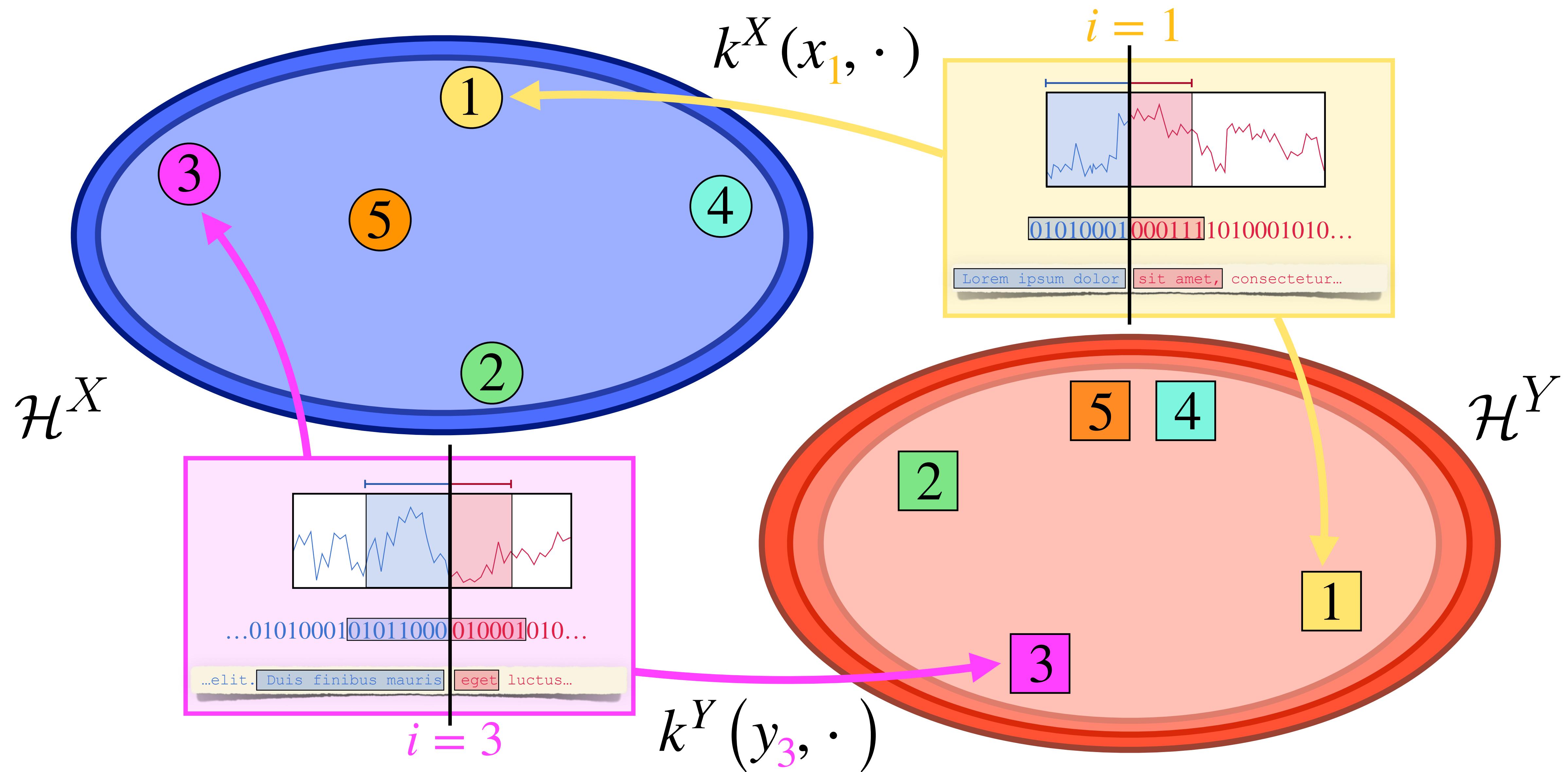




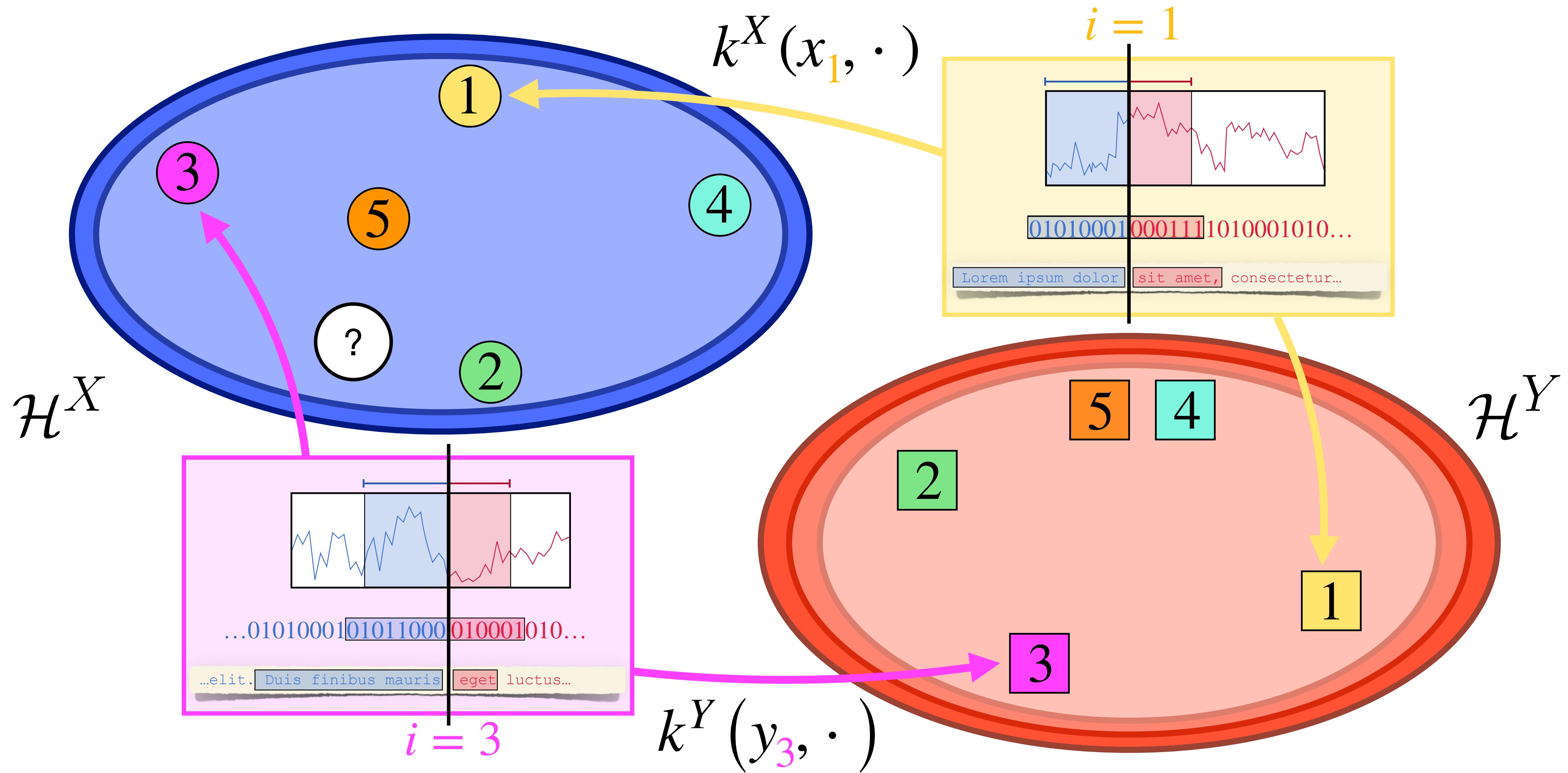
# Kernel Causal States Algorithm



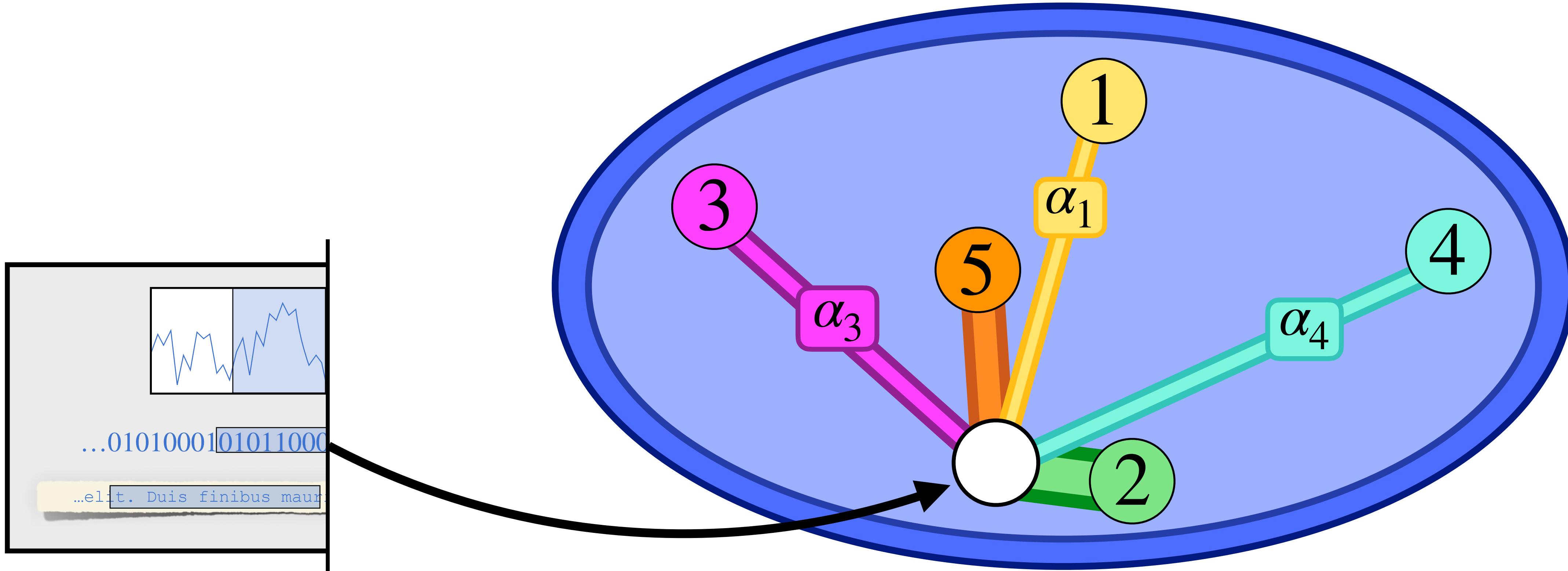
# Embed Pasts and Futures



# Embed Pasts and Futures

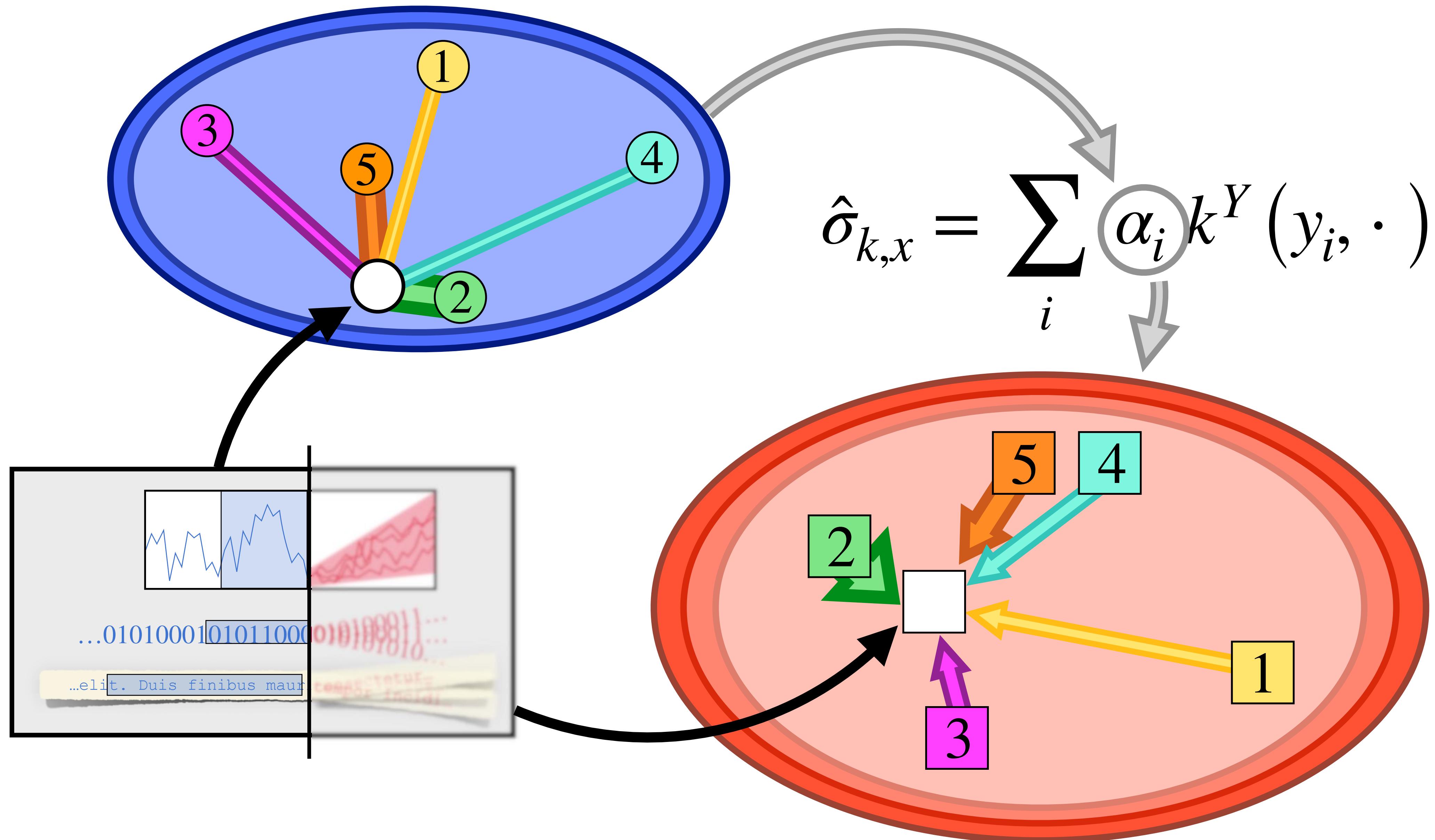


# Calculate Proximity-Based Weights



$$k^X(x, x_i) \uparrow \implies \alpha_i \uparrow$$

# Construct Future Distribution



# Embed Causal State



$$\hat{\sigma}_{k,x} = \sum_i \alpha_i k^Y(y_i, \cdot)$$

