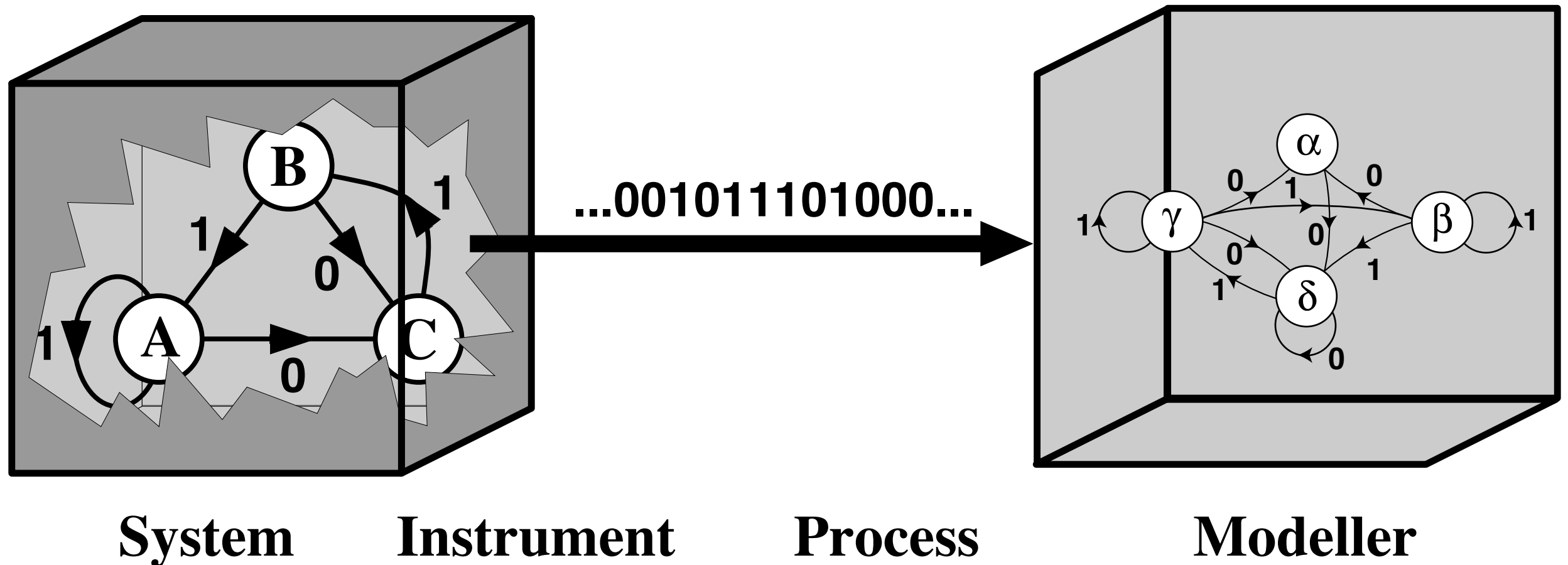


The Learning Channel

The Learning Channel:



Central questions:
What are the states?
What is the dynamic?

The Learning Channel ...

The Prediction Game

Rules:

1. I give you a data stream (an observed past sequence).
2. You predict its future.
3. You give a model (states & transitions) describing the process.

The Learning Channel ...

The Prediction Game ...

Process I:

The Learning Channel ...

The Prediction Game ...

Process I:

Past: ... 111111111111

The Learning Channel ...

The Prediction Game ...

Process I:

Past: ... 111111111111

Your prediction is?

The Learning Channel ...

The Prediction Game ...

Process I:

Past: ... 111111111111

Your prediction is?

Future: 111111111111 ...

The Learning Channel ...

The Prediction Game ...

Process I:

Past: ... 111111111111

Your prediction is?

Future: 111111111111 ...

Your model (states & dynamic) is?

The Learning Channel ...

The Prediction Game ...

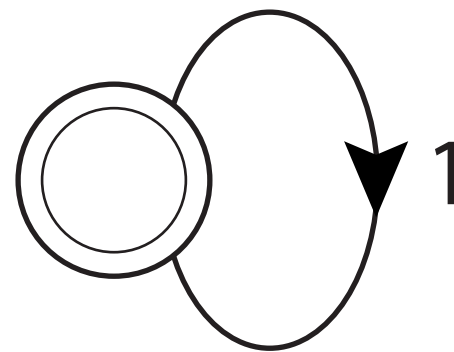
Process I:

Past: ... 111111111111

Your prediction is?

Future: 111111111111...

Your model (states & dynamic) is?



The Learning Channel ...

The Prediction Game ...

Process II:

The Learning Channel ...

The Prediction Game ...

Process II:

Past: ... 10110010001101110

The Learning Channel ...

The Prediction Game ...

Process II:

Past: ... 10110010001101110

Your prediction is?

The Learning Channel ...

The Prediction Game ...

Process II:

Past: ... 10110010001101110

Your prediction is?

Analysis: All words of length L occur & equally often

The Learning Channel ...

The Prediction Game ...

Process II:

Past: ... 10110010001101110

Your prediction is?

Analysis: All words of length L occur & equally often

Future: Well, anything can happen, how about?

The Learning Channel ...

The Prediction Game ...

Process II:

Past: ... 10110010001101110

Your prediction is?

Analysis: All words of length L occur & equally often

Future: Well, anything can happen, how about?

01010111010001101 ...

The Learning Channel ...

The Prediction Game ...

Process II:

Past: ... 10110010001101110

Your prediction is?

Analysis: All words of length L occur & equally often

Future: Well, anything can happen, how about?

01010111010001101 ...

Your model is?

The Learning Channel ...

The Prediction Game ...

Process II:

Past: ... 10110010001101110

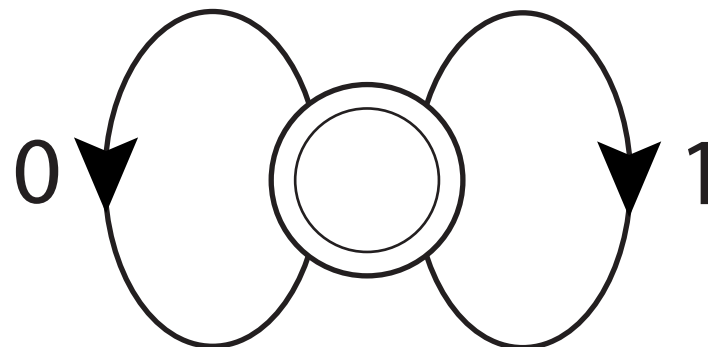
Your prediction is?

Analysis: All words of length L occur & equally often

Future: Well, anything can happen, how about?

01010111010001101 ...

Your model is?



The Learning Channel ...

The Prediction Game ...

Process III:

The Learning Channel ...

The Prediction Game ...

Process III:

Past: ... 10101010101010

The Learning Channel ...

The Prediction Game ...

Process III:

Past: ... 1010101010101010

Your prediction is?

The Learning Channel ...

The Prediction Game ...

Process III:

Past: ... 1010101010101010

Your prediction is?

Future: 101010101010101...

The Learning Channel ...

The Prediction Game ...

Process III:

Past: ... 1010101010101010

Your prediction is?

Future: 101010101010101 ...

Your model is?

The Learning Channel ...

The Prediction Game ...

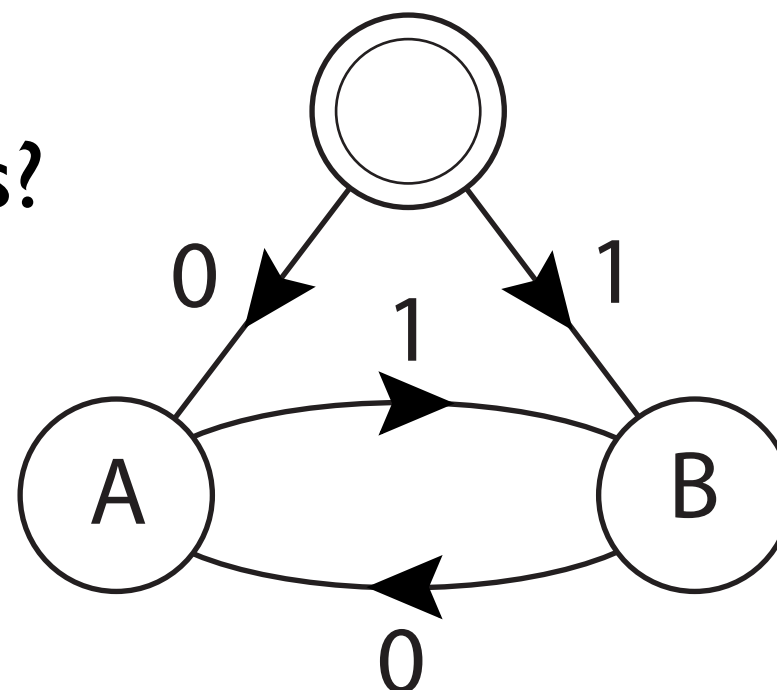
Process III:

Past: ... 10101010101010

Your prediction is?

Future: 1010101010101...

Your model is?



The Learning Channel ...

Goal:

Predict the future \vec{S}
using information from the past \overleftarrow{S}

But what “information” to use?

We want to find the effective “states”
and the dynamic (state-to-state mapping)

How to define “states”, if they are hidden?

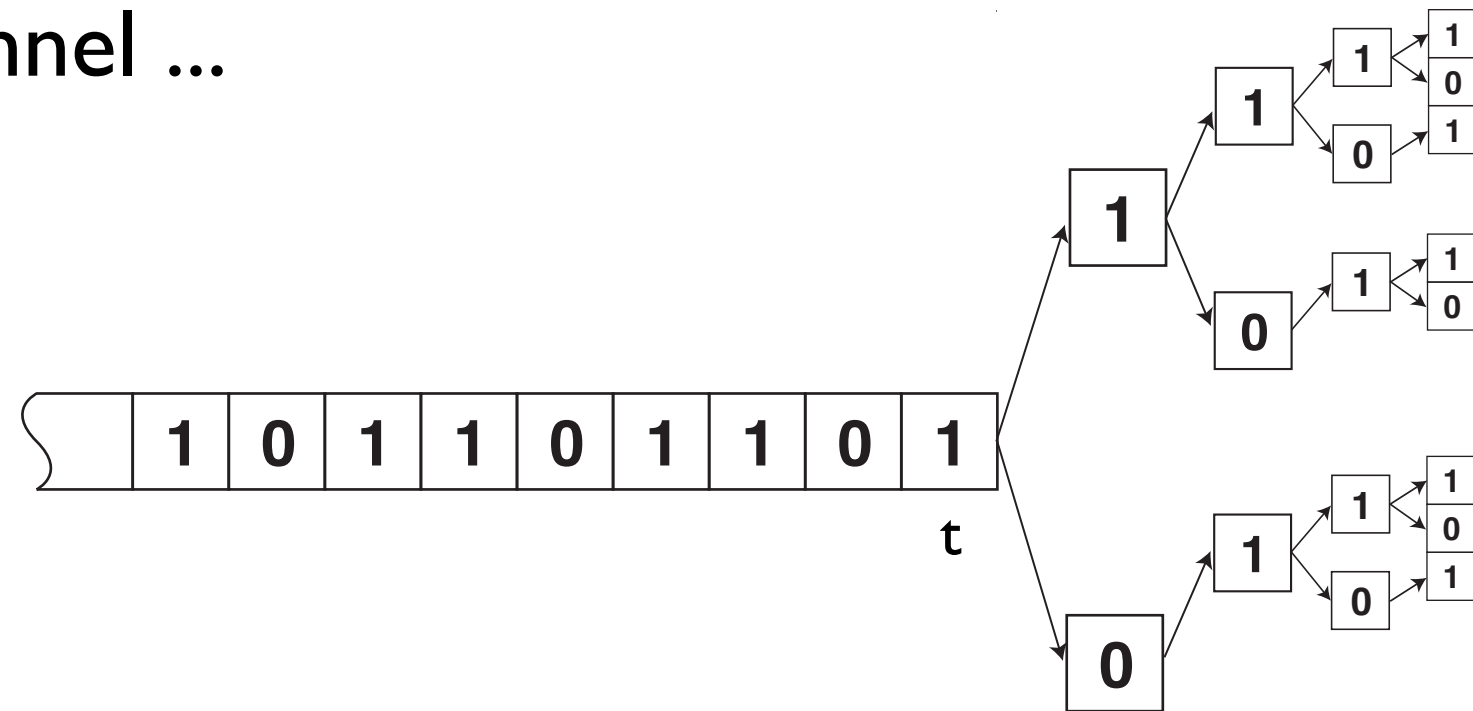
All we have are sequences of observations
Over some measurement alphabet \mathcal{A}
These symbols only indirectly reflect the hidden states

The Learning Channel ...

Effective States:

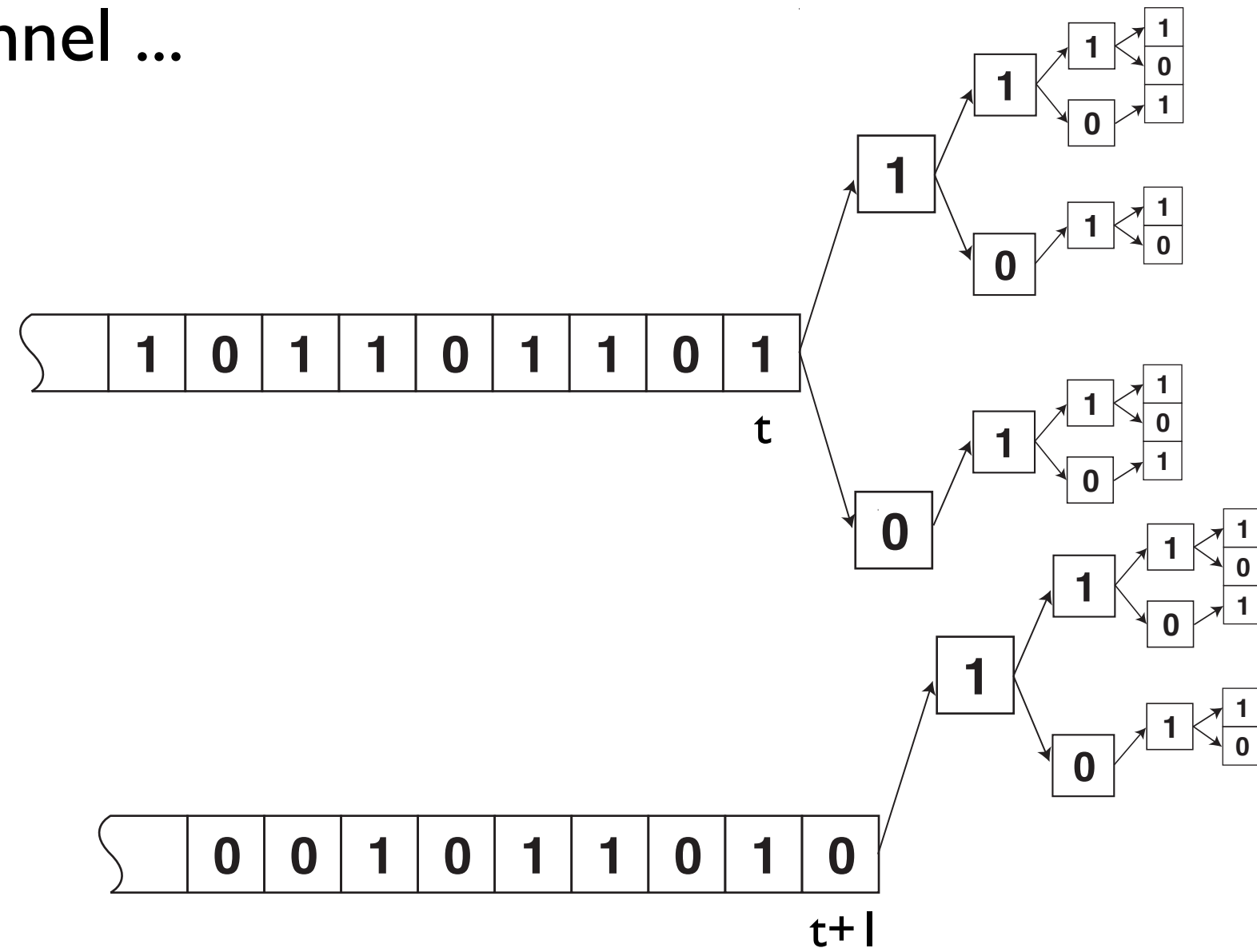
The Learning Channel ...

Effective States:



The Learning Channel ...

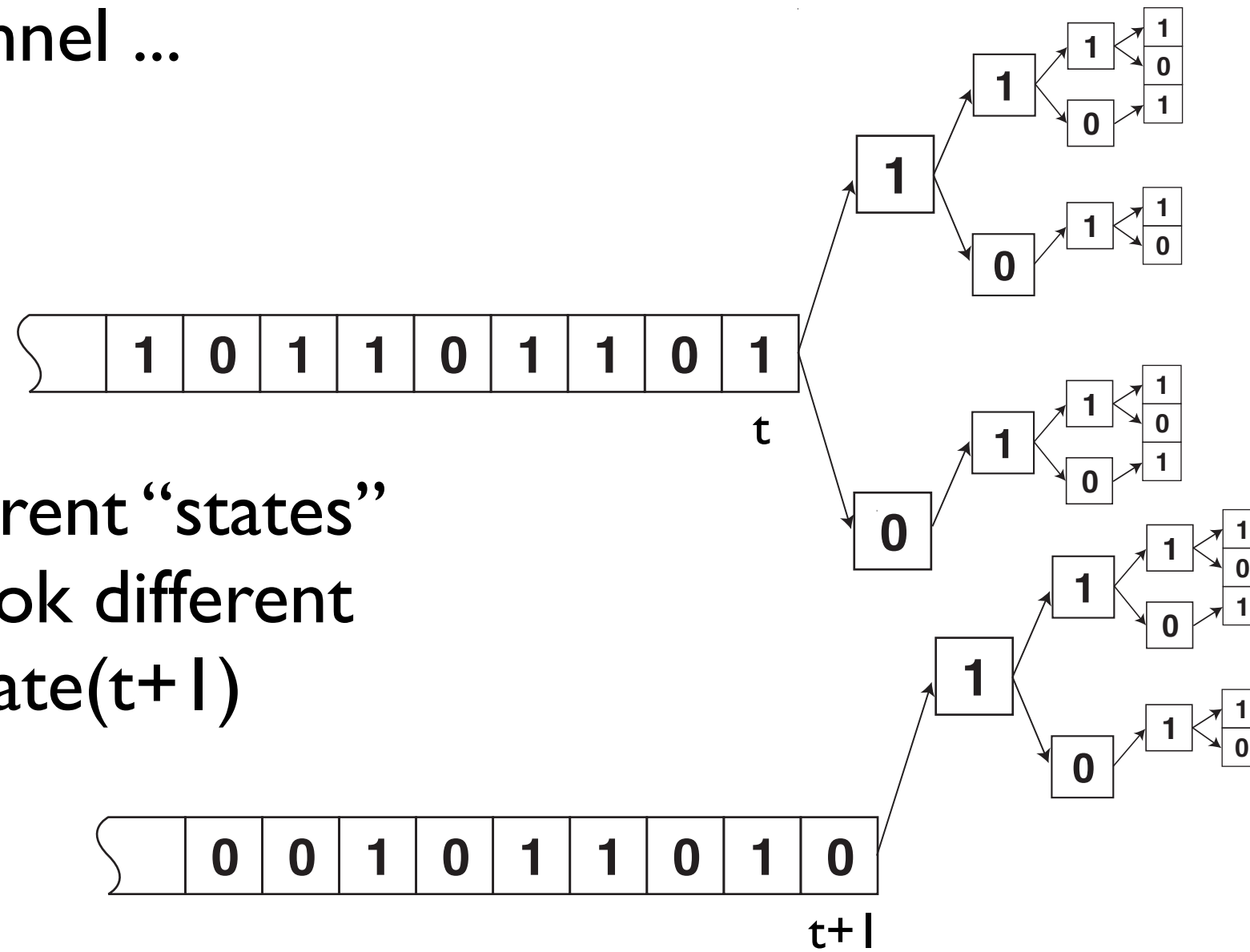
Effective States:



The Learning Channel ...

Effective States:

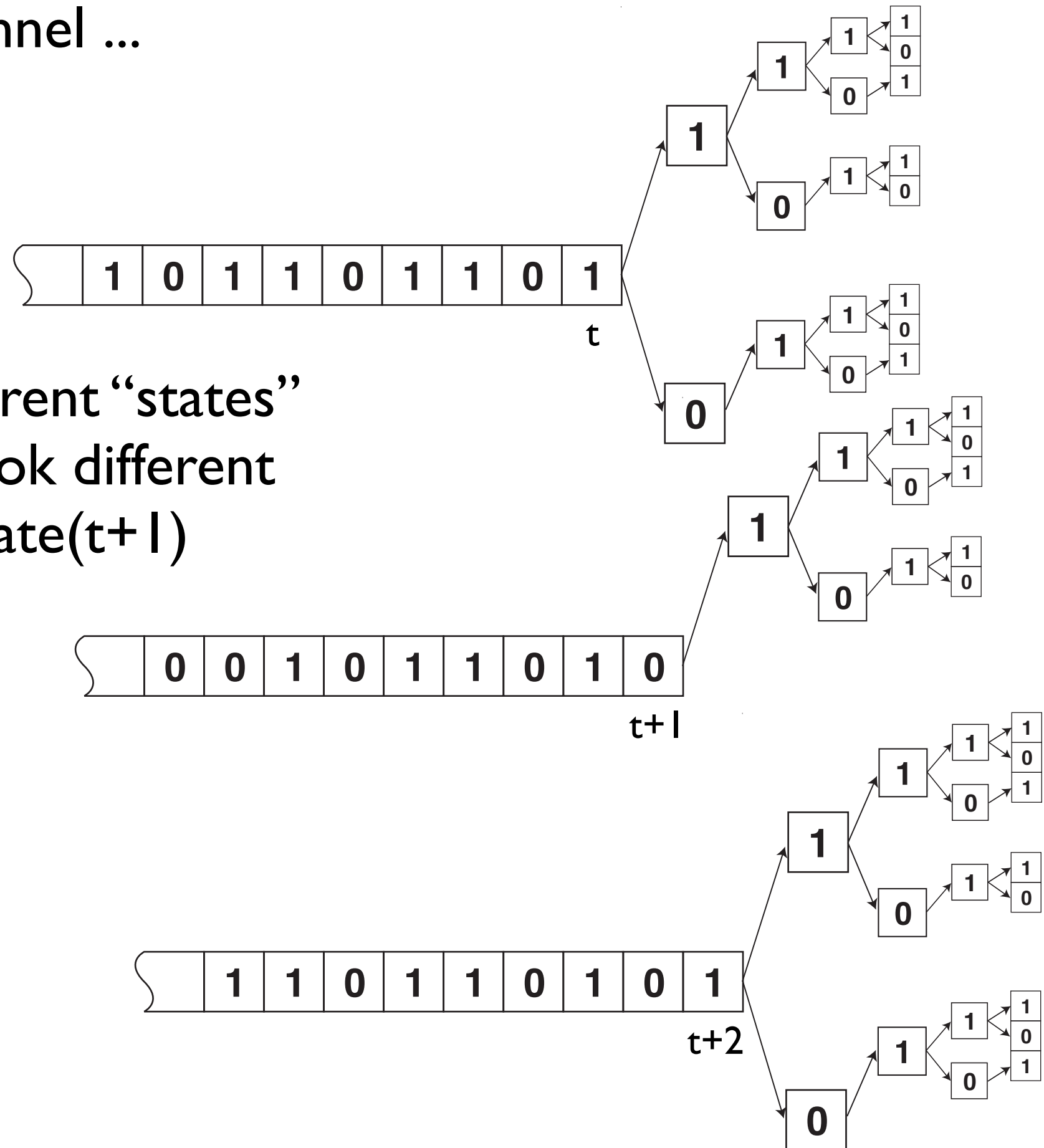
Process is in different “states”
when futures look different
 $\text{State}(t) \approx \text{State}(t+1)$



The Learning Channel ...

Effective States:

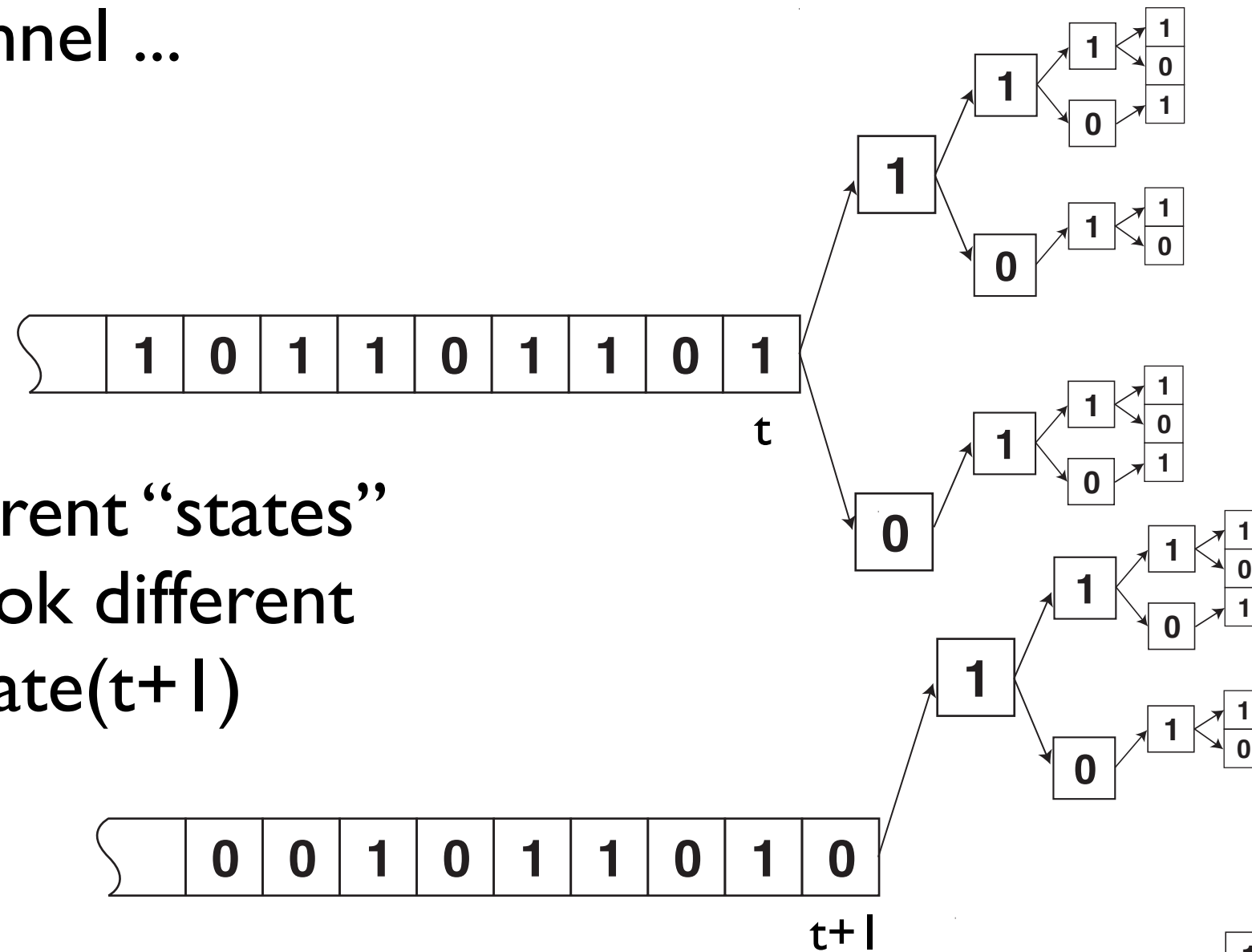
Process is in different “states”
when futures look different
 $\text{State}(t) \approx \text{State}(t+1)$



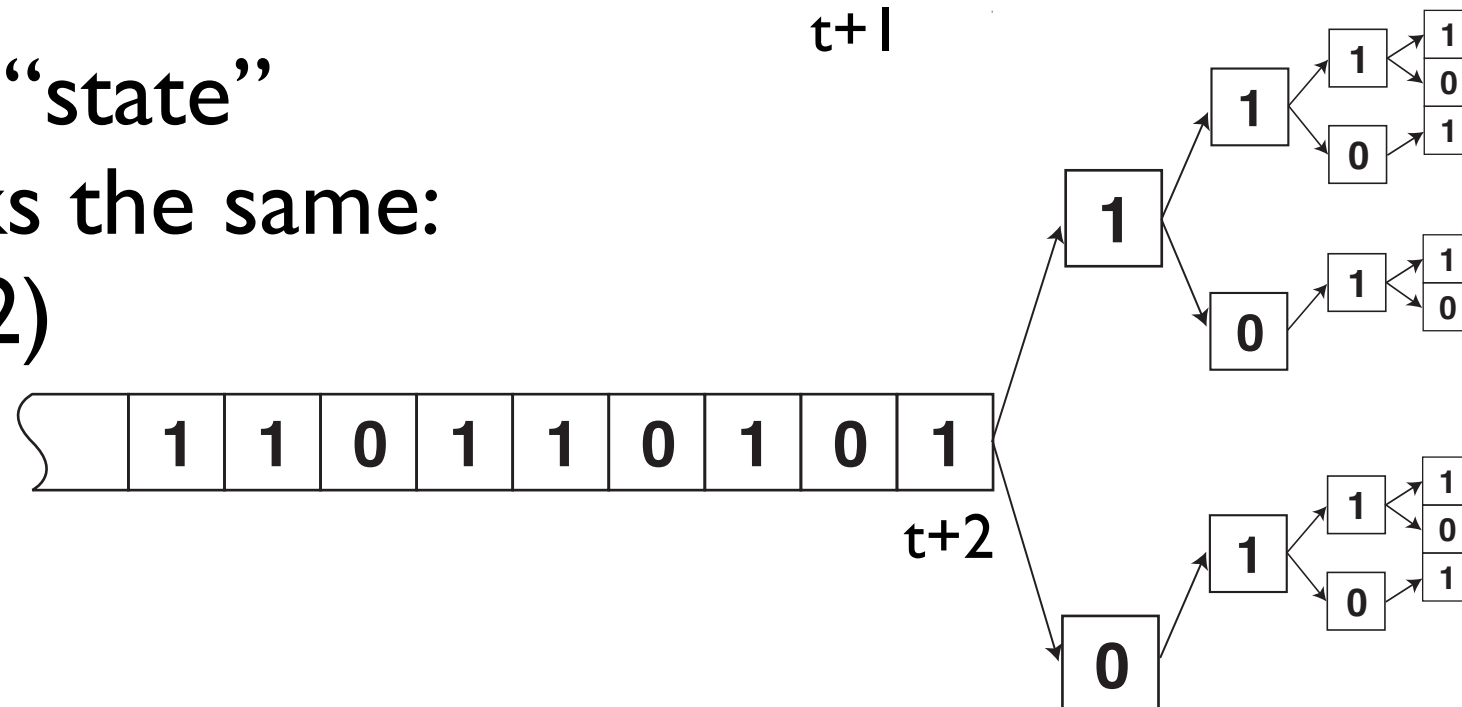
The Learning Channel ...

Effective States:

Process is in different “states”
when futures look different
 $\text{State}(t) \not\sim \text{State}(t+1)$



Process is in the same “state”
when the future looks the same:
 $\text{State}(t) \sim \text{State}(t+2)$



The Learning Channel ...

Effective for what?

What's a prediction?

A mapping from the past to the future.

Process $\Pr(\overleftrightarrow{S}) : \overleftrightarrow{S} = \overleftarrow{S} \overrightarrow{S}$

Future: \overrightarrow{S}^L

Particular past: \overleftarrow{s}

Future Morph: $\Pr(\overrightarrow{S}^L | \overleftarrow{s})$ (the most general mapping)

Refined goal:

Predict as much about the future \overrightarrow{S} ,
using as little of the past \overleftarrow{S} as possible.

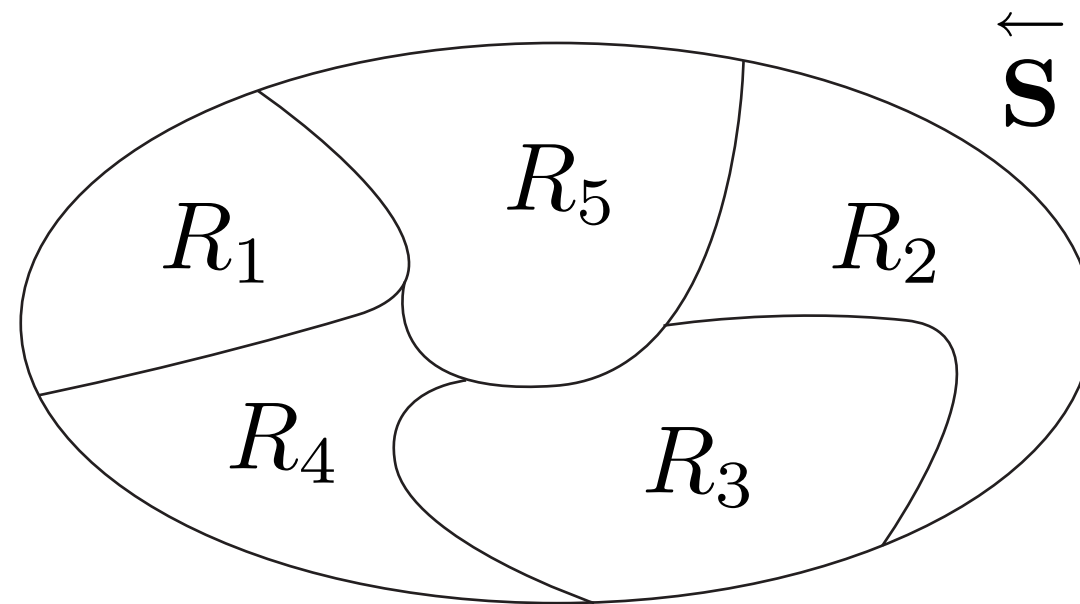
The Learning Channel ...

Space of Histories ...

Histories leading to the same predictions are equivalent.

Effective States = **Partitions of History**:

$$R = \{R_i : R_i \cap R_j = \emptyset, \overleftarrow{\mathbf{S}} = \bigcup_i R_i\}$$



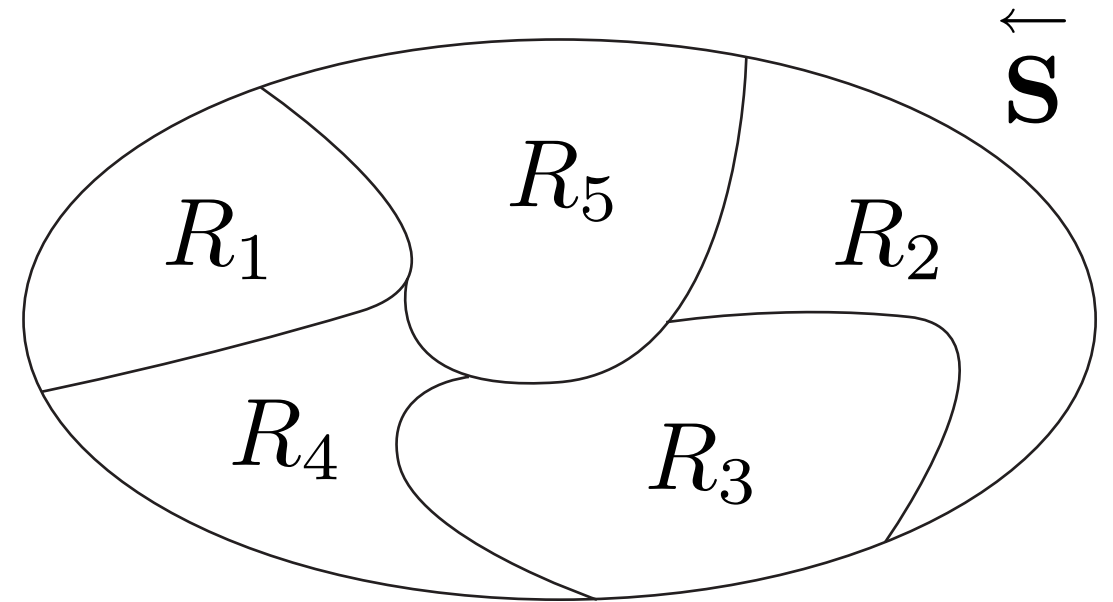
The Learning Channel ...

Space of Histories ...

Map from histories to partition elements:

$$\eta : \overleftarrow{\mathbf{S}} \rightarrow R$$

$$\eta(\overleftarrow{s}) = R_i$$



Random variable:

$$R = \eta(\overleftarrow{S})$$

Distribution over Effective States:

$$\Pr(R = R_i) = \sum_{\overleftarrow{s} : \eta(\overleftarrow{s}) = R_i} \Pr(\overleftarrow{s})$$

The Learning Channel ...

How Effective are the Effective States?

Effective Prediction Error: Given a candidate partition R

$$H[\vec{S}^L | R]$$

Uncertainty about future given effective states

Effective Prediction Error Rate:

$$h_\mu(R) = \lim_{L \rightarrow \infty} \frac{H[\vec{S}^L | R]}{L}$$

Entropy rate given effective states

The Learning Channel ...

How Effective are the Effective States?

Statistical Complexity of the Effective States:

$$C_\mu(R) = H[R] = H(\text{Pr}(R))$$

Interpretations:

Uncertainty in state.

Shannon information one gains when told effective state.

Model “size” $\propto \log_2(\text{number of states})$

Historical memory used by R .

The Learning Channel ...

Goals Restated:

Question 1:

Can we find effective states that give good predictions?

$$H[\vec{S}^{\rightarrow L} | R] = H[\vec{S}^{\rightarrow L} | \vec{S}^{\leftarrow}]$$

or

$$h_{\mu}(R) = h_{\mu}$$

Question 2:

Can we find the smallest such set?

$$\min C_{\mu}(R)$$

The Learning Channel ...

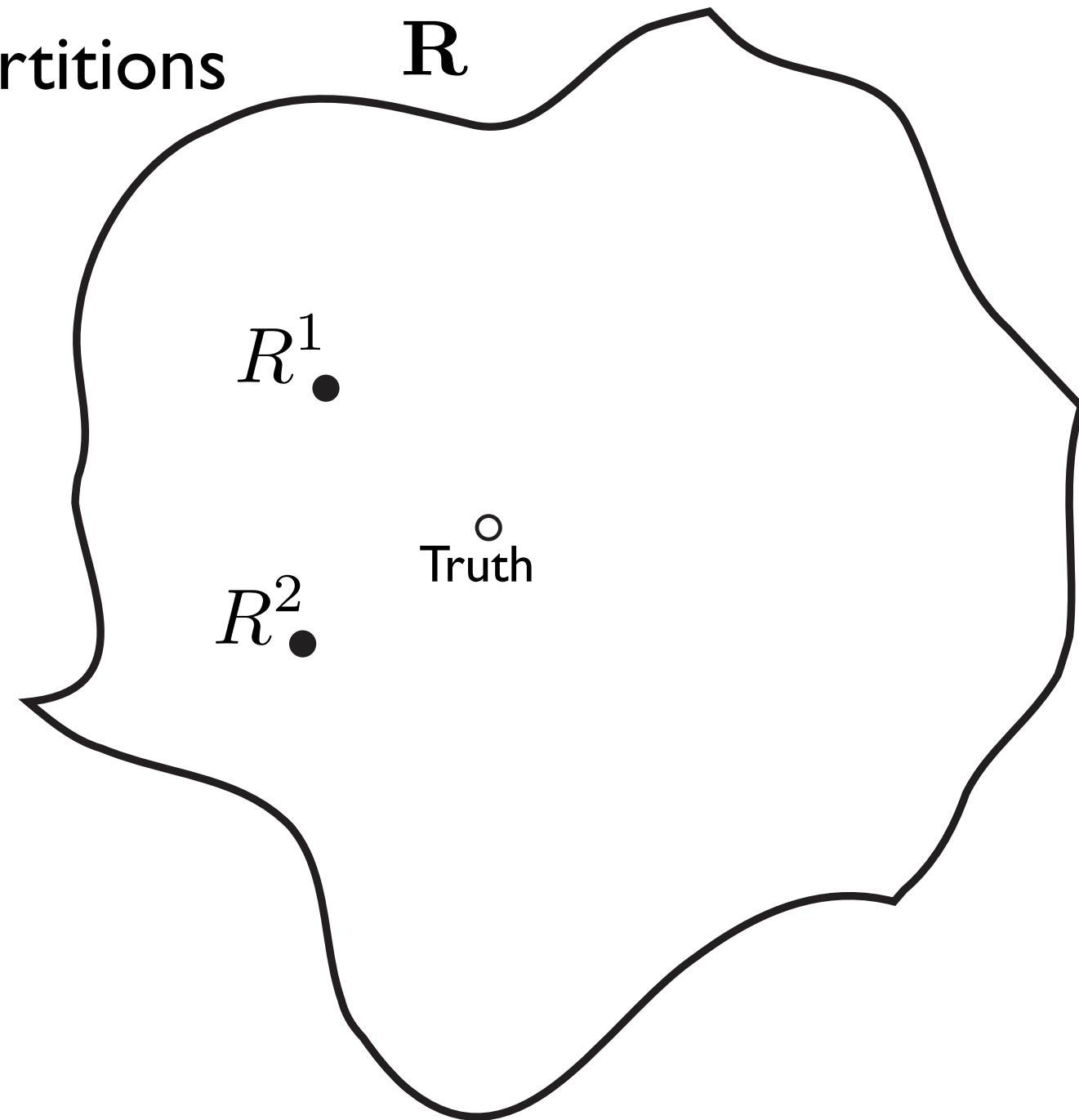
Occam's Pool: The Space of Models

Model = Partition of History Space

Model Space \mathbf{R} = Space of all partitions

Rival Models:

$$R_1, R_2 \in \mathbf{R}$$



The Learning Channel ...

Causal States:

Causal State:

Set of pasts with same morph $\Pr(\vec{S} \mid \overleftarrow{s})$.

Set of histories that lead to same predictions.

Predictive equivalence relation:

$$\overleftarrow{s}' \sim \overleftarrow{s}'' \iff \Pr(\vec{S} \mid \overleftarrow{S} = \overleftarrow{s}') = \Pr(\vec{S} \mid \overleftarrow{S} = \overleftarrow{s}'')$$

$$\overleftarrow{s}', \overleftarrow{s}'' \in \overleftarrow{\mathbf{S}}$$

The Learning Channel ...

Causal State Components

Causal State = Pasts with same morph: $\Pr(\vec{S} \mid \overleftarrow{s})$

$$\mathcal{S} = \{ \overleftarrow{s}' : \overleftarrow{s}' \sim \overleftarrow{s} \}$$

Set of causal states:

$$\mathcal{S} = \overleftarrow{\mathbf{S}} / \sim = \{ \mathcal{S}_0, \mathcal{S}_1, \mathcal{S}_2, \dots \}$$

Causal state map:

$$\epsilon : \overleftarrow{\mathbf{S}} \rightarrow \mathcal{S}$$

$$\epsilon(\overleftarrow{s}) = \{ \overleftarrow{s}' : \overleftarrow{s}' \sim \overleftarrow{s} \}$$

The Learning Channel ...

Causal States ...

We've answered the first part of the modeling goal:

We have the effective states!

Now,

What is the dynamic?

The Learning Channel ...

Causal State Dynamic ...

Causal-state Filtering:

$$\begin{aligned}\overleftrightarrow{s} &= \dots s_{-3} \quad s_{-2} \quad s_{-1} \quad s_0 \quad s_1 \quad s_2 \quad s_3 \quad \dots \\ \epsilon(\overleftrightarrow{s}) &= \dots \epsilon(\overleftarrow{s}_{-3}) \epsilon(\overleftarrow{s}_{-2}) \epsilon(\overleftarrow{s}_{-1}) \epsilon(\overleftarrow{s}_0) \epsilon(\overleftarrow{s}_1) \epsilon(\overleftarrow{s}_2) \epsilon(\overleftarrow{s}_3) \dots \\ \overleftrightarrow{\mathcal{S}} &= \dots \mathcal{S}_{t=-3} \mathcal{S}_{t=-2} \mathcal{S}_{t=-1} \mathcal{S}_{t=0} \mathcal{S}_{t=1} \mathcal{S}_{t=2} \mathcal{S}_{t=3} \dots\end{aligned}$$

Causal-state process:

$$\Pr(\overleftrightarrow{\mathcal{S}})$$

The Learning Channel ...

Causal State Dynamic ...

Conditional transition probability:

$$\begin{aligned} T_{ij}^{(s)} &= \Pr(\mathcal{S}_j, s | \mathcal{S}_i) \\ &= \Pr\left(\mathcal{S} = \epsilon(\overleftarrow{s} s) | \mathcal{S} = \epsilon(\overleftarrow{s})\right) \end{aligned}$$

State-to-State Transitions:

$$\{T_{ij}^{(s)} : s \in \mathcal{A}, i, j = 0, 1, \dots, |\mathcal{S}|\}$$

The ϵ -Machine ...

Process \Rightarrow Predictive equivalence $\Rightarrow \epsilon$ - Machine

$$\text{Pr}(\vec{S}) \Rightarrow \vec{S} / \sim \Rightarrow \epsilon - \text{Machine}$$

$$\mathcal{M} = \left\{ \mathcal{S}, \{T^{(s)}, s \in \mathcal{A}\} \right\}$$

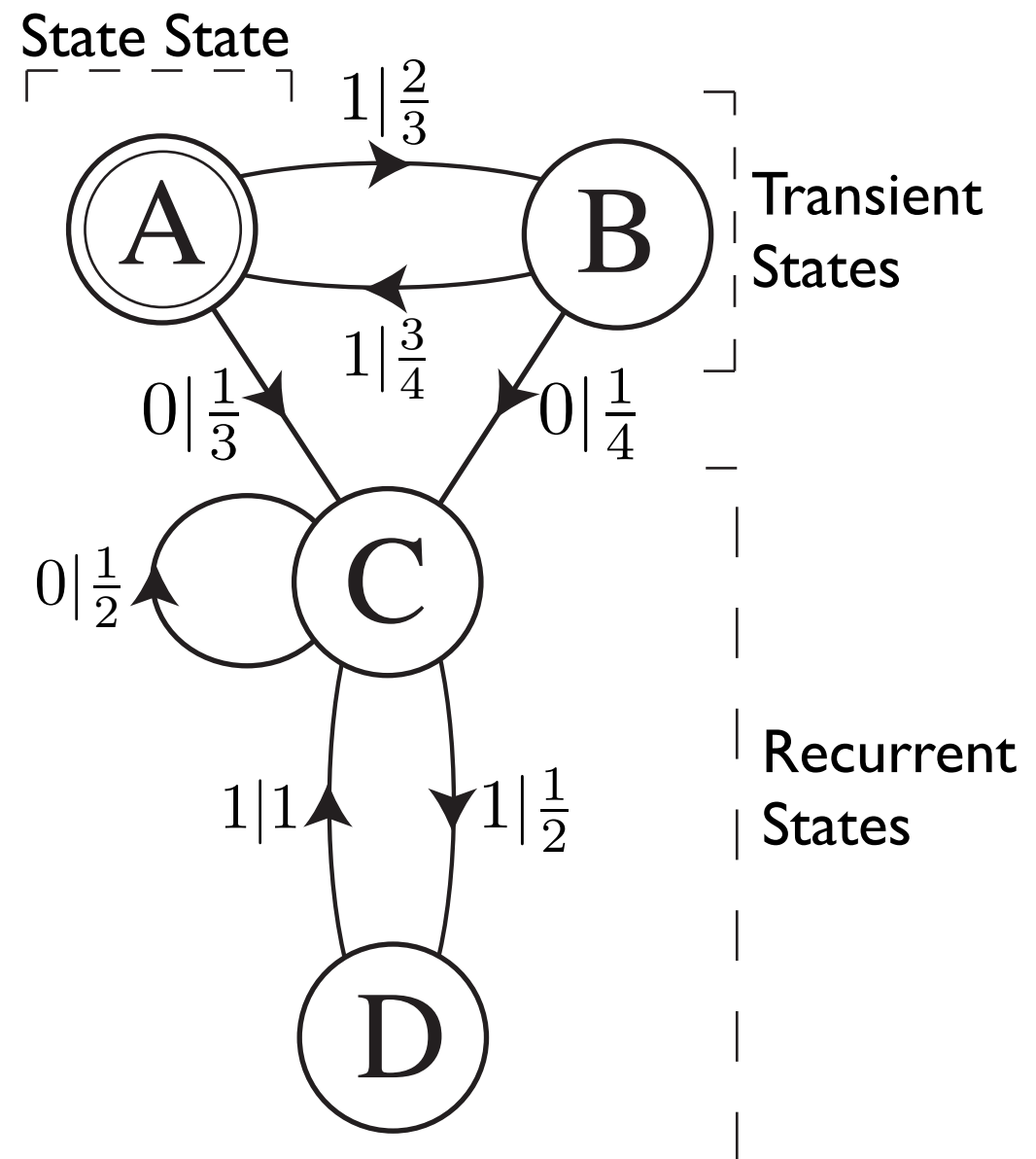
Unique Start State:

$$\mathcal{S}_0 = [\lambda]$$

$$\text{Pr}(\mathcal{S}_0, \mathcal{S}_1, \mathcal{S}_2, \dots) = (1, 0, 0, \dots)$$

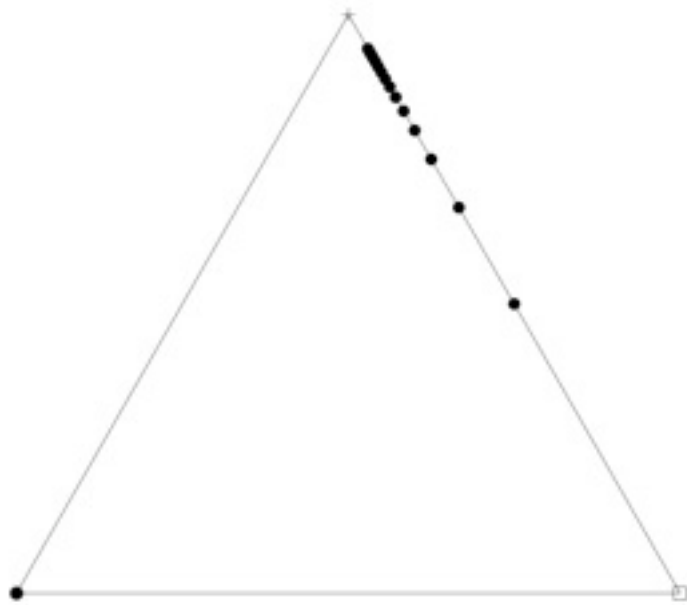
Transient States

Recurrent States

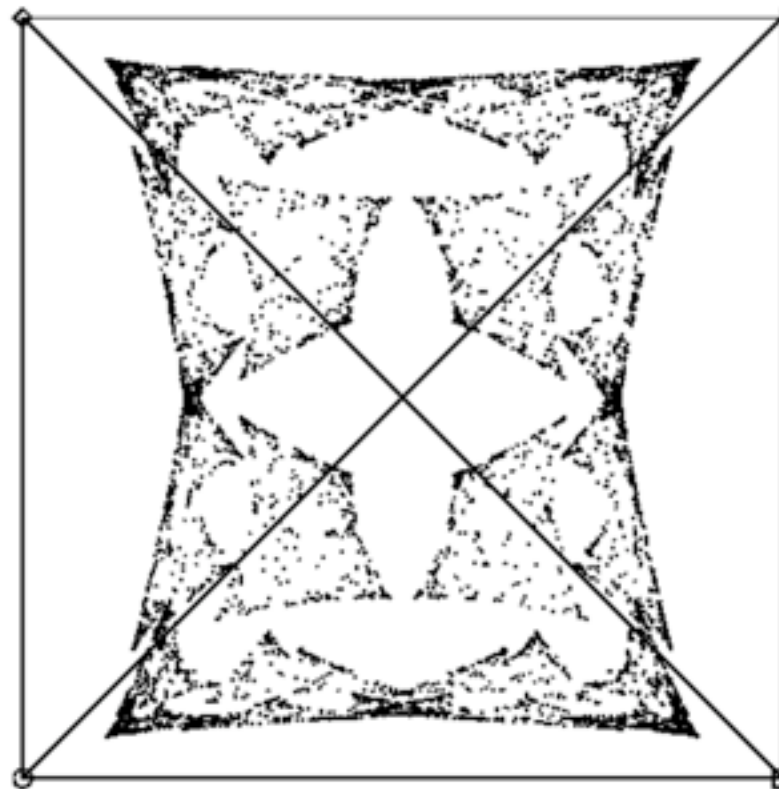


The Learning Channel ...

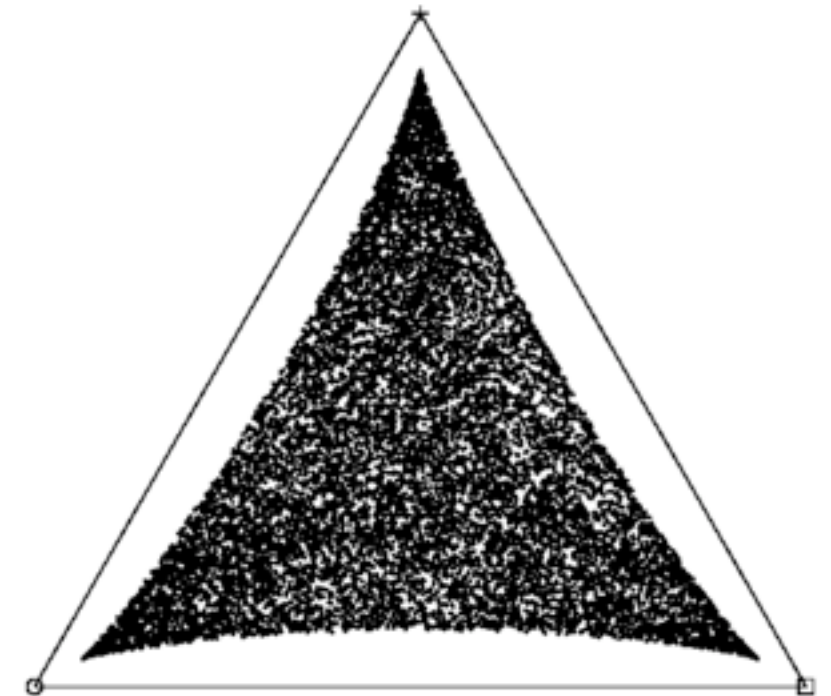
The ϵ -Machine of a Process ...



**Denumerable
Causal States**

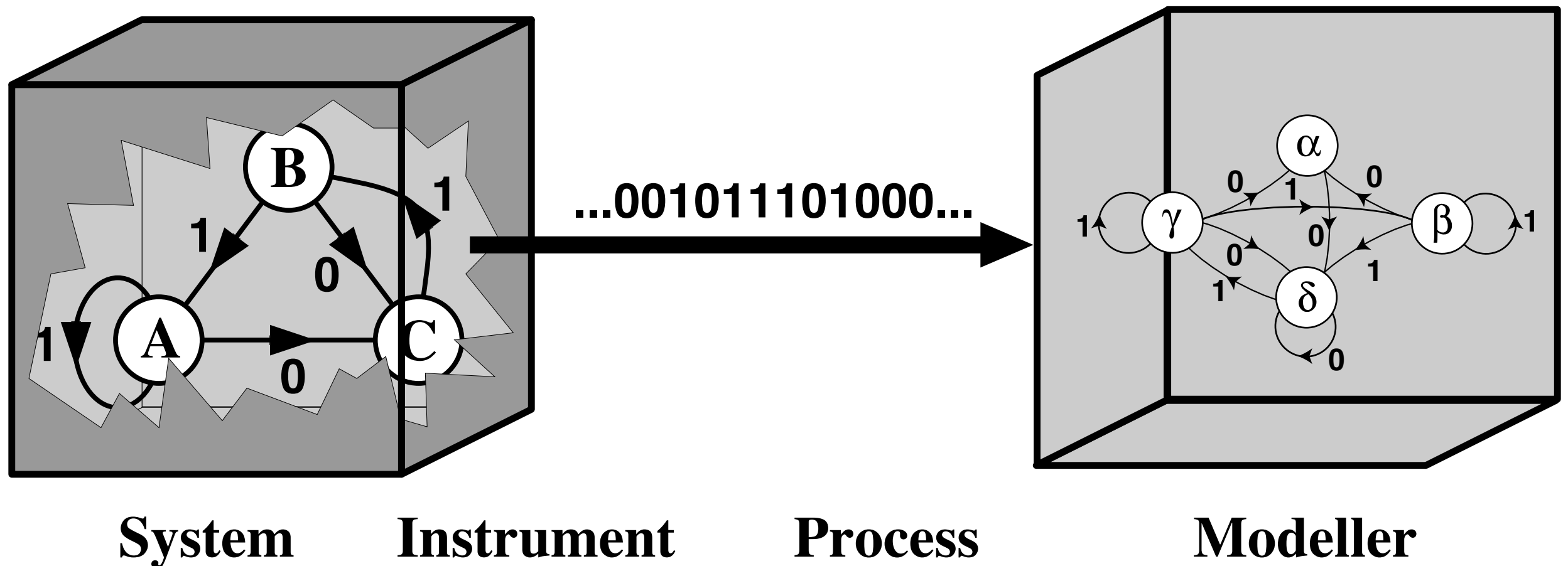


Fractal



Continuous

The Learning Channel:



Central questions:

What are the states? Causal States

What is the dynamic? The ϵ -Machine

The ϵ -Machine ...

A **Model** of a Process $\Pr(\vec{S})$:

ϵ -Machine reproduces the process's word distribution:

$$\Pr(s^1), \Pr(s^2), \Pr(s^3), \dots$$

$$s^L = s_1 s_2 \dots s_L \quad \mathcal{S}(t=0) = \mathcal{S}_0$$

$$\begin{aligned} \Pr(s^L) = & \Pr(\mathcal{S}_0) \Pr(\mathcal{S}_0 \rightarrow_{s=s_1} \mathcal{S}(1)) \Pr(\mathcal{S}(1) \rightarrow_{s=s_2} \mathcal{S}(2)) \\ & \dots \Pr(\mathcal{S}(L-1) \rightarrow_{s=s_L} \mathcal{S}(L)) \end{aligned}$$

Initially, $\Pr(\mathcal{S}_0) = 1$.

$$\Pr(s^L) = \prod_{l=1}^L T_{i=\epsilon(s^{l-1}), j=\epsilon(s^l)}^{(s^l)}$$

The ϵ -Machine ...

Causal shielding:

Past and future are independent given causal state

$$\text{Process: } \Pr(\vec{S}) = \Pr(\overleftarrow{S} \overrightarrow{S})$$

$$\Pr(\overleftarrow{S} \overrightarrow{S} | \mathcal{S}) = \Pr(\overleftarrow{S} | \mathcal{S}) \Pr(\overrightarrow{S} | \mathcal{S})$$

Causal states shield past & future from each other.

Similar to states of a Markov chain, but for hidden processes.

The ϵ -Machine ...

ϵ Ms are **Unifilar**: $(\mathcal{S}_t, s) \rightarrow$ unique \mathcal{S}_{t+1}

(in automata theory, “deterministic”)

That is:

(1) $\mathcal{S}_i \in \mathcal{S}$, $s \in \mathcal{A}$, at most one $\mathcal{S}_j \in \mathcal{S}$:

$$\overleftarrow{s} \in \mathcal{S}_i \Rightarrow \overleftarrow{s} s \in \mathcal{S}_j$$

(2) If there is a next causal state j :

$$\mathcal{S}_{k \neq j} \in \mathcal{S} \Rightarrow T_{ik}^{(s)} = 0$$

(3) If there is not:

$$T_{ij}^{(s)} = 0$$

The ϵ -Machine ...

Unifilarity ...

Consequence:

Unifilarity: 1-1 map between state-sequences & symbol-sequences.

Entropy rate expression requires this 1-1 mapping.

Can use ϵM to calculate entropy rate h_μ .
(Any unifilar presentation will do.)

The ϵ -Machine ...

ϵ Ms are **Optimal Predictors**:

Compared to any rival effective states R :

$$H \left[\overrightarrow{S}^L \mid R \right] \geq H \left[\overrightarrow{S}^L \mid \mathcal{S} \right]$$

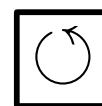
Proof sketch: $H \left[\overrightarrow{S}^L \mid \mathcal{S} \right] = H \left[\overrightarrow{S}^L \mid \overleftarrow{s} \in \mathcal{S} \right]$ (Causal equiv. rel'n)

$$= H \left[\overrightarrow{S}^L \mid \overleftarrow{s} \right]$$

$$\leq H \left[\overrightarrow{S}^L \mid R \right]$$

$$R = \eta(\overleftarrow{s})$$

(Data processing inequality)



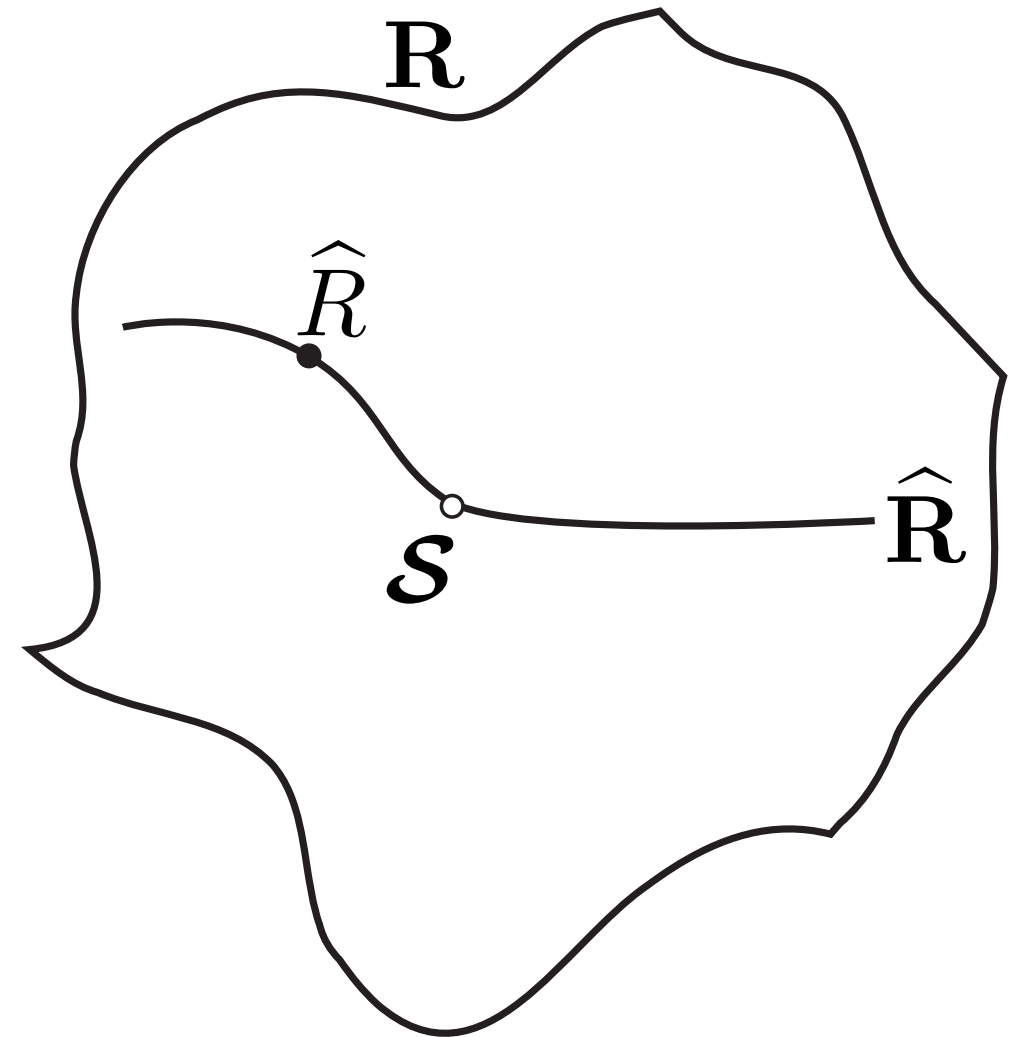
The ϵ -Machine ...

Prescient Rivals $\hat{\mathbf{R}}$:

Alternative models that are optimal predictors

$$\hat{R} \in \hat{\mathbf{R}}$$

$$H[\vec{S}^L | \hat{R}] = H[\vec{S}^L | \mathcal{S}]$$



(Prescient rivals are sufficient statistics for process's future.)

The ϵ -Machine ...

Minimal Statistical Complexity:

For all prescient rivals, ϵM is the smallest:

$$C_\mu(\hat{R}) \geq C_\mu(\mathcal{S})$$

Proof sketch:

(1) Prescient rivals are refinements, so

$$\exists g : \mathcal{S} = g(\hat{R})$$

(2) But

$$H[f(X)] \leq H[X] \Rightarrow H[\mathcal{S}] = H[g(\hat{R})] \leq H[\hat{R}]$$

(3) So $C_\mu \leq H[\hat{R}]$



The ϵ -Machine ...

Minimal Statistical Complexity ...

Consequence:

- (1) C_μ measures historical information process stores.
- (2) This would not be true, if not minimal representation.

Remarks:

- (1) Causal states contain every difference (in past) that makes a difference (to future) (Bateson “information”)
- (2) Causal states are sufficient statistics for the future.

The ϵ -Machine ...

Summary:

ϵM :

- (1) Optimal predictor: Lower prediction error than any rival.
- (2) Minimal size: Smallest of the prescient rivals.
- (3) Unique: Smallest, optimal, unifilar predictor is equivalent.
- (4) Model of the process: Reproduces all of process's statistics.
- (5) Causal shielding: Renders process's future independent of past.

The ϵ -Machine ...

Dynamical system's **intrinsic computation**:

- (1) How much of past does process store?
- (2) In what architecture is that information stored?
- (3) How is stored information used to produce future behavior?

The ϵ -Machine ...

Dynamical system's **intrinsic computation**:

(1) How much of past does process store?

$$C_\mu$$

(2) In what architecture is that information stored?

(3) How is stored information used to produce future behavior?

The ϵ -Machine ...

Dynamical system's **intrinsic computation**:

(1) How much of past does process store?

$$C_\mu$$

(2) In what architecture is that information stored?

$$\left\{ \mathcal{S}, \{T^{(s)}, s \in \mathcal{A}\} \right\}$$

(3) How is stored information used to produce future behavior?

The ϵ -Machine ...

Dynamical system's **intrinsic computation**:

(1) How much of past does process store?

$$C_\mu$$

(2) In what architecture is that information stored?

$$\left\{ \mathcal{S}, \{T^{(s)}, s \in \mathcal{A}\} \right\}$$

(3) How is stored information used to produce future behavior?

$$h_\mu$$