

# The Learning Channel

Reading for this lecture: *CMR* articles

*BOAC*: “Between Order and Chaos”,  
Nature Physics **8** (January 2012) 17-24.

*CMPPSS*: (Sections I and II only).

*CAO*: “*Chance and Order*”, Stanislaw Lem,  
New Yorker **59** (1984) 88-98.

*ROIC*: “*Revealing Order in the Chaos*”, Mark Buchanan,  
New Scientist, 26 February 2005;  
available at [csc.ucdavis.edu/~chaos/news/](http://csc.ucdavis.edu/~chaos/news/).

# PHY 256B Topics

Effective States & Dynamic

$\varepsilon$ -Machines

Measures of Complexity

Irreversibility

Crypticity

Meaning & Measurement Semantics

Information Measures Redux

Directional Computational Mechanics

Intrinsic Semantics of Information

# PHY 256B Topics ...

Topical options:

Structurally Complex Materials

Information Thermodynamics

Cellular Automata Computational Mechanics

Rate Distortion Theory: Statistical Mechanics of Causal Inference

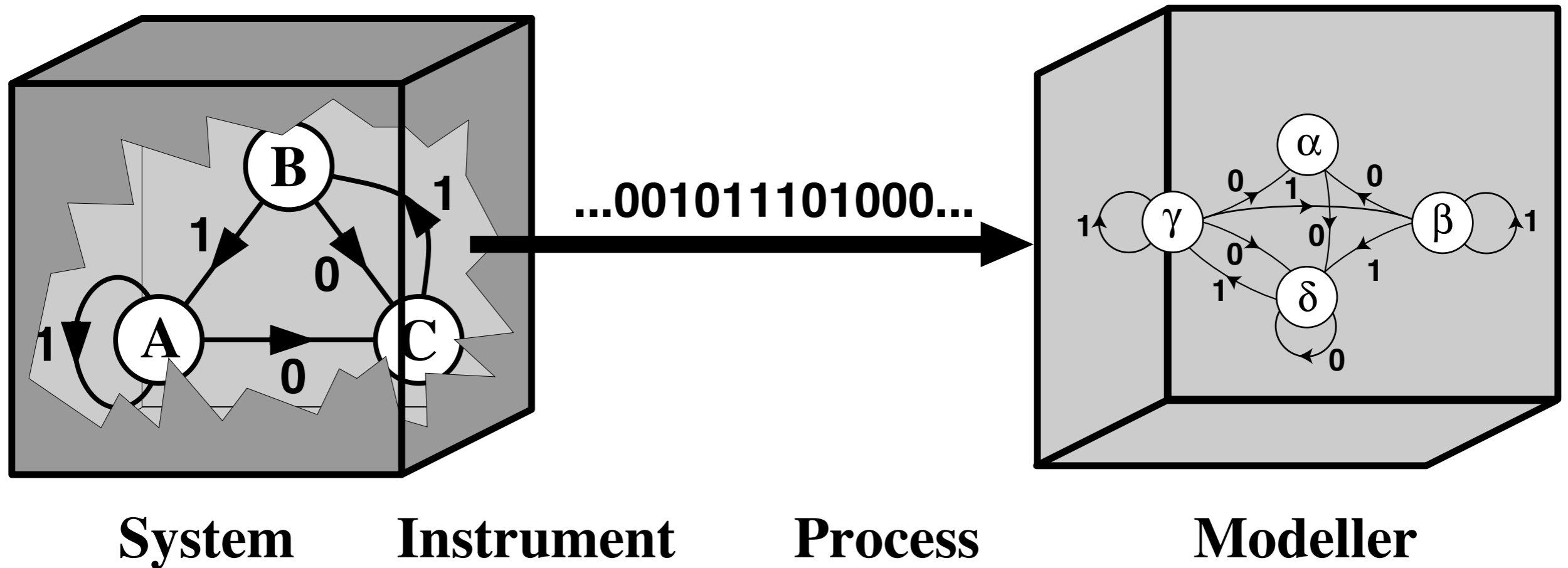
Computation at Phase Transitions and Hierarchical  $\epsilon$ -Machines

Quantum Information & Dynamics

Evolution and Self-Organization

...

# 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 ...

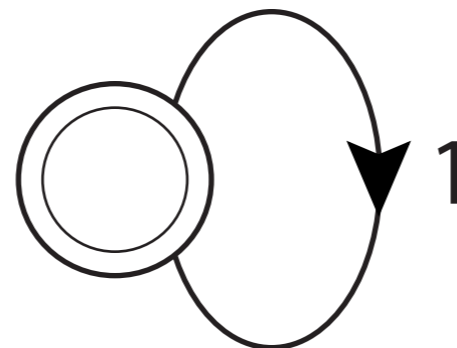
Process I:

Past: ... 111111111111

Your prediction is?

Future: 111111111111 ...

Your model (states & dynamic) is?



# The Learning Channel ...

## The Prediction Game ...

### Process II:

Past: ... 10110010001101110

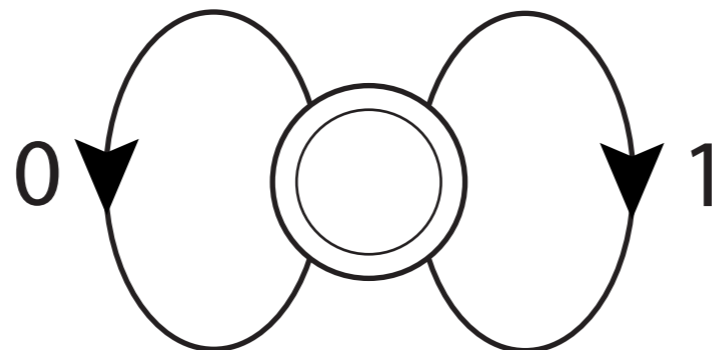
Your prediction is?

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 ...

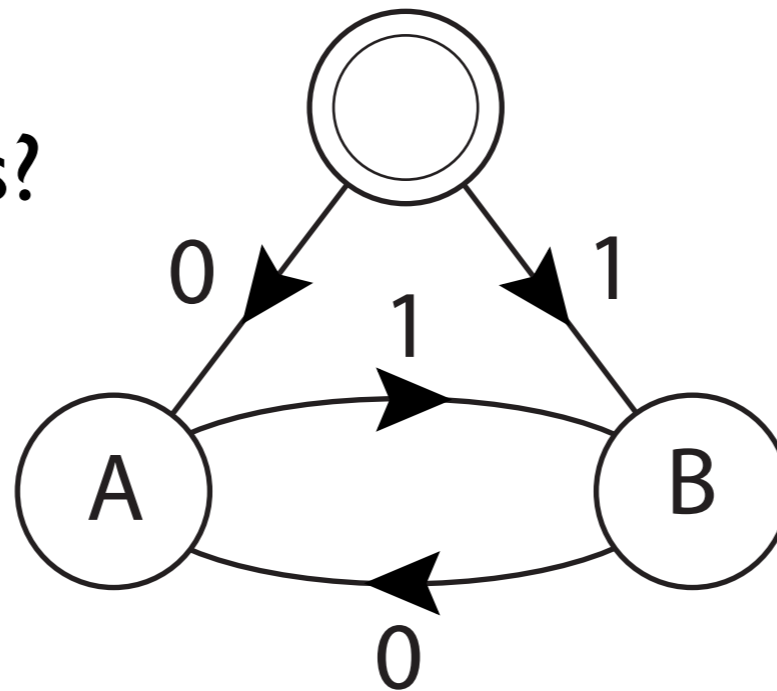
### Process III:

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?

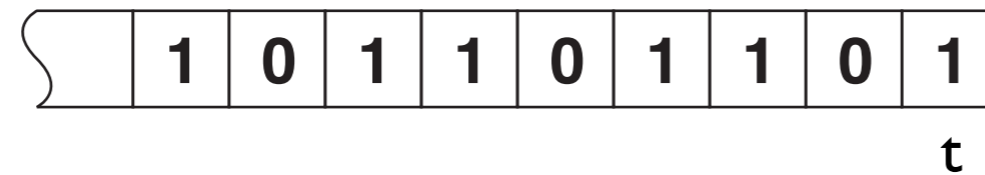
All we have are sequences of observations

Over some measurement alphabet  $\mathcal{A}$

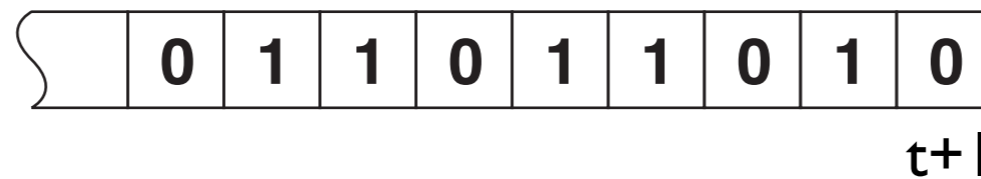
These symbols only indirectly reflect the hidden states

# The Learning Channel ...

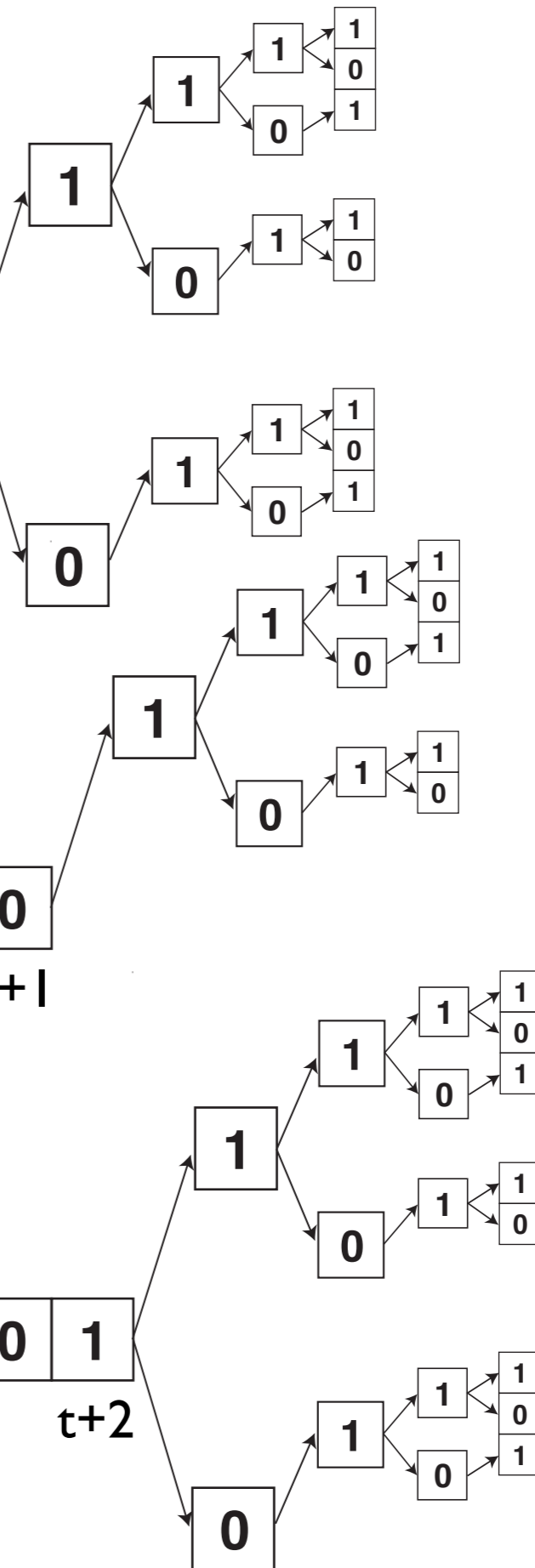
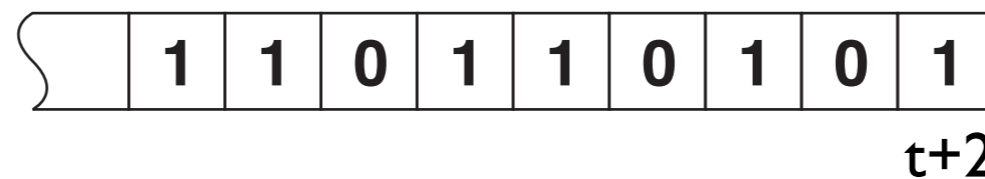
Effective States:



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



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?

Prediction!

Refined Goal: Find states that are effective for prediction.

What are the “predictive states” in the measurements?

Simple, but key observation:

Histories leading to the same predictions are equivalent.

What are these predictions?

What are these groups of histories?

# 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)

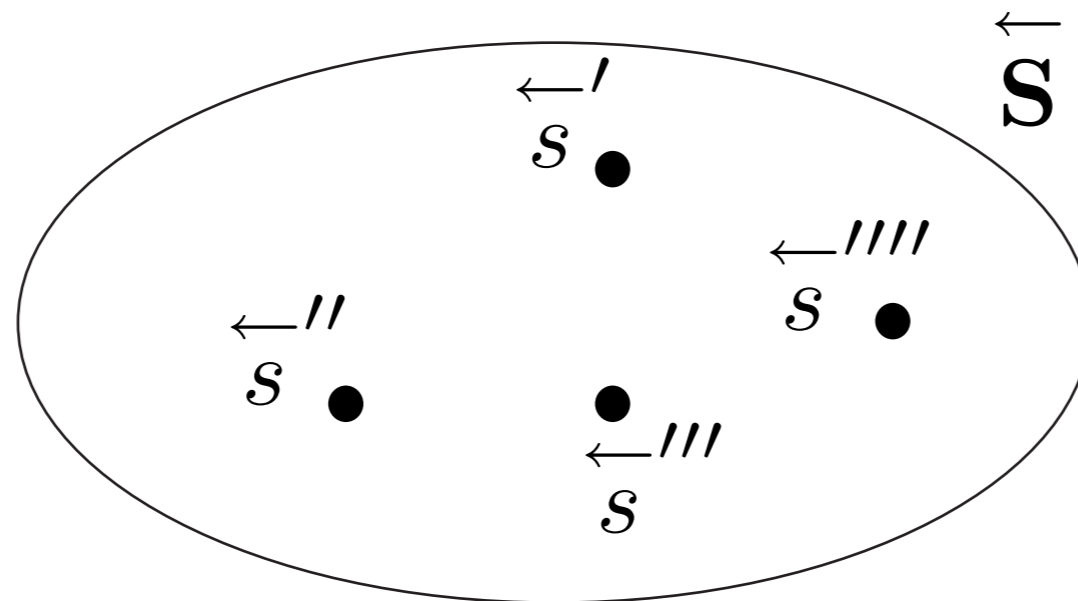
Refined goal:

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

# The Learning Channel ...

## Space of Histories:

$$\overleftarrow{\mathbf{S}} = \mathcal{A}^{\mathbb{Z}^-} = \{\dots s_{-3}s_{-2}s_{-1} : s_i \in \mathcal{A}, i = \dots, -3, -2, -1\}$$



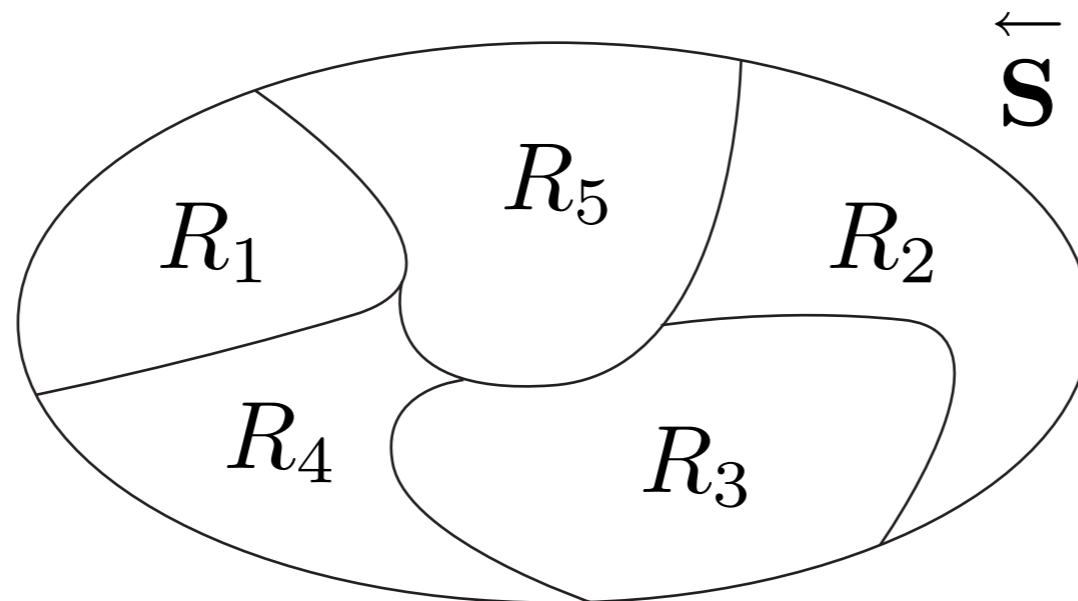
# The Learning Channel ...

## Space of Histories ...

Histories leading to the same predictions are equivalent.

Effective States = **Partitions of History**:

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



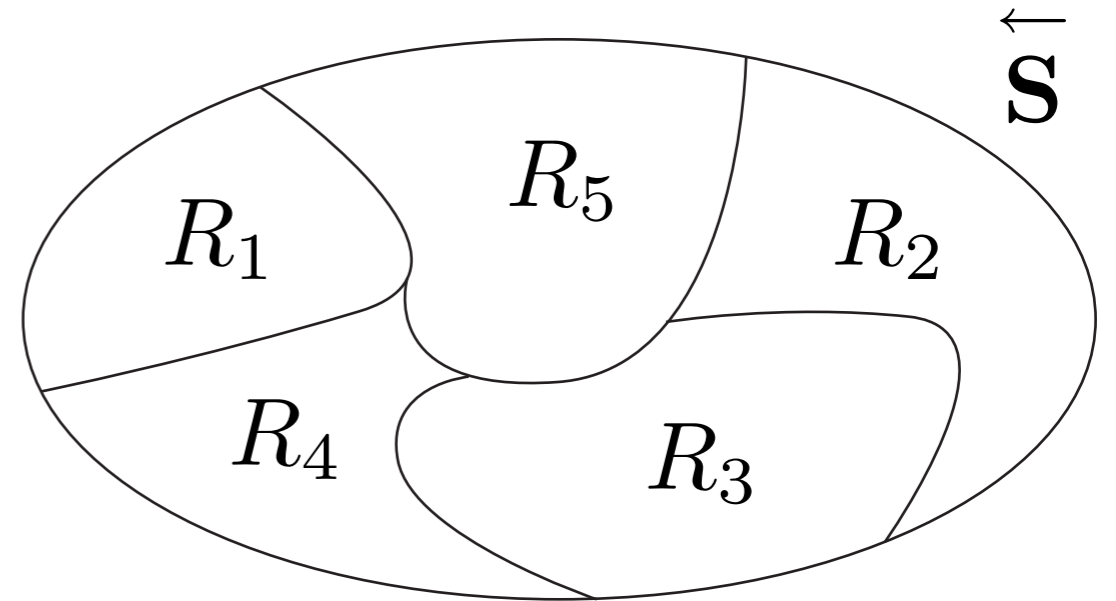
# The Learning Channel ...

## 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:

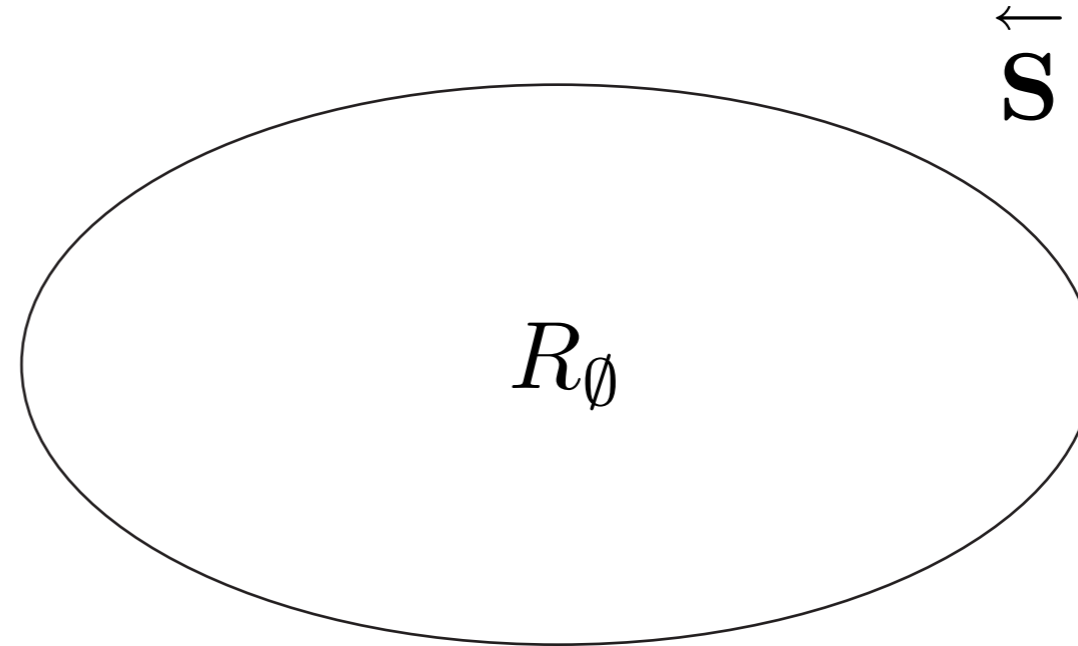
$$\Pr(R = R_i) = \sum_{\overleftarrow{s} : \eta(\overleftarrow{s}) = R_i} \Pr(\overleftarrow{s})$$

# The Learning Channel ...

## Space of Histories ...

### Null Model:

$$R_{\emptyset} = \{ \overleftarrow{s} : \overleftarrow{s} \in \mathcal{A}^{\mathbb{Z}^-} \}$$



Single State:  $R_{\emptyset}$

Single Morph:

$$\begin{aligned} \Pr(\overrightarrow{S} | R_{\emptyset}) &= \sum_{\overleftarrow{s} \in \mathcal{A}^{\mathbb{Z}^-}} \Pr(\overrightarrow{S} | \overleftarrow{s}) \Pr(\overleftarrow{s}) \\ &= \Pr(\overrightarrow{S}) \end{aligned}$$

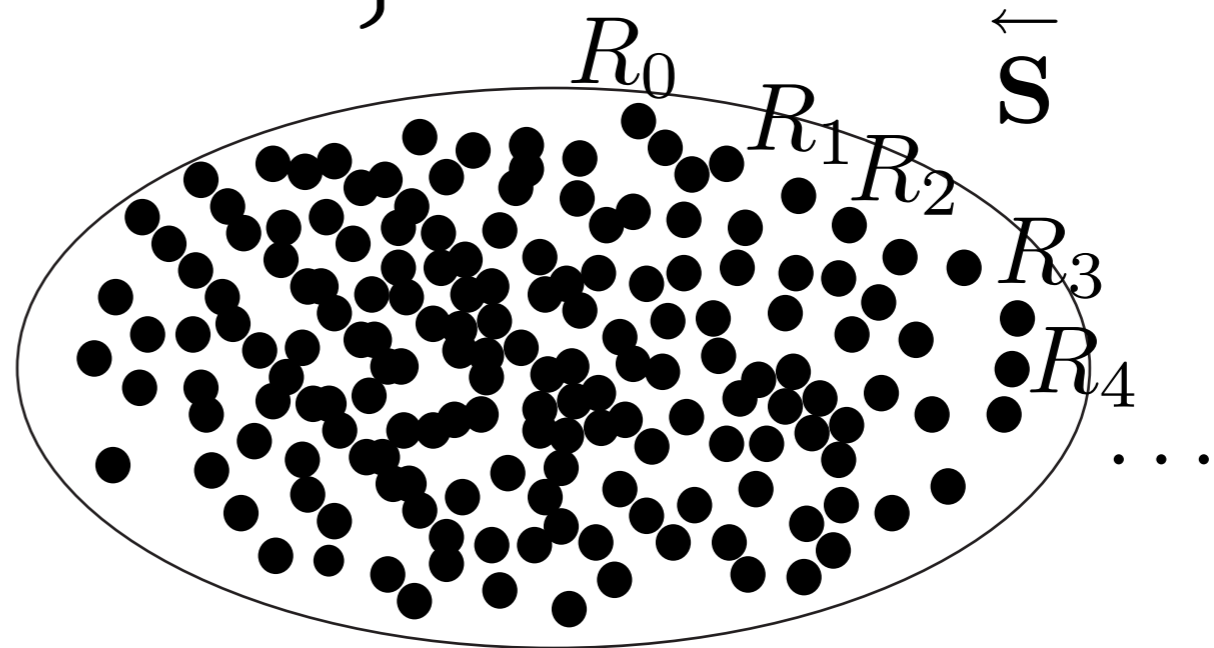


# The Learning Channel ...

## Space of Histories ...

### Every-History-Is-Precious Model:

$$R_\infty = \left\{ R_{\overleftarrow{s}} = \left\{ \overleftarrow{s} \right\} : \overleftarrow{s} \in \mathcal{A}^{\mathbb{Z}^-} \right\}$$



Each past is a state:

$$R_{\overleftarrow{s}} = \left\{ \overleftarrow{s} \right\}$$

Each past has its own future morph:

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

# 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:**

$$h_\mu(R) = H[\vec{S}^1 | R]$$

Entropy rate given effective states

# The Learning Channel ...

## How Effective are the Effective States?

### Effective Prediction Error ...

#### Bounds:

$$h_{\mu}(R) \leq \log_2 |\mathcal{A}|$$

$$h_{\mu}(R_{\emptyset}) = \log_2 |\mathcal{A}|$$

$$h_{\mu}(R_{\infty}) = h_{\mu}$$

# The Learning Channel ...

## How Effective are the Effective States?

### Effective Prediction Error ...

#### Limits on Prediction:

$$\begin{aligned} H[\overset{\rightarrow}{S}^L | R] &= H[\overset{\rightarrow}{S}^L | \eta(\overset{\leftarrow}{S})] \\ &\geq H[\overset{\rightarrow}{S}^L | \overset{\leftarrow}{S}] \end{aligned}$$

A Markov chain:

$$\overset{\leftarrow}{S} \rightarrow \eta(\overset{\leftarrow}{S}) \rightarrow R$$

(Data Processing Inequality)

Models can do no better than to use histories.

That is,  $h_\mu(R) \geq h_\mu$ .

In particular,  $h_\mu(R = \overset{\leftarrow}{S}) = h_\mu$

# The Learning Channel ...

## How Effective are the Effective States ...

### Prescience:

$$\Pi(R) = \log_2 |\mathcal{A}| - h_\mu(R)$$

### Cases:

Null model says nothing about process:

$$\Pi(R_\emptyset) = 0$$

Upper bounded by Total Predictability:

$$\Pi(R) \leq |\mathbf{G}|$$

# The Learning Channel ...

How Effective are the Effective States ...

Refined goal: Find states  $R$  such that  $h_\mu(R) = h_\mu$ .

Solution:  $h_\mu(R_\infty)$  ... rather verbose!

# 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?

$$H[\overset{\rightarrow}{S} | R] = H[\overset{\rightarrow}{S} | \overset{\leftarrow}{S}]$$

or

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

### Question 2:

Can we find the smallest such set?

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



# The Learning Channel ...

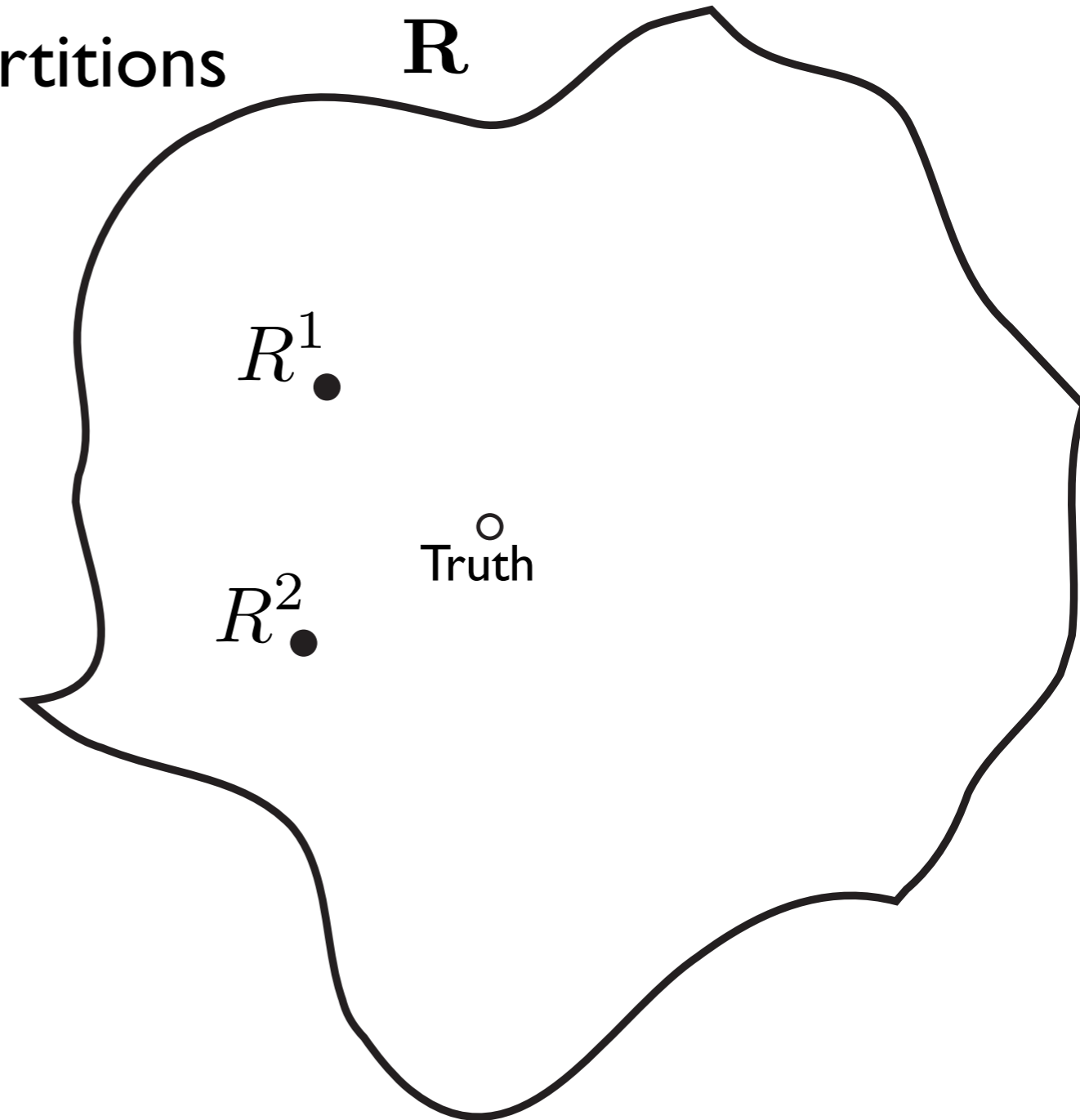
## Occam's Pool: The Space of Models

Model = Partition of History Space

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

### Rival Models:

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



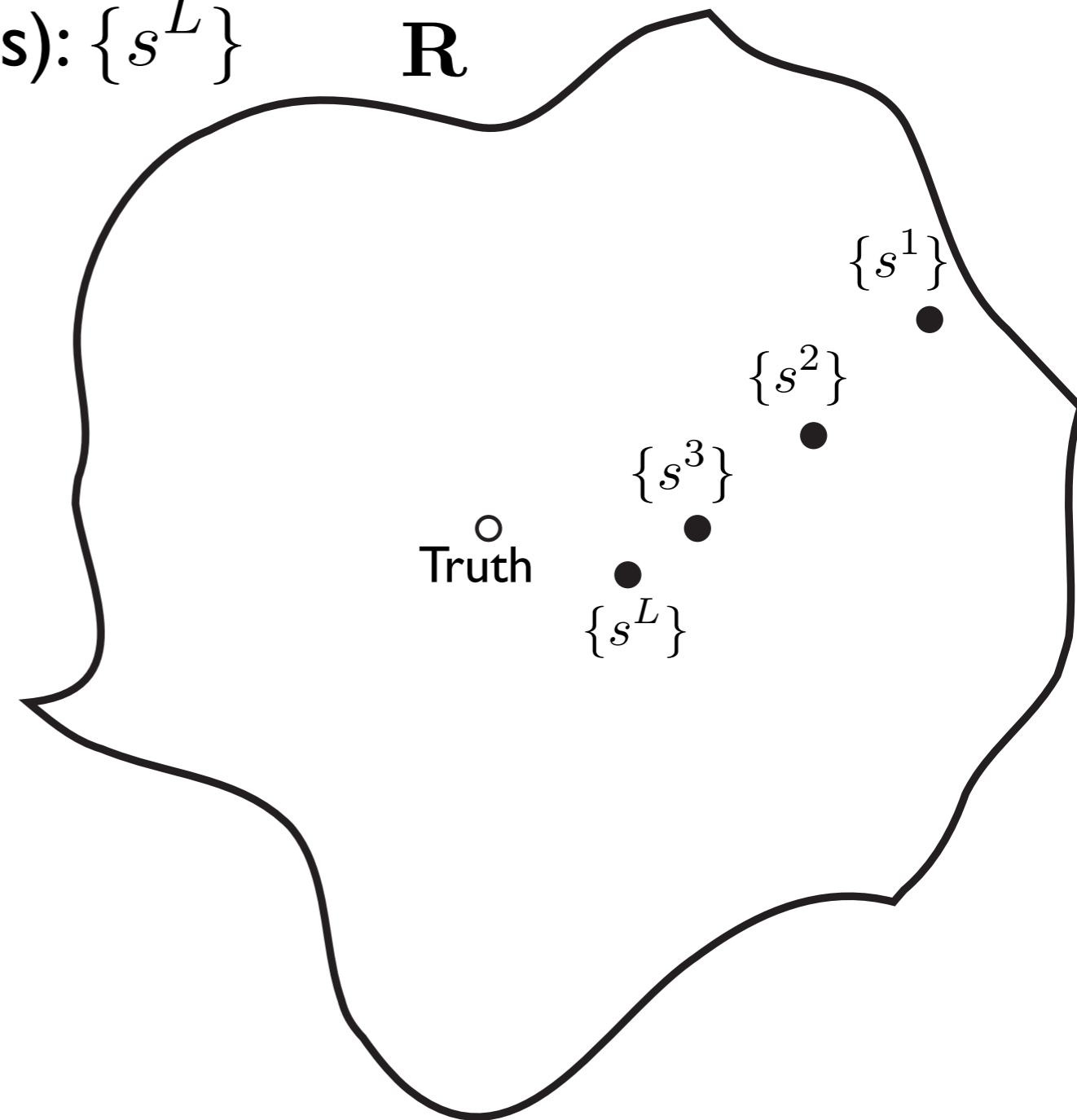
The Learning Channel ...

Occam's Pool: The Space of Models ...

A familiar (and parametrized) class of models:

Word models (aka histograms):  $\{s^L\}$

Partition history space into histories that share last  $L$  symbols.



# The Learning Channel ...

Causal States: Goals finally addressed

## Causal State:

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

Set of histories that lead to same predictions.

## Predictive (or causal) 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 States ...

### Equivalence Relation:

$\sim$  is a **relation** on histories:

Specifies subsets of  $\overleftarrow{\mathbf{S}} \times \overleftarrow{\mathbf{S}} = \{(\overleftarrow{s}', \overleftarrow{s}'') : \overleftarrow{s}', \overleftarrow{s}'' \in \overleftarrow{\mathbf{S}}\}$

$\sim$  is an **equivalence relation** on histories:

(a) Reflexive:  $\overleftarrow{s} \sim \overleftarrow{s}, \forall \overleftarrow{s} \in \overleftarrow{\mathbf{S}}$

(b) Symmetric:  $\overleftarrow{s}' \sim \overleftarrow{s}'' \Rightarrow \overleftarrow{s}'' \sim \overleftarrow{s}'$

(c) Transitive:  $\overleftarrow{s}' \sim \overleftarrow{s}'' \ \& \ \overleftarrow{s}'' \sim \overleftarrow{s}''' \Rightarrow \overleftarrow{s}' \sim \overleftarrow{s}'''$

# The Learning Channel ...

## Causal States ...

Equivalence class of  $\overleftarrow{s} \in \overleftarrow{\mathbf{S}}$ :

$$[\overleftarrow{s}] = \{ \overleftarrow{s}' \in \overleftarrow{\mathbf{S}} : \overleftarrow{s}' \sim \overleftarrow{s} \}$$

Set of all equivalence classes:

$$\overleftarrow{\mathbf{S}} / \sim$$

# The Learning Channel ...

## Causal States ...

### Induced Partition of Histories:

~ induces a partition of history space:

$$(a) [\overleftarrow{s}] \neq \emptyset$$

$$(b) \bigcup_i [\overleftarrow{s}]_i = \overleftarrow{\mathbf{S}}$$

$$(c) [\overleftarrow{s}]_i \cap [\overleftarrow{s}]_j = \emptyset, i \neq j$$

# The Learning Channel ...

## Causal State Components

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

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

Set of causal states:

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

Partition of histories:

$$\overleftarrow{\mathbf{S}} = \bigcup_i \mathcal{S}_i$$

$$\mathcal{S}_i \cap \mathcal{S}_j = \emptyset, i \neq j$$

# The Learning Channel ...

## Causal State Components ...

### Causal state map:

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

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

### Random variable:

$$\mathcal{S} = \epsilon(\overleftarrow{S})$$



# The Learning Channel ...

## Causal States ...

### Causal state morph:

$$\Pr \left( \overset{\rightarrow L}{S} \mid \mathcal{S} \right)$$

$$L = 1, 2, \dots, \forall s^L, \overleftarrow{s}$$

$$\Pr \left( \overset{\rightarrow L}{S} = s^L \mid \mathcal{S} = \epsilon(\overleftarrow{s}) \right) = \Pr \left( \overset{\rightarrow L}{S} = s^L \mid \overleftarrow{s} \right)$$

# The Learning Channel ...

## Causal States ...

Practical morph equality:

$$\Pr \left( \vec{S}^L = s^L | (s')^K \right) = \Pr \left( \vec{S}^L = s^L | (s'')^K \right)$$

$$\forall s^L \in \mathcal{A}^L, s^K \in \mathcal{A}^L \text{ and } K, L = 1, 2, 3, \dots$$

# 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:

Have history:

$$\overleftarrow{s}' = \dots s_{-3} s_{-2} s_{-1}$$

And so in state  $\mathcal{S}_i = \epsilon(\overleftarrow{s}')$

Observe symbol:  $s \in \mathcal{A}$

Have a new history:

$$\overleftarrow{s}'' = \overleftarrow{s}' s$$

$$\overleftarrow{s}'' = \dots s_{-2} s_{-1} s$$

Now in state  $\mathcal{S}_j = \epsilon(\overleftarrow{s}'')$

Transition:  $\mathcal{S}_i \xrightarrow{s} \mathcal{S}_j$

# The Learning Channel ...

## Causal State Dynamic ...

### Causal-state **filtering**:

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

### 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 Learning Channel ...

The  $\epsilon$ -Machine of a Process:

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

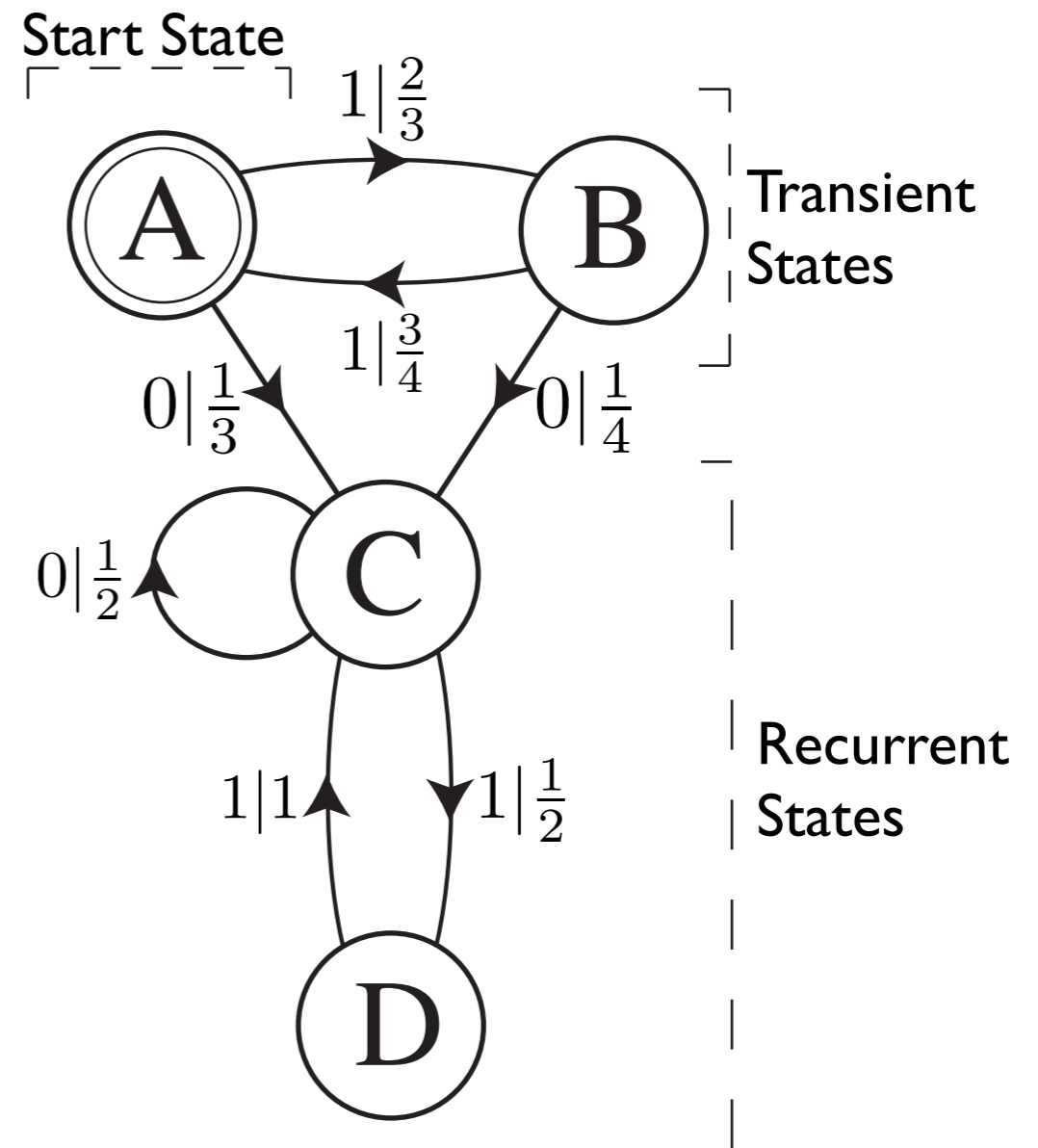
A type of hidden Markov model

# The Learning Channel ...

## The $\epsilon$ -Machine of a Process ...

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

For example ...





# The Learning Channel ...

## The $\epsilon$ -Machine ...

**Unique Start State:** Condition of total ignorance

**Null symbol:**  $\lambda$

No measurements made:

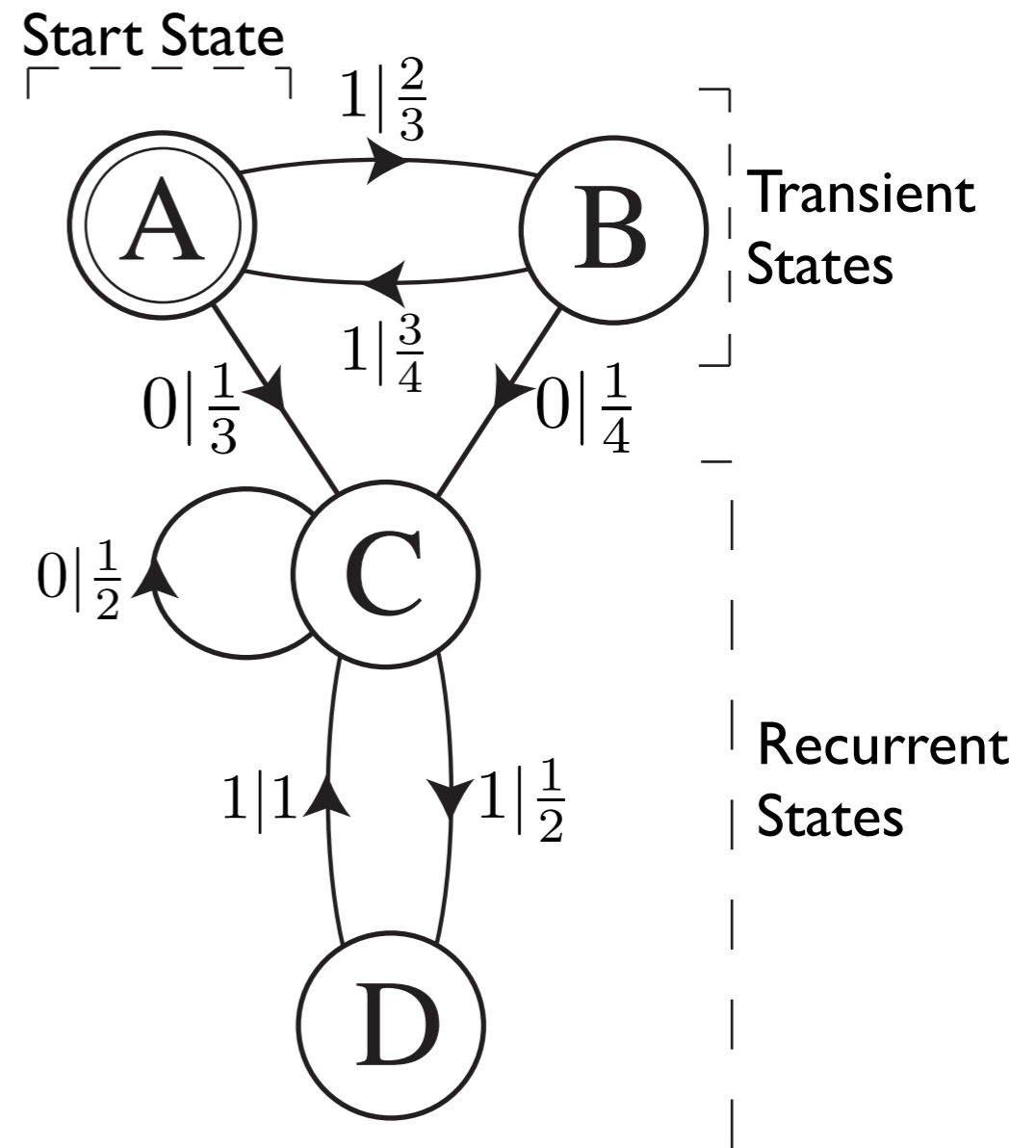
$$\overleftarrow{s} = \lambda$$

Start state:

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

**Start state distribution:**

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



# The Learning Channel ...

## The $\epsilon$ -Machine ...

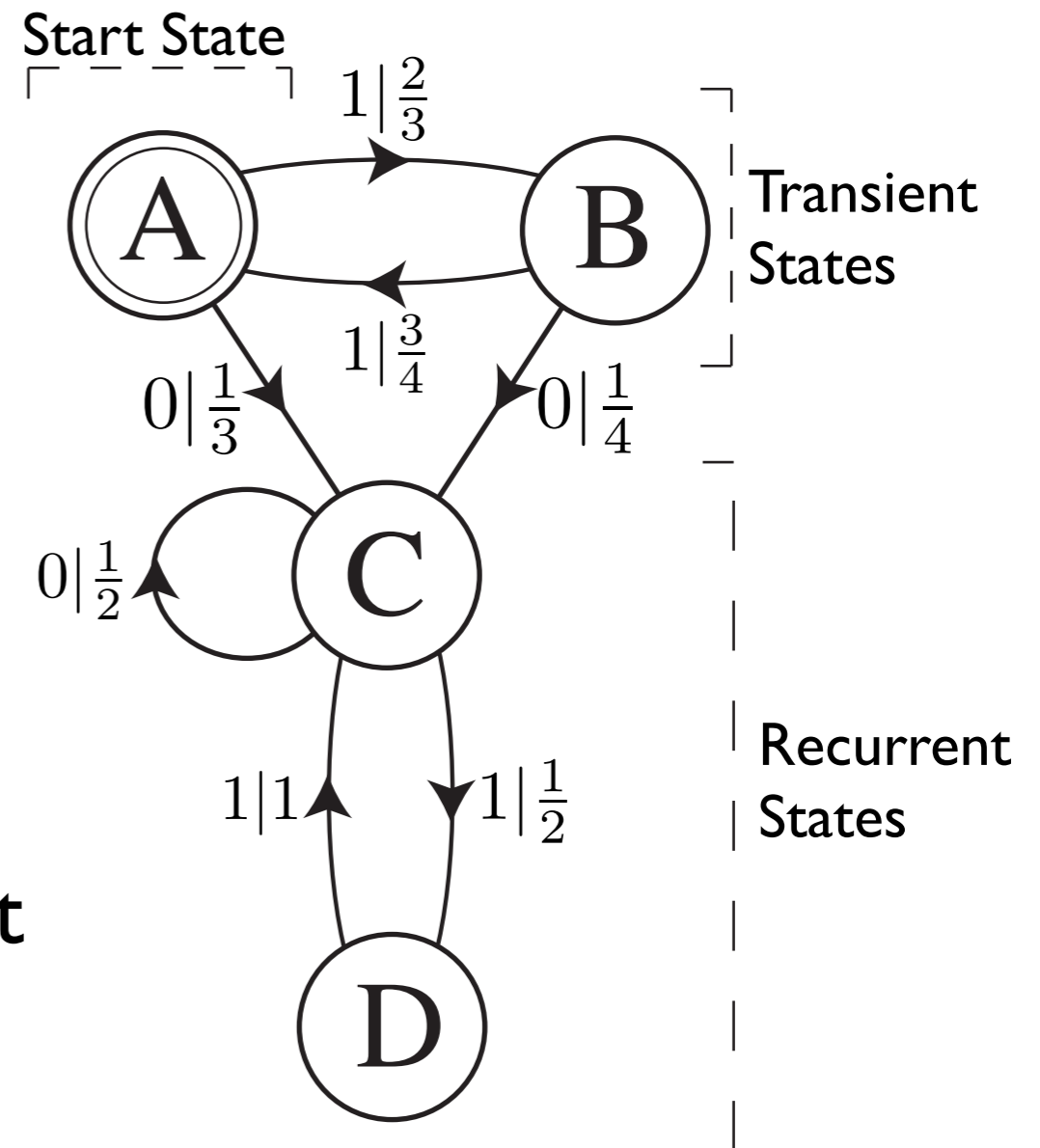
### Transient States:

How one comes to know  
process's recurrent state

### Recurrent States:

Stationary process:

Only one recurrent component



# The Learning Channel ...

## Computational Mechanics in Python: CMPy ... “campy”

A Python package for computational mechanics calculations:

- **Probability theory:**
  - Distributions: Joint, Conditional, Marginal, ...
- **Information theory:**
  - Elementary: Entropy, Conditional entropy, Mutual information, ...
  - Processes: Block entropy, Excess entropy, Transient information, ...
  - Graphics
- **Computational Mechanics:**
  - Causal state distribution
  - Entropy rate
  - Statistical complexity
  - Excess entropy
  - Causal irreversibility
  - Crypticity
  - Mixed state presentations
  - Graphics
  - Statistical inference: Subtree reconstruction, State-splitting, Optimal causal inference, ...
  - Data generation
- **Symbolic and numerical calculations**

# The Learning Channel ...

Reading for next lecture:

**Python installation:**

<http://csc.ucdavis.edu/~chaos/courses/nlp/Software/PythonInstall.html>

**Python tutorials:**

<http://csc.ucdavis.edu/~chaos/courses/nlp/Software/PythonProgramming.html>

**CMPy documentation:**

<http://cmpy.csc.ucdavis.edu/cmpy/>