

# $\epsilon M$ Reconstruction

# Machine Reconstruction ...

$\in M$  **Reconstruction:**

Any method to go from process  $\mathcal{P} \sim \text{Pr}(\vec{S})$  to its  $\in M$

(1) Analytical: Given model, equations of motion, description, ...

(2) Statistical inference: Given samples of  $\mathcal{P}$

(i) Subtree Reconstruction: Time or spacetime data to  $\in M$

(ii) State-splitting (CSSR): Time or spacetime data to  $\in M$

(iii) Spectral (eMSR): Power spectra to  $\in M$

(iv) Optimal Causal Inference: Time or spacetime data to  $\in M$

# Machine Reconstruction ...

How to reconstruct an  $\epsilon\mathcal{M}$ : **Subtree algorithm**

Given: Word distributions  $\Pr(s^D)$ ,  $D = 1, 2, 3, \dots$

Steps:

- (1) Form depth- $D$  parse tree.
- (2) Calculate node-to-node transition probabilities.
- (3) Causal states: Find morphs  $\Pr(\overrightarrow{s}^L \mid \overleftarrow{s}^K)$  as subtrees.
- (4) Label tree nodes with morph (causal state) names.
- (5) Extract state-to-state transitions from parse tree.
- (6) Assemble into  $\epsilon\mathcal{M}$ :  $\mathcal{M} = \left\{ \mathcal{S}, \{T^{(s)}, s \in \mathcal{A}\} \right\}$ .

Algorithm parameters:  $D, L, K$

# Machine Reconstruction ...

How to reconstruct an  $\epsilon M$ ...

Examples:

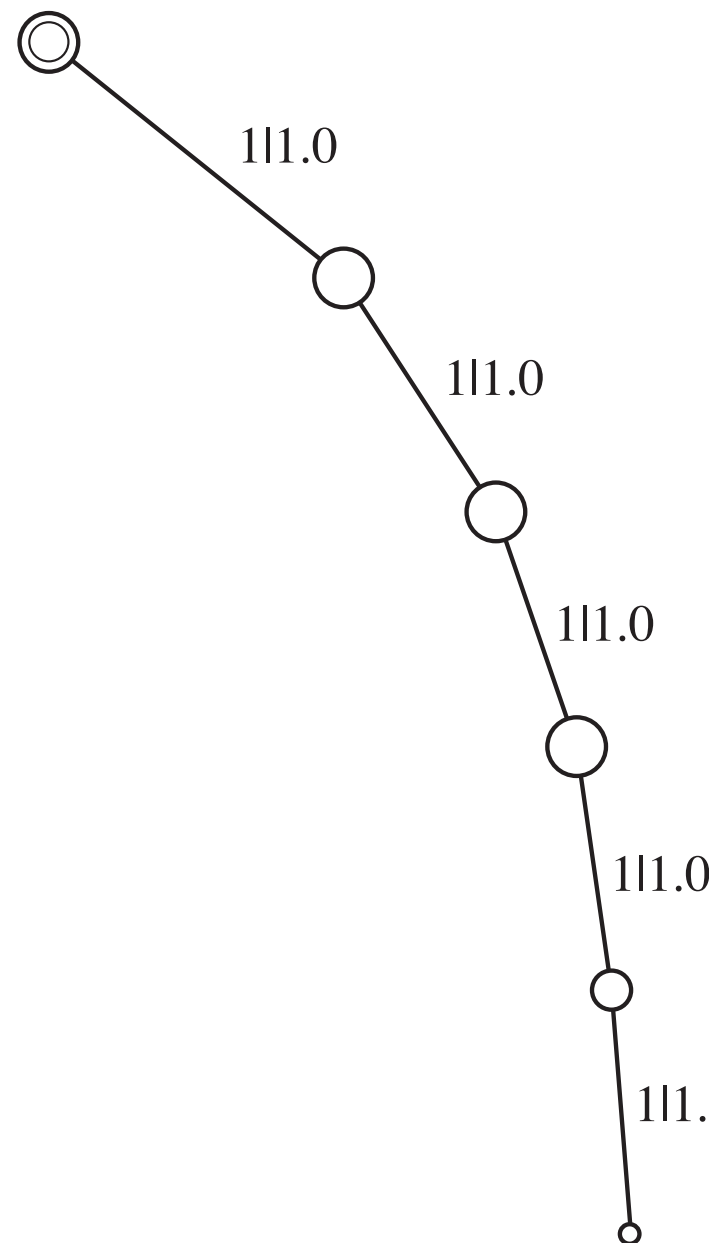
- (1) Period-1
- (2) Fair Coin
- (3) Period-2
- (4) Golden Mean Process
- (5) Even Process

# Machine Reconstruction ...

Examples (back to the Prediction Game):

Period-I: ...1111111111111111

Parse Tree  $D = 5$

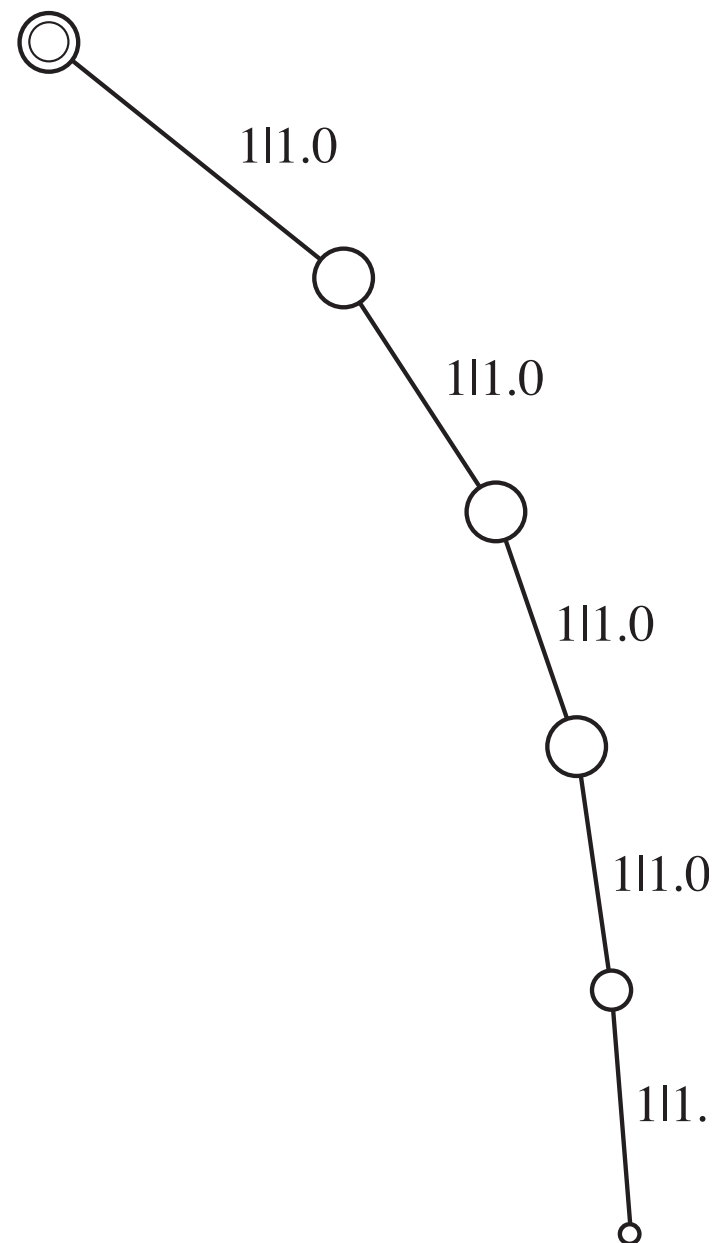


# Machine Reconstruction ...

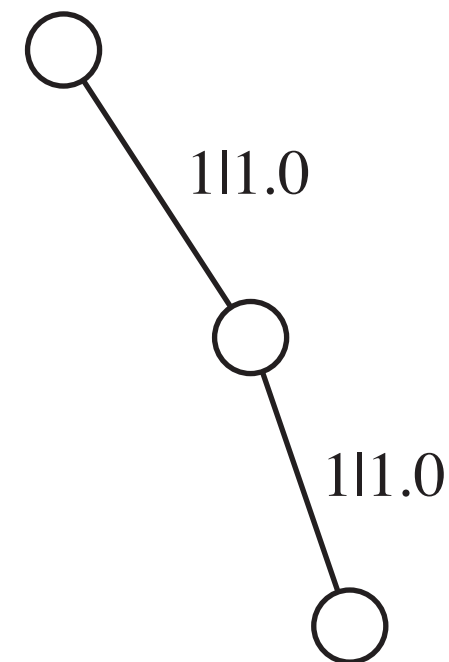
Examples (back to the Prediction Game):

Period-I: ...111111111111111

Parse Tree  $D = 5$



Morph  $L = 2$



# Machine Reconstruction ...

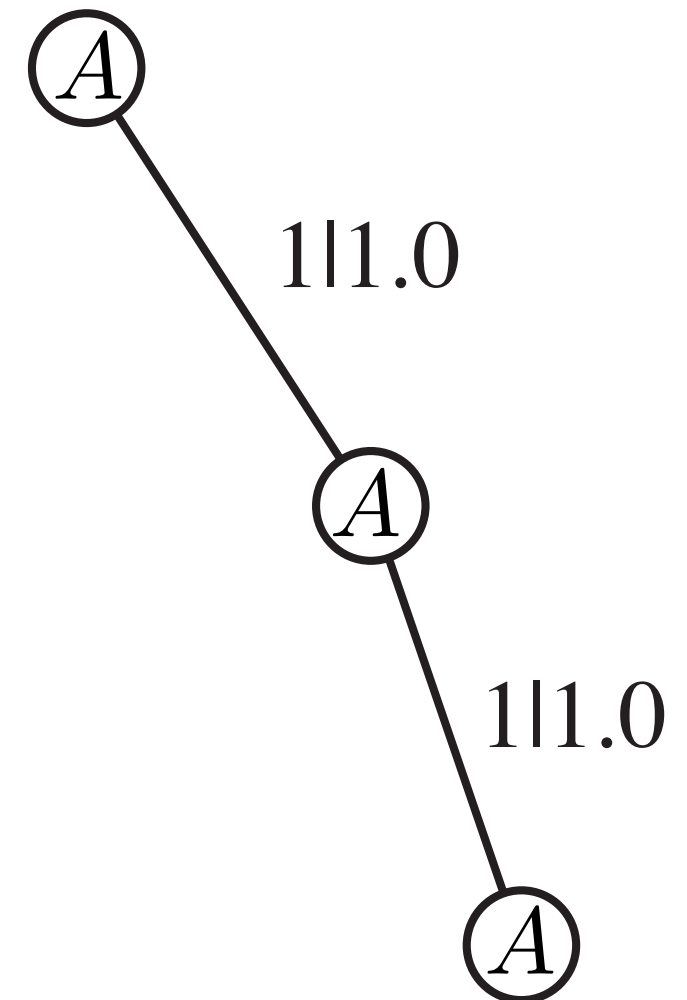
## Examples (back to the Prediction Game) ...

Period-I: ...11111111111111

Space of histories: A single point.

One future morph:  $\{1^+\}$

Morph A distribution:  $\Pr(\vec{S}^L | \overleftarrow{s}) = 1$



# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

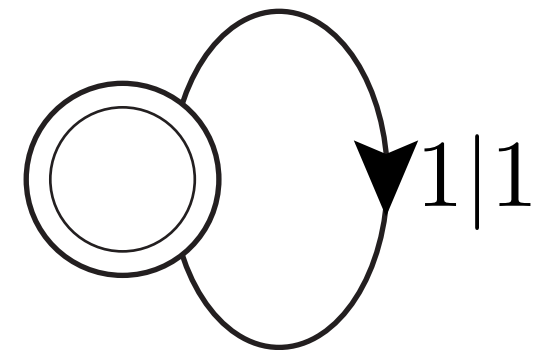
### Period-I ...

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

$$\mathcal{S} = \{\mathcal{S}_0 = \dots 111111\}$$

$$T^{(0)} = (0)$$

$$T^{(1)} = (1)$$





# Machine Reconstruction ...

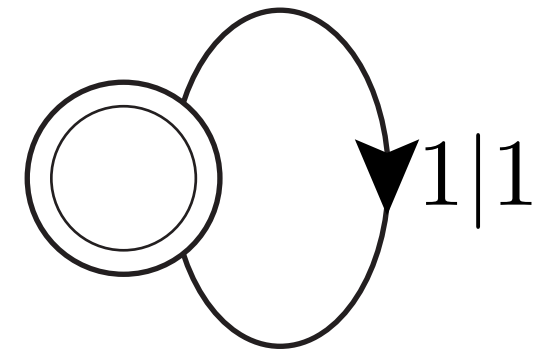
## Examples (back to the Prediction Game) ...

### Period-I ...

Causal state distribution:  $p_{\mathcal{S}} = (1)$

Entropy Rate:  $h_{\mu} = 0$  bits per symbol

Statistical Complexity:  $C_{\mu} = 0$  bits

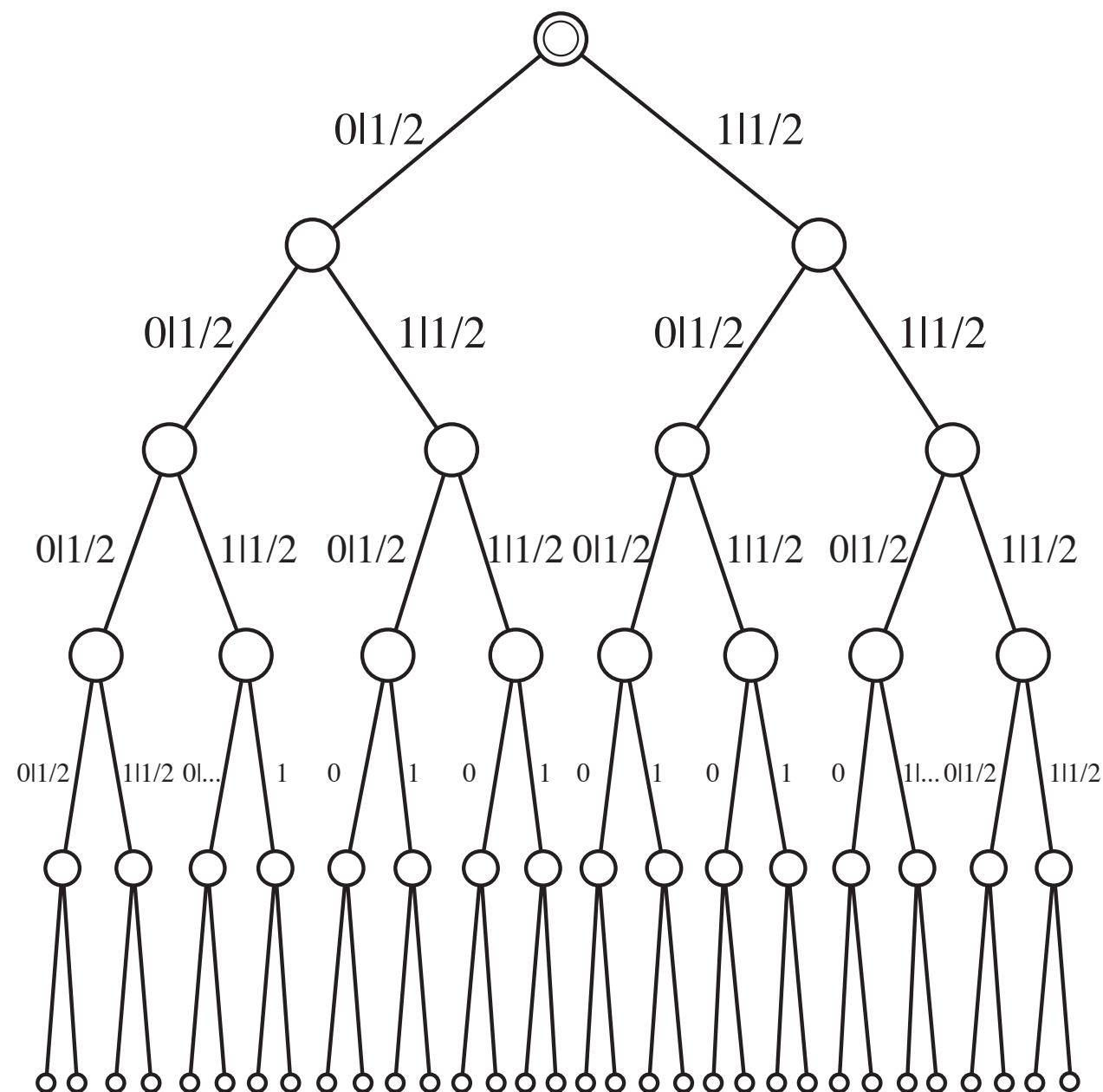


# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

Fair Coin: ... 0101001110001101

Parse Tree  $D = 5$





# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

### Fair Coin ...

Space of histories:  $\mathbf{S}^{\leftarrow K} = \mathcal{A}^K$

One future morph:  $\mathcal{A}^L$

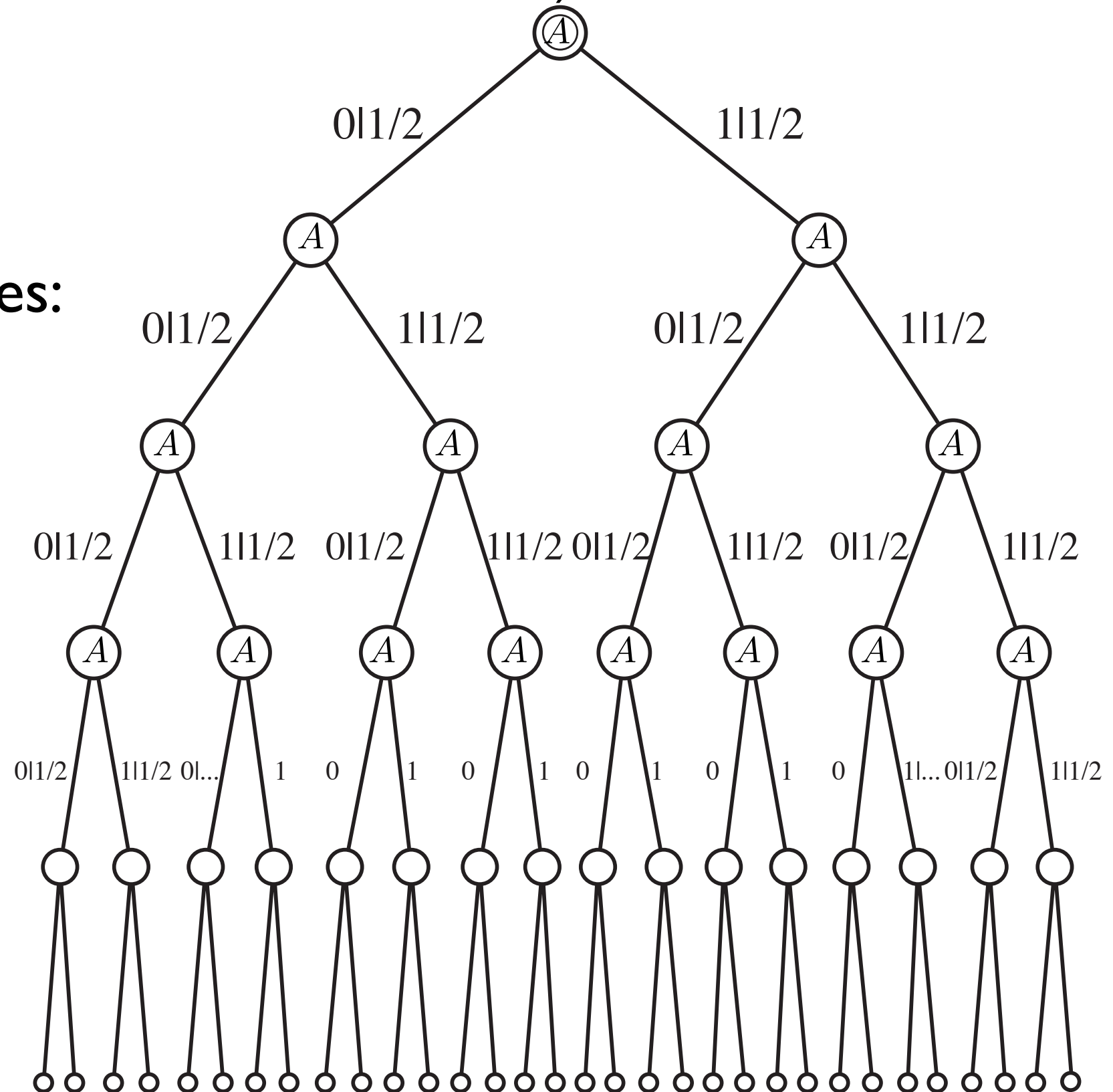
Morph A distribution:  $\Pr(\vec{S}^L \mid \overleftarrow{s}) = 2^{-L}$

# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

Fair Coin ...

Label tree nodes  
with state names:



# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

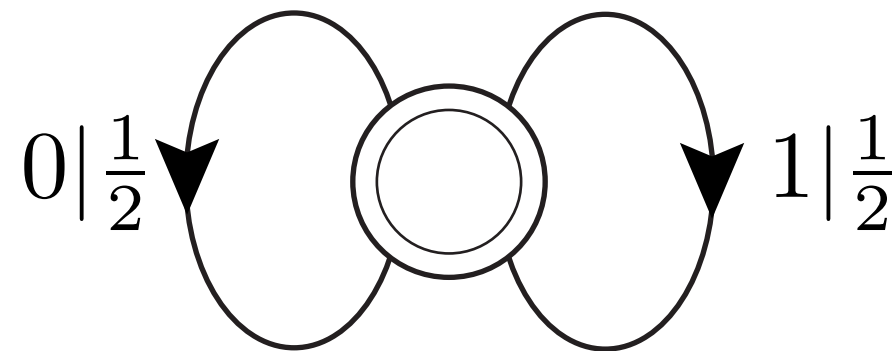
### Fair Coin ...

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

$$\mathcal{S} = \{\mathcal{S}_0 = \mathcal{A}^L\}$$

$$T^{(0)} = \left(\frac{1}{2}\right)$$

$$T^{(1)} = \left(\frac{1}{2}\right)$$



# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

### Fair Coin ...

Causal state distribution:  $p_{\mathcal{S}} = (1)$

Entropy Rate:  $h_{\mu} = 1$  bit per symbol

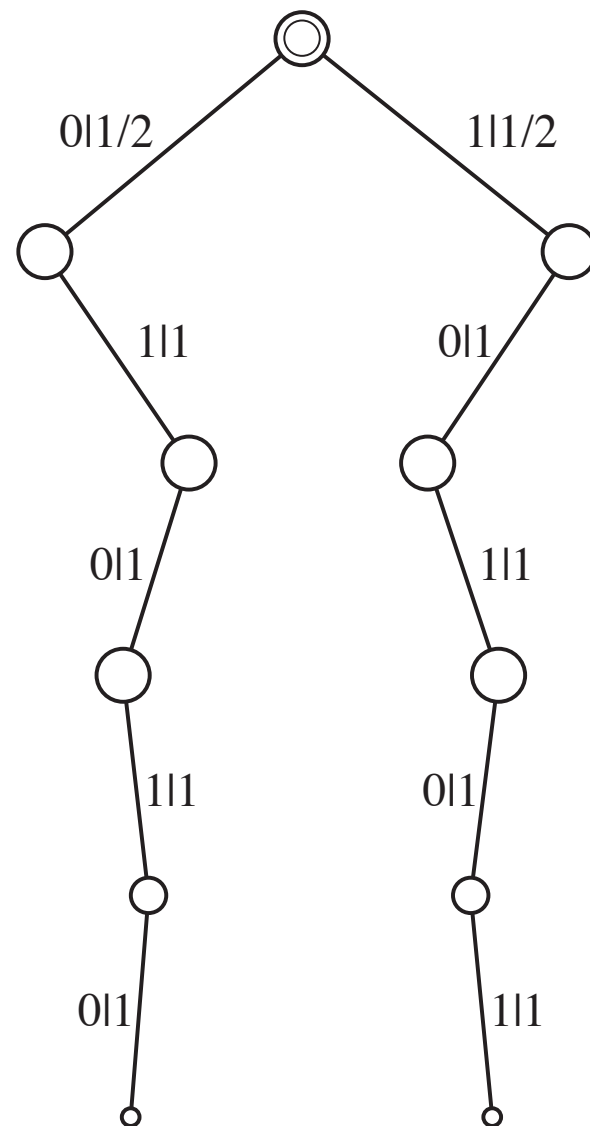
Statistical Complexity:  $C_{\mu} = 0$  bits

# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

# Period-2 Process: ...010101010101

# Parse Tree D = 5







# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

Period-2 Process: ... 010101010101

Space of histories:  $\overleftarrow{S} = \{ \overleftarrow{s}_0 = \dots 101010, \overleftarrow{s}_1 = \dots 010101 \}$

Future morphs:  $\{ \overrightarrow{S}_1 | \lambda \} = \{ 101010 \dots, 010101 \dots \}$

$\{ \overrightarrow{S}_1 | 0 \} = \{ 101010 \dots \}$

$\{ \overrightarrow{S}_1 | 1 \} = \{ 010101 \dots \}$

$\{ \overrightarrow{S}_1 | 10 \} = \{ 101010 \dots \}$

$\{ \overrightarrow{S}_1 | 01 \} = \{ 010101 \dots \}$

Morph distributions:

$$\Pr(0|\lambda) = \frac{1}{2} \quad \Pr(1|0) = 1 \quad \Pr(0|0) = 0$$

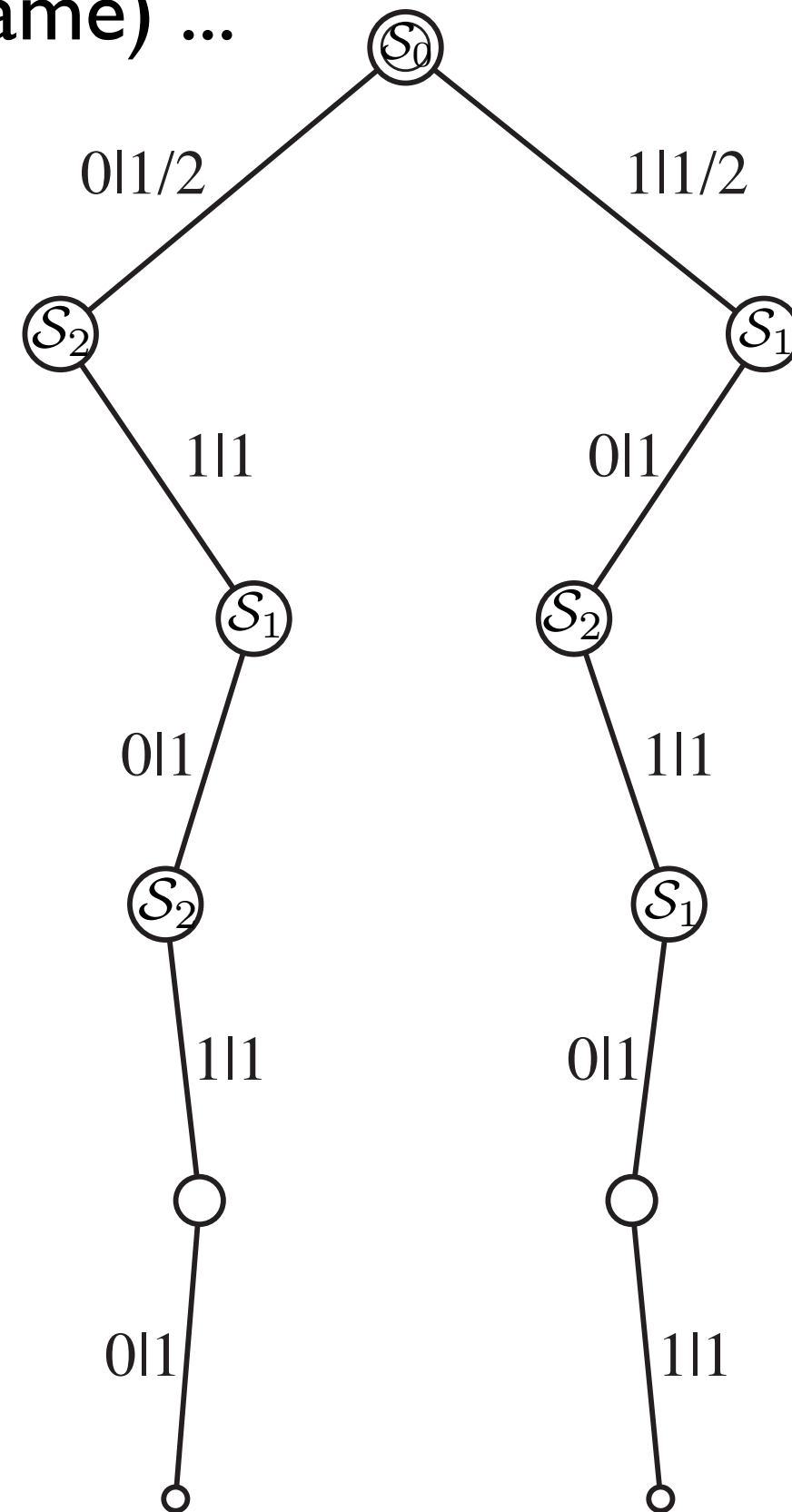
$$\Pr(1|\lambda) = \frac{1}{2} \quad \Pr(1|1) = 0 \quad \Pr(0|1) = 1$$

# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

### Period-2 Process ...

Label tree nodes:



# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

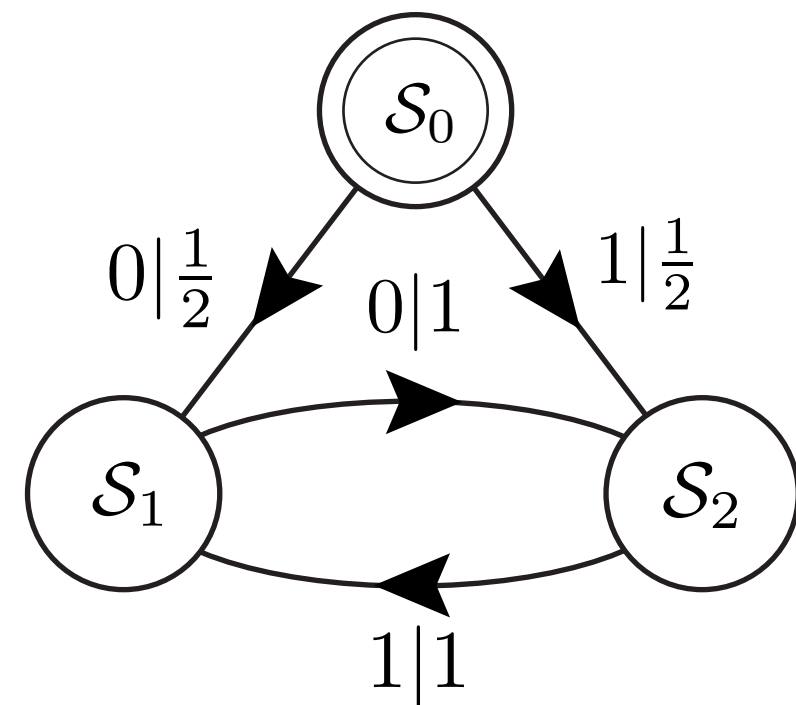
### Period-2 Process ...

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

$$\mathcal{S} = \{\mathcal{S}_0 = \{\dots 0101, \dots 1010\}, \mathcal{S}_1 = \{\dots 1010\}, \mathcal{S}_2 = \{\dots 0101\}\}$$

$$T^{(0)} = \begin{pmatrix} 0 & \frac{1}{2} & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}$$

$$T^{(1)} = \begin{pmatrix} 0 & 0 & \frac{1}{2} \\ 0 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$$



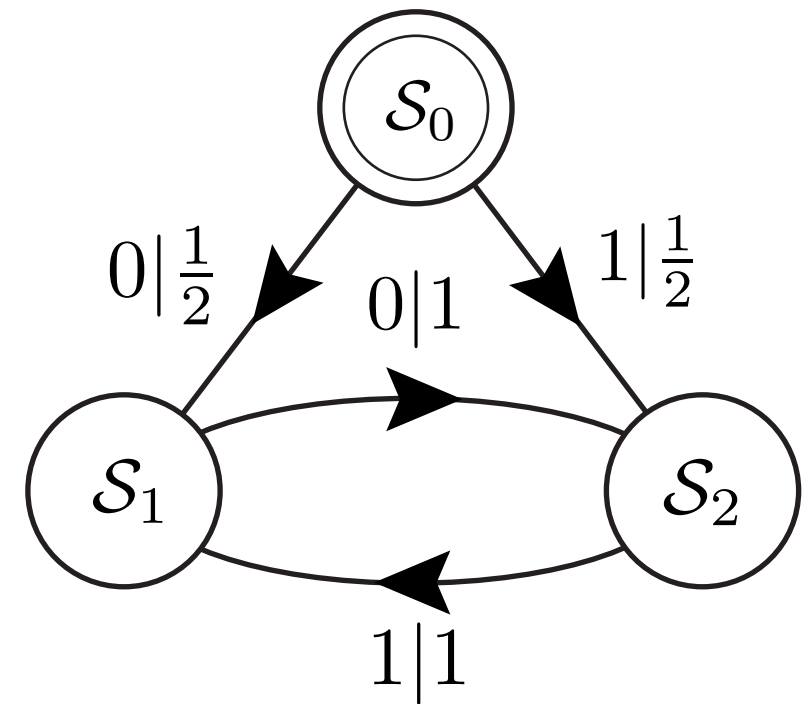
# Machine Reconstruction ...

## Examples (back to the Prediction Game) ...

### Period-2 Process ...

Causal State Distribution:

$$p_{\mathcal{S}} = \left(0, \frac{1}{2}, \frac{1}{2}\right)$$



Entropy rate:  $h_{\mu} = 0$  bits per symbol

Statistical complexity:  $C_{\mu} = 1$  bit





# Machine Reconstruction ...

## Examples ...

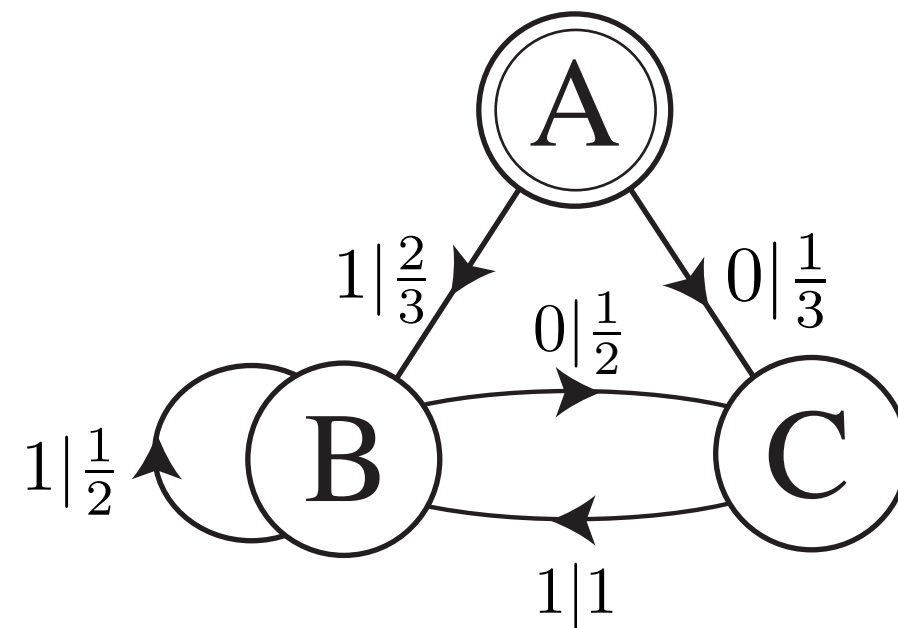
### Golden Mean Process ...

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

$$\mathcal{S} = \{A, B, C\}$$

$$T^{(0)} = \begin{pmatrix} 0 & 0 & 1/3 \\ 0 & 0 & 1/2 \\ 0 & 0 & 0 \end{pmatrix}$$

$$T^{(1)} = \begin{pmatrix} 0 & 2/3 & 0 \\ 0 & 1/2 & 0 \\ 0 & 1 & 0 \end{pmatrix}$$





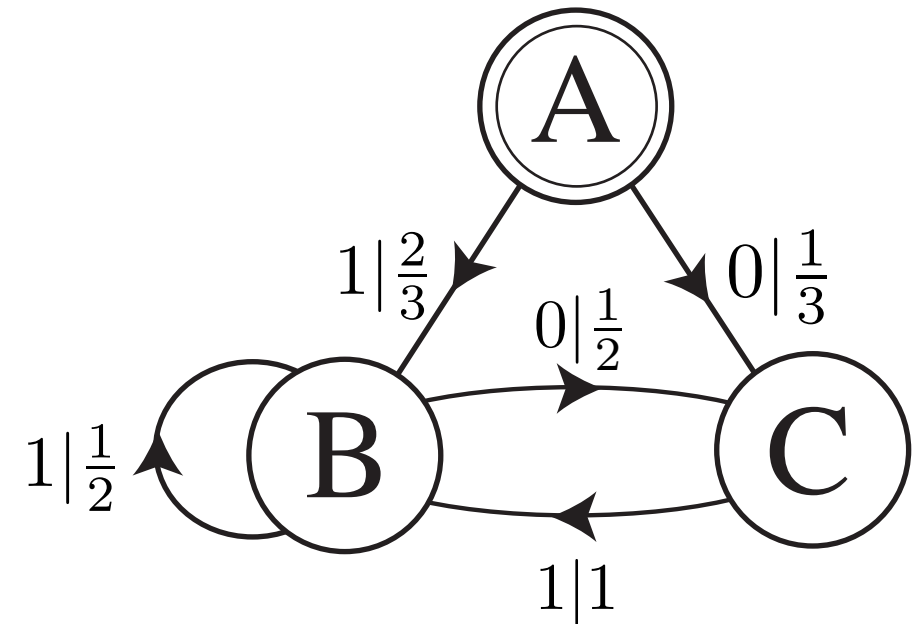
# Machine Reconstruction ...

## Examples ...

### Golden Mean Process ...

Causal State Distribution:

$$p_{\mathcal{S}} = \left(0, \frac{2}{3}, \frac{1}{3}\right)$$



Entropy rate:  $h_{\mu} = \frac{2}{3}$  bits per symbol

Statistical complexity:  $C_{\mu} = H\left(\frac{2}{3}\right)$  bits

# Machine Reconstruction ...

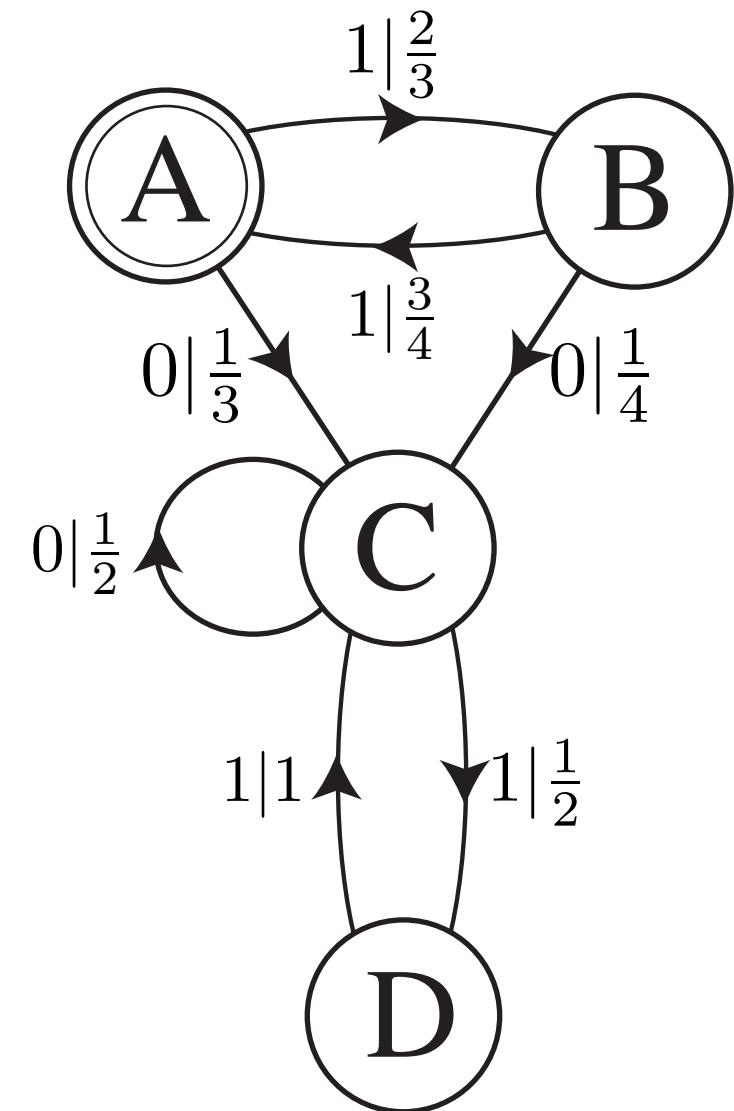
## Examples ...

Even Process:

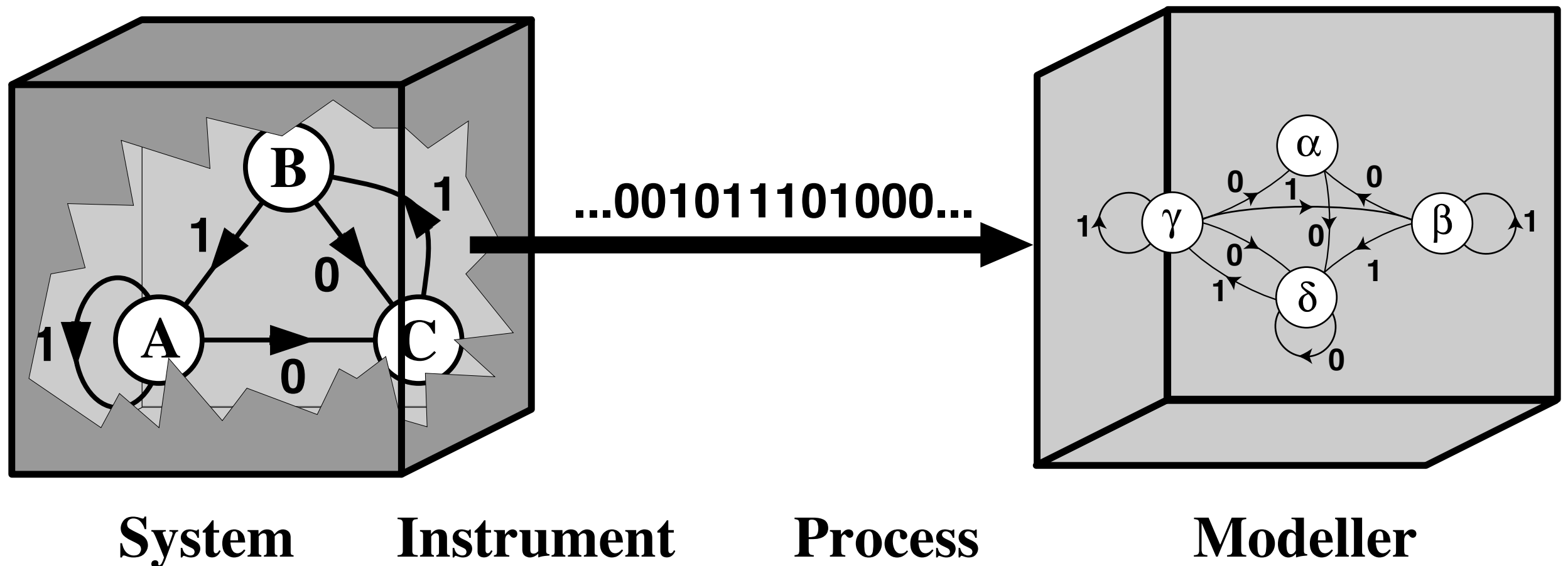
$$\mathcal{M} = \left\{ \mathcal{S}, \{T^{(s)}, s \in \mathcal{A}\} \right\}$$

Entropy rate:  $h_\mu = \frac{2}{3}$  bits per symbol

Statistical complexity:  $C_\mu = H(\frac{2}{3})$  bits



# The Learning Channel:



Central questions:

What are the states? Causal States

What is the dynamic? The  $\epsilon$ -Machine

# The $\epsilon$ -Machine ...

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

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

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

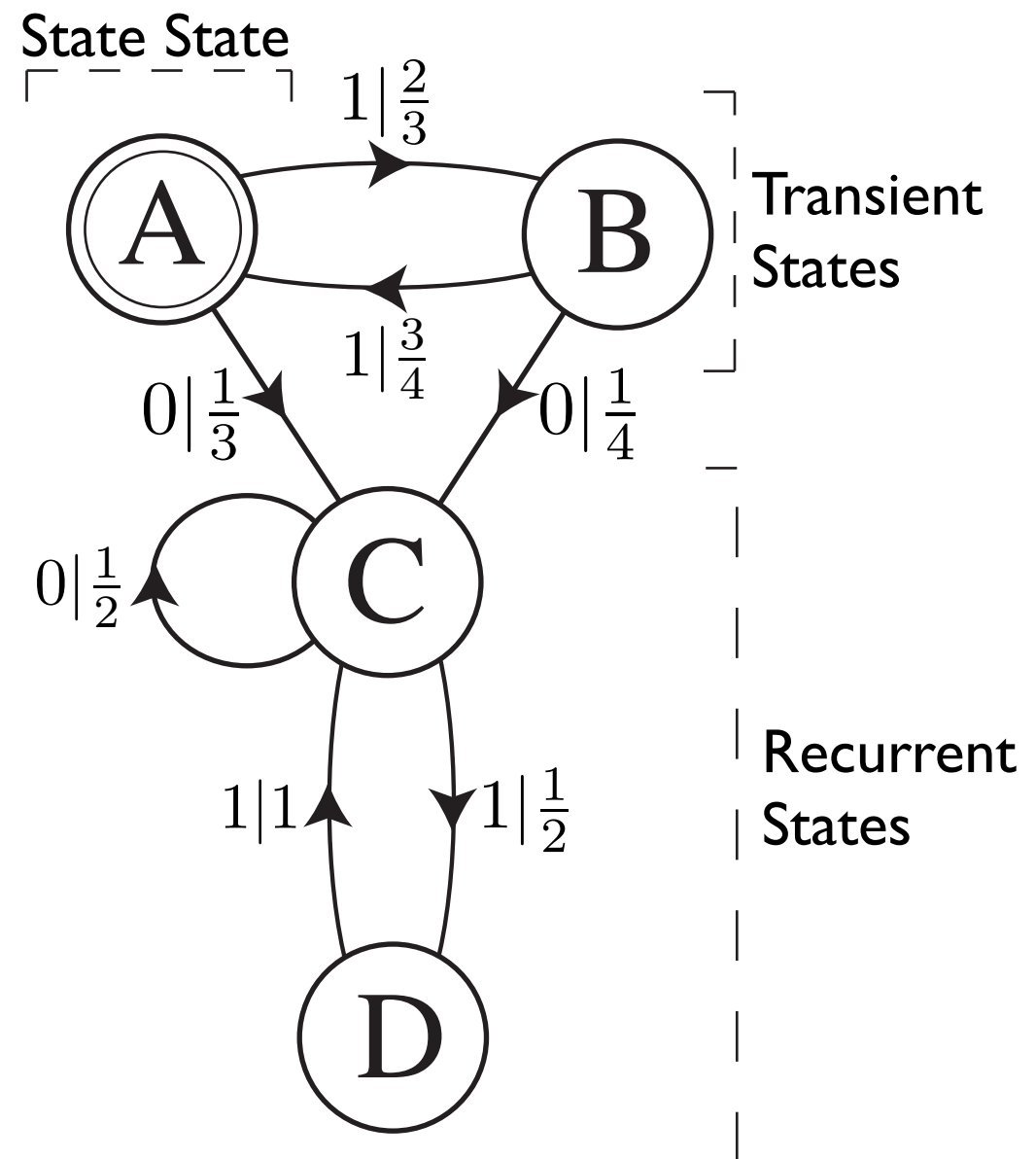
Unique Start State:

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

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

Transient States

Recurrent States



# The $\epsilon$ -Machine ...

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

$\epsilon$ -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 $\epsilon$ -Machine ...

A Model of a Process  $\text{Pr}(\vec{S})$  ...

Calculate word distribution from recurrent states:  $\mathcal{S}_i \in \mathcal{S}_{\text{recurrent}}$

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

Get  $\langle \pi | = (p_{\mathcal{S}_1}, p_{\mathcal{S}_2}, \dots)$  from  $T = \sum_{s \in \mathcal{A}} T^{(s)}$

Then

$$\text{Pr}(s) = \langle \pi | T^{(s)} | 1 \rangle \quad | 1 \rangle = \begin{pmatrix} 1 \\ 1 \\ \vdots \end{pmatrix}$$

$$\text{Pr}(s_0 s_1) = \langle \pi | T^{(s_0)} T^{(s_1)} | 1 \rangle$$

...

$$\text{Pr}(s^L) = \langle \pi | T^{(s^L)} | 1 \rangle$$

$$T^{(s^L)} = T^{(s_0)} T^{(s_1)} \dots T^{(s_{L-1})}$$

# The $\epsilon$ -Machine ...

## Properties:

Conditional independence of future & past

Unifilar

Markovian

Optimal predictor

Minimal size

Unique

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

**Proof sketch:**

$$\begin{aligned}\Pr(\overleftrightarrow{S} \mid \mathcal{S}) &= \Pr(\overleftarrow{S} \overrightarrow{S} \mid \mathcal{S}) \\ &= \Pr(\overrightarrow{S} \mid \overleftarrow{S}, \mathcal{S}) \Pr(\overleftarrow{S} \mid \mathcal{S})\end{aligned}$$

**Will show:**  $\Pr(\overrightarrow{S} \mid \overleftarrow{S}, \mathcal{S}) = \Pr(\overrightarrow{S} \mid \mathcal{S})$

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

$$\Pr\left(\overrightarrow{S} \mid \overleftarrow{S} = \overleftarrow{s}', \mathcal{S} = \epsilon(\overleftarrow{s})\right) = \Pr\left(\overrightarrow{S} \mid \overleftarrow{S} = \overleftarrow{s}\right) \quad (\overleftarrow{s}' \in [\overleftarrow{s}])$$

**But, also,**  $\Pr\left(\overrightarrow{S} \mid \overleftarrow{S} = \overleftarrow{s}\right) = \Pr\left(\overrightarrow{S} \mid \mathcal{S} = \epsilon(\overleftarrow{s})\right)$  (Causal equiv. rel'n)

$$\text{So, } \Pr\left(\overrightarrow{S} \mid \overleftarrow{S} = \overleftarrow{s}, \mathcal{S} = \sigma\right) = \Pr\left(\overrightarrow{S} \mid \mathcal{S} = \sigma\right) \quad \square$$

# The $\epsilon$ -Machine ...

$\epsilon$ 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 $\epsilon$ -Machine ...

## Unifilarity ...

Proof sketch:

Must show  $\overleftarrow{s} \sim \overleftarrow{s}' \Rightarrow \overleftarrow{s}s \sim \overleftarrow{s}'s$

Futures with symbol prefixed:  $sF$   $F \subset \mathcal{A}^\infty$

$$\begin{aligned}
 \overleftarrow{s} \sim \overleftarrow{s}' &\Rightarrow \Pr\left(\overrightarrow{S} \in sF \mid \overleftarrow{S} = \overleftarrow{s}\right) = \Pr\left(\overrightarrow{S} \in sF \mid \overleftarrow{S} = \overleftarrow{s}'\right) \\
 &\Pr\left(\overrightarrow{S}^1 = s, \overrightarrow{S}_1 \in F \mid \overleftarrow{S} = \overleftarrow{s}\right) = \Pr\left(\overrightarrow{S}^1 = s, \overrightarrow{S}_1 \in F \mid \overleftarrow{S} = \overleftarrow{s}'\right) \\
 \Pr\left(\overrightarrow{S}_1 \in F \mid \overrightarrow{S}^1 = s, \overleftarrow{S} = \overleftarrow{s}\right) \Pr\left(\overrightarrow{S}^1 = s \mid \overleftarrow{S} = \overleftarrow{s}\right) &= \Pr\left(\overrightarrow{S}_1 \in F \mid \overrightarrow{S}^1 = s, \overleftarrow{S} = \overleftarrow{s}'\right) \Pr\left(\overrightarrow{S}^1 = s \mid \overleftarrow{S} = \overleftarrow{s}'\right) \\
 \Pr\left(\overrightarrow{S}_1 \in F \mid \overleftarrow{S} = \overleftarrow{s}s\right) &= \Pr\left(\overrightarrow{S}_1 \in F \mid \overleftarrow{S} = \overleftarrow{s}'s\right) \\
 &\Rightarrow \overleftarrow{s}s \sim \overleftarrow{s}'s \quad \square
 \end{aligned}$$

(Stationarity and  
by assumption  
 $\Pr(\overrightarrow{S}^1 = s \mid \overleftarrow{s}) = 1$   
 $\Pr(\overrightarrow{S}^1 = s \mid \overleftarrow{s}') = 1$ )

# The $\epsilon$ -Machine ...

## Unifilarity ...

### Consequence:

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

Entropy rate expression requires this 1-1 mapping.

Can (must) use  $\epsilon M$  to calculate entropy rate  $h_\mu$ .

# The $\epsilon$ -Machine ...

$\epsilon$ Ms are **Optimal Predictors**:

Compared to any rival effective states  $R$ :

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

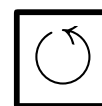
Proof sketch:  $H \left[ \overrightarrow{S}^L \mid \mathcal{S} \right] = H \left[ \overrightarrow{S}^L \mid \overleftarrow{s} \in \mathcal{S} \right]$  (Causal equiv. rel'n)

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

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

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

(Data processing inequality)



# The $\epsilon$ -Machine ...

## $\epsilon$ Ms are Optimal Predictors ...

Lemma:


$$h_\mu(\mathcal{S}) = h_\mu$$

**Proof:**  $h_\mu(\mathcal{S}) = \lim_{L \rightarrow \infty} \frac{1}{L} H \left[ \vec{S}^L | \mathcal{S} \right]$  (Block entropy)

$= \lim_{L \rightarrow \infty} \frac{1}{L} H \left[ \vec{S}^L | \overleftarrow{\mathcal{S}} \right]$  (Causal equiv. rel'n)

$= \lim_{L \rightarrow \infty} \frac{1}{L} L H \left[ S | \overleftarrow{\mathcal{S}} \right]$  (Stationarity)

$= H \left[ S | \overleftarrow{\mathcal{S}} \right]$

$= h_\mu$  

**Corollary:**  $h_\mu(R) \geq h_\mu$

# The $\epsilon$ -Machine ...

$\epsilon$ Ms are Optimal Predictors ...

Corollary (**Maximal Prescience**):  $\Pi(R) \leq \Pi(\mathcal{S})$

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

But:  $\Pi(\mathcal{S}) = \log_2 |\mathcal{A}| - h_\mu = \mathbf{G}$

So:  $\Pi(R) \leq \Pi(\mathcal{S})$   $h_\mu(R) \geq h_\mu$

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.  
(See below.)

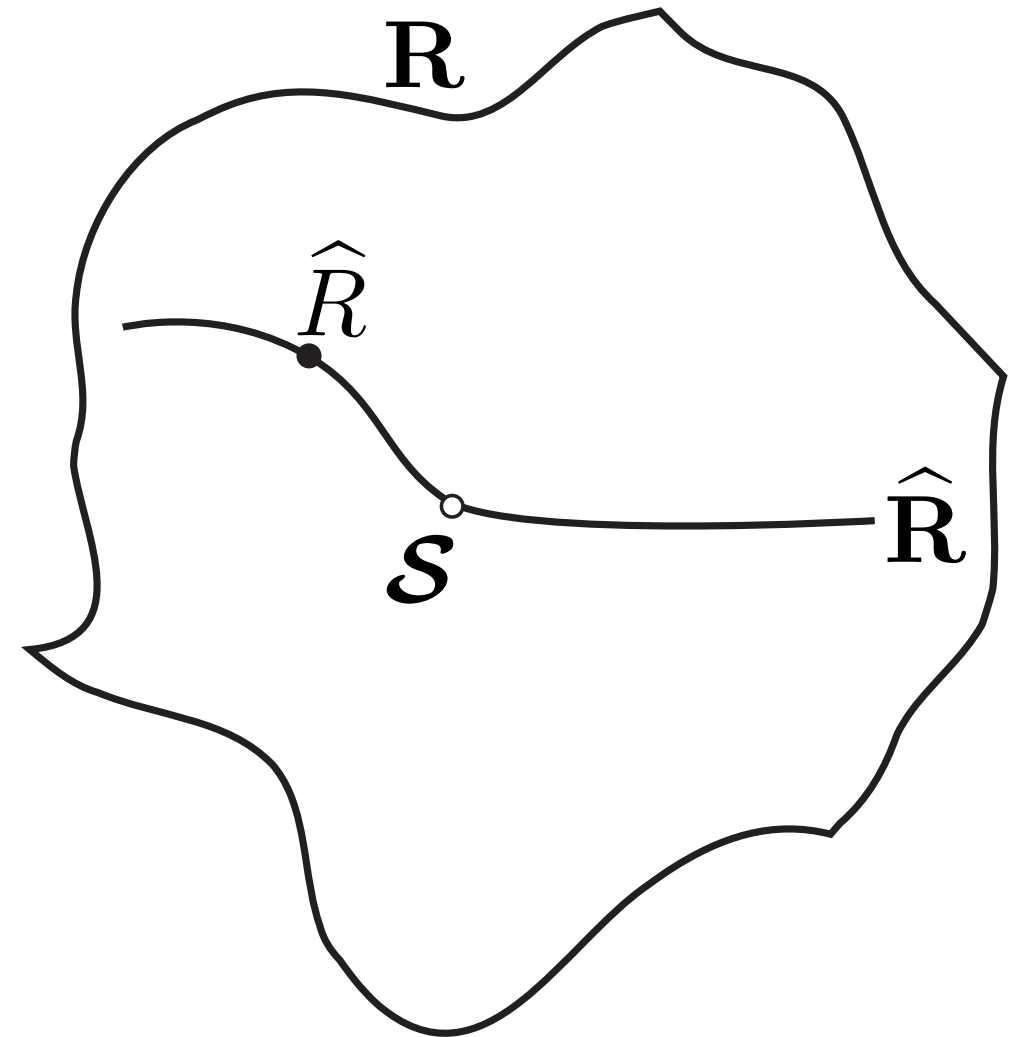
# The $\epsilon$ -Machine ...

## 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 $\epsilon$ -Machine ...

## Minimal Statistical Complexity:

For all prescient rivals,  $\epsilon M$  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 $\epsilon$ -Machine ...

## Minimal Statistical Complexity ...

### Consequence:

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

# The $\epsilon$ -Machine ...

## Summary:

$\epsilon M$ :

- (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 $\epsilon$ -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?