

eM's for Model-Based Reinforcement Learning

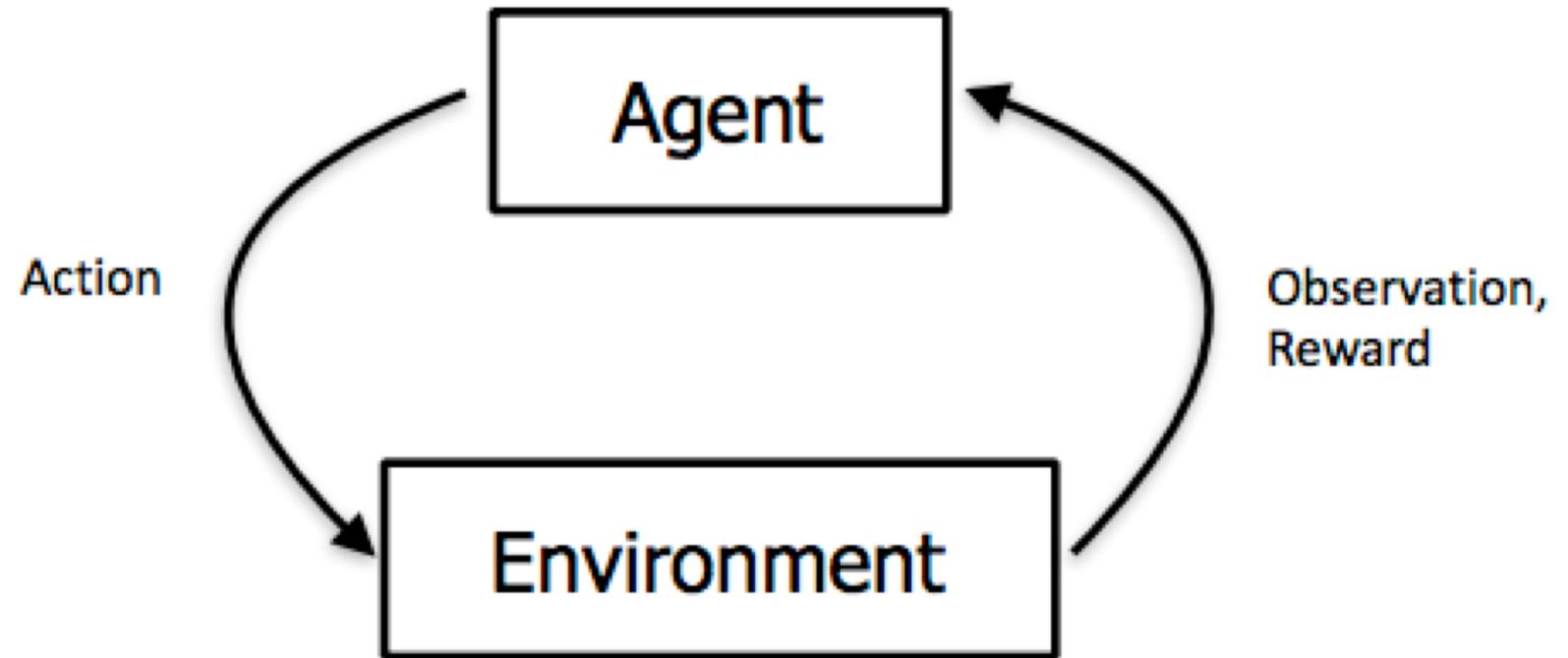
POCI Project, 2021

Sam Holton

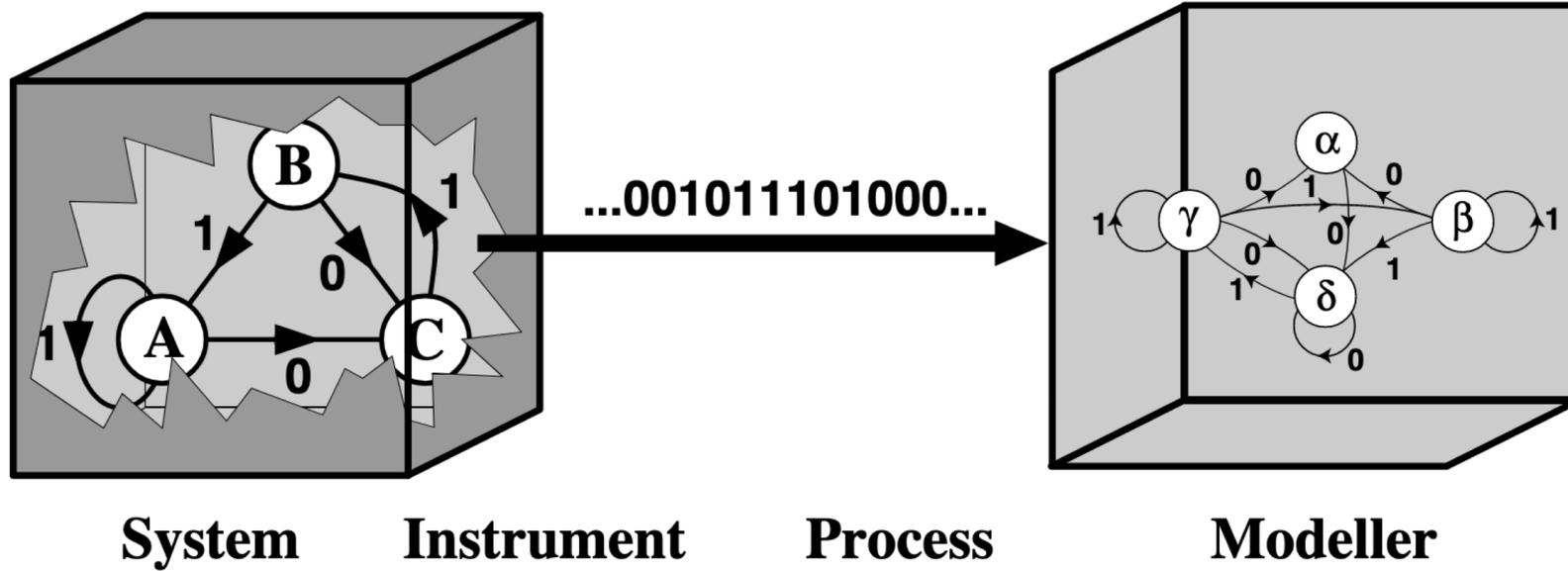
Motivation

Chaos, Quantum Mechanics -> “observer” has active role

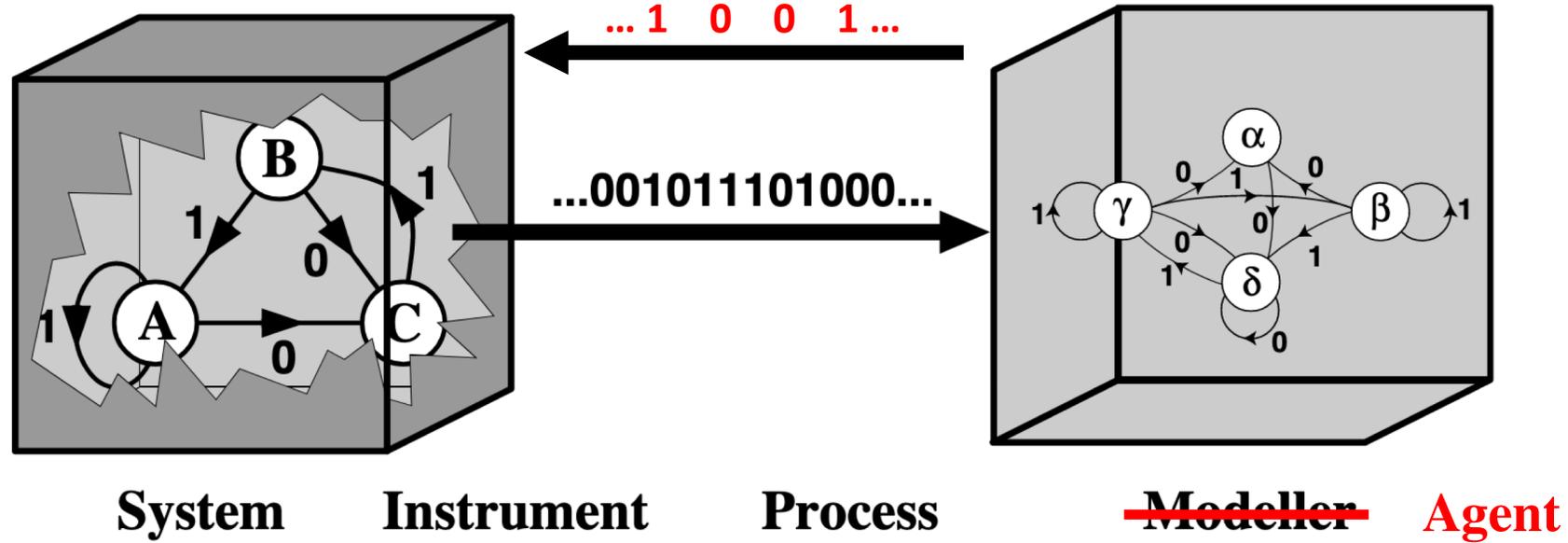
Control the system -> RL



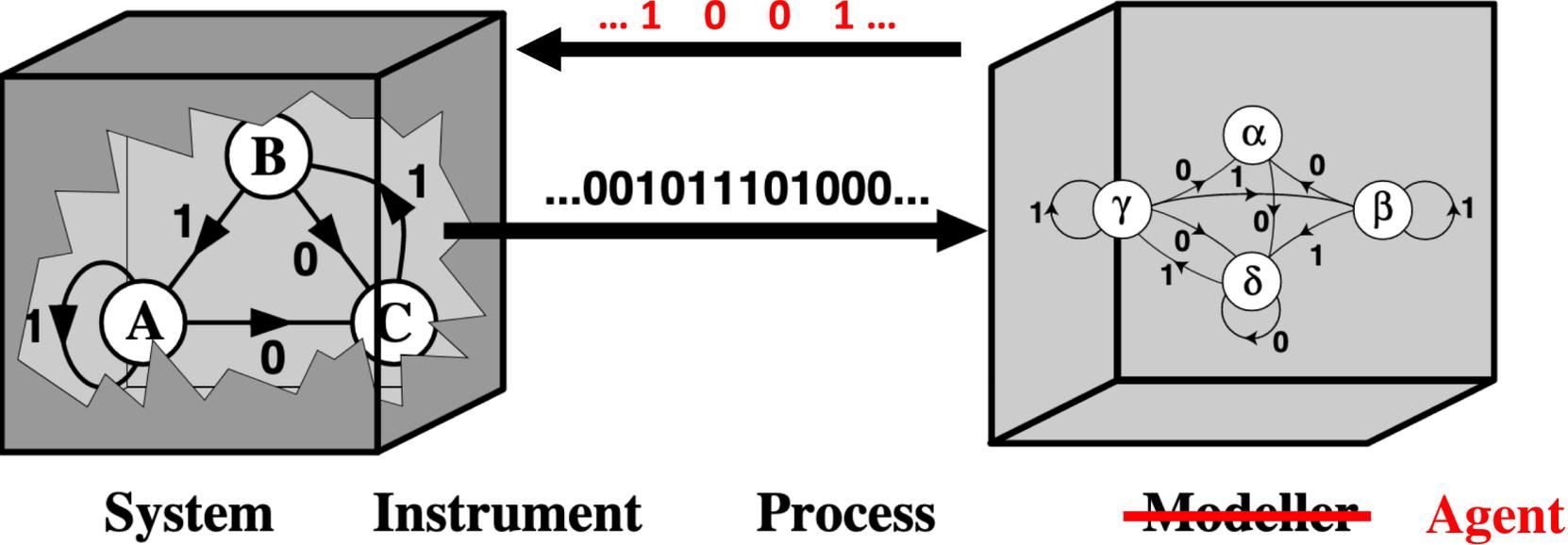
Modelling the RL problem



Modelling the RL problem



Modelling the RL problem



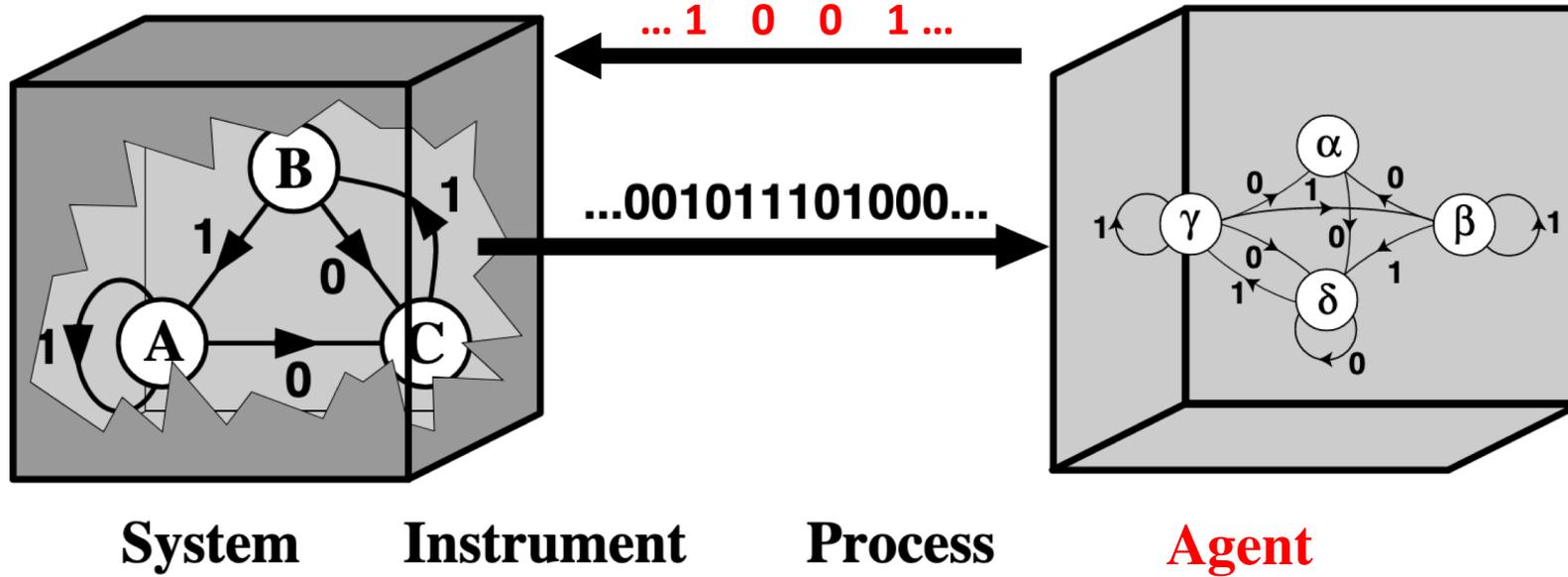
POCI Lecture 22

System History: ... 001011101000 ...

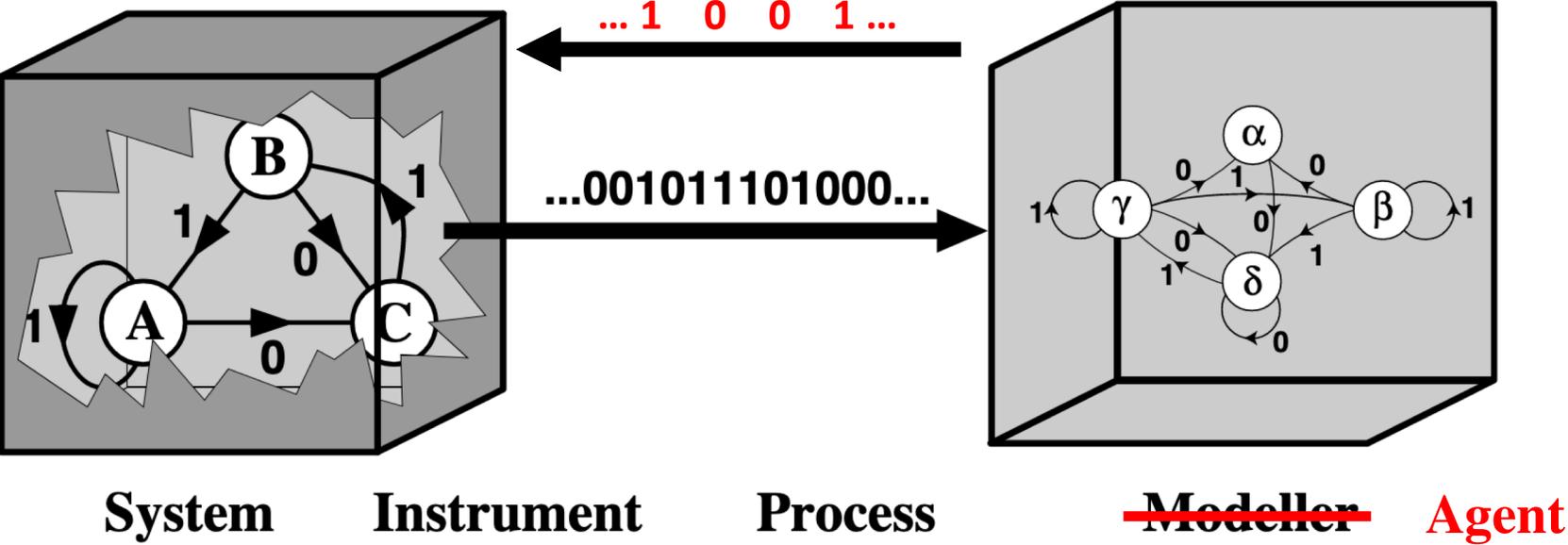


System + Agent History: ... 001**1**011**0**101000**01** ...

Modelling the RL problem



Modelling the RL problem



POCI Lecture 22

System History: ... 001011101000 ...



System + Agent History: ... 0011011010100001 ...

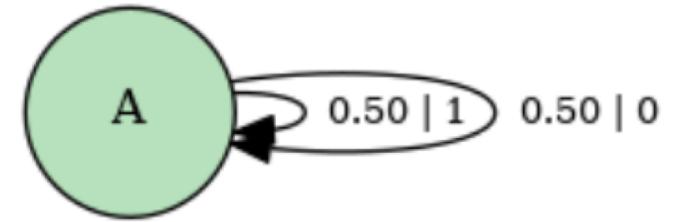
1. Initial random policy
2. Observe output -> Infer eM
3. Choose new policy

Simple Example: Rain or Shine

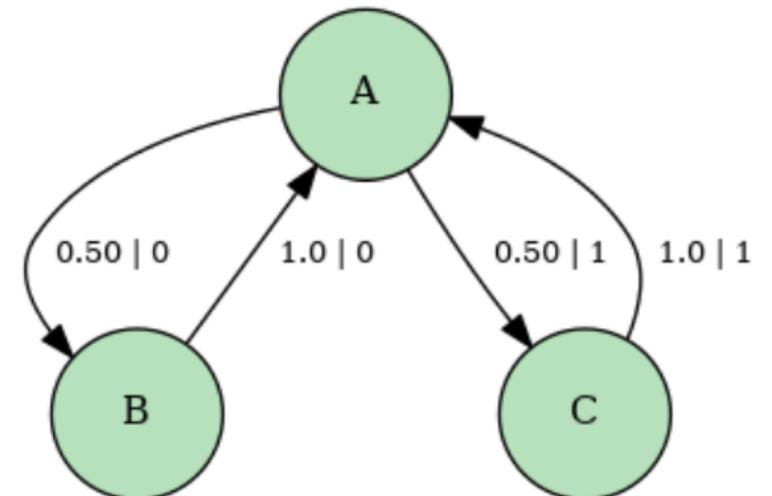
- Random weather
- Agent chooses 2nd bit e.g. 0**1** 1**1** 1**1** 0**0**
- Want sunglasses when sunny (11) and raincoat when rainy (00)

<u>Word</u>	<u>Result</u>
00	Raincoat when raining -> 😊
01	Sunglasses when raining -> ☹️
10	Raincoat when sunny -> ☹️
11	Sunglasses when sunny -> 😊

Random Policy



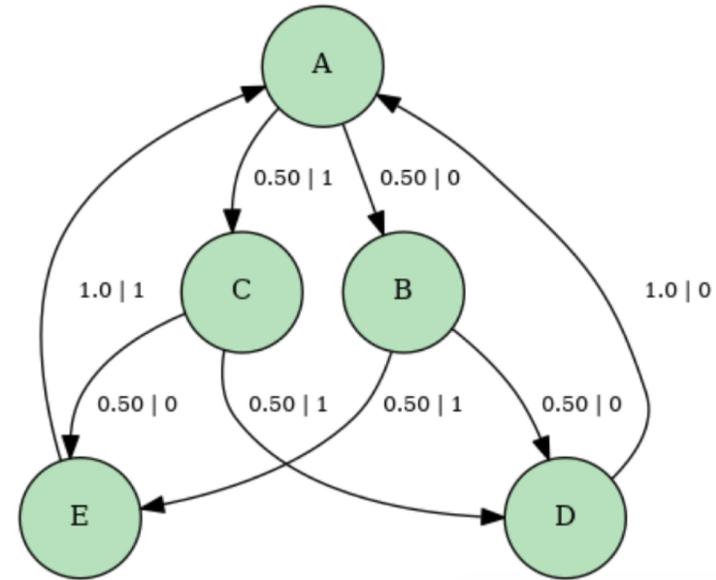
Optimal Policy



Controlling RRXOR

Agent chooses
2nd bit in RRXOR

e.g. 0**1**1

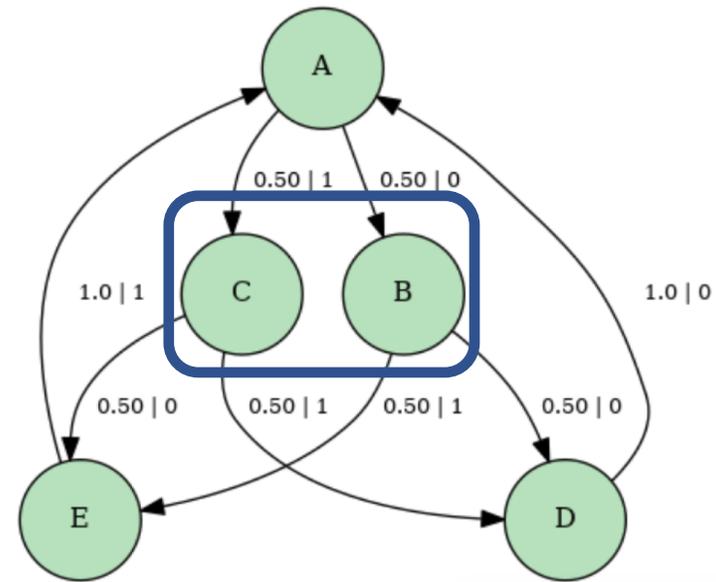


Random
Policy

Controlling RRXOR

Agent chooses
2nd bit in RRXOR

e.g. 0**1**1

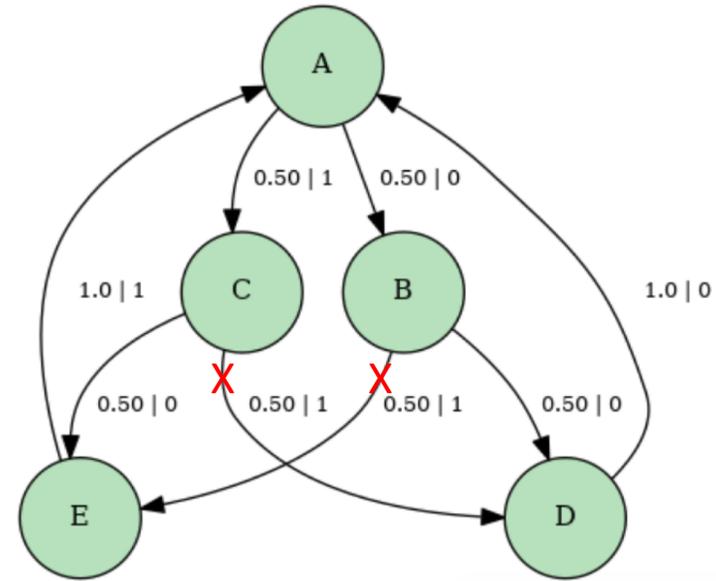


Random
Policy

Controlling RRXOR

Agent chooses
2nd bit in RRXOR

e.g. 0**1**1

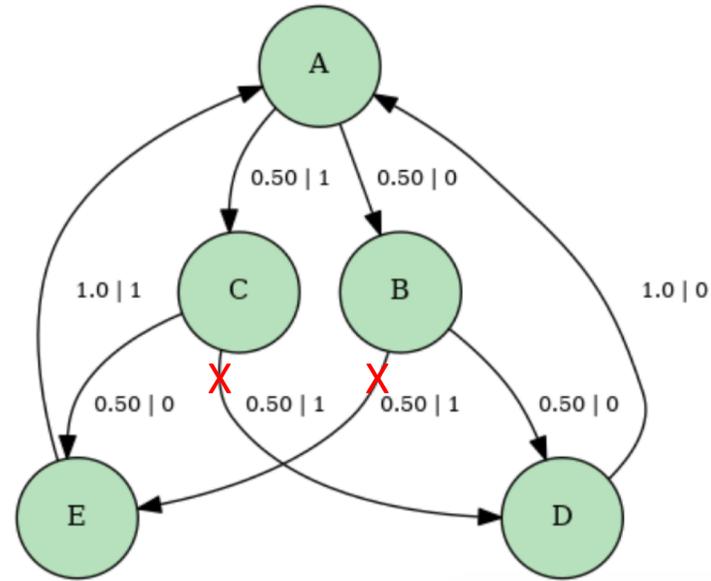


Controlling RRXOR

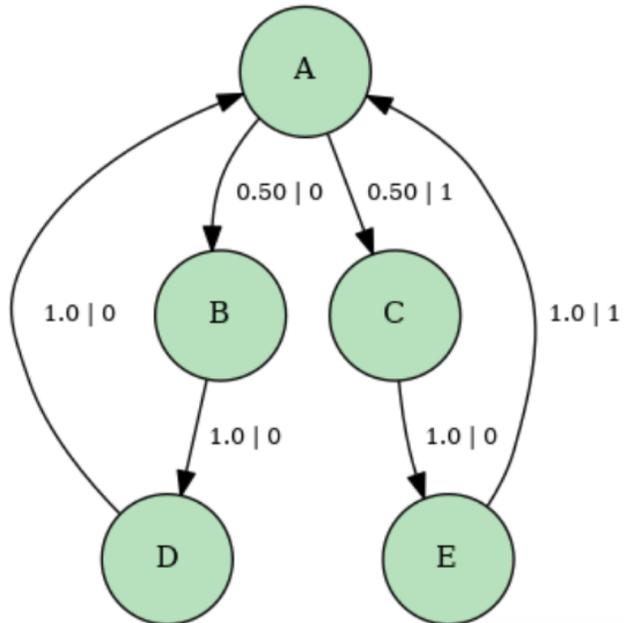
Agent chooses
2nd bit in RRXOR

e.g. 0**1**1

Random
Policy



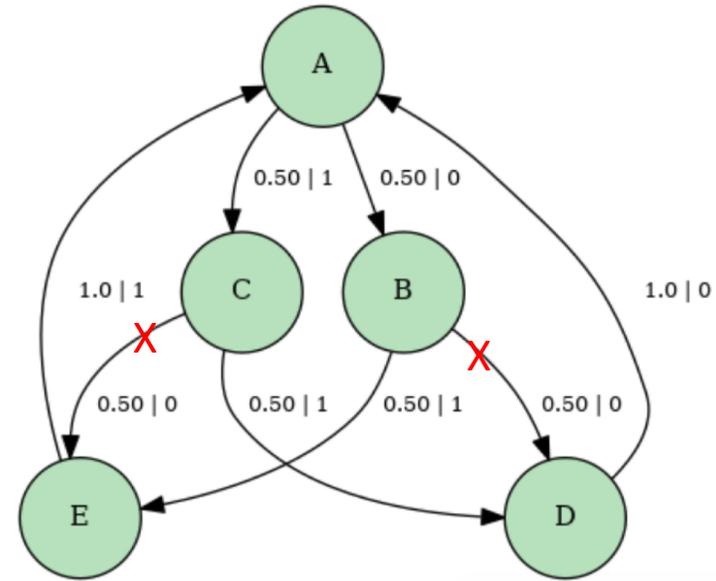
All-0's
Policy



Controlling RRXOR

Agent chooses
2nd bit in RRXOR

e.g. 0**1**1

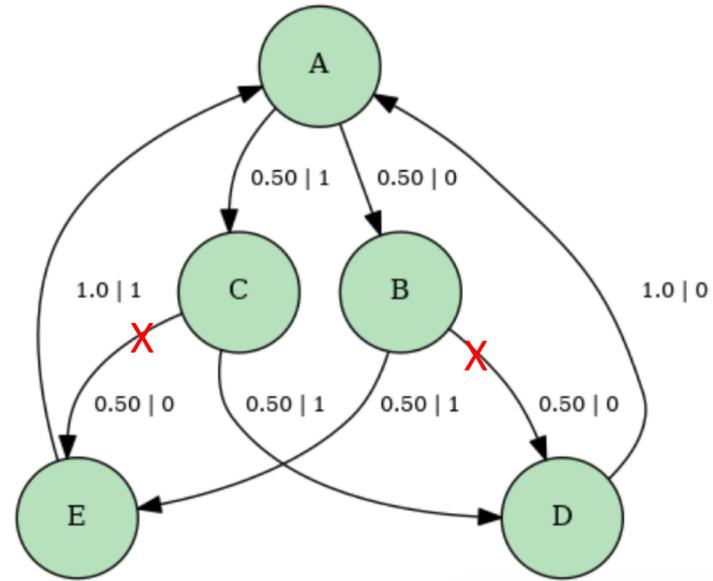


Random
Policy

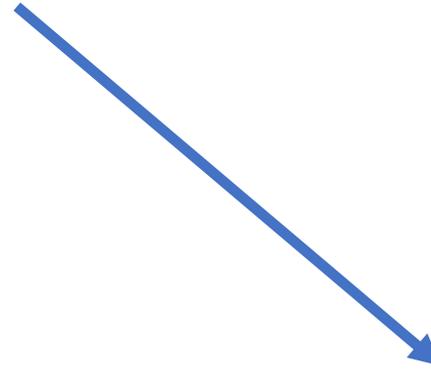
Controlling RRXOR

Agent chooses
2nd bit in RRXOR

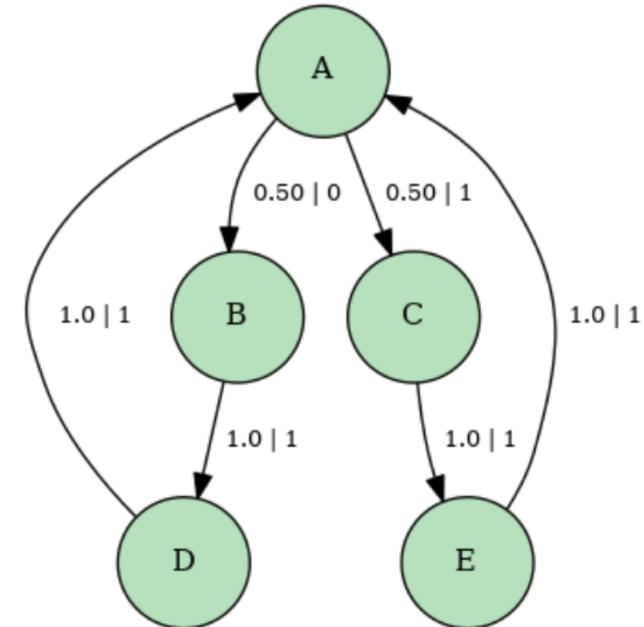
e.g. 0**1**1



Random
Policy



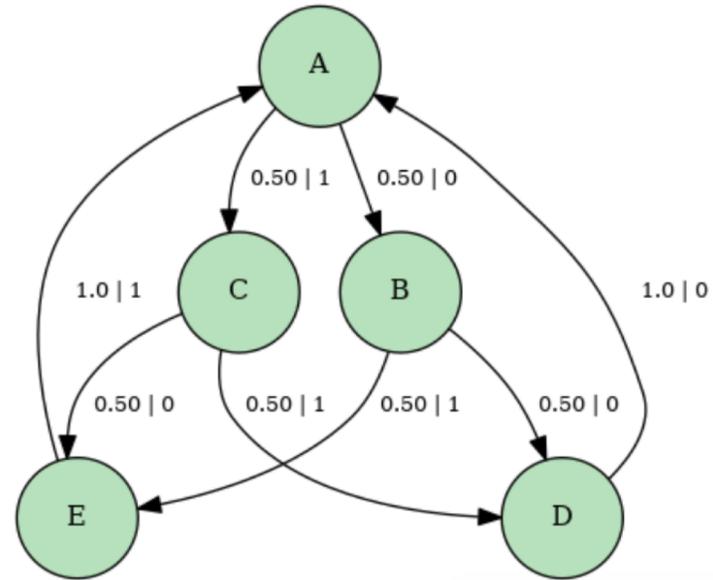
All-1's
Policy



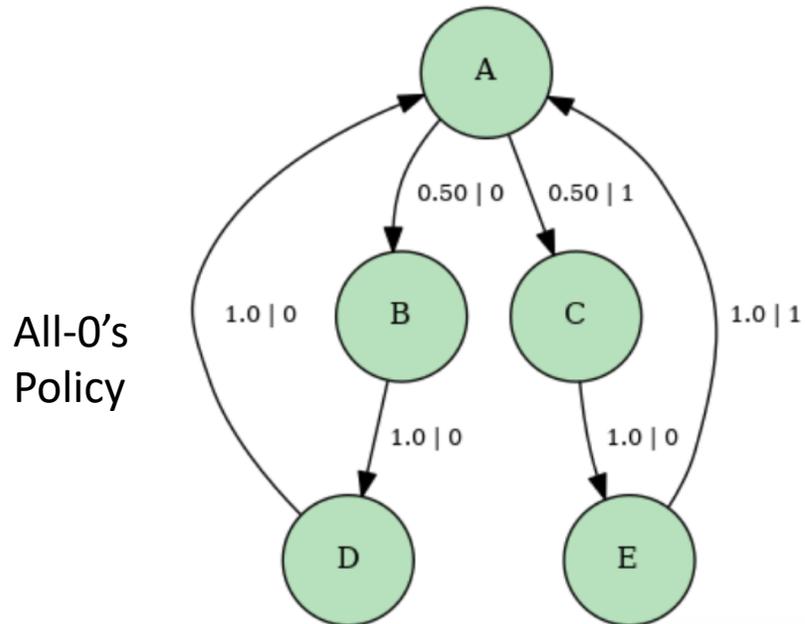
Controlling RRXOR

Agent chooses
2nd bit in RRXOR

e.g. 0**1**1

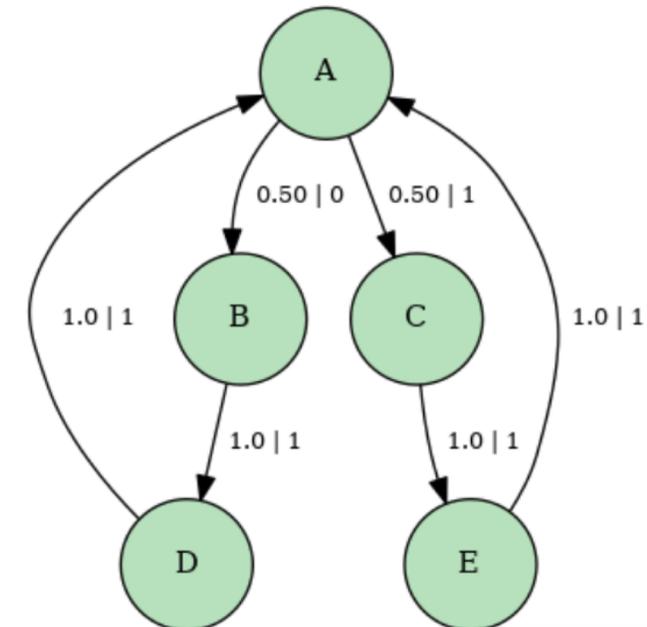


Random
Policy



All-0's
Policy

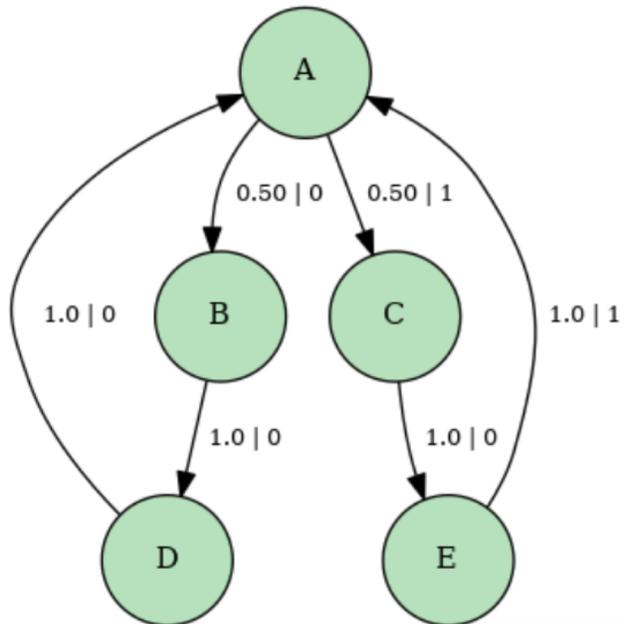
All-1's
Policy



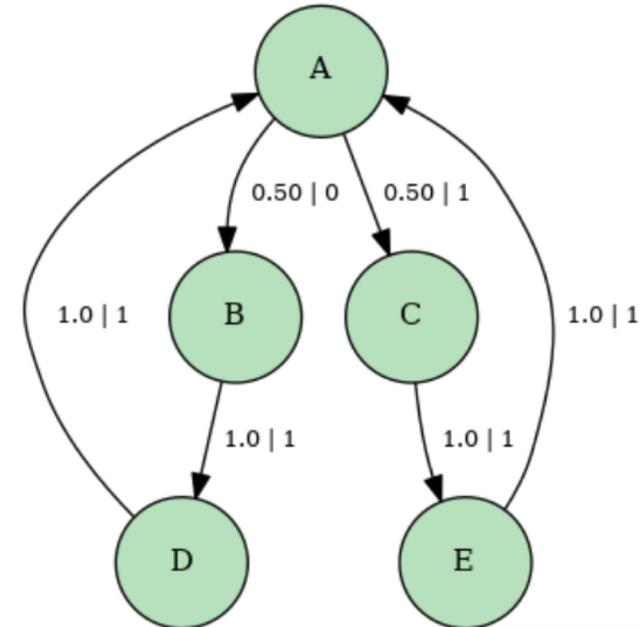
Picking a policy

- Choose output at controlled causal states
- Policy \rightarrow eM \rightarrow word distribution \rightarrow utility

All-0's
Policy

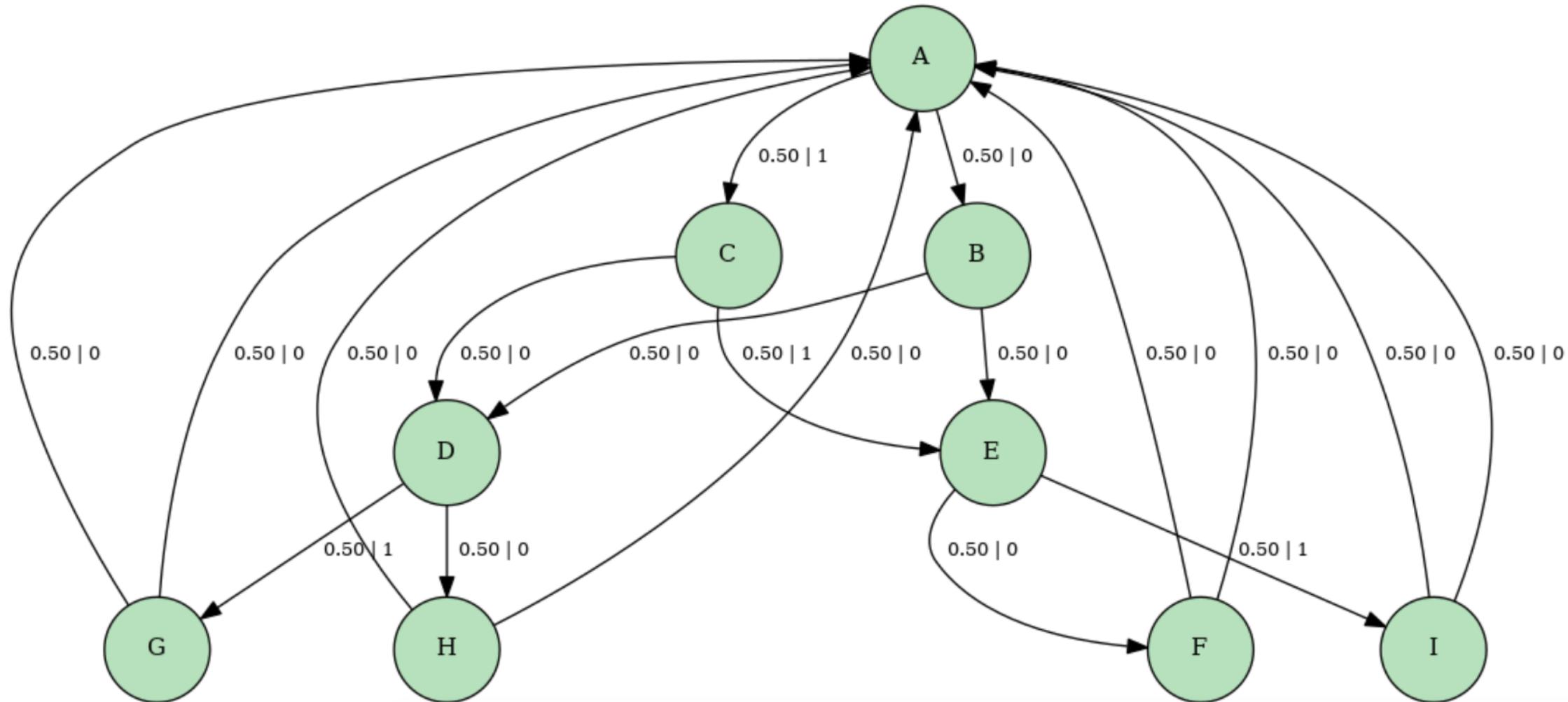


All-1's
Policy



Larger System

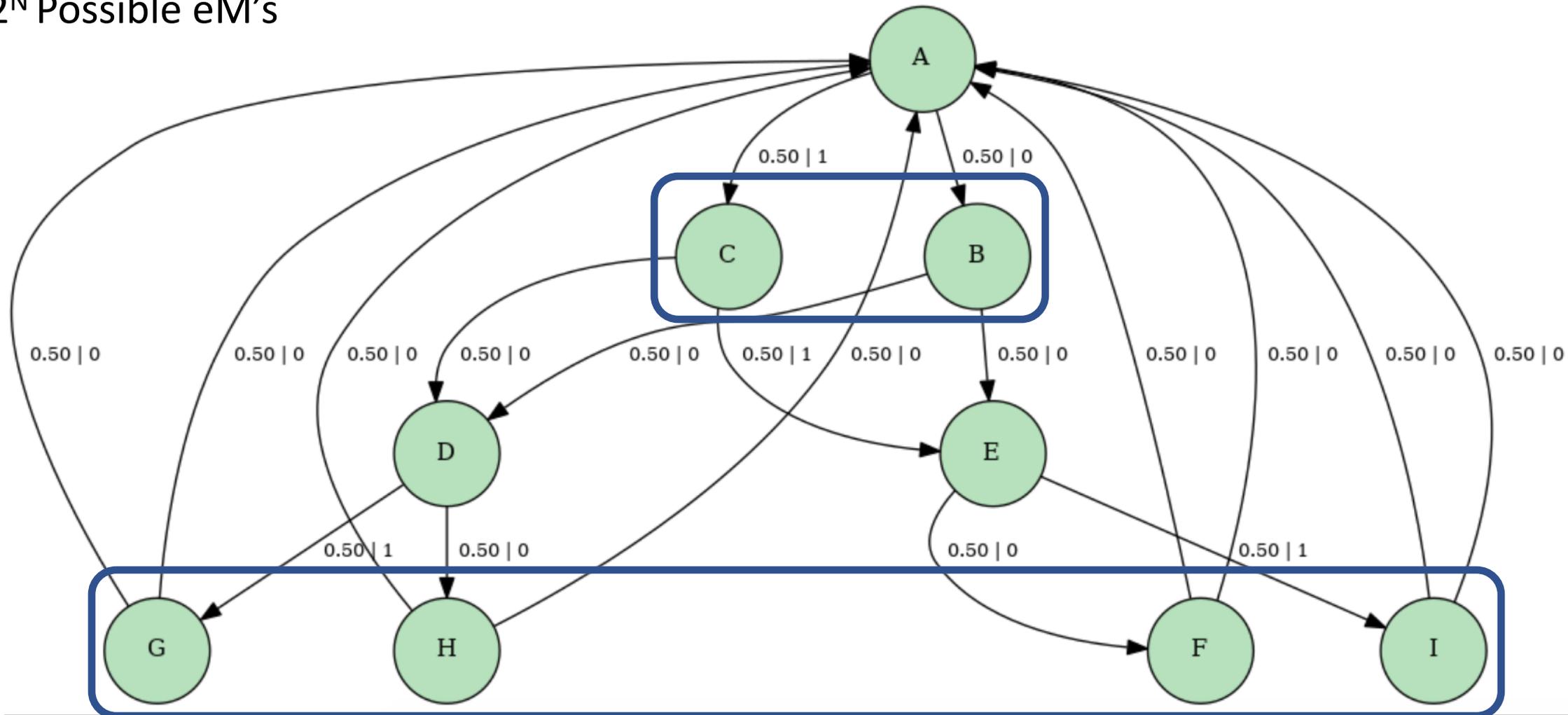
Random-**Controlled**-Random-**Controlled**



Larger System

Random-**Controlled**-Random-**Controlled**

Problem: 2^N Possible eM's

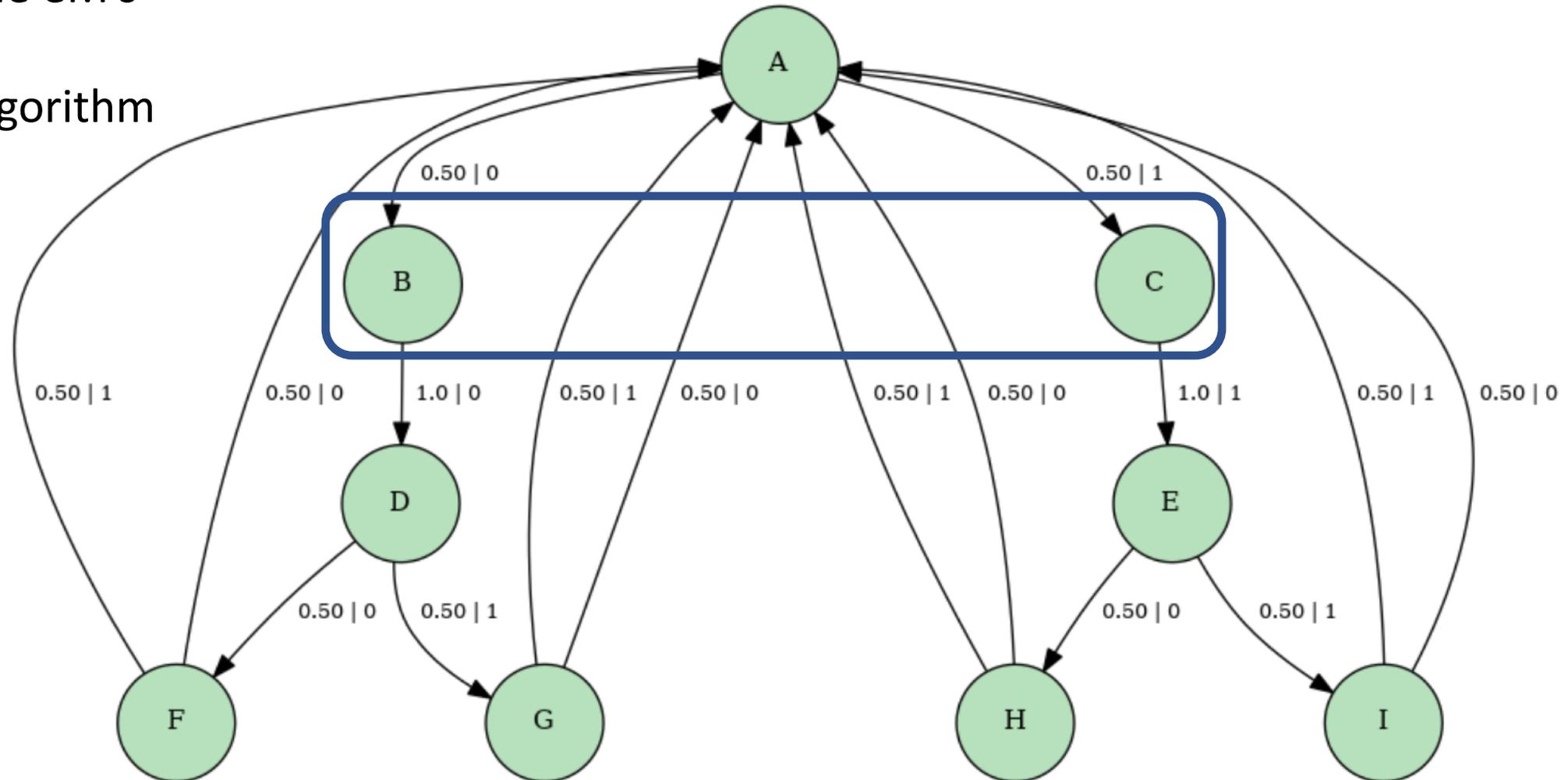


Larger System

Random-**Controlled**-Random-**Controlled**

Problem: 2^N Possible eM's

Solution: greedy algorithm

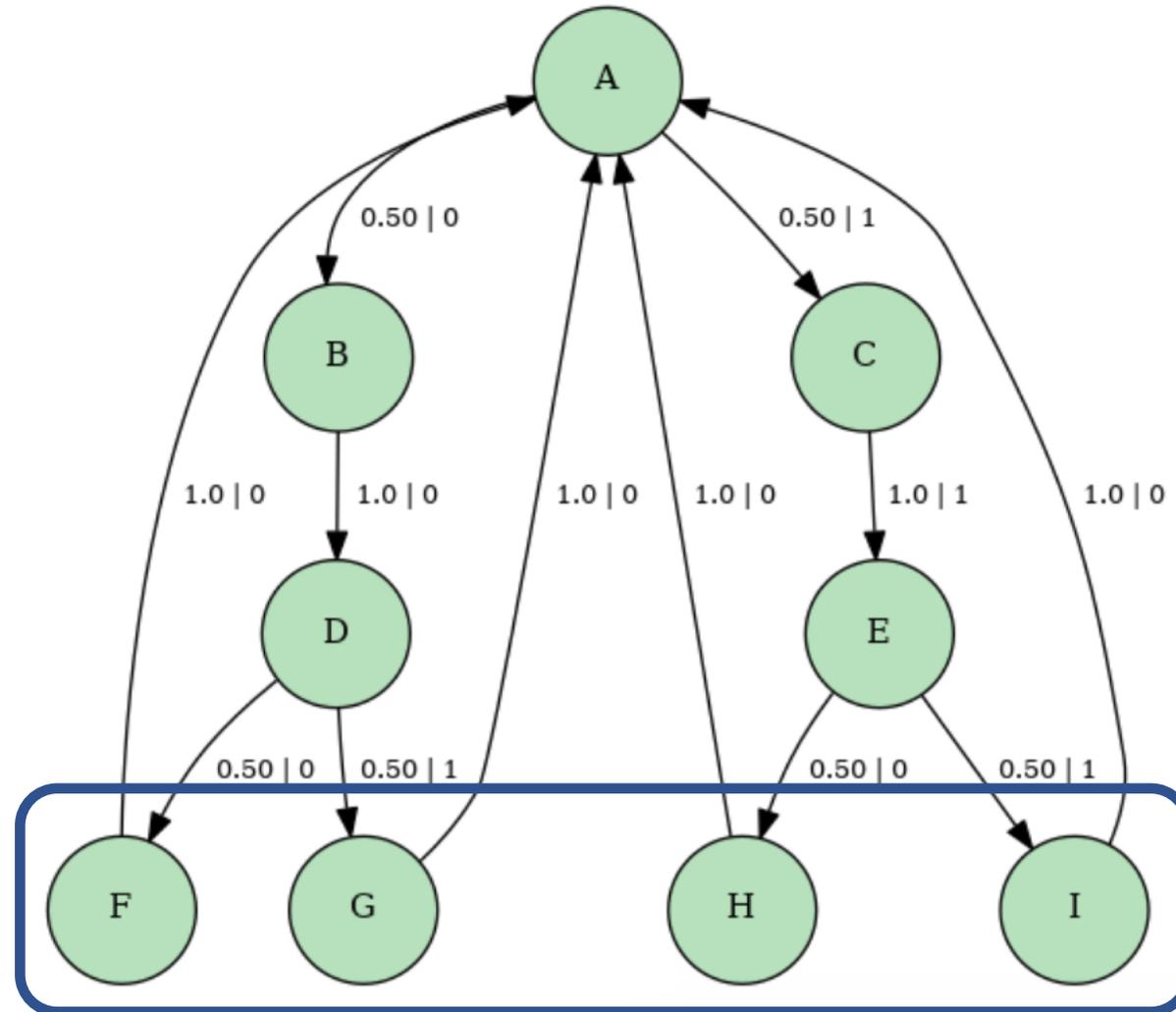


Larger System

Random-**Controlled**-Random-**Controlled**

Problem: 2^N Possible eM's

Solution: greedy algorithm



Future Directions

- Bigger eM's, full Bayesian approach, explore-exploit etc.
- Relationship between controller/system info measures -> Good Regulator Theorem?
- Interacting agents