

## Overview

- 1 Nondeterminism
- 2 Markov Decision Processes
- 3 Probabilities in MDPs
- 4 Policies
  - Finite-memory policies
- 5 Summary

## Overview

- 1 Nondeterminism
- 2 Markov Decision Processes
- 3 Probabilities in MDPs
- 4 Policies
  - Finite-memory policies
- 5 Summary

## Modeling and Verification of Probabilistic Systems

Joost-Pieter Katoen

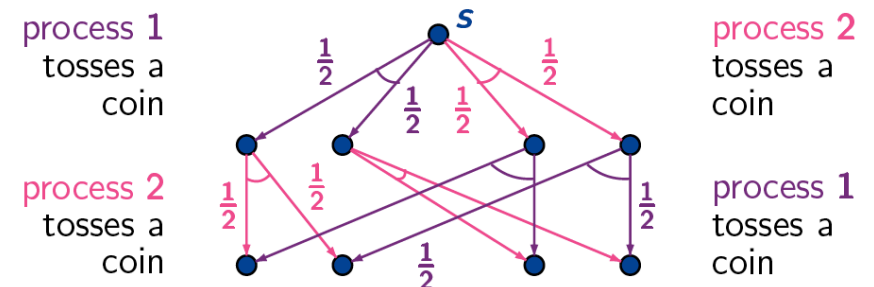
Lehrstuhl für Informatik 2  
Software Modeling and Verification Group

<http://moves.rwth-aachen.de/teaching/ws-1819/movep18/>

November 13, 2018

## Randomness and concurrency

Markov chains are not appropriate for modeling randomized **distributed** systems, since they cannot adequately model the interleaving behavior of the concurrent processes.



## Nondeterminism

### The use of nondeterminism

- ▶ **Concurrency** – scheduling of parallel components
  - ▶ in randomised distributed algorithms, several components run partly autonomously and interact asynchronously
- ▶ **Abstraction**
  - ▶ partition state space of a DTMC in similar (but not bisimilar) states
  - ▶ replace probabilistic branching by a nondeterministic choice
- ▶ **Unknown environments**
  - ▶ interaction with unknown environment
  - ▶ example: security in which the environment is an unknown adversary

### Beware

Nondeterminism is not the same as a uniform distribution!

## Overview

- 1 Nondeterminism
- 2 Markov Decision Processes
- 3 Probabilities in MDPs
- 4 Policies
  - Finite-memory policies
- 5 Summary

## Markov decision process (MDP)

### Markov decision processes

- ▶ In MDPs, **both** nondeterministic and probabilistic choices coexist.
- ▶ MDPs are transition systems in which in any state a nondeterministic choice between probability distributions exists.
- ▶ Once a probability distribution has been chosen nondeterministically, the next state is selected probabilistically—as in DTMCs.
- ▶ Any MC is thus an MDP in which in any state the probability distribution is uniquely determined.

Randomized distributed algorithms are typically appropriately modeled by MDPs, as probabilities affect just a small part of the algorithm and nondeterminism is used to model concurrency between processes by means of interleaving.

## Markov decision process (MDP)

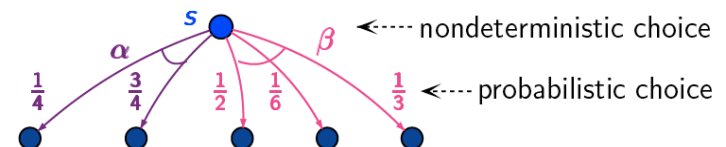
### Markov decision process

An MDP  $\mathcal{M}$  is a tuple  $(S, Act, \mathbf{P}, \iota_{\text{init}}, AP, L)$  where

- ▶  $S$  is a countable set of states with initial distribution  $\iota_{\text{init}} : S \rightarrow [0, 1]$
- ▶  $Act$  is a finite set of actions
- ▶  $\mathbf{P} : S \times Act \times S \rightarrow [0, 1]$ , transition probability function such that:

$$\text{for all } s \in S \text{ and } \alpha \in Act : \sum_{s' \in S} \mathbf{P}(s, \alpha, s') \in \{0, 1\}$$

- ▶  $AP$  is a set of atomic propositions and labeling  $L : S \rightarrow 2^{AP}$ .



## Markov decision process (MDP)

### Markov decision process

An MDP  $\mathcal{M}$  is a tuple  $(S, Act, \mathbf{P}, \iota_{init}, AP, L)$  where

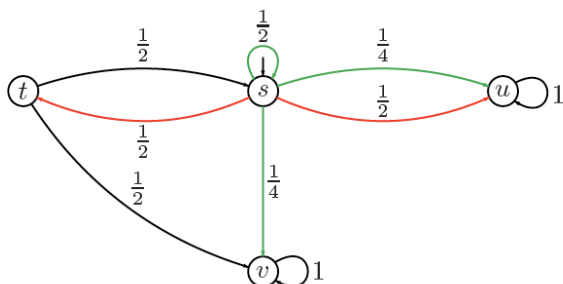
- ▶  $S, \iota_{init} : S \rightarrow [0, 1]$ ,  $AP$  and  $L$  are as before, i.e., as for DTMCs, and
- ▶  $Act$  is a finite set of actions
- ▶  $\mathbf{P} : S \times Act \times S \rightarrow [0, 1]$ , transition probability function such that:

$$\text{for all } s \in S \text{ and } \alpha \in Act : \sum_{s' \in S} \mathbf{P}(s, \alpha, s') \in \{0, 1\}$$

### Enabled actions

Let  $Act(s) = \{\alpha \in Act \mid \exists s' \in S. \mathbf{P}(s, \alpha, s') > 0\}$  be the set of enabled actions in state  $s$ . We require  $Act(s) \neq \emptyset$  for any state  $s$ .

## An example MDP



- ▶ Initial distribution:  $\iota_{init}(s) = 1$  and  $\iota_{init}(t) = \iota_{init}(u) = \iota_{init}(v) = 0$
- ▶ Set of enabled actions in state  $s$  is  $Act(s) = \{\alpha, \beta\}$  where
  - ▶  $\mathbf{P}(s, \alpha, s) = \frac{1}{2}$ ,  $\mathbf{P}(s, \alpha, t) = 0$  and  $\mathbf{P}(s, \alpha, u) = \mathbf{P}(s, \alpha, v) = \frac{1}{4}$
  - ▶  $\mathbf{P}(s, \beta, s) = \mathbf{P}(s, \beta, v) = 0$ , and  $\mathbf{P}(s, \beta, t) = \mathbf{P}(s, \beta, u) = \frac{1}{2}$
- ▶  $Act(t) = \{\alpha\}$  with  $\mathbf{P}(t, \alpha, s) = \mathbf{P}(t, \alpha, u) = \frac{1}{2}$  and 0 otherwise

## Markov decision process (MDP)

### Markov decision process

An MDP  $\mathcal{M}$  is a tuple  $(S, Act, \mathbf{P}, \iota_{init}, AP, L)$  where

- ▶  $S, \iota_{init} : S \rightarrow [0, 1]$ ,  $AP$  and  $L$  are as before, i.e., as for DTMCs, and
- ▶  $Act$  is a finite set of actions
- ▶  $\mathbf{P} : S \times Act \times S \rightarrow [0, 1]$ , transition probability function such that:

$$\text{for all } s \in S \text{ and } \alpha \in Act : \sum_{s' \in S} \mathbf{P}(s, \alpha, s') \in \{0, 1\}$$

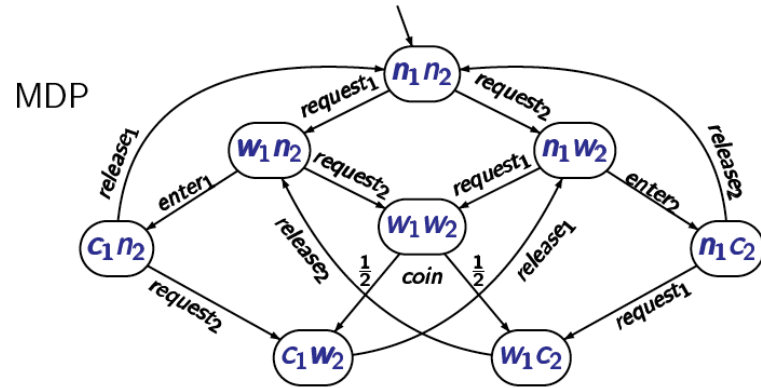
If  $|Act(s)| = 1$  for any state  $s$ , then the nondeterministic choice in any state is over a singleton set. In this case,  $\mathcal{M}$  is a DTMC. Vice versa, a DTMC is an MDP such that  $|Act(s)| = 1$  for all  $s$ .

## Example: randomized mutual exclusion

- 2 concurrent processes  $\mathcal{P}_1, \mathcal{P}_2$  with 3 phases:
  - $n_i$  noncritical actions of process  $\mathcal{P}_i$
  - $w_i$  waiting phase of process  $\mathcal{P}_i$
  - $c_i$  critical section of process  $\mathcal{P}_i$
- competition of both processes are waiting
- resolved by a randomized arbiter who tosses a coin

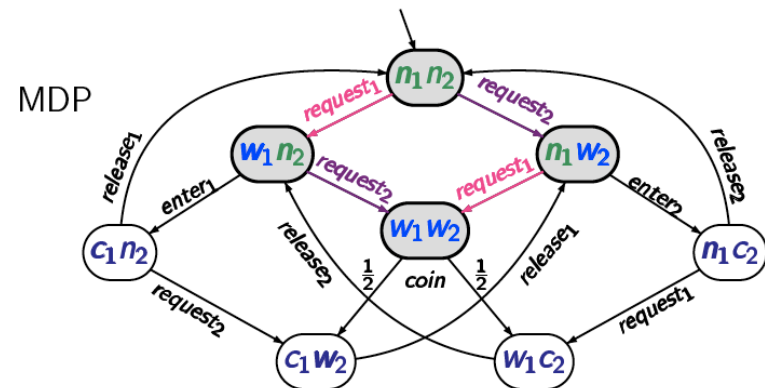
### Randomized mutual exclusion

- interleaving of the request operations
- competition if both processes are waiting
- randomized arbiter tosses a coin if both are waiting



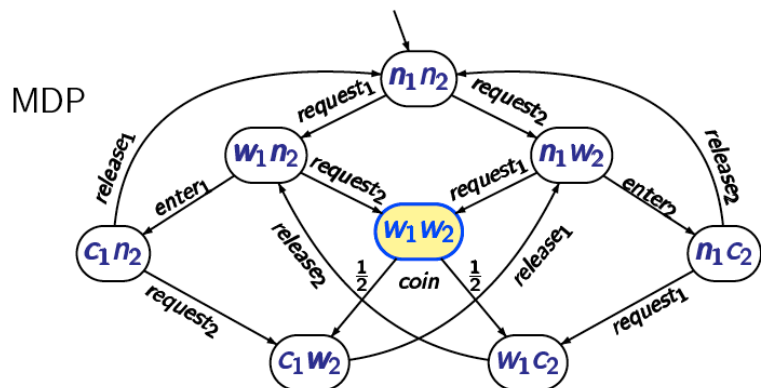
### Randomized mutual exclusion

- interleaving of the request operations
- competition if both processes are waiting
- randomized arbiter tosses a coin if both are waiting



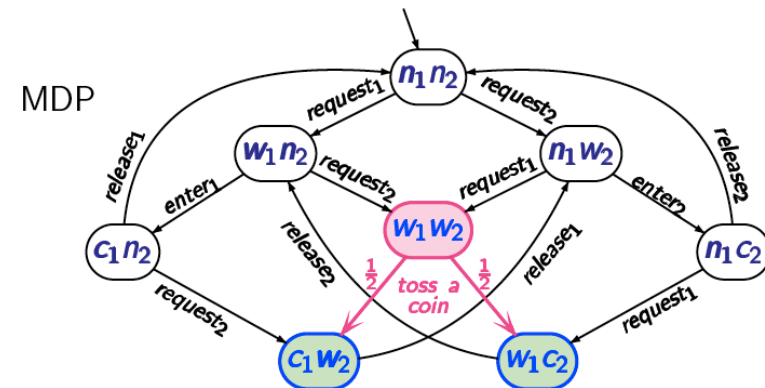
### Randomized mutual exclusion

- interleaving of the request operations
- competition if both processes are waiting
- randomized arbiter tosses a coin if both are waiting



### Randomized mutual exclusion

- interleaving of the request operations
- competition if both processes are waiting
- randomized arbiter tosses a coin if both are waiting



# Intuitive operational behavior

## Intuitive operational MDP behavior

1. A stochastic experiment according to  $\iota_{\text{init}}$  yields starting state  $s_0$  with probability  $\iota_{\text{init}}(s_0) > 0$ .
2. On entering state  $s$ , a non-deterministic choice among  $Act(s)$  determines the next action  $\alpha \in Act(s)$ , say.
3. The next state  $t$  is randomly chosen with probability  $\mathbf{P}(s, \alpha, t)$ .
4. If  $t$  is the unique  $\alpha$ -successor of  $s$ , then almost surely  $t$  is the successor after selecting  $\alpha$ , i.e.,  $\mathbf{P}(s, \alpha, t) = 1$ .
5. Continue with step 2.

# Paths in an MDP

## State graph

The *state graph* of MDP  $\mathcal{M}$  is a digraph  $G = (V, E)$  with  $V$  are the states of  $\mathcal{M}$ , and  $(s, s') \in E$  iff  $\mathbf{P}(s, \alpha, s') > 0$  for some  $\alpha \in Act$ .

## Paths

An infinite *path* in an MDP  $\mathcal{M} = (S, Act, \mathbf{P}, \iota_{\text{init}}, AP, L)$  is an infinite sequence  $s_0 \alpha_1 s_1 \alpha_2 s_2 \alpha_3 \dots \in (S \times Act)^\omega$ , written as

$$\pi = s_0 \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} s_2 \xrightarrow{\alpha_3} \dots,$$

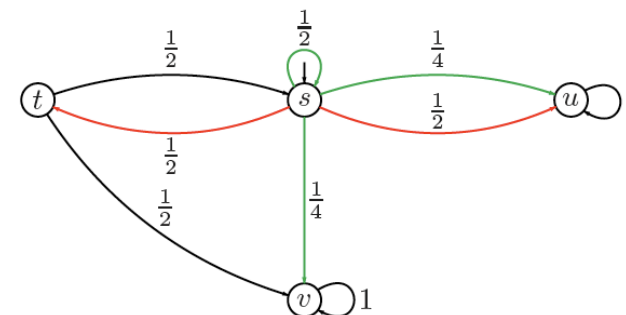
such that  $\mathbf{P}(s_i, \alpha_{i+1}, s_{i+1}) > 0$  for all  $i \geq 0$ . Any finite prefix of  $\pi$  that ends in a state is a *finite path*.

Let  $Paths(\mathcal{M})$  denote the set of paths in  $\mathcal{M}$ , and  $Paths^*(\mathcal{M})$  the set of finite prefixes thereof.

# Overview

- 1 Nondeterminism
- 2 Markov Decision Processes
- 3 Probabilities in MDPs
- 4 Policies
  - Finite-memory policies
- 5 Summary

# Paths in MDPs



$$s \xrightarrow{\alpha} s \xrightarrow{\alpha} s \xrightarrow{\beta} t \xrightarrow{\alpha} s \xrightarrow{\beta} u \dots$$

$$s \xrightarrow{\beta} t \xrightarrow{\alpha} s \xrightarrow{\beta} t \xrightarrow{\alpha} s \dots$$

# Probabilities in MDPs

- ▶ For DTMCs, a set of infinite paths is equipped with a  $\sigma$ -algebra and a probability measure that reflects the intuitive notion of probabilities for paths.
- ▶ Due to the presence of nondeterminism, MDPs are not augmented with a unique probability measure.
- ▶ Example: suppose we have two coins: a fair one, and a biased one, say  $\frac{1}{6}$  for heads and  $\frac{5}{6}$  for tails. We select nondeterministically one of the coins, and are interested in the probability of obtaining tails. This, however, is **not** specified! This also applies if we select one of the two coins repeatedly.
- ▶ Reasoning about probabilities of sets of paths of an MDP relies on the **resolution of nondeterminism**. This resolution is performed by a **policy**.<sup>1</sup> A policy chooses in any state  $s$  one of the actions  $\alpha \in Act(s)$ .

<sup>1</sup>Also called scheduler, strategy or adversary.

# Overview

- 1 Nondeterminism
- 2 Markov Decision Processes
- 3 Probabilities in MDPs
- 4 Policies
  - Finite-memory policies
- 5 Summary

# Policies

## Policy

Let  $\mathcal{M} = (S, Act, \mathbf{P}, \iota_{init}, AP, L)$  be an MDP. A *policy* for  $\mathcal{M}$  is a function  $\mathfrak{G} : S^+ \rightarrow Act$  such that  $\mathfrak{G}(s_0 s_1 \dots s_n) \in Act(s_n)$  for all  $s_0 s_1 \dots s_n \in S^+$ .

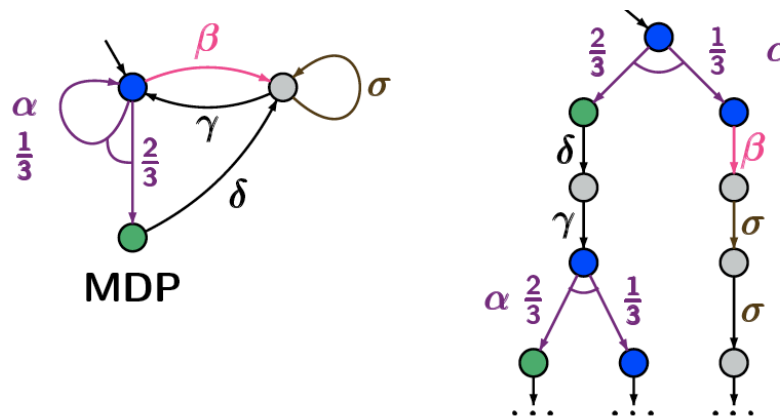
The path

$$\pi = s_0 \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} s_2 \xrightarrow{\alpha_3} \dots$$

is called a  $\mathfrak{G}$ -path if  $\alpha_i = \mathfrak{G}(s_0 \dots s_{i-1})$  for all  $i > 0$ .

For any scheduler, the actions are omitted from the *history*  $s_0 s_1 \dots s_n$ . This is not a restriction as for any sequence  $s_0 s_1 \dots s_n$  the relevant actions  $\alpha_i$  are given by  $\alpha_{i+1} = \mathfrak{G}(s_0 s_1 \dots s_i)$ . Hence, the scheduled action sequence can be constructed from prefixes of the path at hand.

# Induced Markov chain



Each policy induces an infinite DTMC. States are finite prefixes of paths in the MDP.

## Induced DTMC of an MDP by a policy

### DTMC of an MDP induced by a policy

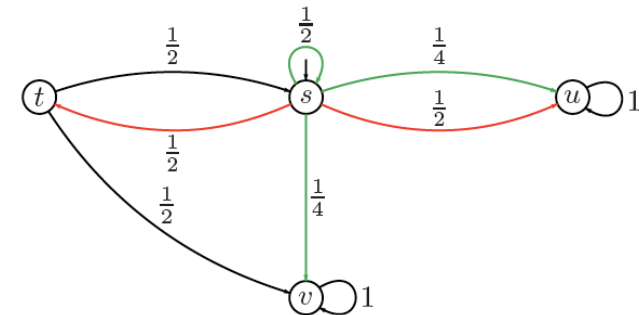
Let  $\mathcal{M} = (S, Act, \mathbf{P}, l_{init}, AP, L)$  be an MDP and  $\mathfrak{G}$  a policy on  $\mathcal{M}$ . The DTMC induced by  $\mathfrak{G}$ , denoted  $\mathcal{M}_{\mathfrak{G}}$ , is given by

$$\mathcal{M}_{\mathfrak{G}} = (S^+, \mathbf{P}_{\mathfrak{G}}, l_{init}, AP, L')$$

where for  $\sigma = s_0 s_1 \dots s_n$ :  $\mathbf{P}_{\mathfrak{G}}(\sigma, \sigma s_{n+1}) = \mathbf{P}(s_n, \mathfrak{G}(\sigma), s_{n+1})$  and  $L'(\sigma) = L(s_n)$ .

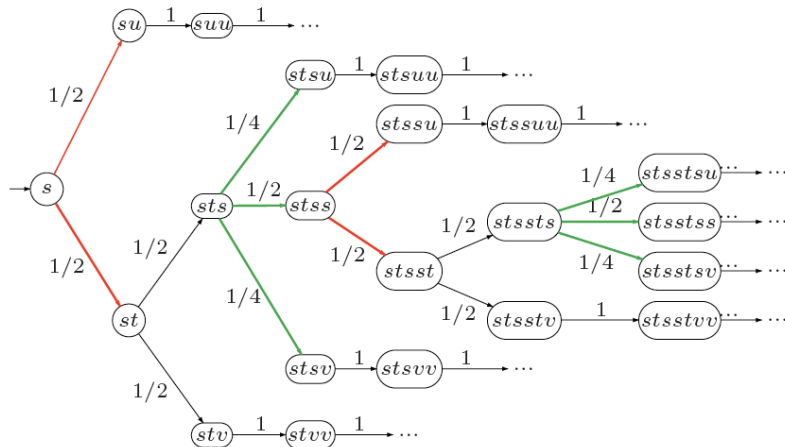
$\mathcal{M}_{\mathfrak{G}}$  is infinite, even if the MDP  $\mathcal{M}$  is finite. Intuitively, state  $s_0 s_1 \dots s_n$  of DTMC  $\mathcal{M}_{\mathfrak{G}}$  represents the configuration where the MDP  $\mathcal{M}$  is in state  $s_n$  and  $s_0 s_1 \dots s_{n-1}$  stands for the history. Since policy  $\mathfrak{G}$  might select different actions for finite paths that end in the same state  $s$ , a policy as defined above is also referred to as *history-dependent*.

## Example MDP



Consider a policy that alternates between selecting **red** and **green**, starting with **red**.

## Example induced DTMC



Induced DTMC for a policy that alternates between selecting **red** and **green**.

## MDP paths versus paths in the induced DTMC

There is a one-to-one correspondence between the  $\mathfrak{G}$ -paths of the MDP  $\mathcal{M}$  and the paths in the Markov chain  $\mathcal{M}_{\mathfrak{G}}$ .

For  $\mathfrak{G}$ -path  $\pi = s_0 \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} \dots$ , the corresponding path in DTMC  $\mathcal{M}_{\mathfrak{G}}$  is:

$$\pi^{\mathfrak{G}} = \hat{\pi}_0 \hat{\pi}_1 \hat{\pi}_2 \dots \quad \text{where } \hat{\pi}_n = s_0 s_1 \dots s_n.$$

Vice versa, for a path  $\hat{\pi}_0 \hat{\pi}_1 \hat{\pi}_2 \dots$  in the DTMC  $\mathcal{M}_{\mathfrak{G}}$ ,  $\hat{\pi}_0 = s_0$  for some state  $s_0$  such that  $l_{init}(s_0) > 0$  and, for each  $n > 0$ ,  $\hat{\pi}_n = \hat{\pi}_{n-1} s_n$  for some state  $s_n$  in the MDP  $\mathcal{M}$  such that  $\mathbf{P}(s_{n-1}, \mathfrak{G}(\hat{\pi}_{n-1}), s_n) > 0$ . Hence:

$$s_0 \xrightarrow{\mathfrak{G}(\hat{\pi}_0)} s_1 \xrightarrow{\mathfrak{G}(\hat{\pi}_1)} s_2 \xrightarrow{\mathfrak{G}(\hat{\pi}_2)} \dots$$

is a  $\mathfrak{G}$ -path in  $\mathcal{M}$ .

## Probability measure on MDP

### Probability measure on MDP

Let  $Pr_{\mathfrak{G}}^{\mathcal{M}}$ , or simply  $Pr^{\mathfrak{G}}$ , denote the probability measure  $Pr^{\mathcal{M}_{\mathfrak{G}}}$  associated with the DTMC  $\mathcal{M}_{\mathfrak{G}}$ .

This measure is the basis for associating probabilities with events in the MDP  $\mathcal{M}$ . Let, e.g.,  $P \subseteq (2^{AP})^{\omega}$  be an  $\omega$ -regular property. Then  $Pr^{\mathfrak{G}}(P)$  is defined as:

$$Pr^{\mathfrak{G}}(P) = Pr^{\mathcal{M}_{\mathfrak{G}}}(P) = Pr_{\mathcal{M}_{\mathfrak{G}}}\{\pi \in Paths(\mathcal{M}_{\mathfrak{G}}) \mid trace(\pi) \in P\}.$$

Similarly, for fixed state  $s$  of  $\mathcal{M}$ , which is considered as the unique starting state,

$$Pr^{\mathfrak{G}}(s \models P) = Pr_s^{\mathcal{M}_{\mathfrak{G}}}\{\pi \in Paths(s) \mid trace(\pi) \in P\}$$

where we identify the paths in  $\mathcal{M}_{\mathfrak{G}}$  with the corresponding  $\mathfrak{G}$ -paths in  $\mathcal{M}$ .

## Finite-memory policies

- ▶ *Finite-memory policies* (shortly: fm-policies) are a generalisation of positional policies.
- ▶ The behavior of an fm-policy is described by a deterministic finite automaton (DFA).
- ▶ The selection of the action to be performed in the MDP  $\mathcal{M}$  depends on the current state of  $\mathcal{M}$  (as before) and the current state (called *mode*) of the policy, i.e., the DFA.

## Positional policy

### Positional policy

Let  $\mathcal{M}$  be an MDP with state space  $S$ . Policy  $\mathfrak{G}$  on  $\mathcal{M}$  is *positional* (or: *memoryless*) iff for each sequence  $s_0 s_1 \dots s_n$  and  $t_0 t_1 \dots t_m \in S^+$  with  $s_n = t_m$ :

$$\mathfrak{G}(s_0 s_1 \dots s_n) = \mathfrak{G}(t_0 t_1 \dots t_m).$$

In this case,  $\mathfrak{G}$  can be viewed as a function  $\mathfrak{G} : S \rightarrow Act$ .

Policy  $\mathfrak{G}$  is positional if it always selects the same action in a given state. This choice is independent of what has happened in the history, i.e., which path led to the current state.

## Finite-memory policy

### Finite-memory policy

Let  $\mathcal{M}$  be an MDP with state space  $S$  and action set  $Act$ .

A *finite-memory policy*  $\mathfrak{G}$  for  $\mathcal{M}$  is a tuple  $\mathfrak{G} = (Q, act, \Delta, start)$  with:

- ▶  $Q$  is a finite set of **modes**,
- ▶  $\Delta : Q \times S \rightarrow Q$  is the **transition function**,
- ▶  $act : Q \times S \rightarrow Act$  is a function that selects an action  $act(q, s) \in Act(s)$  for any mode  $q \in Q$  and state  $s \in S$  of  $\mathcal{M}$ ,
- ▶  $start : S \rightarrow Q$  is a function that selects a **starting mode** for state  $s \in S$ .



## An MDP under a finite-memory policy

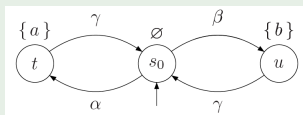
The behavior of an MDP  $\mathcal{M}$  under fm-policy  $\mathfrak{G} = (Q, act, \Delta, start)$  is:

- ▶ Initially, a starting state  $s_0$  is randomly determined according to the initial distribution  $\iota_{init}$ , i.e.,  $\iota_{init}(s_0) > 0$ .
- ▶ The fm-policy  $\mathfrak{G}$  initializes its DFA to the mode  $q_0 = start(s_0) \in Q$ .
- ▶ If  $\mathcal{M}$  is in state  $s$  and the current mode of  $\mathfrak{G}$  is  $q$ , then the decision of  $\mathfrak{G}$ , i.e., the selected action, is  $\alpha = act(q, s) \in Act(s)$ .
- ▶ The policy changes to mode  $\Delta(q, s)$ , while  $\mathcal{M}$  performs the selected action  $\alpha$  and randomly moves to the next state according to the distribution  $\mathbf{P}(s, \alpha, \cdot)$ .

## Positional versus fm-policies

### Positional policies are insufficient for $\omega$ -regular properties

Consider the MDP:



Positional policy  $\mathfrak{G}_\alpha$  always chooses  $\alpha$  in state  $s_0$

Positional policy  $\mathfrak{G}_\beta$  always chooses  $\beta$  in state  $s_0$ . Then:

$$Pr_{\mathfrak{G}_\alpha}(s_0 \models \diamond a \wedge \diamond b) = Pr_{\mathfrak{G}_\beta}(s_0 \models \diamond a \wedge \diamond b) = 0.$$

Now consider fm-policy  $\mathfrak{G}_{\alpha\beta}$  which alternates between selecting  $\alpha$  and  $\beta$ .

Then:  $Pr_{\mathfrak{G}_{\alpha\beta}}(s_0 \models \diamond a \wedge \diamond b) = 1$ .

Thus, the class of positional policies is insufficiently powerful to characterise minimal (or maximal) probabilities for  $\omega$ -regular properties.

## Finite-memory policies

### Relation fm-policy to definition policy

An fm-policy  $\mathfrak{G} = (Q, act, \Delta, start)$  is identified with policy,  $\mathfrak{G}' : Paths^* \rightarrow Act$  which is defined as follows.

1. For the starting state  $s_0$ , let  $\mathfrak{G}'(s_0) = act(start(s_0), s_0)$ .
2. For path fragment  $\hat{\pi} = s_0 s_1 \dots s_n$  let

$$\mathfrak{G}'(\hat{\pi}) = act(q_n, s_n)$$

where  $q_0 = start(s_0)$  and  $q_{i+1} = \Delta(q_i, s_i)$  for  $0 \leq i \leq n$ .

Positional policies can be considered as fm-policies with just a single mode.

## Other kinds of policies

- ▶ **Counting** policies that base their decision on the number of visits to a state, or the length of the history (i.e., number of visits to all states)
- ▶ **Partial-observation** policies that base their decision on the trace  $L(s_0) \dots L(s_n)$  of the history  $s_0 \dots s_n$ .
- ▶ **Randomised** policies. This is applicable to all (deterministic) policies. For instance, a randomised positional policy  $\mathfrak{G} : S \rightarrow Dist(Act)$ , where  $Dist(X)$  is the set of probability distributions on  $X$ , such that  $\mathfrak{G}(s)(\alpha) > 0$  iff  $\alpha \in Act(s)$ . Similar can be done for fm-policies and history-dependent policies etc..
- ▶ There is a **strict hierarchy** of policies, showing their expressiveness (black board).

# Overview

- 1 Nondeterminism
- 2 Markov Decision Processes
- 3 Probabilities in MDPs
- 4 Policies
  - Finite-memory policies
- 5 **Summary**

# Summary

## Important points

1. An MDP is a model exhibiting non-determinism and probabilities.
2. Non-determinism is important for e.g., randomized distributed algorithms.
3. Policies are functions that select enabled actions in states.
4. A policy on an MDP induces an infinite DTMC, even if the MDP is finite.
5. Probability measures on MDP paths are defined using induced DTMC paths.
6. A positional policy selects in a state always the same action.