# Formal Approaches to Decision-Making under Uncertainty

## Lecture 2-2: Markov Decision Processes
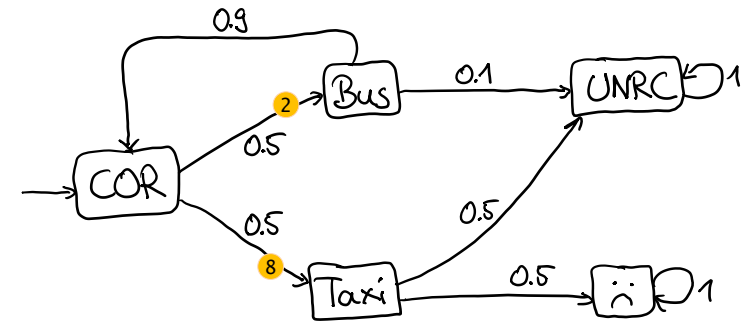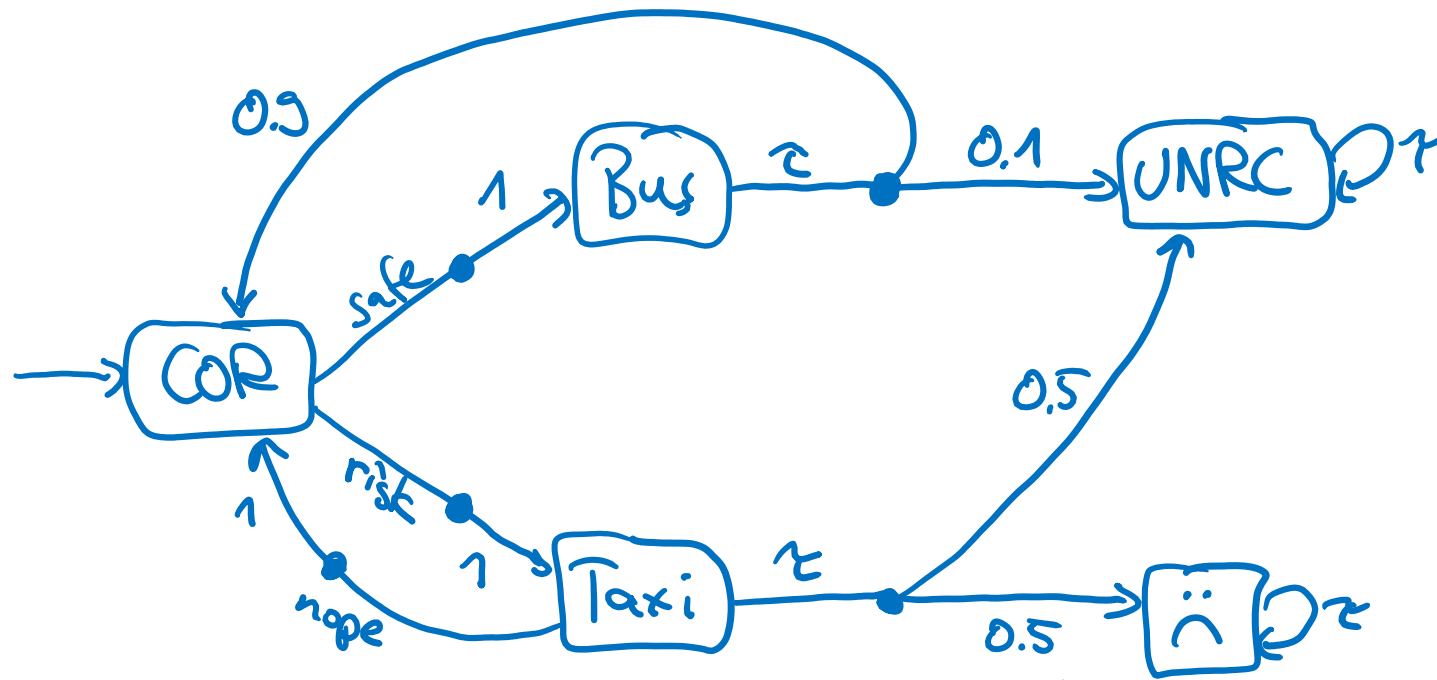
Arnd Hartmanns

Formal Methods and Tools

**UNIVERSITY OF TWENTE**

# Markov Decision Processes

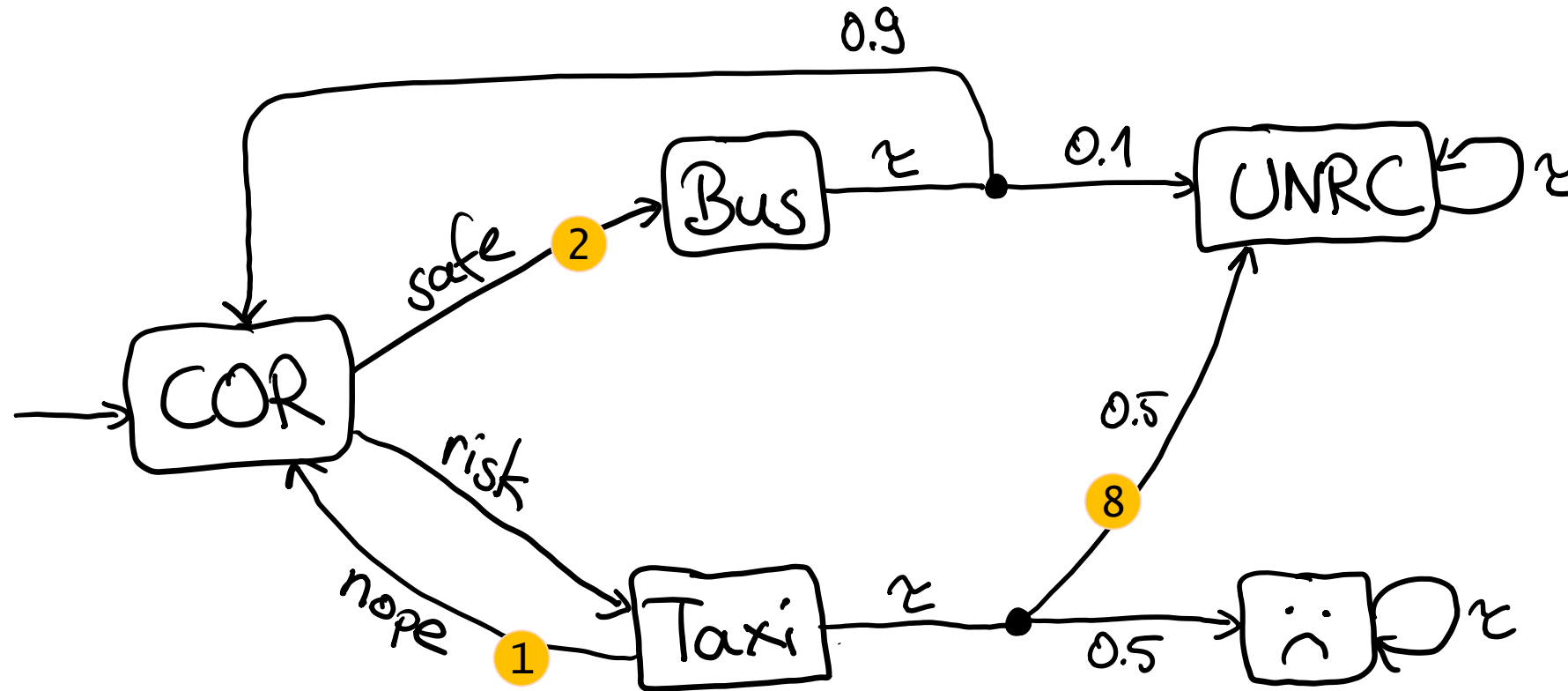A Markov decision process **MDP**
from Córdoba airport to Rio Cuarto:

# Markov Decision Processes

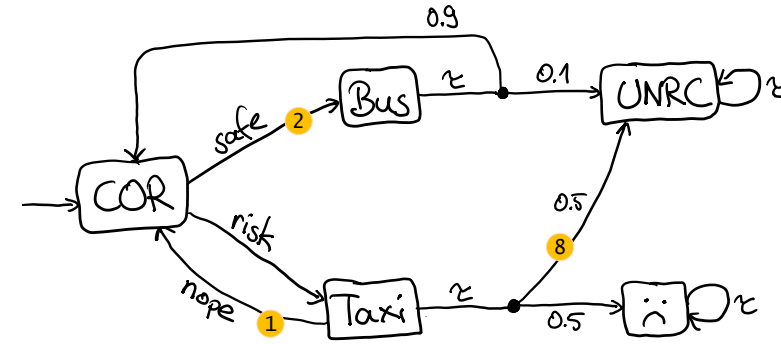A Markov decision process **MDP**
from Córdoba airport to Rio Cuarto:

DTMC: $P(\Diamond\{UNRC\}) = ?$

here: $P_{max}(\Diamond\{UNRC\}) = ?$

An **MDP** is a 4-tuple $M = \langle S, A, T, s_I \rangle$ where

$S$ are the states with initial state $s_I$,

$A$ is the finite set of action names,

$T: S \to 2^{A \times \mathrm{Dist}(S)}$ is the

   nondeterministic transition function.

DTMC: $T: S \to \mathrm{Dist}(s)$

Shorthand for transitions: write $s \xrightarrow{a} \mu$ if $\langle a, \mu \rangle \in T(s)$.

state   transition   next state   ...

Paths: $s_0\, a_0 \mu_0\, s_1\, a_1 \mu_1\, \ldots$

$T(\text{Taxi}) = \{\ \langle \text{nope}, \{\text{COR} \mapsto 1\}\rangle,$
$\langle \tau, \{\text{UNRC} \mapsto 0.5, \ddot{\frown} \mapsto 0.5\}\rangle\ \}$

A **reward structure** maps *branches* of transitions to reward values:

$$R: S \times \big(A \times \mathrm{Dist}(S)\big) \times S \to \mathbb{R}$$

DTMC: $R: S \times S \to \mathbb{R}$

## Typical restrictions:

– no deadlocks: $\forall s \in S: |T(s)| \geq 1$

– all branch probabilities in $\mathbb{Q}$

– action determinism (in ML and planning):

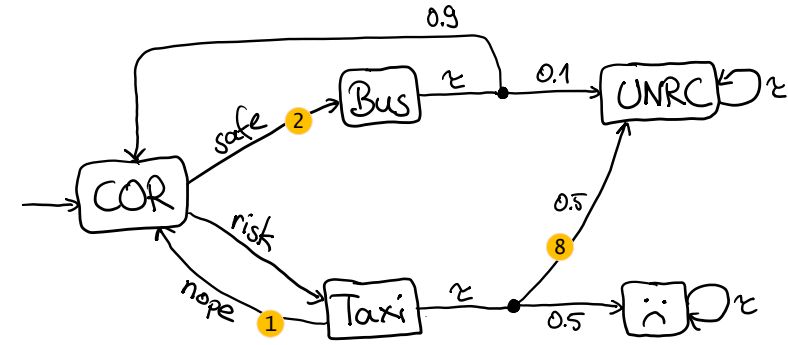$$s \xrightarrow{a} \mu_1 \land s \xrightarrow{a} \mu_2 \Rightarrow \mu_1 = \mu_2$$

– discrete-time Markov chain (DTMC):

no nondeterminism $\rightarrow \forall s \in S: |T(s)| = 1$

– rewards in $[0, \infty)$ or $\mathbb{N}$

– state rewards:
    same reward on every incoming branch to a state

strategies, policies, adversaries

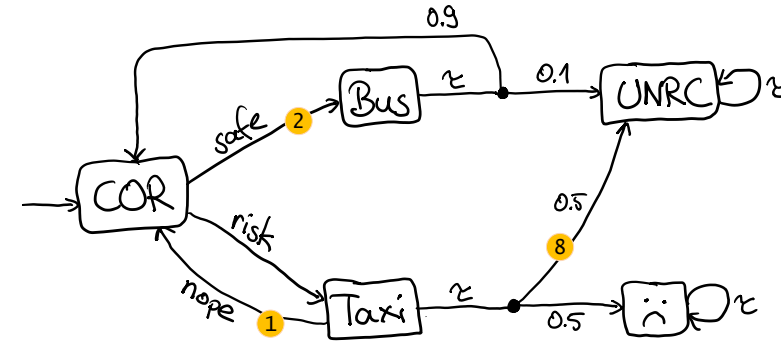**Schedulers** resolve all nondeterminism; e.g.

$$\mathcal{S}: S \to A \times \text{Dist}(S)$$
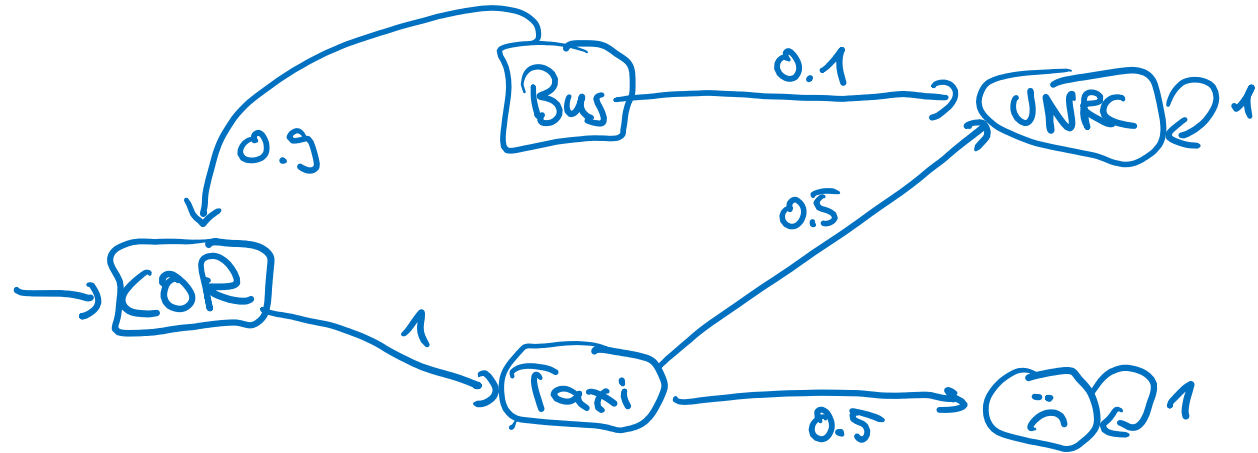
= memoryless deterministic scheduler

Example:

$S(COR) = \langle risk, \{Taxi \mapsto 1\} \rangle$

$S(Bus) = \tau$    $S(Taxi) = \tau$

simplify: risk ] definition of one of many possible schedulers

Induced DTMC $M|_{\mathcal{S}}$:

# Unbounded probabilistic reachability:

$$P_{opt}(\lozenge\, G)$$

for $opt \in \{\max, \min\}$ and $G \subseteq S$:

induced DTMC of $\mathcal{S}$

$$P_{\max}(\lozenge\, G) = \sup_{\mathcal{S}} \mathbb{P}_{M|_{\mathcal{S}}}(\{s_0 \ldots \in \text{Paths}(M) \mid \exists i \colon s_i \in G\})$$

$$P_{\min}(\lozenge\, G) = \inf_{\mathcal{S}} \mathbb{P}_{M|_{\mathcal{S}}}(\{s_0 \ldots \in \text{Paths}(M) \mid \exists i \colon s_i \in G\})$$

Examples:

$$P_{\max}(\lozenge\, \{\ddot{n}\}) = 0.5$$
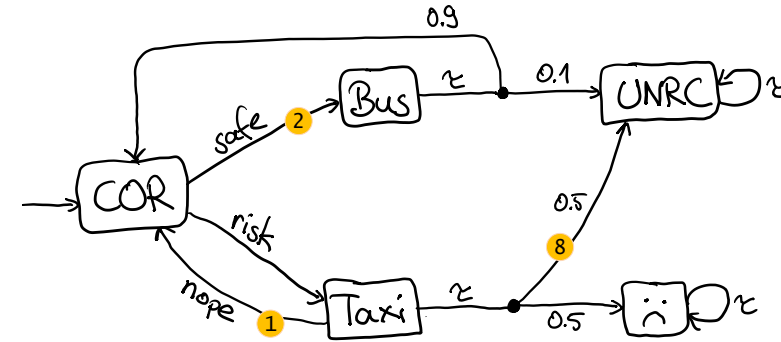
$$P_{\min}(\lozenge\, \{\ddot{n}\}) = 0$$

# Step-bounded probabilistic reachability:

$$P_{opt}(\Diamond^{\leq b} G)$$

for $opt \in \{\max, \min\}$ and $G \subseteq S$:

$$P_{\max}(\Diamond^{\leq b} G) = \sup_{\mathcal{S}} \mathbb{P}_{M|_{\mathcal{S}}}(\{s_0 \ldots \in \text{Paths}(M) \mid \exists i < b : s_i \in G\})$$

$$P_{\min}(\Diamond^{\leq b} G) = \inf_{\mathcal{S}} \mathbb{P}_{M|_{\mathcal{S}}}(\{s_0 \ldots \in \text{Paths}(M) \mid \exists i < b : s_i \in G\})$$

Example:

$$P_{\max}(\Diamond^{\leq 2} \{UNRC\}) = 0.5$$
$$\underline{\phantom{P_{\max}}} \; {}^{\leq 3} \; \underline{\phantom{xxx}} = 0.5$$
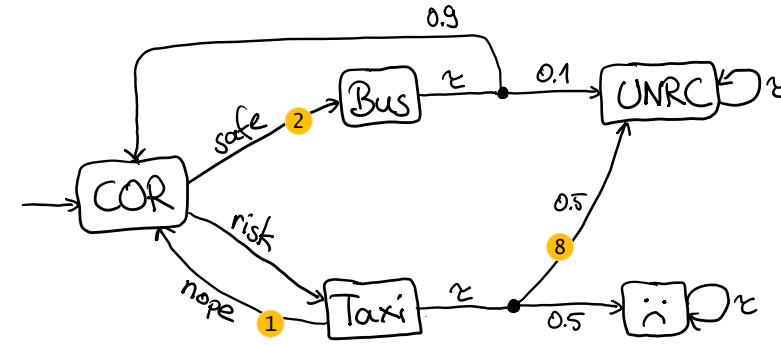$$\underline{\phantom{P_{\max}}} \; {}^{\leq 4} \; \underline{\phantom{xxx}} = 0.1 + 0.9 \cdot 0.5 = 0.55$$

## Step-bounded probabilistic reachability:

$$P_{opt}(\diamond^{\leq b} G)$$

for $opt \in \{\max, \min\}$ and $G \subseteq S$:

$$P_{\max}(\diamond G) = \sup_{\mathcal{S}} \mathbb{P}_{M|_{\mathcal{S}}}(\{s_0 \ldots \in \mathrm{Paths}(M) \mid \exists i < b : s_i \in G\})$$

$$P_{\min}(\diamond G) = \inf_{\mathcal{S}} \mathbb{P}_{M|_{\mathcal{S}}}(\{s_0 \ldots \in \mathrm{Paths}(M) \mid \exists i < b : s_i \in G\})$$

Complication: need step-positional schedulers

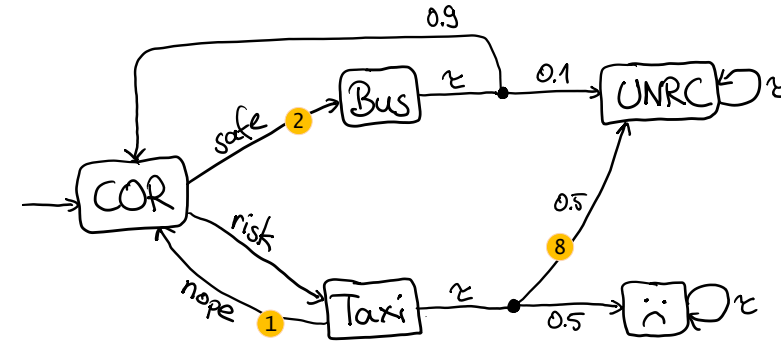$$\mathcal{S}_{st} \colon S \times \mathbb{N} \to A \times \mathrm{Dist}(S)$$

## Step-bounded probabilistic reachability:

$$\mathrm{P}_{opt}(\diamond^{\leq b} G)$$

for $opt \in \{\max, \min\}$ and $G \subseteq S$:

$$\mathrm{P}_{\max}(\diamond G) = \sup_{\mathcal{S}_{st}} \mathbb{P}_{M|_{\mathcal{S}_{st}}}(\{s_0 \ldots \in \mathrm{Paths}(M) \mid \exists i < b : s_i \in G\})$$

$$\mathrm{P}_{\min}(\diamond G) = \inf_{\mathcal{S}_{st}} \mathbb{P}_{M|_{\mathcal{S}_{st}}}(\{s_0 \ldots \in \mathrm{Paths}(M) \mid \exists i < b : s_i \in G\})$$

Complication: need step-positional schedulers

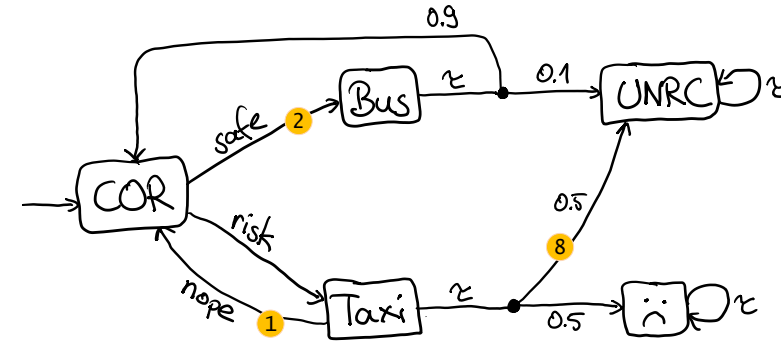$$\mathcal{S}_{st} : S \times \mathbb{N} \to A \times \mathrm{Dist}(S)$$

**Expected** accumulated **reward** to reach a goal:

$$R_{opt}(\diamond\, G)$$

for $opt \in \{\max, \min\}$ and $G \subseteq S$:

$$R_{\max}(\diamond\, G) = \sup_{\mathcal{S}} \mathbb{E}\left(rew_{\diamond G}^{M|_{\mathcal{S}}}\right)$$

$$R_{\min}(\diamond\, G) = \inf_{\mathcal{S}} \mathbb{E}\left(rew_{\diamond G}^{M|_{\mathcal{S}}}\right)$$

Example:

$$R_{min}(\diamond\{UNRC\}) = \min\{\, \infty,\ \infty,\ \underline{2 + 0.9 \cdot 2 + 0.9^2 \cdot 2 + \ldots}\,\}$$

$$R_{max}(\diamond\{UNRC\}) = \max\{\ \underline{\quad\quad}\ \text{"}\ \underline{\quad\quad}\ \} = \infty$$

**Expected** accumulated **reward** to reach a goal:

$$R_{opt}(\diamond G)$$

for $opt \in \{\max, \min\}$ and $G \subseteq S$:

$$R_{\max}(\diamond G) = \sup_{\mathcal{S}} \mathbb{E}\left(rew_{\diamond G}^{M|_{\mathcal{S}}}\right)$$

$$R_{\min}(\diamond G) = \inf_{\mathcal{S}} \mathbb{E}\left(rew_{\diamond G}^{M|_{\mathcal{S}}}\right)$$

$\rightarrow R_{\max}(\diamond G) = \infty$ if $\quad P_{min}(\diamond G) < 1$

$R_{\min}(\diamond G) = \infty$ if $\quad P_{max}(\diamond G) < 1$