

# Long-Term Values in Markov Decision Processes, (Co)Algebraically

Frank Feys & Helle Hvid Hansen & Larry Moss

TU Delft

TU Delft

Indiana University

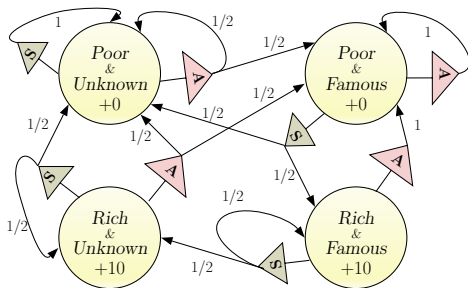
April 14, 2018

CMCS'18, Thessaloniki, Greece

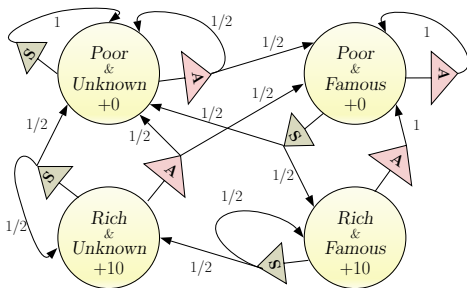
- Markov decision processes (MDPs) are state-based models of sequential decision-making under uncertainty
- Applications:
  - Planning
  - Reinforcement learning
  - Insurance and finance
  - ...
- We restrict to *finite, discrete time-homogeneous, infinite-horizon* MDPs, with the *discounting criterion*

## INTRO – MDPs: probab., state-based systems, with rewards 3/18

Example: A start-up company needs to decide to **Advertise** or **Save** money



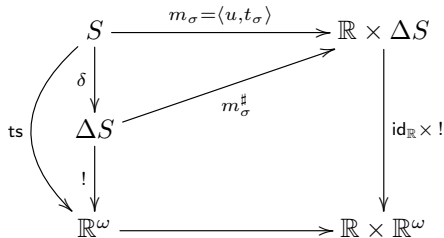
Example: A start-up company needs to decide to **Advertise** or **Save** money



A **Markov Decision Process (MDP)** is coalgebra  $m = \langle u, t \rangle: S \rightarrow \mathbb{R} \times (\Delta S)^A$  where

- $S$  is finite set of *states* and  $A$  is finite set of *actions*
- $u: S \rightarrow \mathbb{R}$  is *reward* map
- $t: S \rightarrow (\Delta S)^A$  is *transition* map (where  $\Delta S$  is set of prob. distr. on  $S$ )

A **policy**  $\sigma$  is a map  $\sigma: S \rightarrow A$



Given  $m = \langle u, t \rangle: S \rightarrow \mathbb{R} \times (\Delta S)^A$   
and policy  $\sigma: S \rightarrow A$

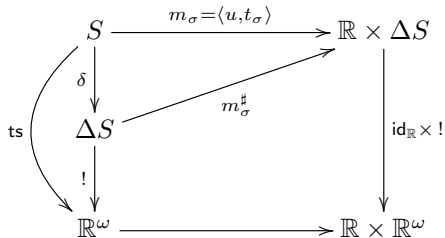
$$t_\sigma \stackrel{\text{def}}{=} t(s)(\sigma(s))$$

$\rightsquigarrow m_\sigma^\#$  by **determinization**

[Jacobs, Silva, Sokolova]

$\text{ts}(s) = (r_0^\sigma(s), r_1^\sigma(s), r_2^\sigma(s), \dots)$  is **trace semantics** of  $m_\sigma$

( $r_n^\sigma(s)$  is expected reward at time  $n$ , starting from  $s$ )



Given  $m = \langle u, t \rangle: S \rightarrow \mathbb{R} \times (\Delta S)^A$   
and policy  $\sigma: S \rightarrow A$

$$t_\sigma \stackrel{\text{def}}{=} t(s)(\sigma(s))$$

$\rightsquigarrow m_\sigma^\sharp$  by **determinization**

[Jacobs, Silva, Sokolova]

$ts(s) = (r_0^\sigma(s), r_1^\sigma(s), r_2^\sigma(s), \dots)$  is **trace semantics** of  $m_\sigma$

( $r_n^\sigma(s)$  is expected reward at time  $n$ , starting from  $s$ )

**Discounting criterion:** letting  $0 \leq \gamma < 1$  be a *discount factor*,  
the **long-term value** of policy  $\sigma$  is  $V^\sigma: S \rightarrow \mathbb{R}$ :

$$V^\sigma(s) = r_0^\sigma(s) + \gamma \cdot r_1^\sigma(s) + \dots + \gamma^n \cdot r_n^\sigma(s) + \dots$$

The **optimal value function**  $V^* : S \rightarrow \mathbb{R}$  of  $m$  in state  $s$  is given by

$$V^*(s) = \max_{\sigma} V^{\sigma}(s)$$

The **optimal value function**  $V^* : S \rightarrow \mathbb{R}$  of  $m$  in state  $s$  is given by

$$V^*(s) = \max_{\sigma} V^{\sigma}(s)$$

We define:

- $\sigma \geq \sigma'$  if for all  $s$ ,  $V^{\sigma}(s) \geq V^{\sigma'}(s)$
- $\sigma$  is **optimal** iff  $\sigma \geq \tau$  for all  $\tau$



The **optimal value function**  $V^* : S \rightarrow \mathbb{R}$  of  $m$  in state  $s$  is given by

$$V^*(s) = \max_{\sigma} V^{\sigma}(s)$$

We define:

- $\sigma \geq \sigma'$  if for all  $s$ ,  $V^{\sigma}(s) \geq V^{\sigma'}(s)$
- $\sigma$  is **optimal** iff  $\sigma \geq \tau$  for all  $\tau$

Classical facts (cf. (Puterman, 2014)):

- Optimal policy **always exists**
- Optimal policies need not be unique
- If  $\sigma$  is optimal, then  $V^{\sigma} = V^*$
- **Stationary** (memoryless), **deterministic** policies suffice

Observations:

- Classic theory uses low-level, analytic methods
- MDPs are coalgebras

**Goal:** to develop coalgebraic methods for reasoning about LTVs

The main **contributions**:

Part 1: Value function  $V^\sigma$  from *b*-corecursive algebras

Part 2: *Coinductive proof* of policy improvement theorem

---

# PART 1:

Value Function Arises from a Universal Property

- Following Bellman, the value function has a natural **recursive** structure:

( $V^\sigma$  from today) = reward today +  $\gamma \cdot$  ( $V^\sigma$  from tomorrow)

$$\boxed{V^\sigma = u + \gamma P_\sigma V^\sigma} \quad (1)$$

- Following Bellman, the value function has a natural **recursive** structure:

( $V^\sigma$  from today) = reward today +  $\gamma \cdot$  ( $V^\sigma$  from tomorrow)

$$\boxed{V^\sigma = u + \gamma P_\sigma V^\sigma} \quad (1)$$

- So,  $V^\sigma$  arises as a **fixpoint** of the operator  $\Psi_\sigma: \mathbb{R}^S \rightarrow \mathbb{R}^S$  given by

$$\Psi_\sigma(v) = u + \gamma P_\sigma v$$

- Following Bellman, the value function has a natural **recursive** structure:

$$(V^\sigma \text{ from today}) = \text{reward today} + \gamma \cdot (V^\sigma \text{ from tomorrow})$$

$$\boxed{V^\sigma = u + \gamma P_\sigma V^\sigma} \quad (1)$$

- So,  $V^\sigma$  arises as a **fixpoint** of the operator  $\Psi_\sigma: \mathbb{R}^S \rightarrow \mathbb{R}^S$  given by

$$\Psi_\sigma(v) = u + \gamma P_\sigma v$$

- Our observation:** we can re-express (1) as  $V^\sigma$  being a **coalgebra-to-algebra morphism**, as in

$$\begin{array}{ccc}
 S & \xrightarrow{m_\sigma = \langle u, t_\sigma \rangle} & \mathbb{R} \times \Delta S \\
 V^\sigma \downarrow & & \downarrow \text{id}_{\mathbb{R}} \times \Delta(V^\sigma) \\
 \mathbb{R} & \xleftarrow{\alpha_\gamma} \mathbb{R} \times \mathbb{R} \xleftarrow{\text{id}_{\mathbb{R}} \times \mathbf{E}} & \mathbb{R} \times \Delta \mathbb{R}
 \end{array}$$

where  $\alpha_\gamma: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  is  $\alpha_\gamma(x, y) = x + \gamma \cdot y$  and  $\mathbf{E}: \Delta \mathbb{R} \rightarrow \mathbb{R}$  is EV

- Recall that a **corecursive algebra** (for functor  $F$ ) is an  $F$ -algebra  $\alpha$  s.t.

$$\begin{array}{ccc} C & \xrightarrow{\forall f} & FC \\ \exists! f^\dagger \downarrow & & \downarrow Ff^\dagger \\ A & \xleftarrow{\alpha} & FA \end{array}$$



- Recall that a **corecursive algebra** (for functor  $F$ ) is an  $F$ -algebra  $\alpha$  s.t.

$$\begin{array}{ccc}
 C & \xrightarrow{\forall f} & FC \\
 \exists! f^\dagger \downarrow & & \downarrow Ff^\dagger \\
 A & \xleftarrow{\alpha} & FA
 \end{array}$$

- Recall:

$$\begin{array}{ccc}
 S & \xrightarrow{m_\sigma = \langle u, t_\sigma \rangle} & \mathbb{R} \times \Delta S \\
 V^\sigma \downarrow & & \downarrow \text{id}_{\mathbb{R}} \times \Delta(V^\sigma) \\
 \mathbb{R} & \xleftarrow{\alpha_\gamma \circ (\text{id}_{\mathbb{R}} \times E)} & \mathbb{R} \times \Delta \mathbb{R}
 \end{array}$$

- Recall that a **corecursive algebra** (for functor  $F$ ) is an  $F$ -algebra  $\alpha$  s.t.

$$\begin{array}{ccc}
 C & \xrightarrow{\forall f} & FC \\
 \exists! f^\dagger \downarrow & & \downarrow Ff^\dagger \\
 A & \xleftarrow{\alpha} & FA
 \end{array}$$

- Recall:

$$\begin{array}{ccc}
 S & \xrightarrow{m_\sigma = \langle u, t_\sigma \rangle} & \mathbb{R} \times \Delta S \\
 V^\sigma \downarrow & & \downarrow \text{id}_{\mathbb{R}} \times \Delta(V^\sigma) \\
 \mathbb{R} & \xleftarrow{\alpha_\gamma \circ (\text{id}_{\mathbb{R}} \times \mathbf{E})} & \mathbb{R} \times \Delta \mathbb{R}
 \end{array}$$

- Question:** is  $\alpha_\gamma \circ (\text{id}_{\mathbb{R}} \times \mathbf{E})$  a corecursive algebra?

- Recall that a **corecursive algebra** (for functor  $F$ ) is an  $F$ -algebra  $\alpha$  s.t.

$$\begin{array}{ccc}
 C & \xrightarrow{\forall f} & FC \\
 \exists! f^\dagger \downarrow & & \downarrow Ff^\dagger \\
 A & \xleftarrow{\alpha} & FA
 \end{array}$$

- Recall:

$$\begin{array}{ccc}
 S & \xrightarrow{m_\sigma = \langle u, t_\sigma \rangle} & \mathbb{R} \times \Delta S \\
 V^\sigma \downarrow & & \downarrow \text{id}_{\mathbb{R}} \times \Delta(V^\sigma) \\
 \mathbb{R} & \xleftarrow{\alpha_\gamma \circ (\text{id}_{\mathbb{R}} \times \mathbf{E})} & \mathbb{R} \times \Delta \mathbb{R}
 \end{array}$$

- Question:** is  $\alpha_\gamma \circ (\text{id}_{\mathbb{R}} \times \mathbf{E})$  a corecursive algebra?
- Consider a more basic question: is algebra  $\alpha_\gamma : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  **corecursive**?  
By (Capretta et al., 2004), this is **equivalent** with  $\alpha_\gamma \circ (\text{id}_{\mathbb{R}} \times \mathbf{E})$  corec.

- Letting  $H = \mathbb{R} \times \text{id}$ ,  $\alpha_\gamma$  is corecursive if

$$\begin{array}{ccc}
 X & \xrightarrow{\forall f} & \mathbb{R} \times X \\
 \exists! f^\dagger \downarrow & & \downarrow Hf^\dagger = \text{id}_{\mathbb{R}} \times f^\dagger \\
 \mathbb{R} & \xleftarrow{\alpha_\gamma} & \mathbb{R} \times \mathbb{R}
 \end{array}$$

- Letting  $H = \mathbb{R} \times \text{id}$ ,  $\alpha_\gamma$  is corecursive if

$$\begin{array}{ccc}
 X & \xrightarrow{\forall f} & \mathbb{R} \times X \\
 \exists! f^\dagger \downarrow & & \downarrow Hf^\dagger = \text{id}_{\mathbb{R}} \times f^\dagger \\
 \mathbb{R} & \xleftarrow{\alpha_\gamma} & \mathbb{R} \times \mathbb{R}
 \end{array}$$

- In particular, let  $X$  be a set of “variables”  $\{x_0, x_1, x_2, \dots\}$ , and fix  $\gamma$ .

$$\left\{ \begin{array}{l} \text{System of lin. eqs.} \\ (3) \quad \boxed{x_n = a_n + \gamma \cdot x_{n+1}} \\ \quad (n = 0, 1, 2, \dots) \\ \text{Solutions to (3)} \end{array} \right\} \begin{array}{l} \iff \text{coalgebra } f: X \rightarrow \mathbb{R} \times X \\ \iff f^\dagger \text{ s.t. } f^\dagger = \alpha_\gamma \cdot (\text{id}_{\mathbb{R}} \times f^\dagger) \cdot f \end{array}$$

- Letting  $H = \mathbb{R} \times \text{id}$ ,  $\alpha_\gamma$  is corecursive if

$$\begin{array}{ccc}
 X & \xrightarrow{\forall f} & \mathbb{R} \times X \\
 \exists! f^\dagger \downarrow & & \downarrow Hf^\dagger = \text{id}_{\mathbb{R}} \times f^\dagger \\
 \mathbb{R} & \xleftarrow{\alpha_\gamma} & \mathbb{R} \times \mathbb{R}
 \end{array}$$

- In particular, let  $X$  be a set of “variables”  $\{x_0, x_1, x_2, \dots\}$ , and fix  $\gamma$ .

$$\left\{ \begin{array}{l} \text{System of lin. eqs.} \\ (3) \quad \boxed{x_n = a_n + \gamma \cdot x_{n+1}} \\ \quad (n = 0, 1, 2, \dots) \end{array} \right\} \iff \text{coalgebra } f: X \rightarrow \mathbb{R} \times X$$

$$\text{Solutions to (3)} \iff f^\dagger \text{ s.t. } f^\dagger = \alpha_\gamma \cdot (\text{id}_{\mathbb{R}} \times f^\dagger) \cdot f$$

- But, (3) has *infinite* number of solutions, even if  $(a_n)_n$  is bounded  
 $\Rightarrow$  answer to earlier question is **NO**

- Letting  $H = \mathbb{R} \times \text{id}$ ,  $\alpha_\gamma$  is corecursive if

$$\begin{array}{ccc}
 X & \xrightarrow{\forall f} & \mathbb{R} \times X \\
 \exists! f^\dagger \downarrow & & \downarrow Hf^\dagger = \text{id}_{\mathbb{R}} \times f^\dagger \\
 \mathbb{R} & \xleftarrow{\alpha_\gamma} & \mathbb{R} \times \mathbb{R}
 \end{array}$$

- In particular, let  $X$  be a set of “variables”  $\{x_0, x_1, x_2, \dots\}$ , and fix  $\gamma$ .

$$\left\{ \begin{array}{l} \text{System of lin. eqs.} \\ (3) \quad \boxed{x_n = a_n + \gamma \cdot x_{n+1}} \\ \quad (n = 0, 1, 2, \dots) \\ \text{Solutions to (3)} \end{array} \right\} \begin{array}{l} \iff \text{coalgebra } f: X \rightarrow \mathbb{R} \times X \\ \iff f^\dagger \text{ s.t. } f^\dagger = \alpha_\gamma \cdot (\text{id}_{\mathbb{R}} \times f^\dagger) \cdot f \end{array}$$

- But, (3) has *infinite* number of solutions, even if  $(a_n)_n$  is bounded  
 $\Rightarrow$  answer to earlier question is **NO**
- However**, if  $(a_n)_n$  is bounded then (3) has a unique **bounded** solution

- To get uniqueness  $\Rightarrow$  incorporate boundedness information



- To get uniqueness  $\Rightarrow$  incorporate **boundedness information**
- **Definition** A *b-category*  $(\mathbf{C}, \mathcal{B})$  is a category  $\mathbf{C}$  with a subcollection of “bounded” morphism  $\mathcal{B}$  s.t.  $(f \in \mathcal{B} \Rightarrow f \circ g \in \mathcal{B})$

- To get uniqueness  $\Rightarrow$  incorporate **boundedness information**
- Definition** A *b-category*  $(\mathcal{C}, \mathcal{B})$  is a category  $\mathcal{C}$  with a subcollection of “bounded” morphism  $\mathcal{B}$  s.t.  $(f \in \mathcal{B} \Rightarrow f \circ g \in \mathcal{B})$
- Definition** Let  $(\mathcal{C}, \mathcal{B})$  *b-category*,  $F$  endofunctor on  $\mathcal{C}$ ,  $\alpha: FA \rightarrow A$  an  $F$ -algebra. Then  $\alpha$  is a *b-corecursive algebra (bca)* if

$$\begin{array}{ccc}
 X & \xrightarrow{\forall f \in \mathcal{B}} & FX \\
 \exists! f^\dagger \in \mathcal{B} \downarrow & & \downarrow Ff^\dagger \\
 A & \xleftarrow{\alpha} & FA
 \end{array}$$

- To get uniqueness  $\Rightarrow$  incorporate **boundedness information**
- Definition** A *b-category*  $(\mathcal{C}, \mathcal{B})$  is a category  $\mathcal{C}$  with a subcollection of “bounded” morphism  $\mathcal{B}$  s.t.  $(f \in \mathcal{B} \Rightarrow f \circ g \in \mathcal{B})$
- Definition** Let  $(\mathcal{C}, \mathcal{B})$  *b-category*,  $F$  endofunctor on  $\mathcal{C}$ ,  $\alpha: FA \rightarrow A$  an  $F$ -algebra. Then  $\alpha$  is a *b-corecursive algebra (bca)* if

$$\begin{array}{ccc}
 X & \xrightarrow{\forall f \in \mathcal{B}} & FX \\
 \exists! f^\dagger \in \mathcal{B} \downarrow & & \downarrow Ff^\dagger \\
 A & \xleftarrow{\alpha} & FA
 \end{array}$$

- Proposition**  $\alpha_\gamma$  is a **bca** in  $(\text{Met}, B)$  for  $H$
- Then by *b-version* of (Capretta et al., 2004) result:  
**Proposition**  $\alpha_\gamma \circ (\mathbb{R} \times \mathbb{E})$  is a **bca** in  $(\text{Met}, B)$  for  $H \circ \Delta$

Part 1: Value function  $V^\sigma$  from *b*-corecursive algebras ✓

Part 2: *Coinductive proof* of policy improvement theorem

---

# PART 2:

## Correctness of Policy Iteration via Contraction (Co)Induction

- Suppose  $\sigma$  is policy  $\rightsquigarrow$  find  $V^\sigma$

- Suppose  $\sigma$  is policy  $\rightsquigarrow$  find  $V^\sigma$
- Now we define a **new policy**  $\sigma'$  by putting for each state  $s$

$$\sigma'(s) = \operatorname{argmax}_{a \in A} \left\{ \sum_{s' \in S} P(s, a, s') V^\sigma(s') \right\}$$

- Suppose  $\sigma$  is policy  $\rightsquigarrow$  find  $V^\sigma$
- Now we define a **new policy**  $\sigma'$  by putting for each state  $s$

$$\sigma'(s) = \operatorname{argmax}_{a \in A} \left\{ \sum_{s' \in S} P(s, a, s') V^\sigma(s') \right\}$$

**Theorem** (Howard, 1960)

- The policy  $\sigma'$  is a **better** policy than  $\sigma$ , i.e.,  $\sigma' \geq \sigma$ .
- If  $\sigma' = \sigma$ , then  $\sigma$  is **optimal**.



- Suppose  $\sigma$  is policy  $\rightsquigarrow$  find  $V^\sigma$
- Now we define a **new policy**  $\sigma'$  by putting for each state  $s$

$$\sigma'(s) = \operatorname{argmax}_{a \in A} \left\{ \sum_{s' \in S} P(s, a, s') V^\sigma(s') \right\}$$

### Theorem (Howard, 1960)

- The policy  $\sigma'$  is a **better** policy than  $\sigma$ , i.e.,  $\sigma' \geq \sigma$ .
  - If  $\sigma' = \sigma$ , then  $\sigma$  is **optimal**.
- 
- **Policy Iteration**: start with *any*  $\sigma$ , iteratively obtain  $\sigma', \sigma'', \sigma''', \dots$ , and continue until there is fixpoint  $\Rightarrow$  this outputs an *optimal policy*

**Policy Improvement Theorem** (Howard, 1960)

$$P_{\sigma'} \cdot V^{\sigma} \geq P_{\sigma} \cdot V^{\sigma} \Rightarrow V^{\sigma'} \geq V^{\sigma}$$

**Policy Improvement Theorem** (Howard, 1960)

$$P_{\sigma'} \cdot V^{\sigma} \geq P_{\sigma} \cdot V^{\sigma} \Rightarrow V^{\sigma'} \geq V^{\sigma}$$

- Recall

$$\sigma'(s) = \operatorname{argmax}_{a \in A} \left\{ \sum_{s' \in S} P(s, a, s') V^{\sigma}(s') \right\}$$

$\Rightarrow$  antecedent holds

**Policy Improvement Theorem** (Howard, 1960)

$$P_{\sigma'} \cdot V^{\sigma} \geq P_{\sigma} \cdot V^{\sigma} \Rightarrow V^{\sigma'} \geq V^{\sigma}$$

- Recall

$$\sigma'(s) = \operatorname{argmax}_{a \in A} \left\{ \sum_{s' \in S} P(s, a, s') V^{\sigma}(s') \right\}$$

$\Rightarrow$  antecedent holds

- New proof using **Contraction (Co)Induction**

### Contraction (Co)Induction Theorem

Let  $M$  be a non-empty, complete, **ordered** (i.e., there is partial order  $\leq$  s.t.  $\uparrow x = \{y \mid x \leq y\}$  and  $\downarrow x = \{y \mid y \leq x\}$  are closed) metric space. If  $f: M \rightarrow M$  is contractive and order-preserving, then the (unique) fixpoint  $x^*$  is

- **least pre-fixpoint** (if  $f(x) \leq x$ , then  $x^* \leq x$ ),
- **greatest post-fixpoint** (if  $x \leq f(x)$ , then  $x \leq x^*$ ).

Cf. *Metric Coinduction* (Kozen & Ruozzi, 2009) and (Denardo, 1967).

### Contraction (Co)Induction Theorem

Let  $M$  be a non-empty, complete, **ordered** (i.e., there is partial order  $\leq$  s.t.  $\uparrow x = \{y \mid x \leq y\}$  and  $\downarrow x = \{y \mid y \leq x\}$  are closed) metric space. If  $f: M \rightarrow M$  is contractive and order-preserving, then the (unique) fixpoint  $x^*$  is

- **least pre-fixpoint** (if  $f(x) \leq x$ , then  $x^* \leq x$ ),
- **greatest post-fixpoint** (if  $x \leq f(x)$ , then  $x \leq x^*$ ).

Cf. *Metric Coinduction* (Kozen & Ruozzi, 2009) and (Denardo, 1967).

**Proof of Policy Improvement** (i.e.,  $P_{\sigma'} \cdot V^\sigma \geq P_\sigma \cdot V^\sigma \Rightarrow V^{\sigma'} \geq V^\sigma$ ).

Apply theorem to  $\Psi_\pi: \mathbb{R}^S \rightarrow \mathbb{R}^S$  (contractive and order-preserving  $\checkmark$ )

$$\Psi_\pi(v) = u + \gamma P_\pi v, \quad \text{and } V^\pi \text{ is its fixpoint.}$$

Thus  $P_{\sigma'} \cdot V^\sigma \geq P_\sigma \cdot V^\sigma \Rightarrow \Psi_{\sigma'}(V^\sigma) \geq \Psi_\sigma(V^\sigma) = V^\sigma \Rightarrow V^{\sigma'} \geq V^\sigma$ . ■

## Contributions:

- Value functions  $V^\sigma$  and  $V^*$  from  $b$ -corecursive algebras
- Coinductive proof of policy improvement theorem

## Future work:

- Generalize the setting (e.g., to stochastic games)
- Make connections with related literature:
  - Combining semantics of computation and game theory (Pavlovic, 2009)
  - Coalgebraic formulation of infinite games (Abramsky & Winschel, 2017)
  - Open games (Hedges, Ghani, Winschel, Zahn, 2018)
- Investigate contraction coinduction further and look for other applications (e.g., in social choice)

## Contributions:

- Value functions  $V^\sigma$  and  $V^*$  from  $b$ -corecursive algebras
- Coinductive proof of policy improvement theorem

## Future work:

- Generalize the setting (e.g., to stochastic games)
- Make connections with related literature:
  - Combining semantics of computation and game theory (Pavlovic, 2009)
  - Coalgebraic formulation of infinite games (Abramsky & Winschel, 2017)
  - Open games (Hedges, Ghani, Winschel, Zahn, 2018)
- Investigate contraction coinduction further and look for other applications (e.g., in social choice)

Thank you!