



Global Decision-Making via Local Economic Transactions



Michael Chang



Sid Kaushik



Matt Weinberg

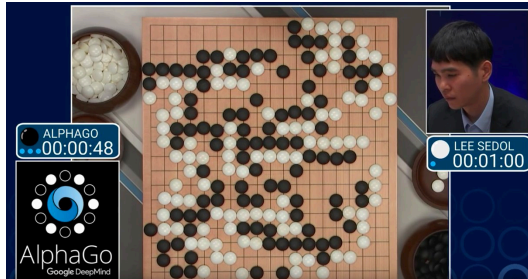


Tom Griffiths

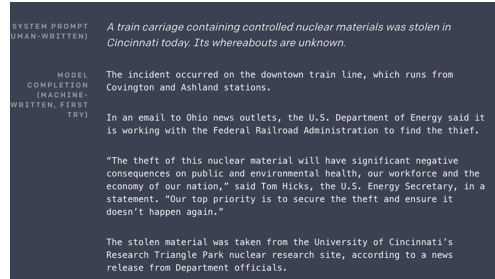


Sergey Levine

Much Success So Far



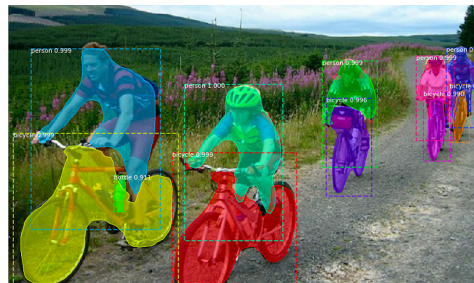
Game Playing
Silver et al. (2016)



Natural Language Processing
Radford et al. (2019)

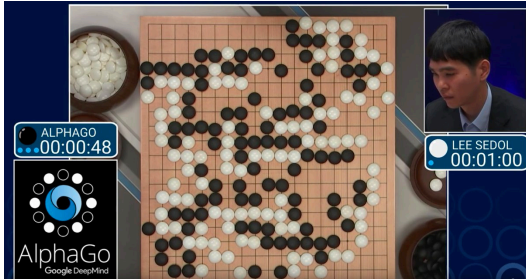


Robotics
Levine et al. (2016)

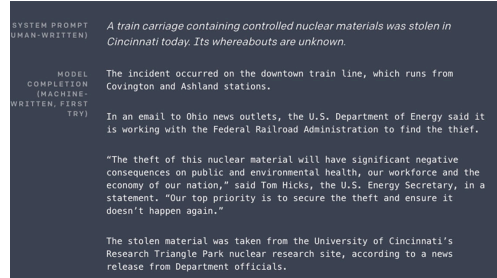


Computer Vision
He et al. (2017)

Much Success So Far: Monolithic Optimization



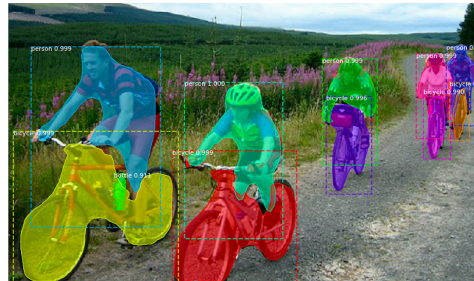
Game Playing
Silver et al. (2016)



Natural Language Processing
Radford et al. (2019)



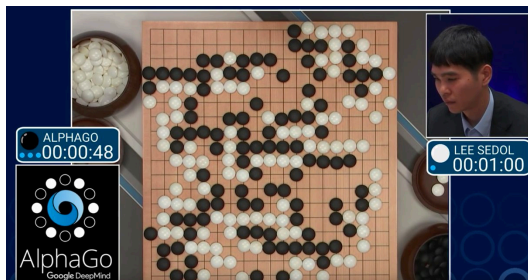
Robotics
Levine et al. (2016)



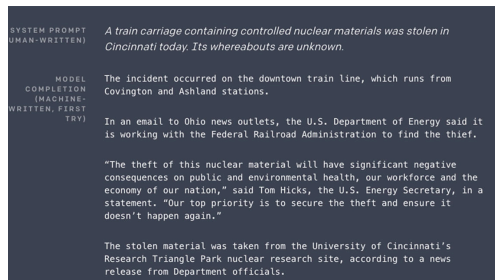
Computer Vision
He et al. (2017)



Much Success So Far: Monolithic Optimization



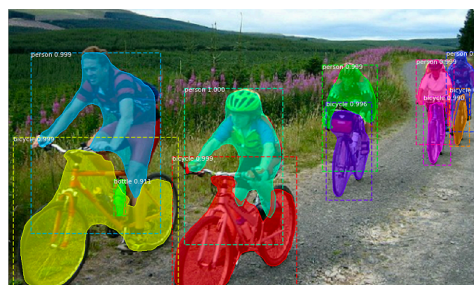
Game Playing
Silver et al. (2016)



Natural Language Processing
Radford et al. (2019)



Robotics
Levine et al. (2016)

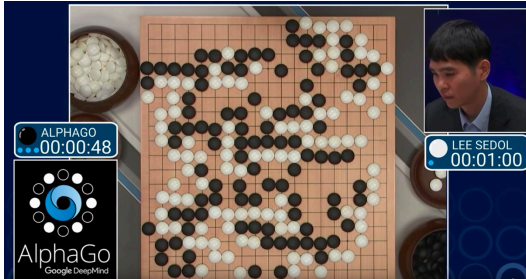


Computer Vision
He et al. (2017)

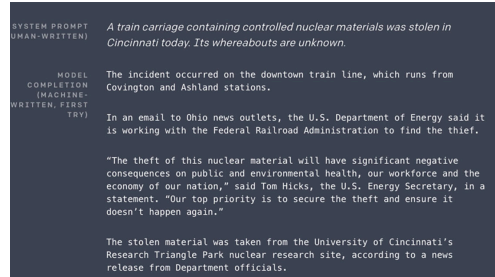


One optimization problem

Much Success So Far: Monolithic Optimization



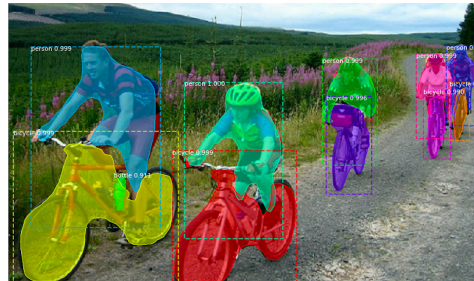
Game Playing
Silver et al. (2016)



Natural Language Processing
Radford et al. (2019)



Robotics
Levine et al. (2016)

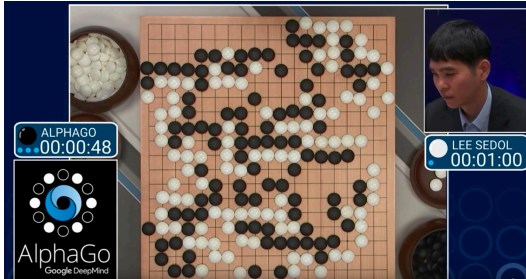


Computer Vision
He et al. (2017)

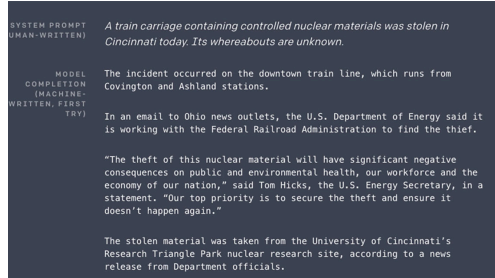


One optimization problem
One agent

Much Success So Far: Monolithic Optimization



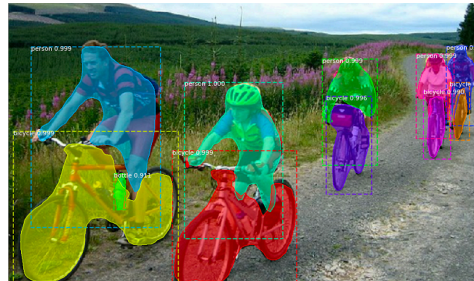
Game Playing
Silver et al. (2016)



Natural Language Processing
Radford et al. (2019)



Robotics
Levine et al. (2016)



Computer Vision
He et al. (2017)



One optimization problem
One agent
One objective

Decentralized Optimization



Corporation



One optimization problem
One agent
One objective

Decentralized Optimization



Corporation



Many optimization problems

Many agents

Many objectives

Decentralized Optimization



Many *local* optimization problems
Many *local* agents
Many *local* objectives



Emergent *global* optimization problem
Emergent *global* agent
Emergent *global* objective

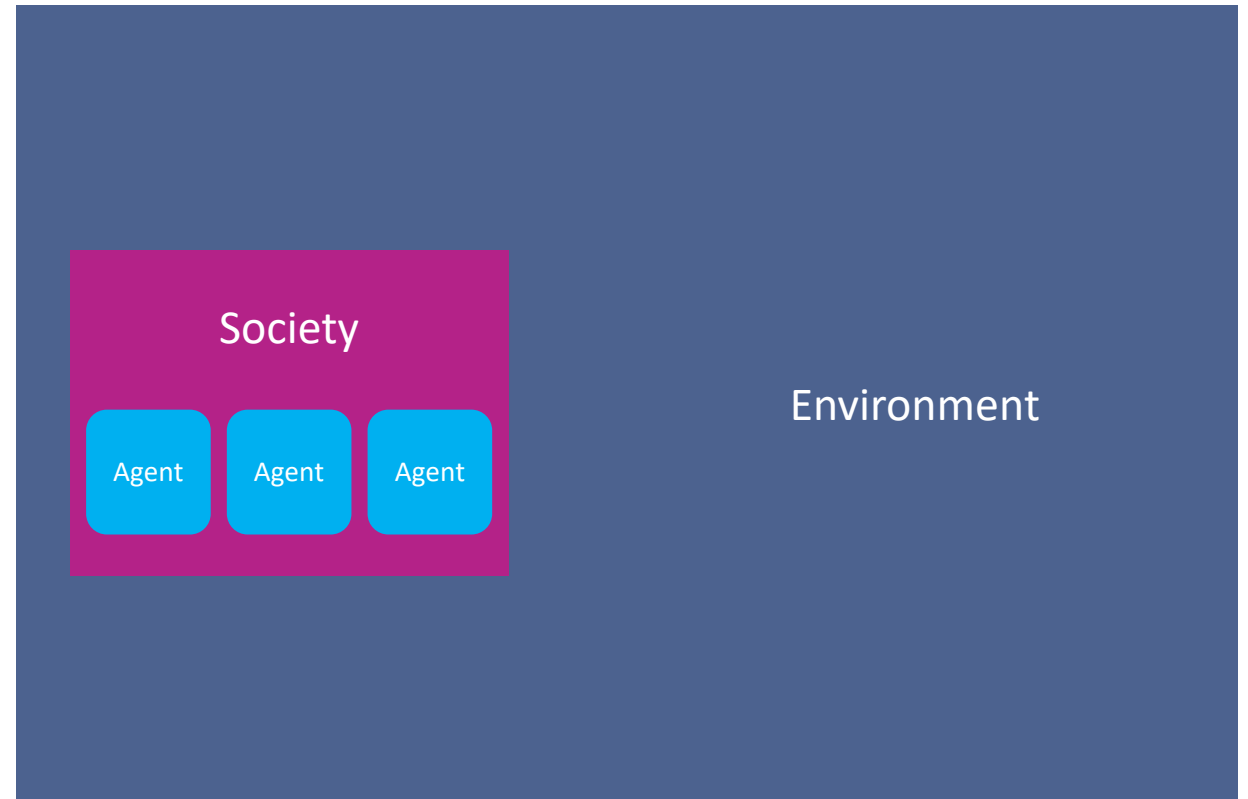
Decentralized Optimization



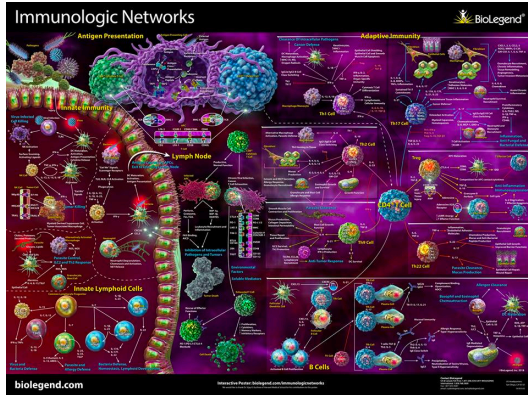
Decentralized Optimization



Decentralized Optimization



Decentralized Optimization



Biological Processes



Ecosystems



Economies



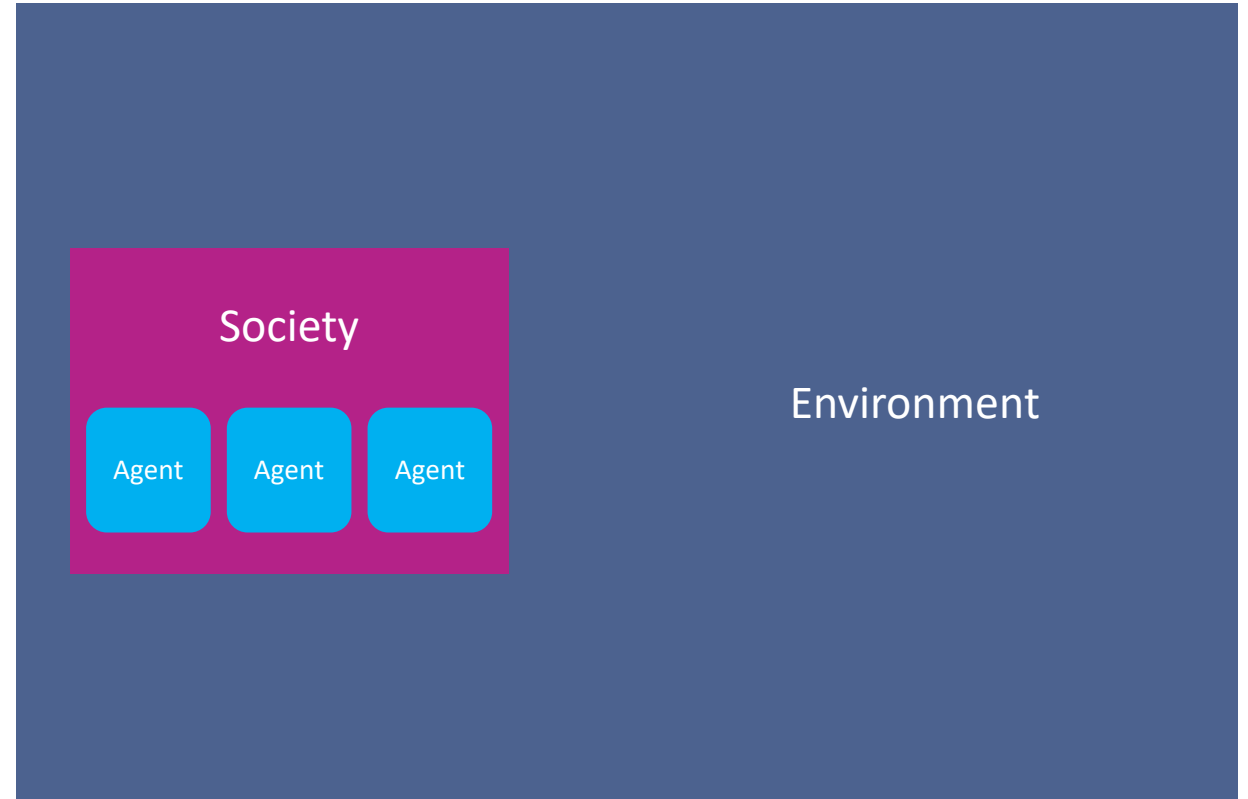
Organizations



Optimization at Two Levels of Abstraction

Challenge

How can we build machine learning algorithms that relate the global level of the society and the local level of the agent?



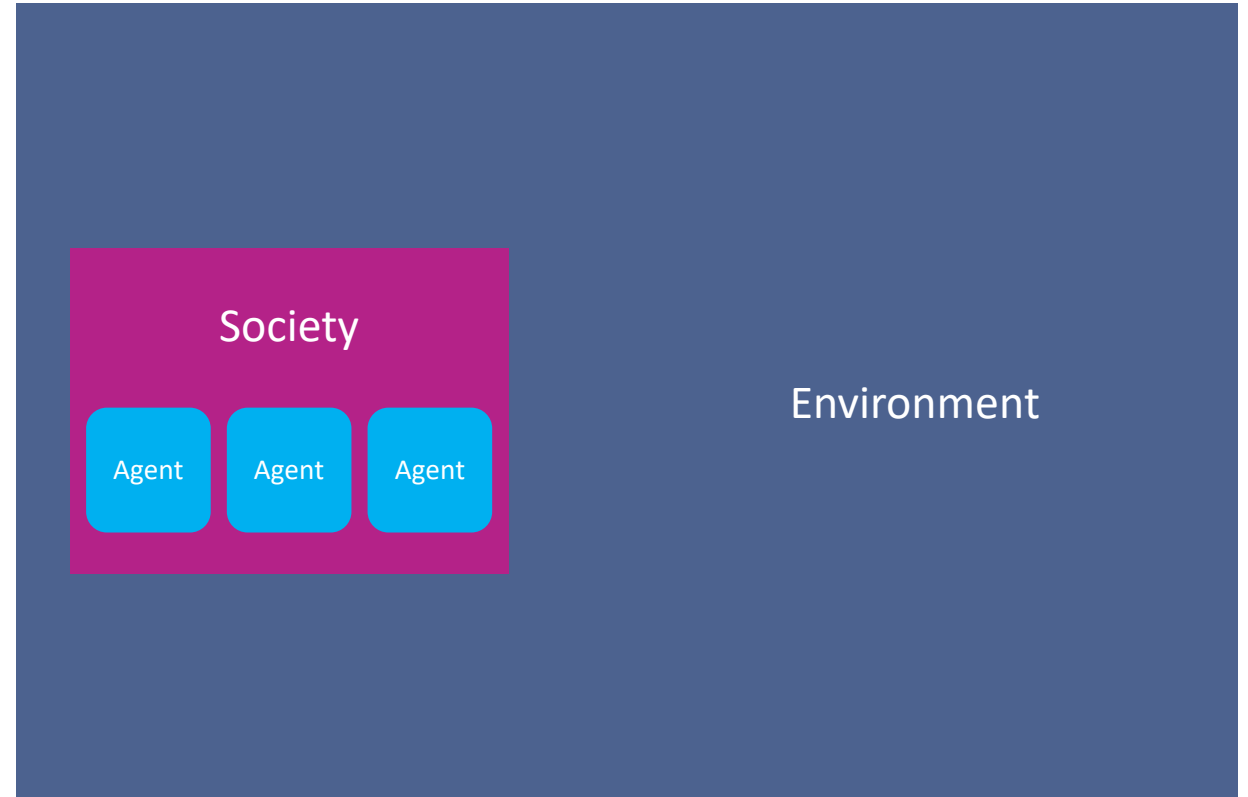
Optimization at Two Levels of Abstraction

Challenge

How can we build machine learning algorithms that relate the global level of the society and the local level of the agent?

Implications

- Enable the design of learning algorithms that are inherently modular



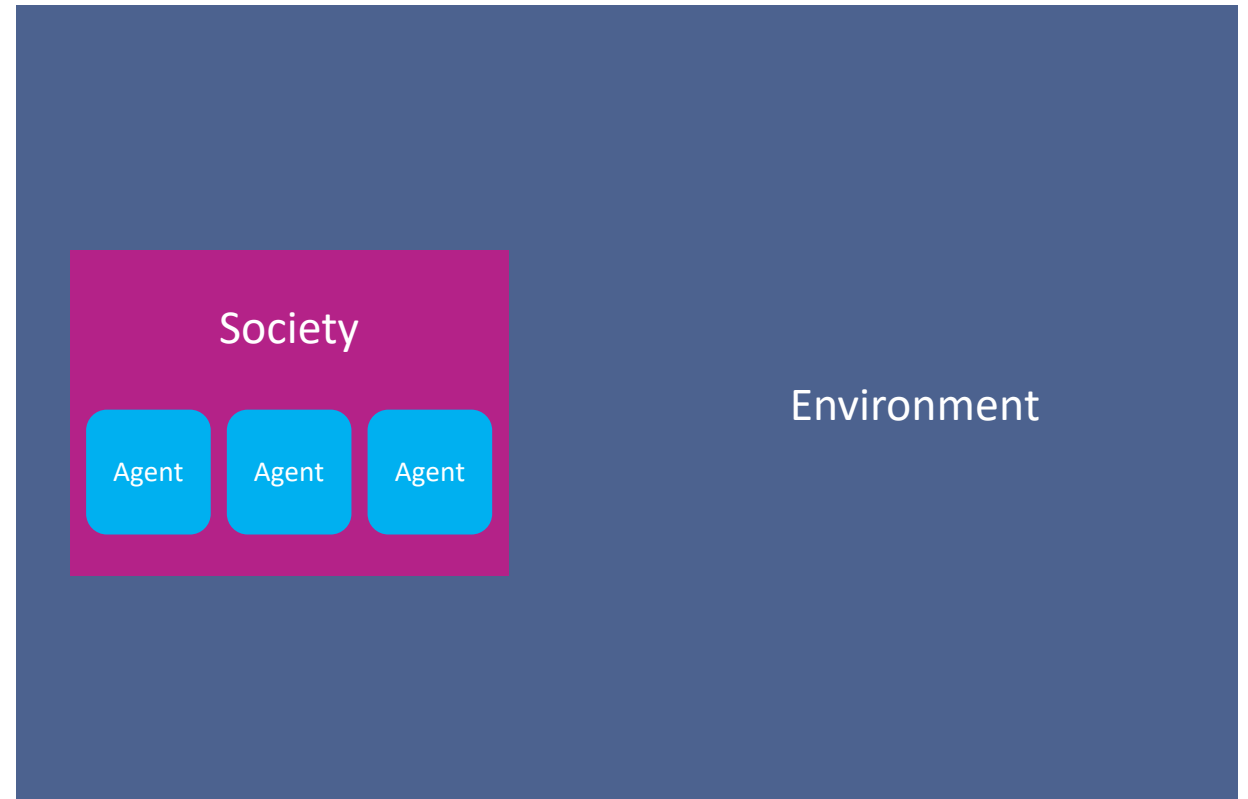
Optimization at Two Levels of Abstraction

Challenge

How can we build machine learning algorithms that relate the global level of the society and the local level of the agent?

Implications

- Enable the design of learning algorithms that are inherently modular
- Provide a recipe for engineering and analyzing a multi-agent system to achieve a desired global outcome



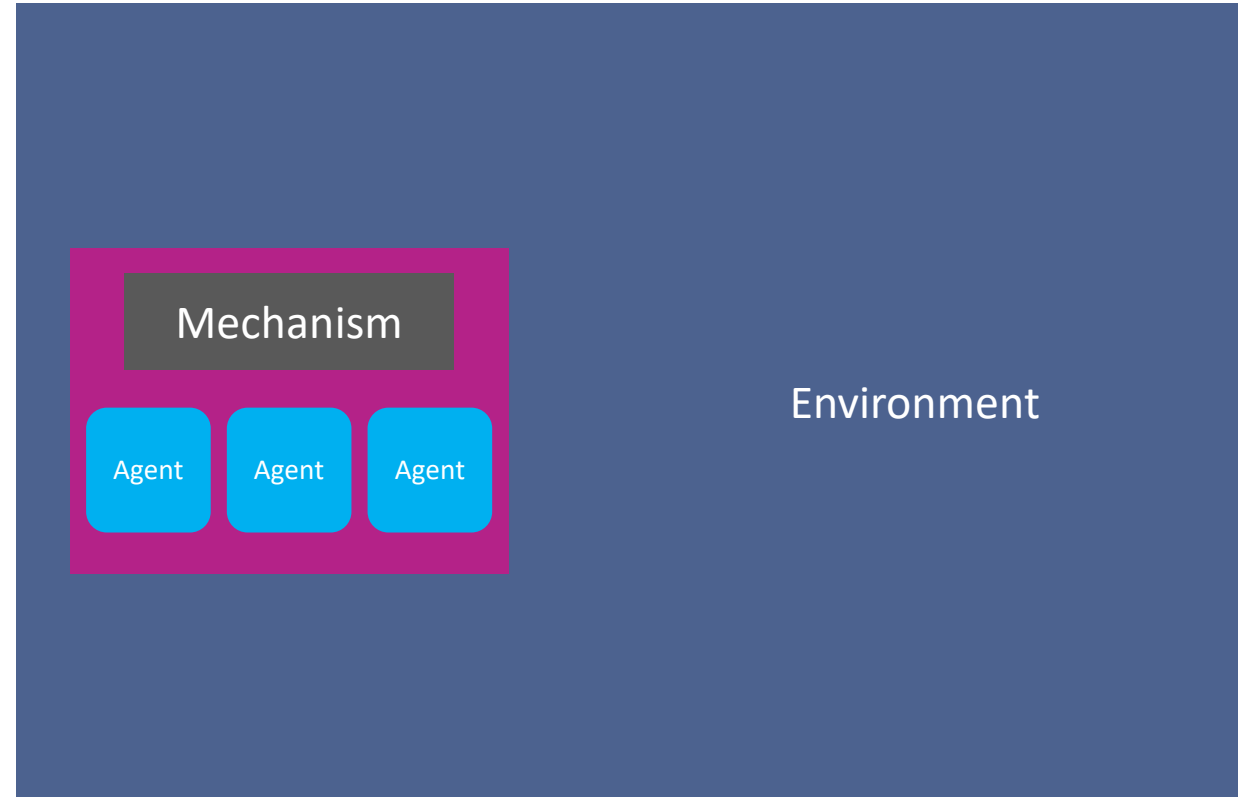
Optimization at Two Levels of Abstraction

Challenge

How can we build machine learning algorithms that relate the global level of the society and the local level of the agent?

Implications

- Enable the design of learning algorithms that are inherently modular
- Provide a recipe for engineering and analyzing a multi-agent system to achieve a desired global outcome



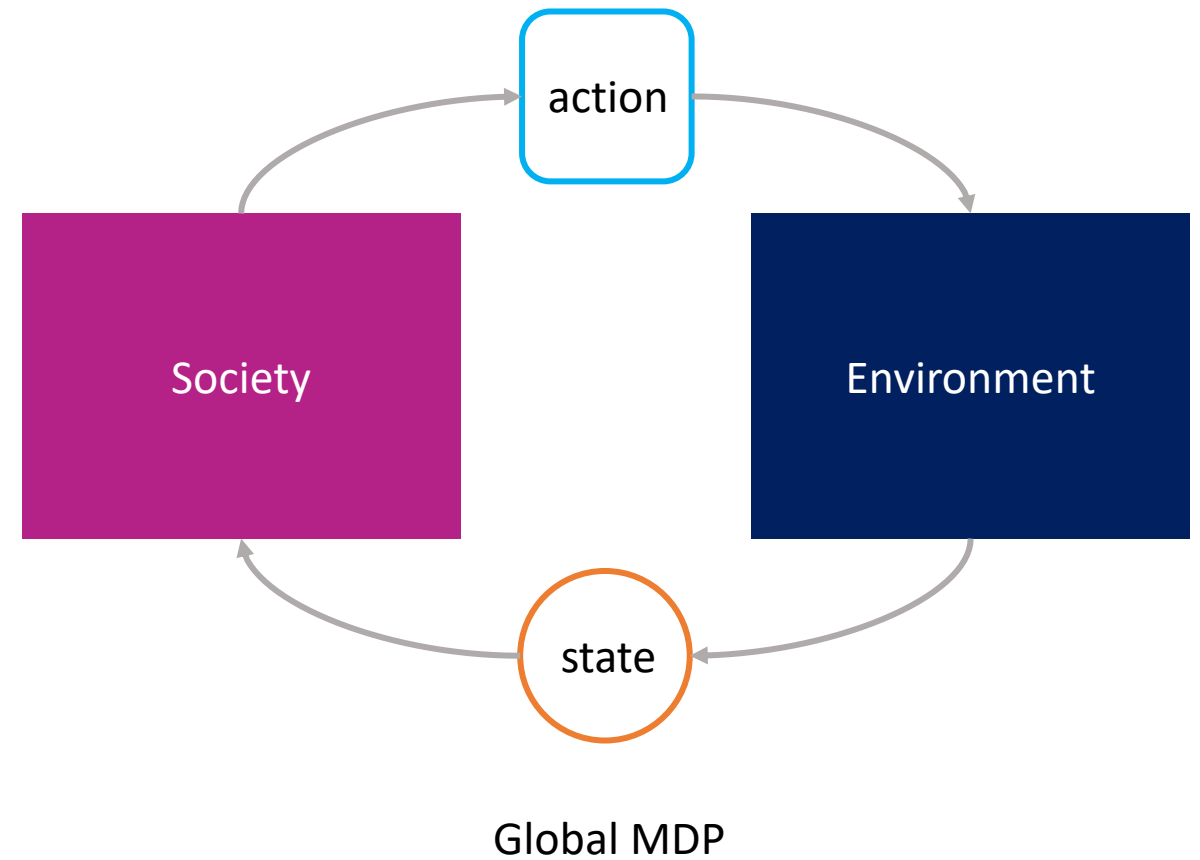
This Paper

This Paper: Assumptions

This Paper: Assumptions

Assumptions

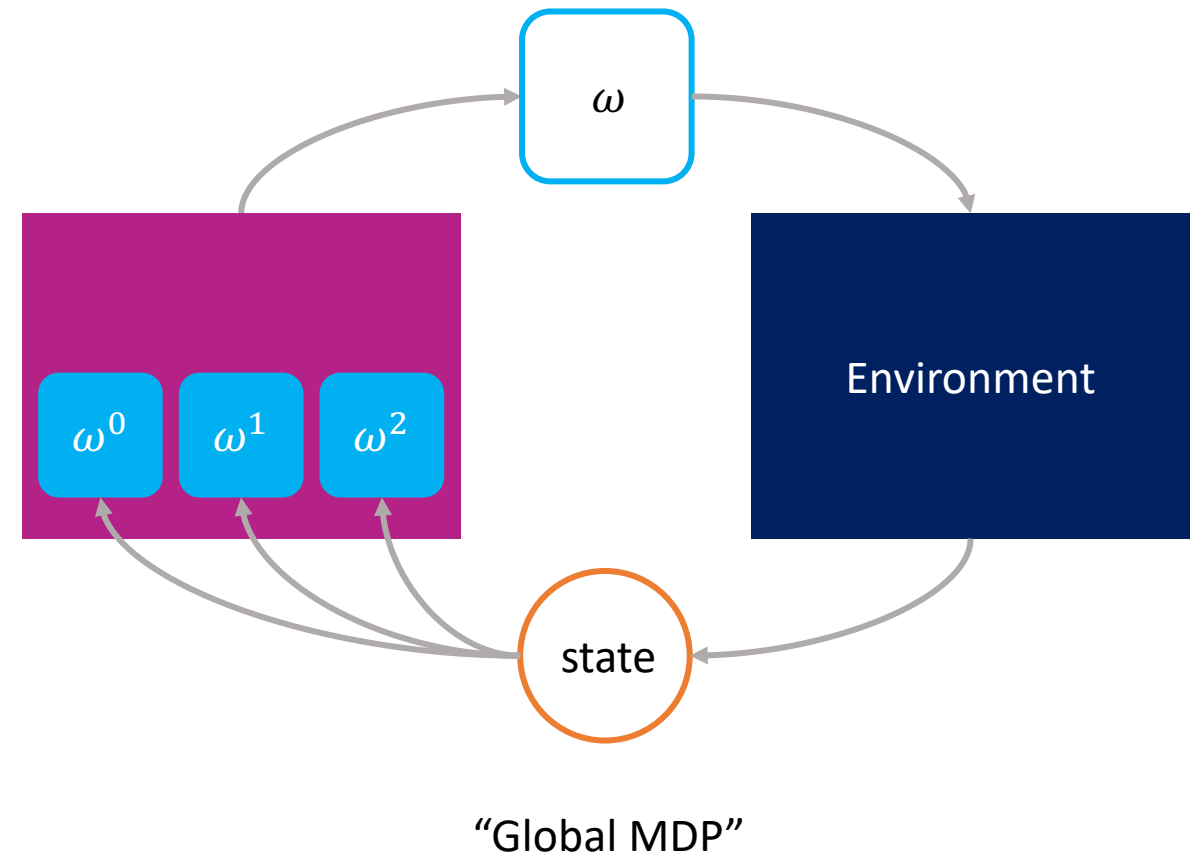
- Sequential decision-making setting



This Paper: Assumptions

Assumptions

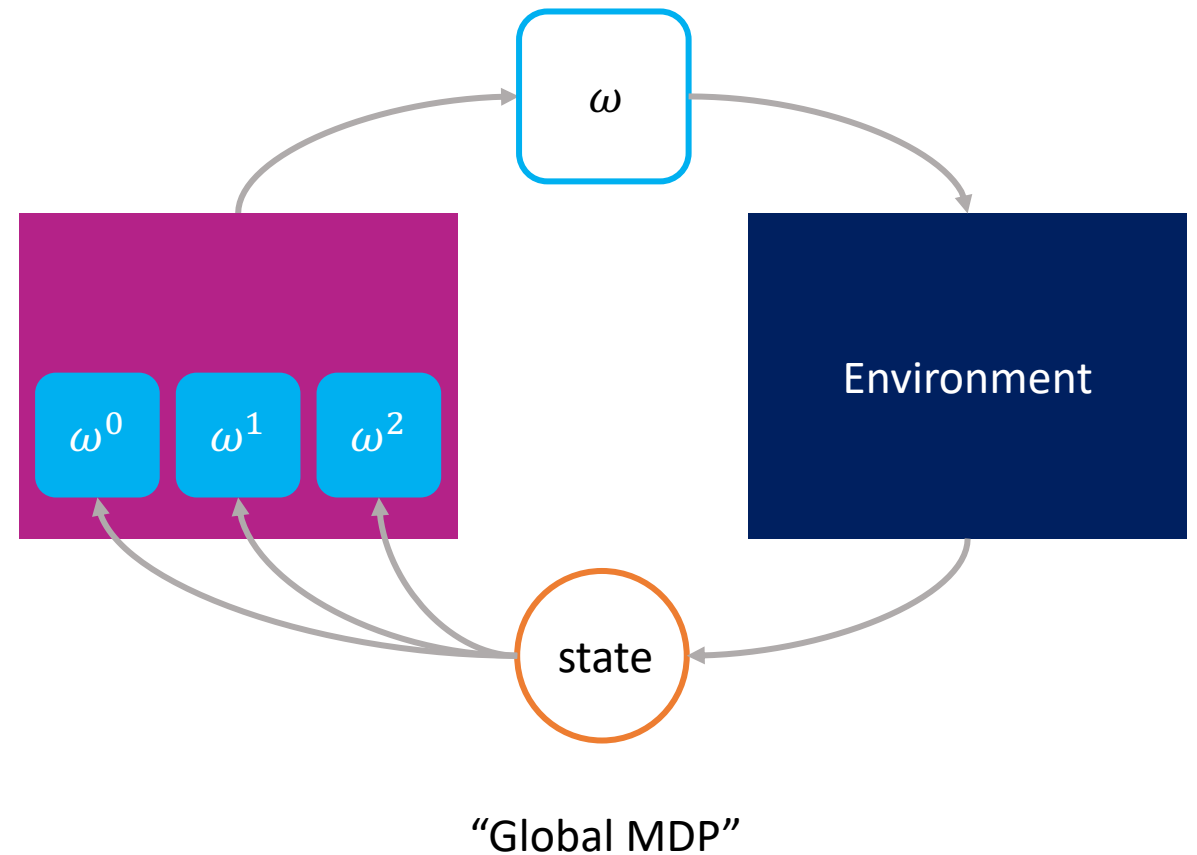
- Sequential decision making setting
- Each agent produces a specialized transformation to the state (e.g. a literal action)



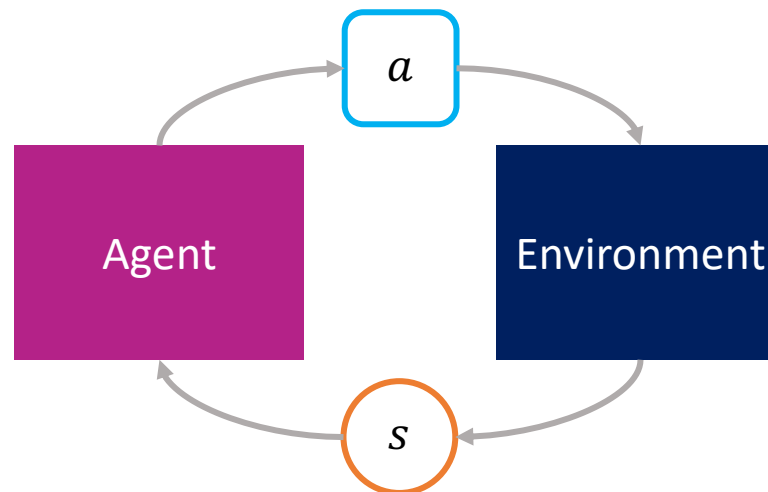
This Paper: Assumptions

Assumptions

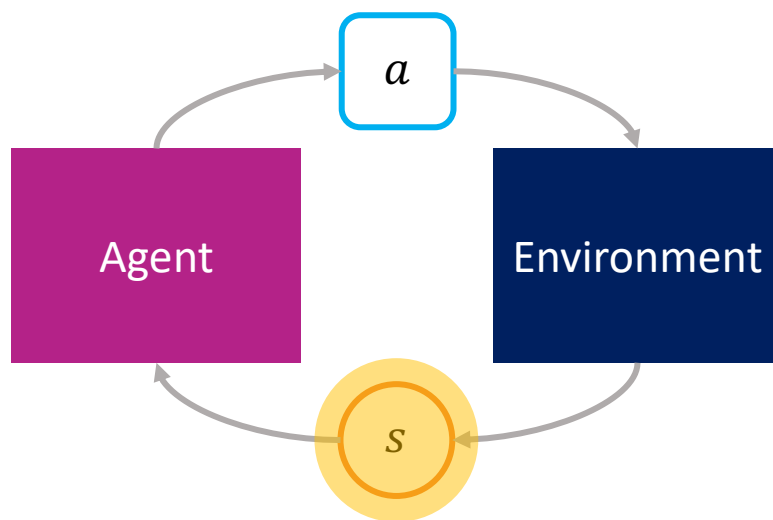
- Sequential decision making setting
- Each agent produces a specialized transformation to the state (e.g. a literal action)
- Only one agent activates at each time step

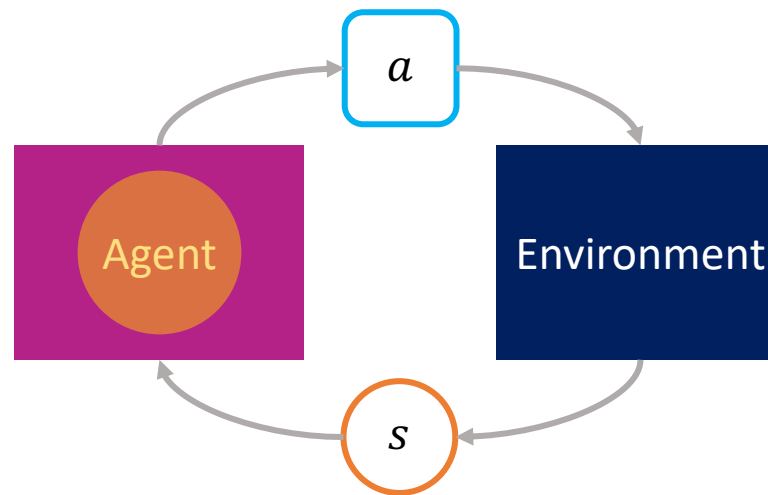
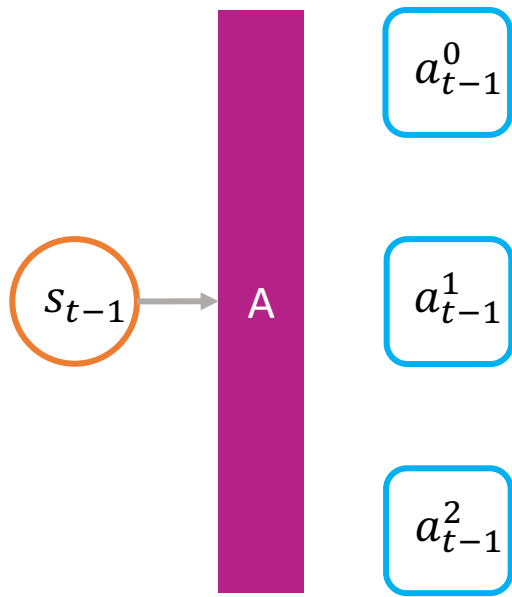


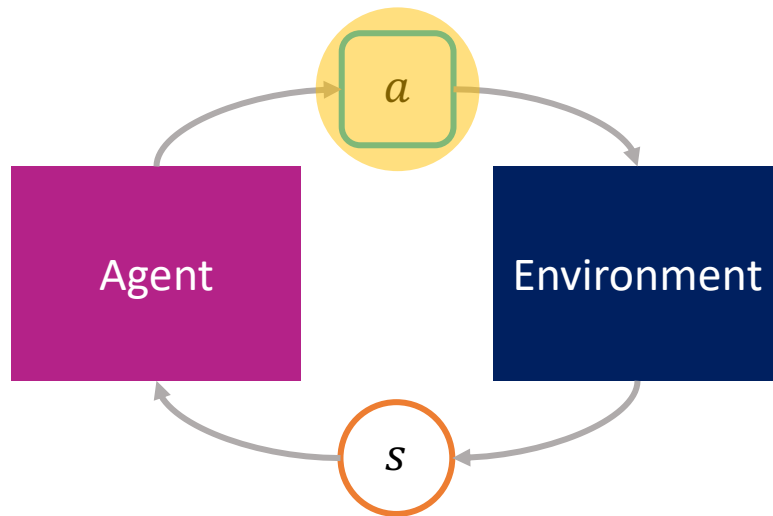
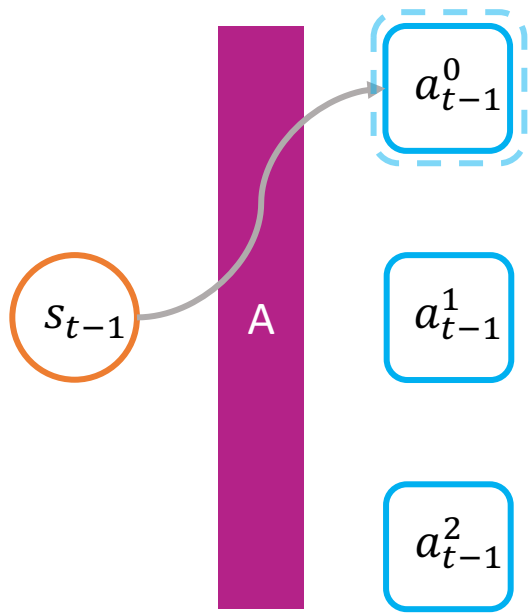
Intuition

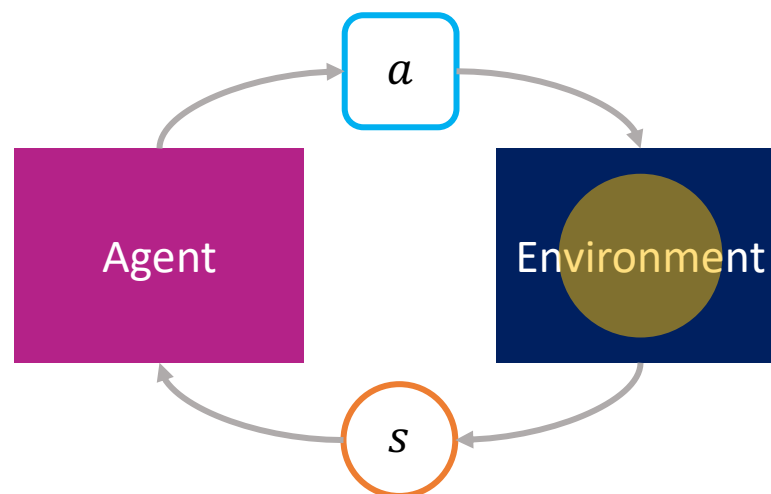
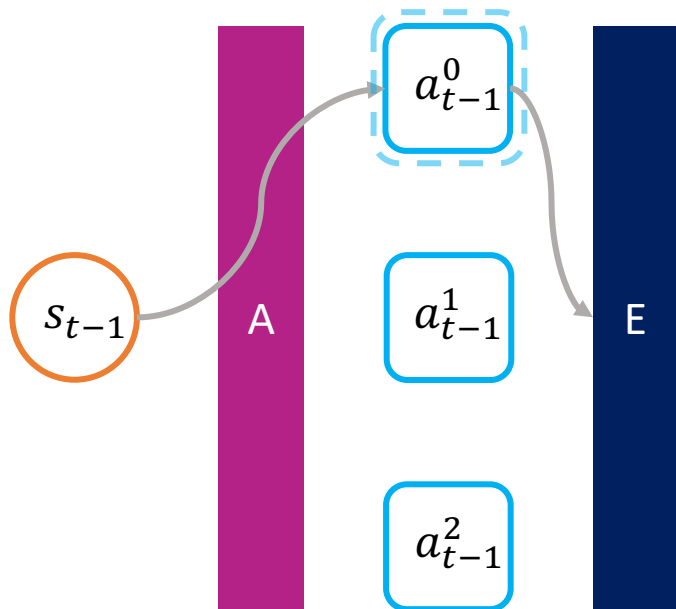


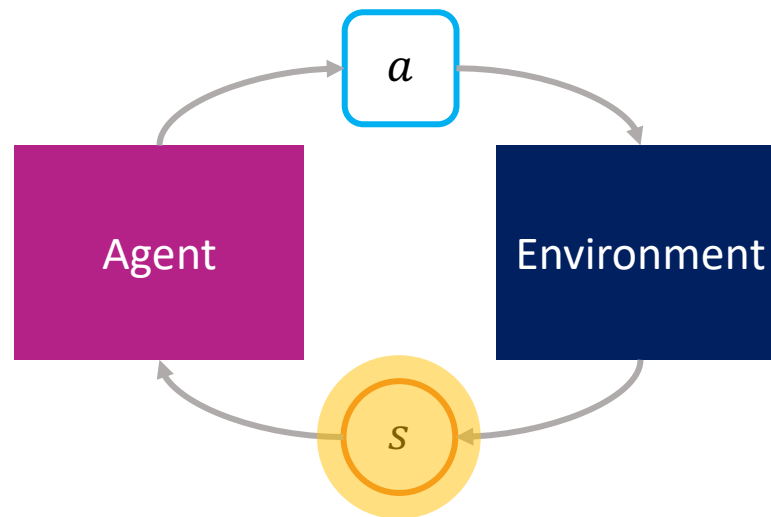
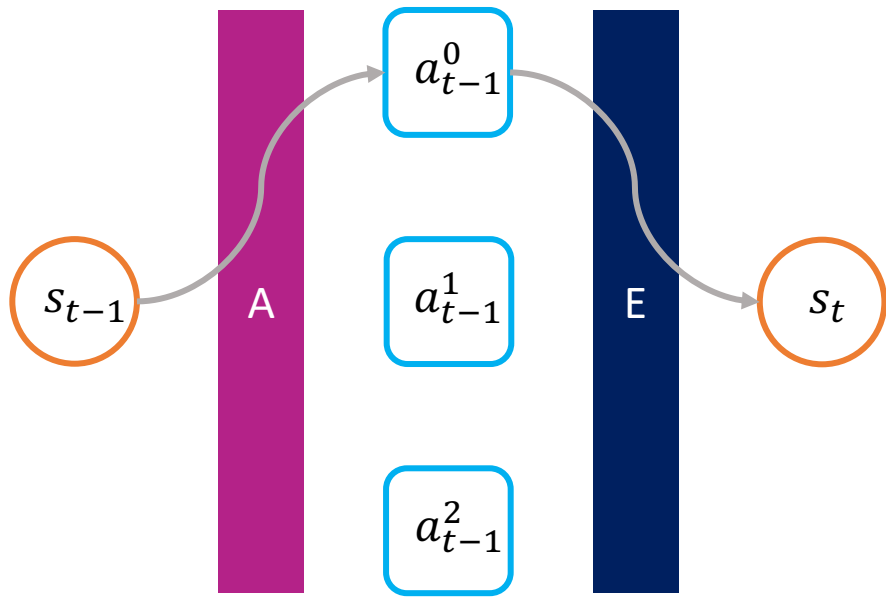
s_{t-1}

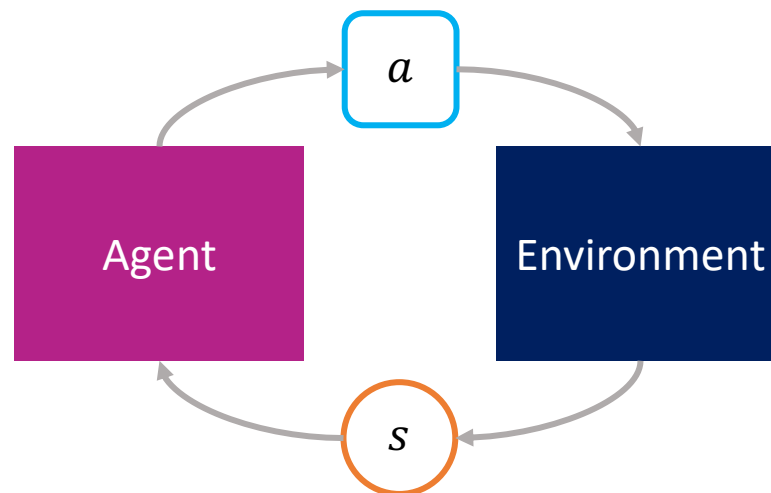
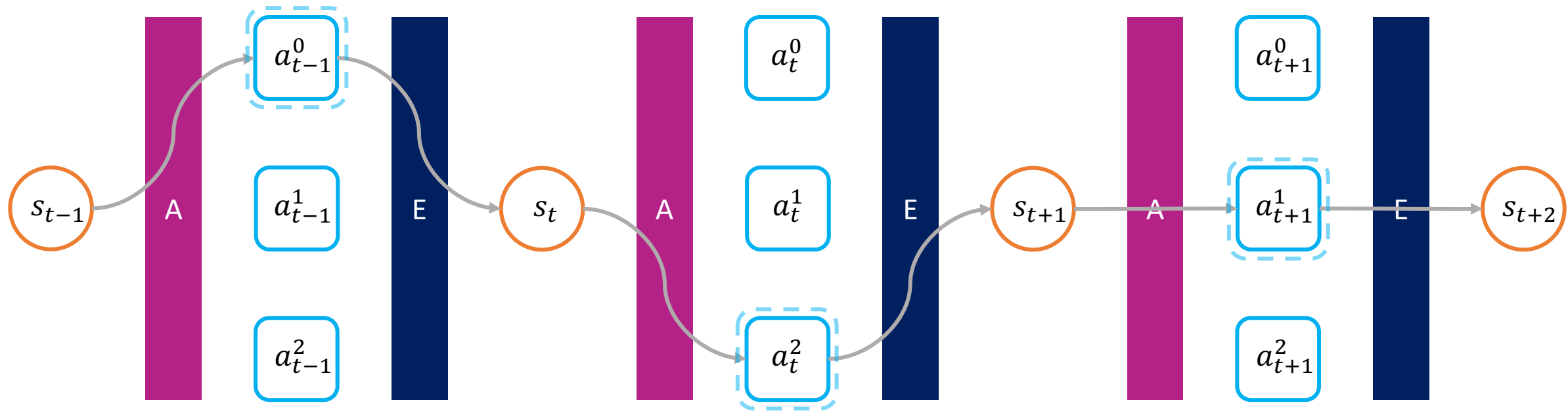


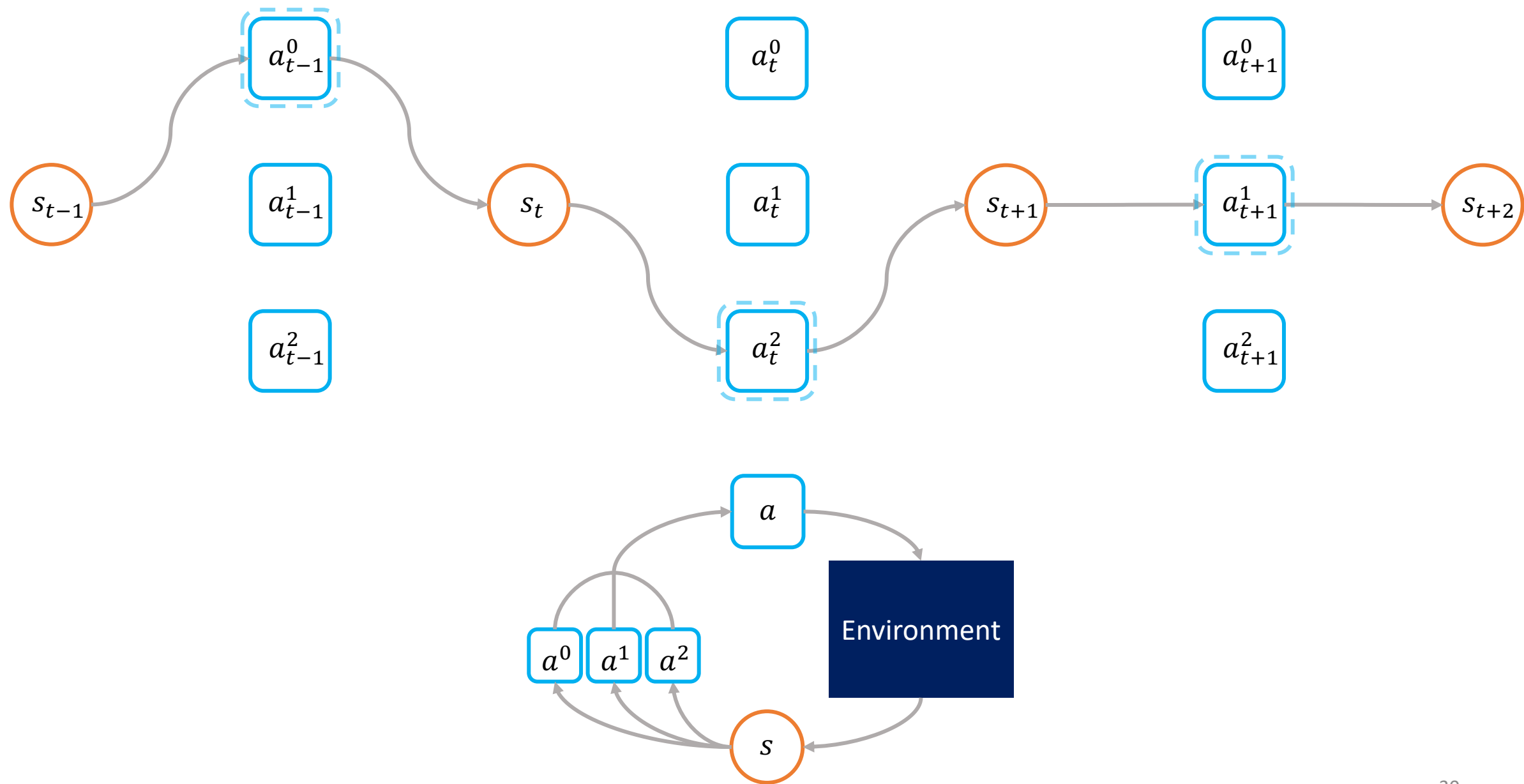


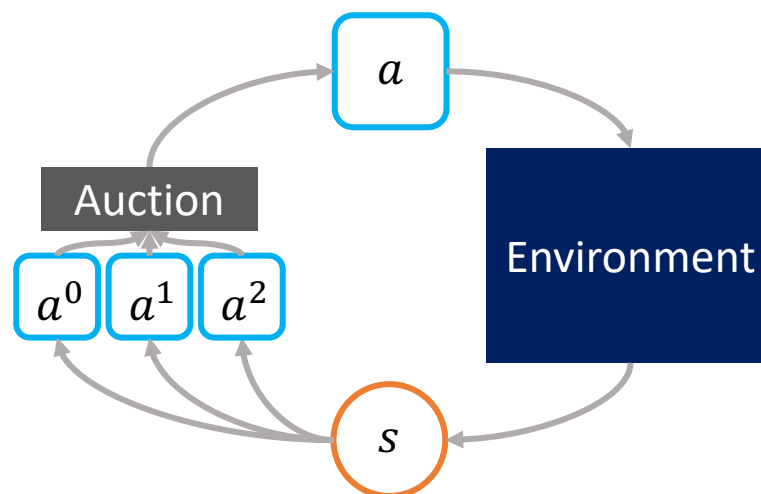
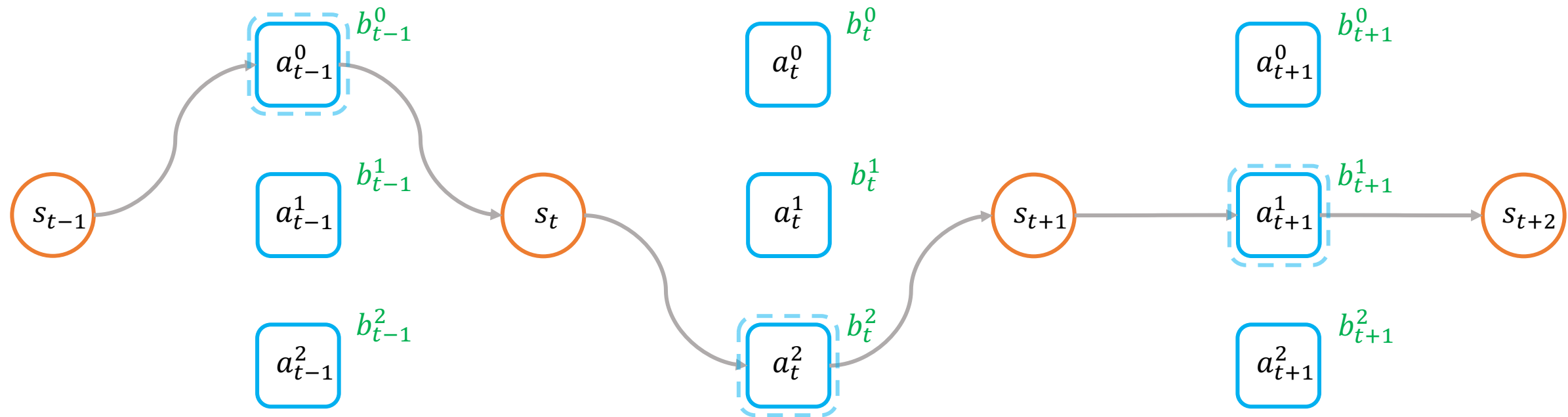


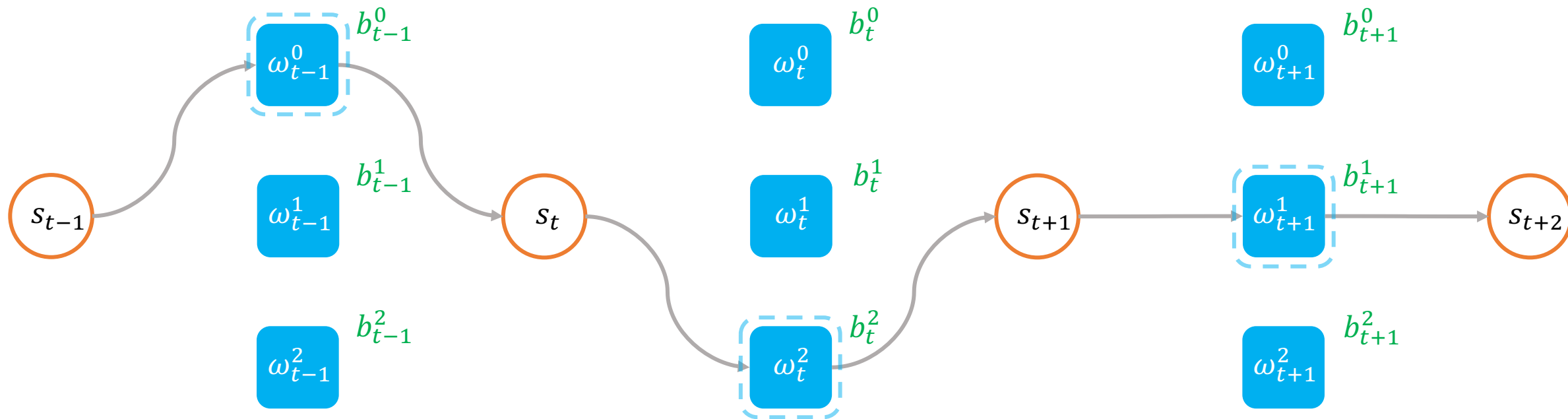








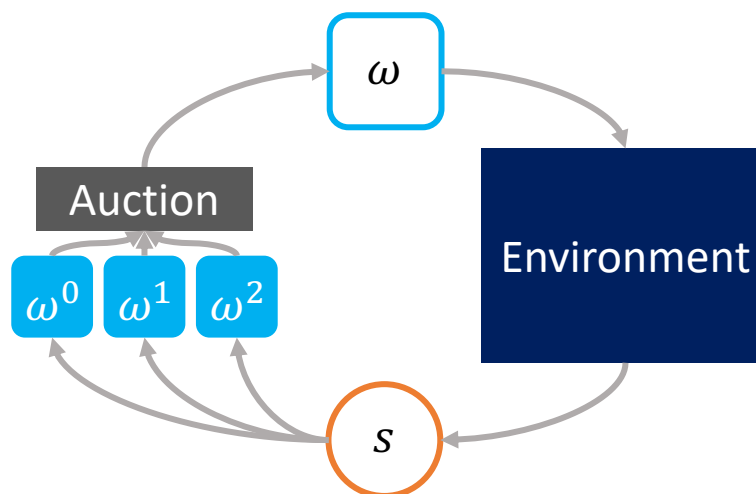


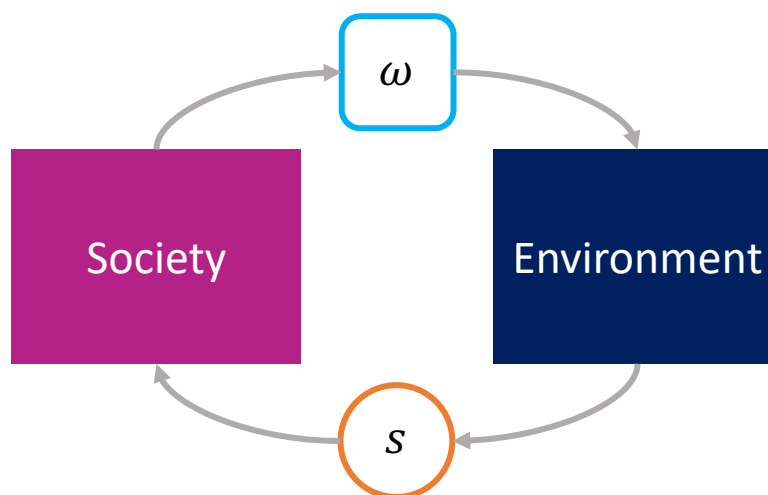
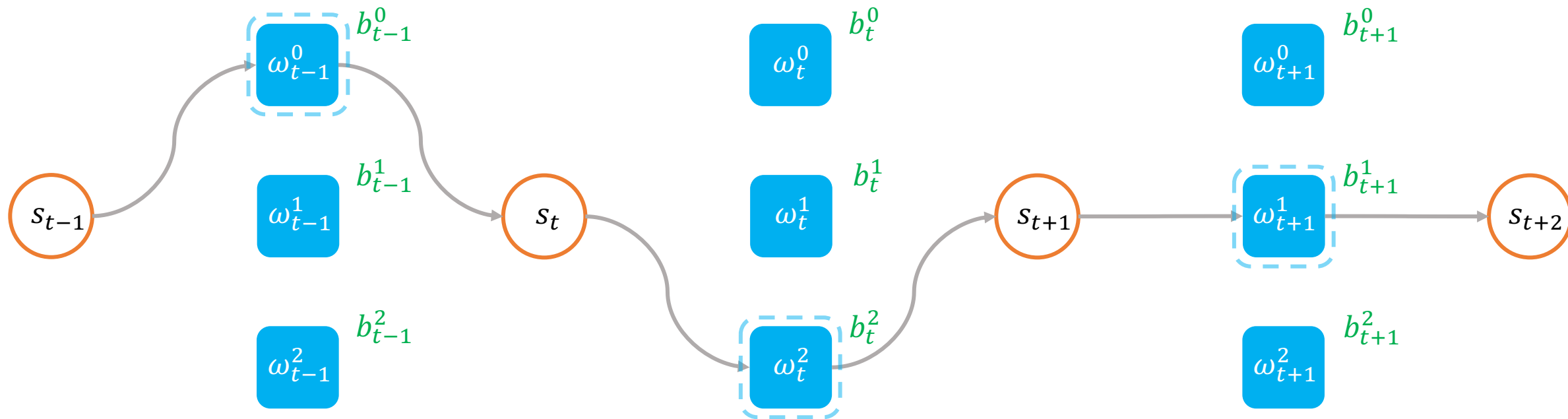


Local Auction

Action Space: bids b

Objective: optimize utility in auction

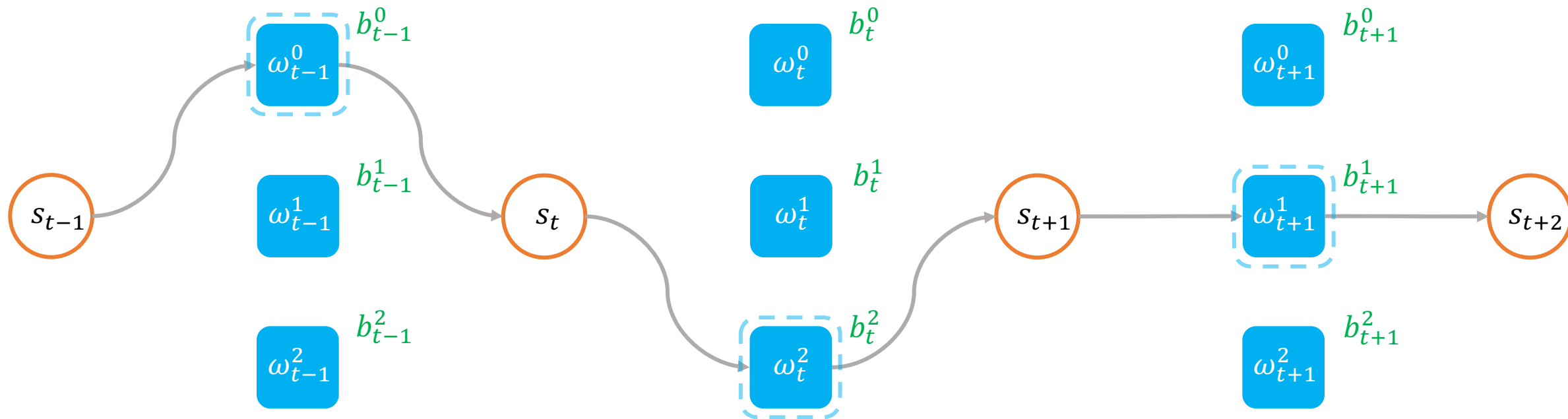




Global MDP

Action Space: agents ω

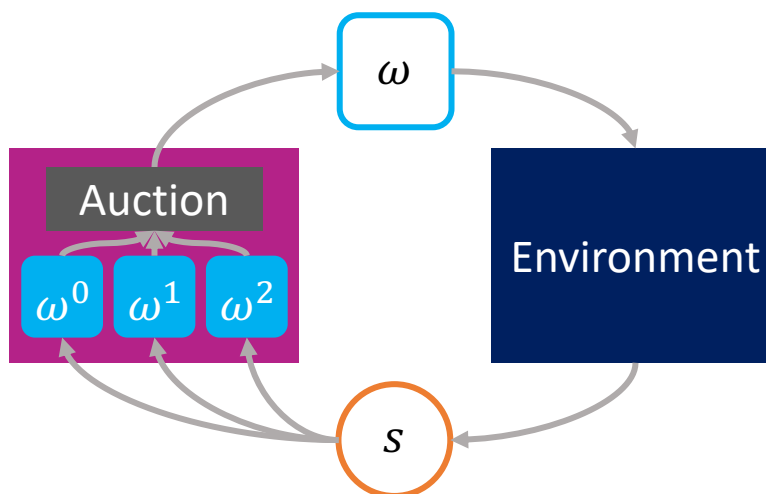
Objective: optimize return in environment



Local Auction

Action Space: bids b

Objective: optimize utility in auction



Global MDP

Action Space: agents ω

Objective: optimize return in environment

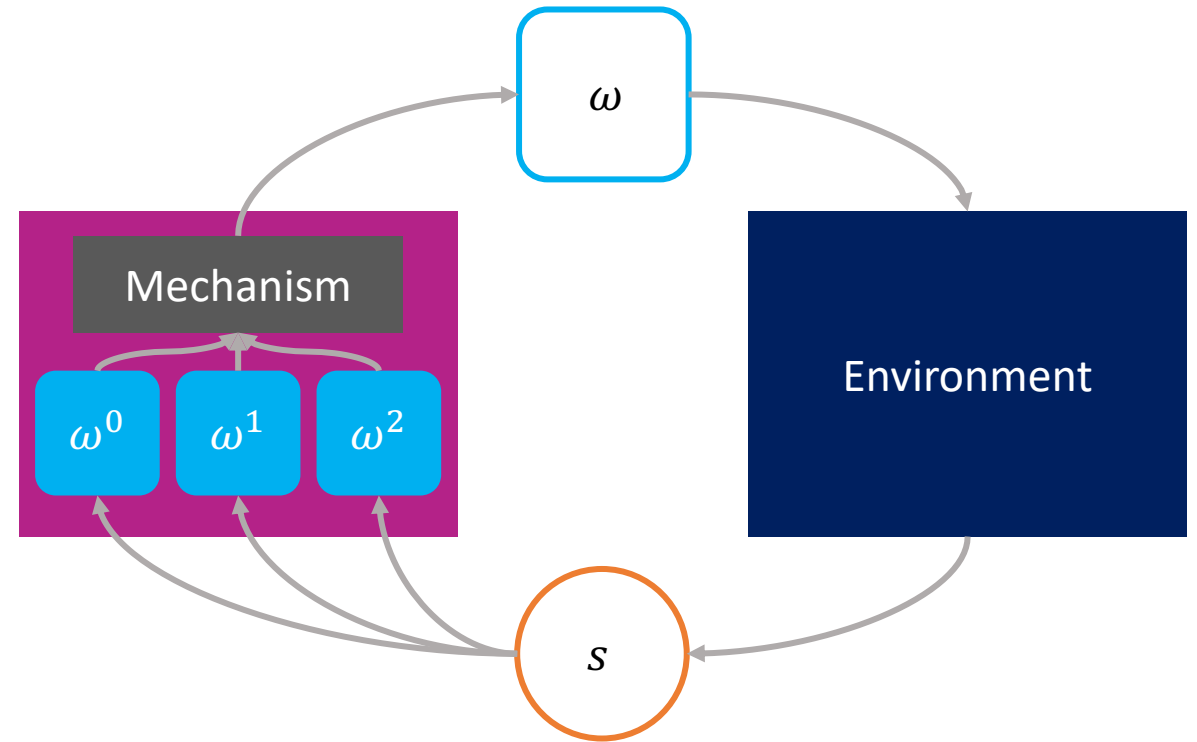
This Paper: Contributions

Assumptions

- Sequential decision making setting
- Each agent produces a specialized transformation to the state (e.g. a literal action)
- Only one agent activates at each time step

Main Contribution

We show that the Vickrey Auction can be adapted to MDPs such that the solution of the global societal objective emerges as a Nash equilibrium strategy profile of the local agents



This Paper: Contributions

Assumptions

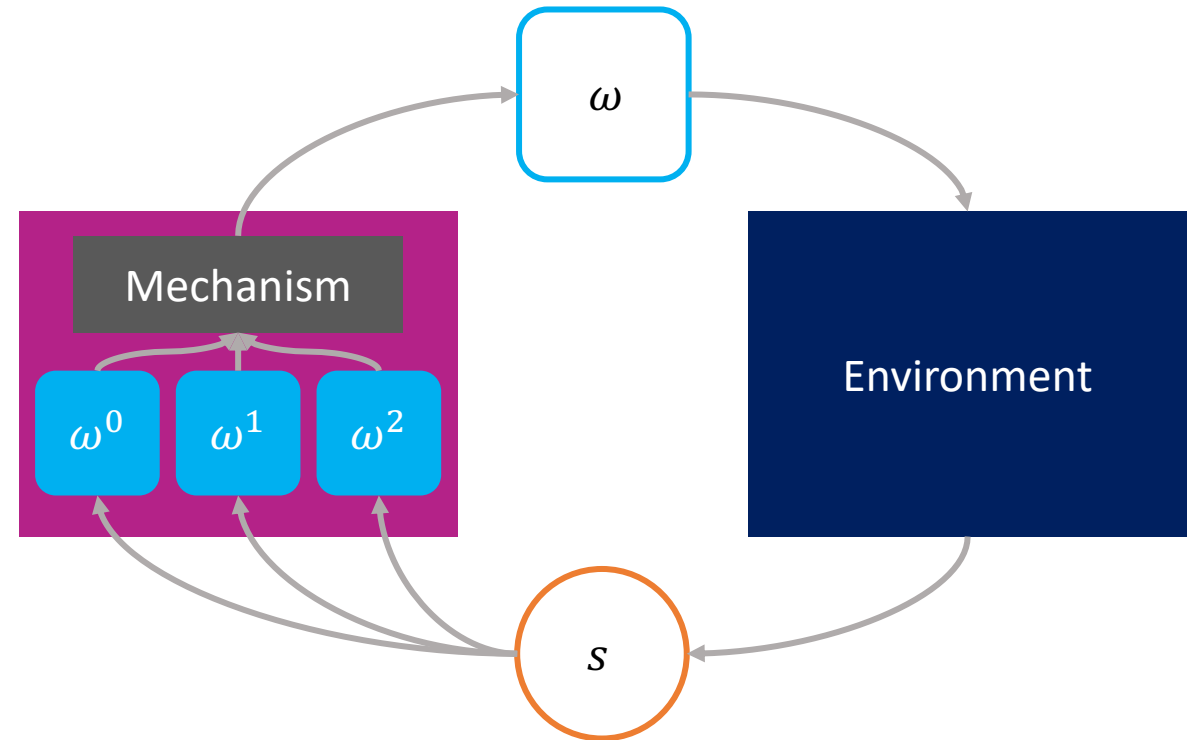
- Sequential decision making setting
- Each agent produces a specialized transformation to the state (e.g. a literal action)
- Only one agent activates at each time step

Main Contribution

We show that the Vickrey Auction can be adapted to MDPs such that the solution of the global societal objective emerges as a Nash equilibrium strategy profile of the local agents

Implication: Bridging Two Levels of Abstraction

- A recipe for translating a global objective of a society into local learning problems for the agents



This Paper: Contributions

Assumptions

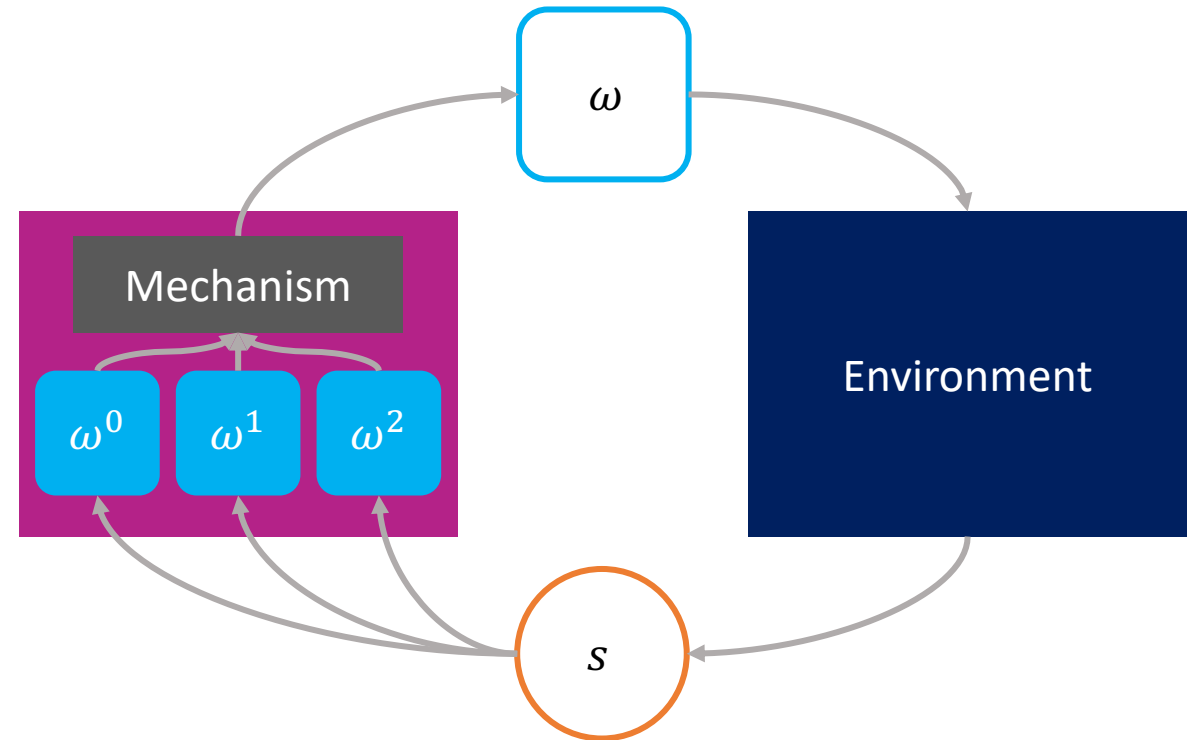
- Sequential decision making setting
- Each agent produces a specialized transformation to the state (e.g. a literal action)
- Only one agent activates at each time step

Main Contribution

We show that the Vickrey Auction can be adapted to MDPs such that the solution of the global societal objective emerges as a Nash equilibrium strategy profile of the local agents

Implication: Bridging Two Levels of Abstraction

- A recipe for translating a global objective of a society into local learning problems for the agents
- A decentralized reinforcement learning algorithm with credit assignment local in space and time



Roadmap

Question

Key Idea

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

For what auction mechanism would these optimal bids be an equilibrium strategy?

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

For what auction mechanism would these optimal bids be an equilibrium strategy?

How can we adapt this auction mechanism for discrete-action MDPs?

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

For what auction mechanism would these optimal bids be an equilibrium strategy?

How can we adapt this auction mechanism for discrete-action MDPs?

How can we avoid suboptimal equilibria?

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

For what auction mechanism would these optimal bids be an equilibrium strategy?

How can we adapt this auction mechanism for discrete-action MDPs?

How can we avoid suboptimal equilibria?

How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

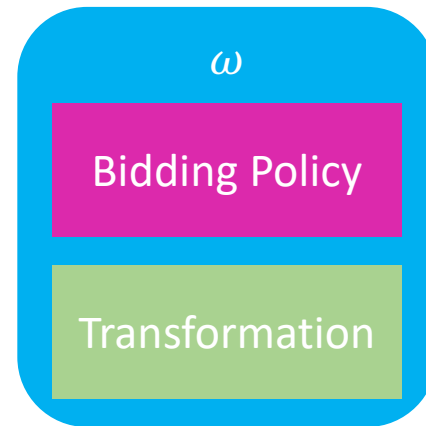
For what auction mechanism would these optimal bids be an equilibrium strategy?

How can we adapt this auction mechanism for discrete-action MDPs?

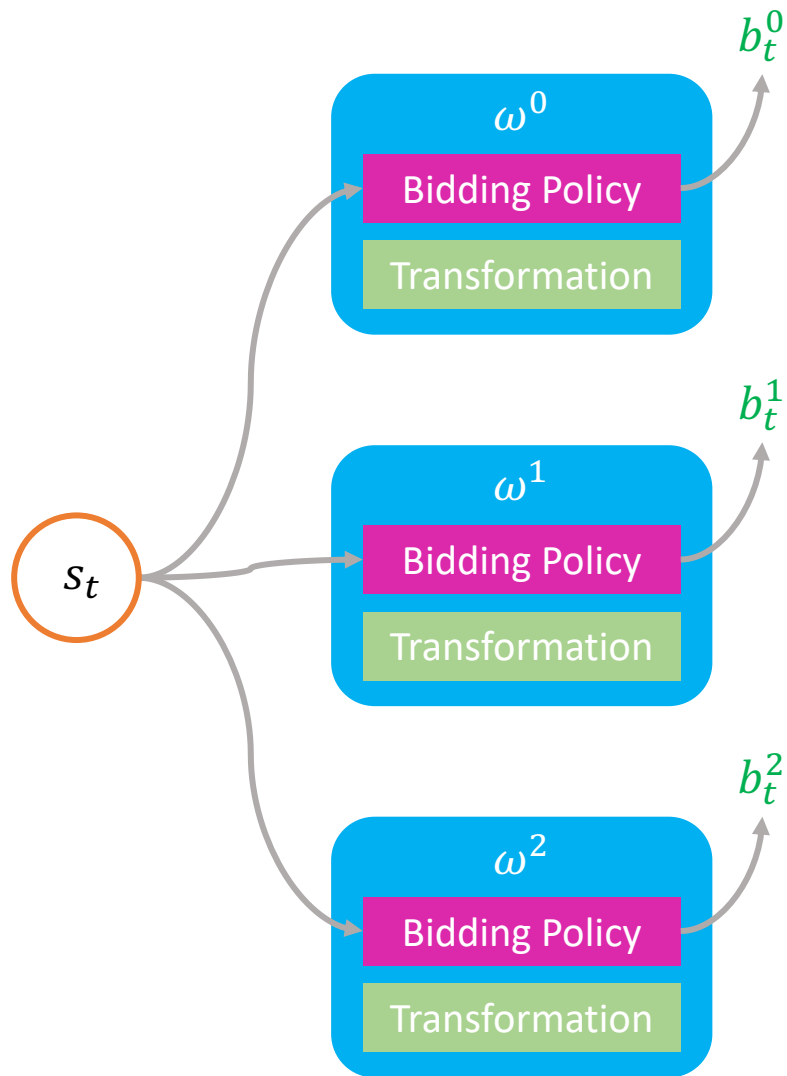
How can we avoid suboptimal equilibria?

How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?

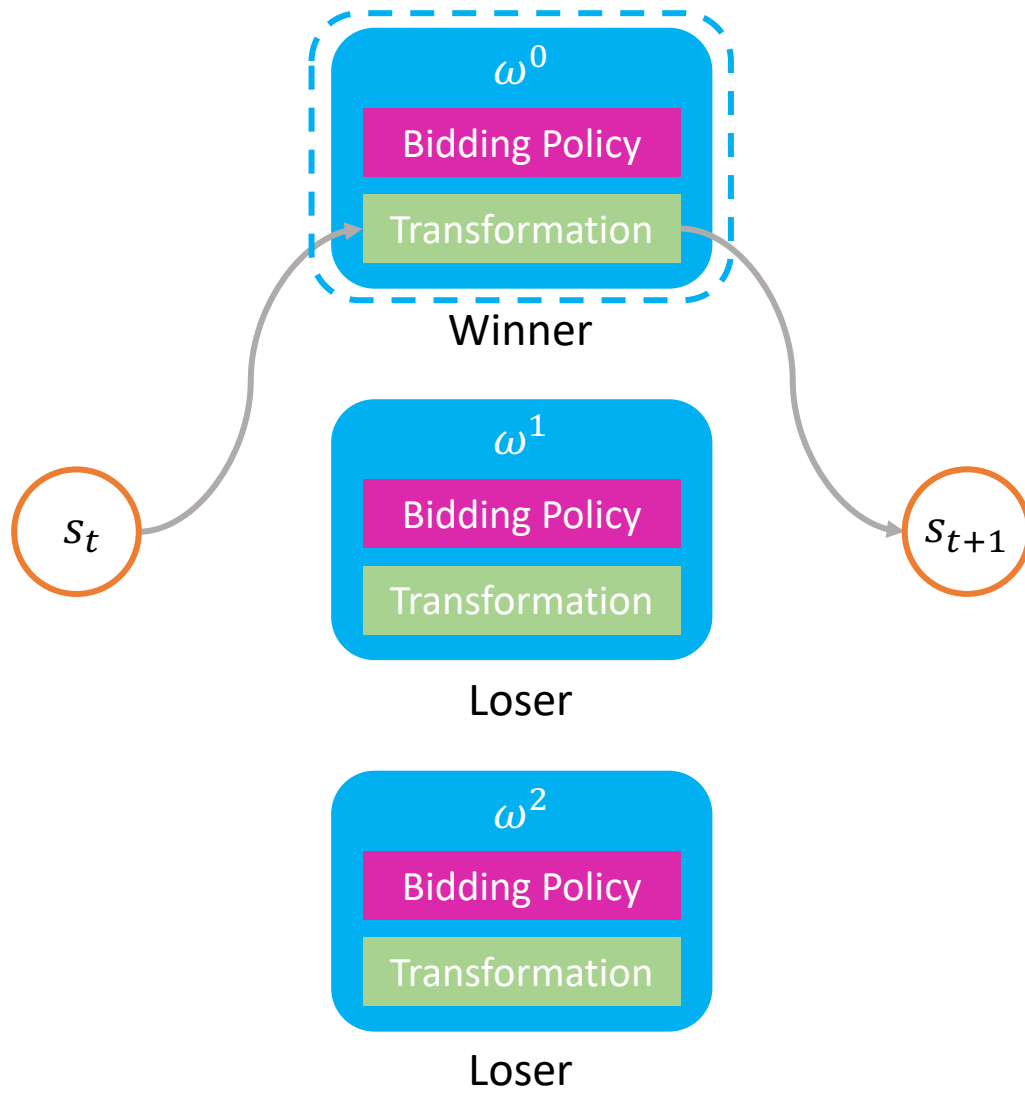
Architecture of an Agent



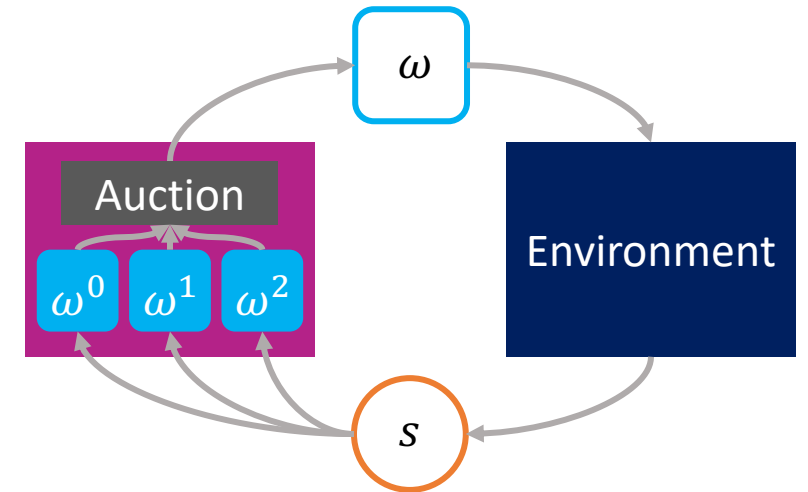
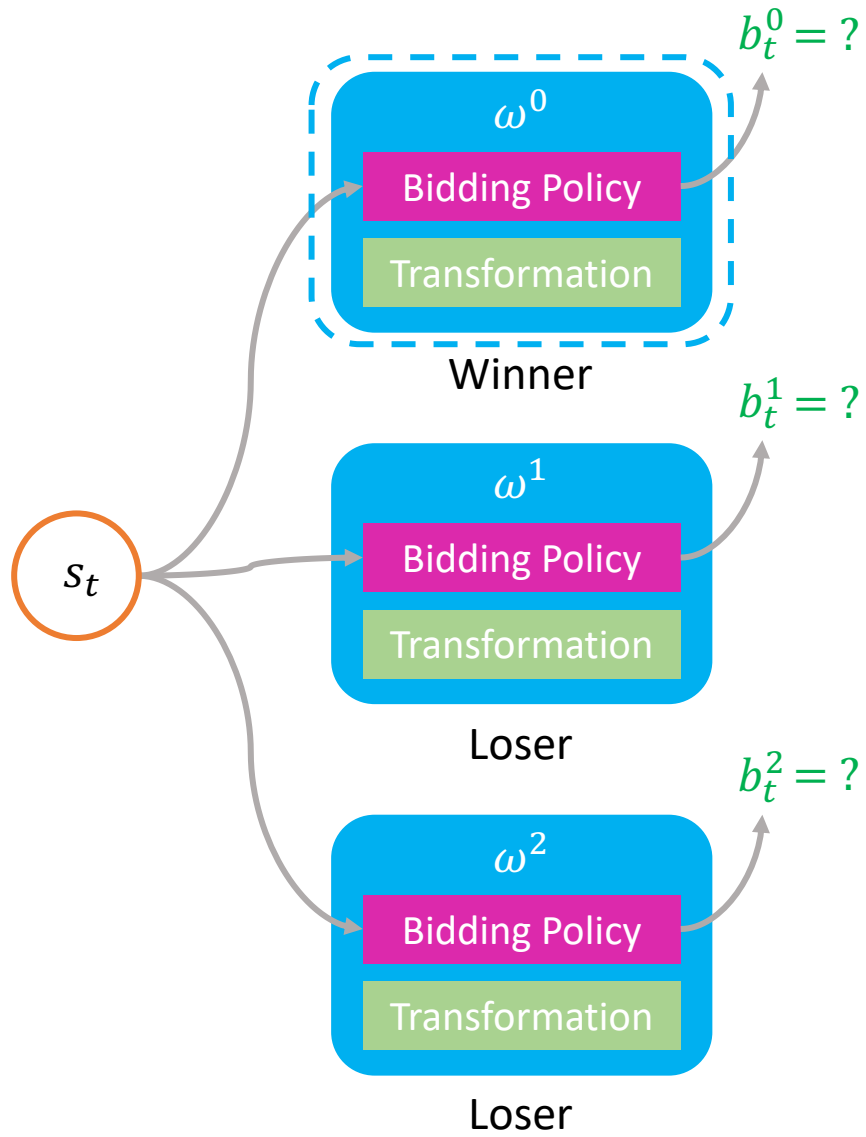
Activating Agents via Auction



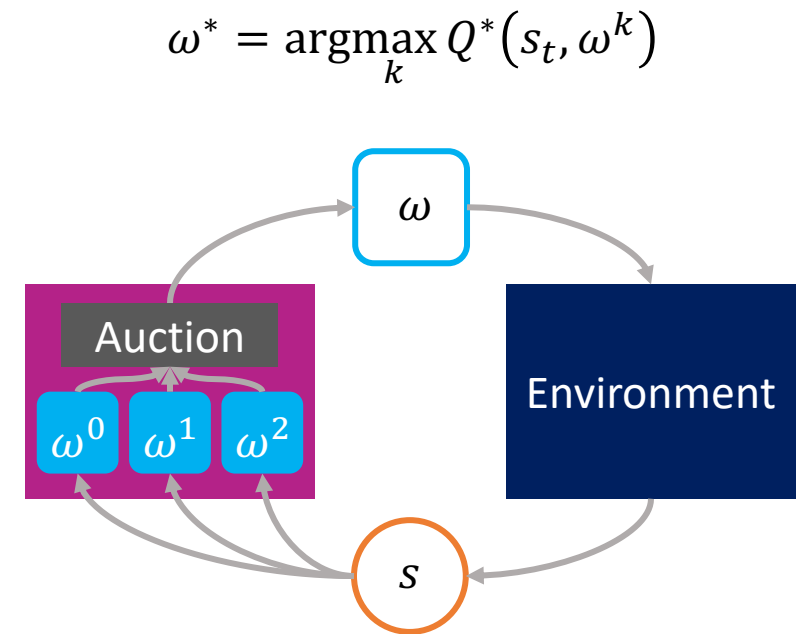
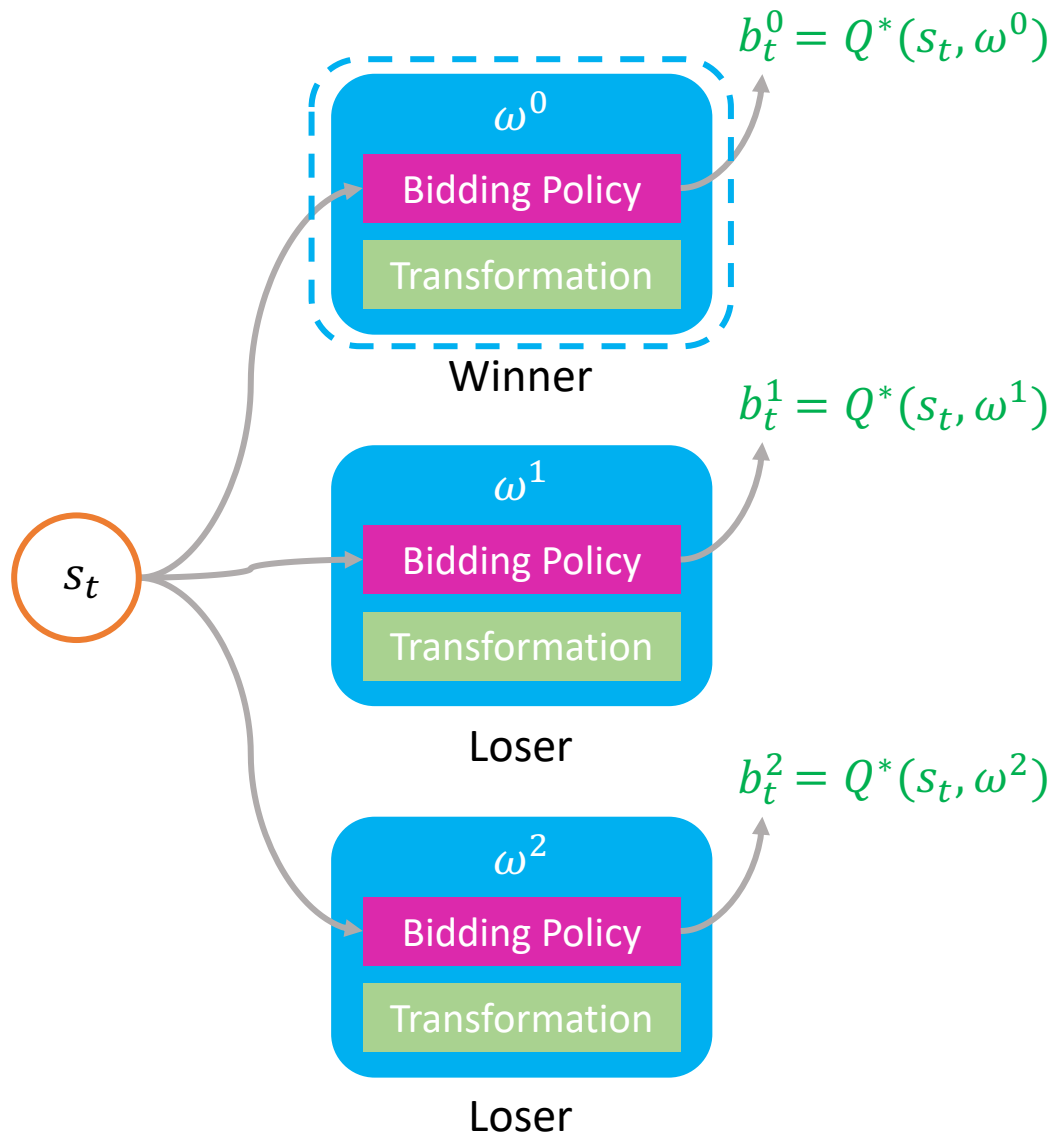
Transforming the State



What should the optimal bids be?



Key Idea: the optimal bid is your optimal Q value



Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

For what auction mechanism would these optimal bids be an equilibrium strategy?

How can we adapt this auction mechanism for discrete-action MDPs?

How can we avoid suboptimal equilibria?

How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

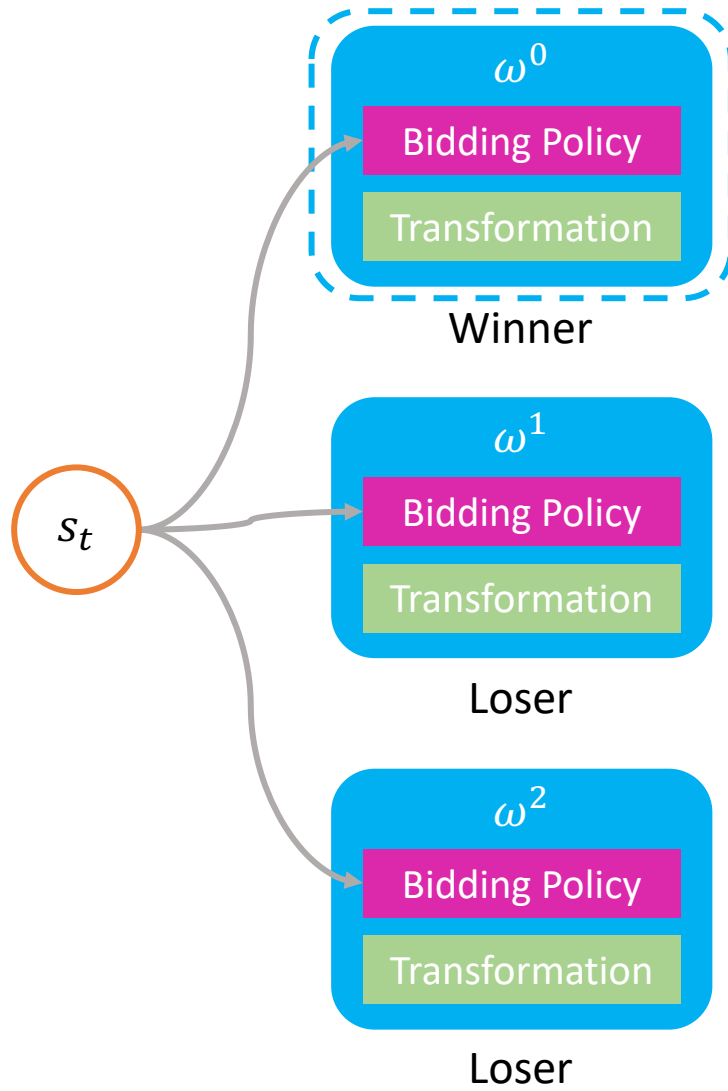
For what auction mechanism would these optimal bids be an equilibrium strategy?

How can we adapt this auction mechanism for discrete-action MDPs?

How can we avoid suboptimal equilibria?

How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?

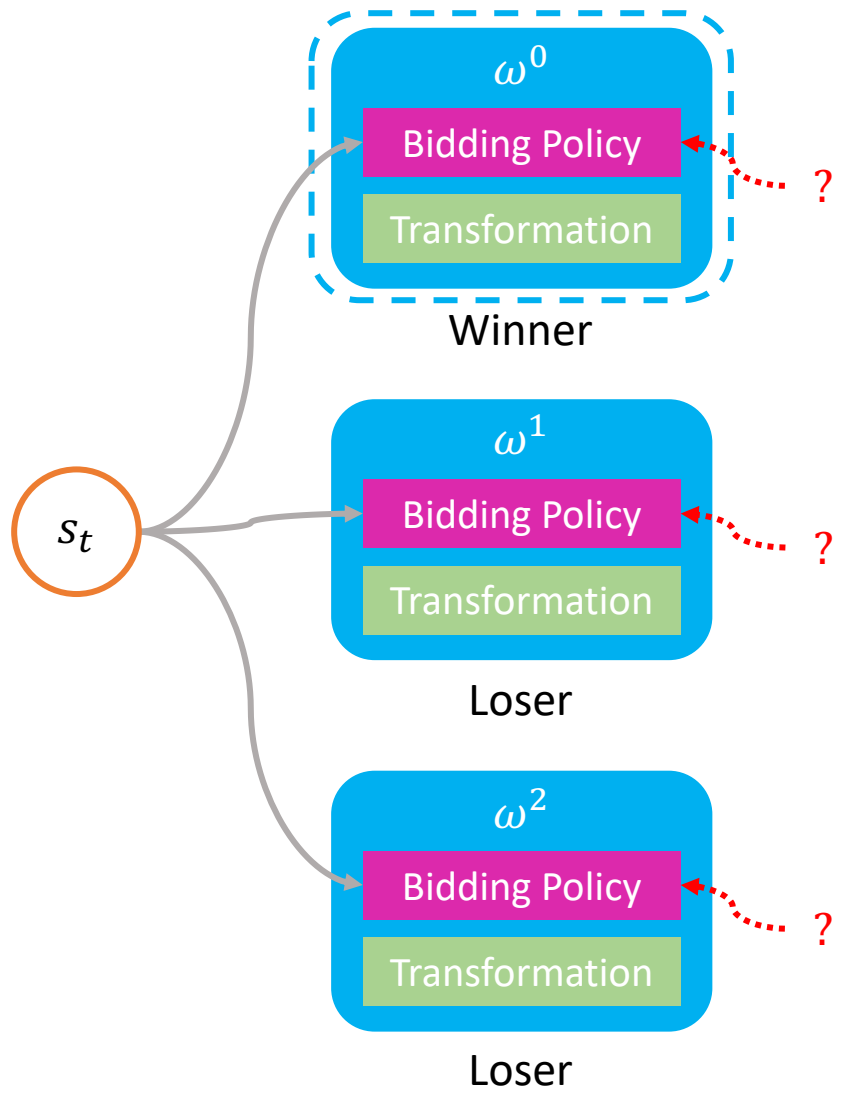
What Should the Auction Mechanism be?



Assume

Each agent ω^k has a valuation $v^k(s_t)$ for state s_t

What Should the Auction Mechanism be?



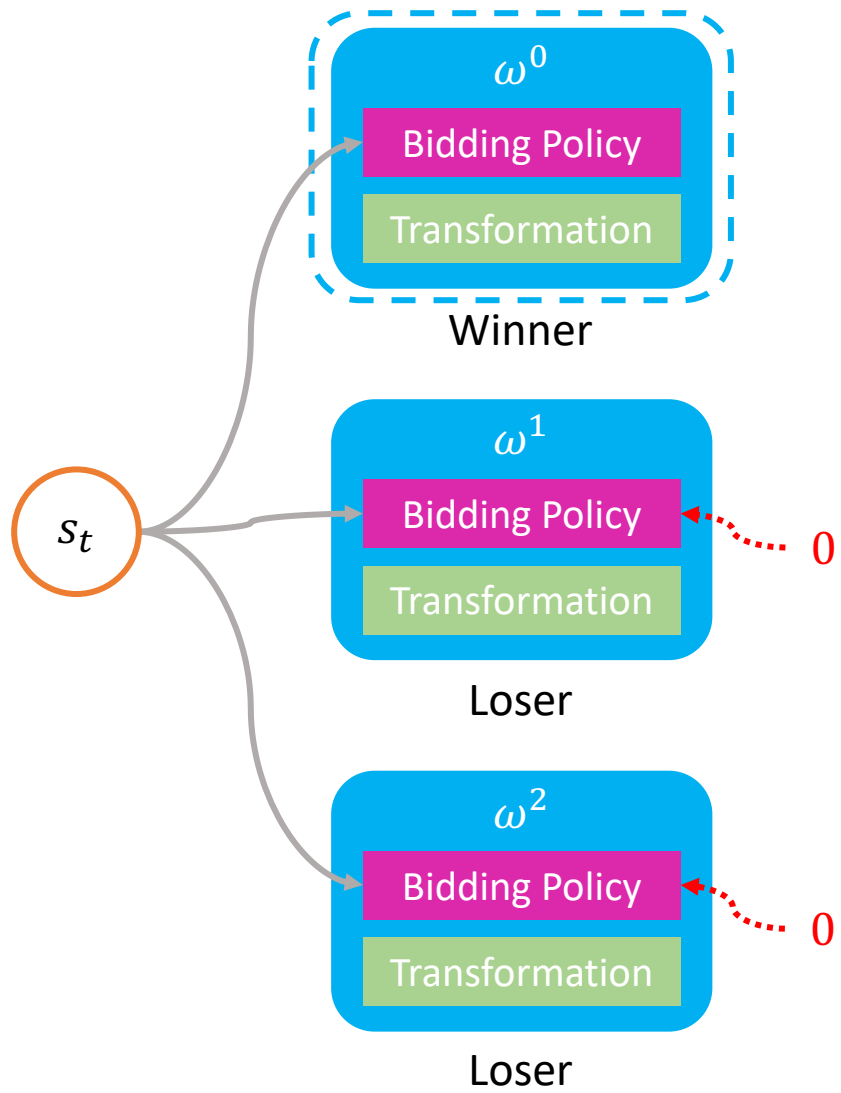
Assume

Each agent ω^k has a valuation $v^k(s_t)$ for state s_t

Question

What should the agents' utilities be?

What Should the Auction Mechanism be?



Assume

Each agent ω^k has a valuation $v^k(s_t)$ for state s_t

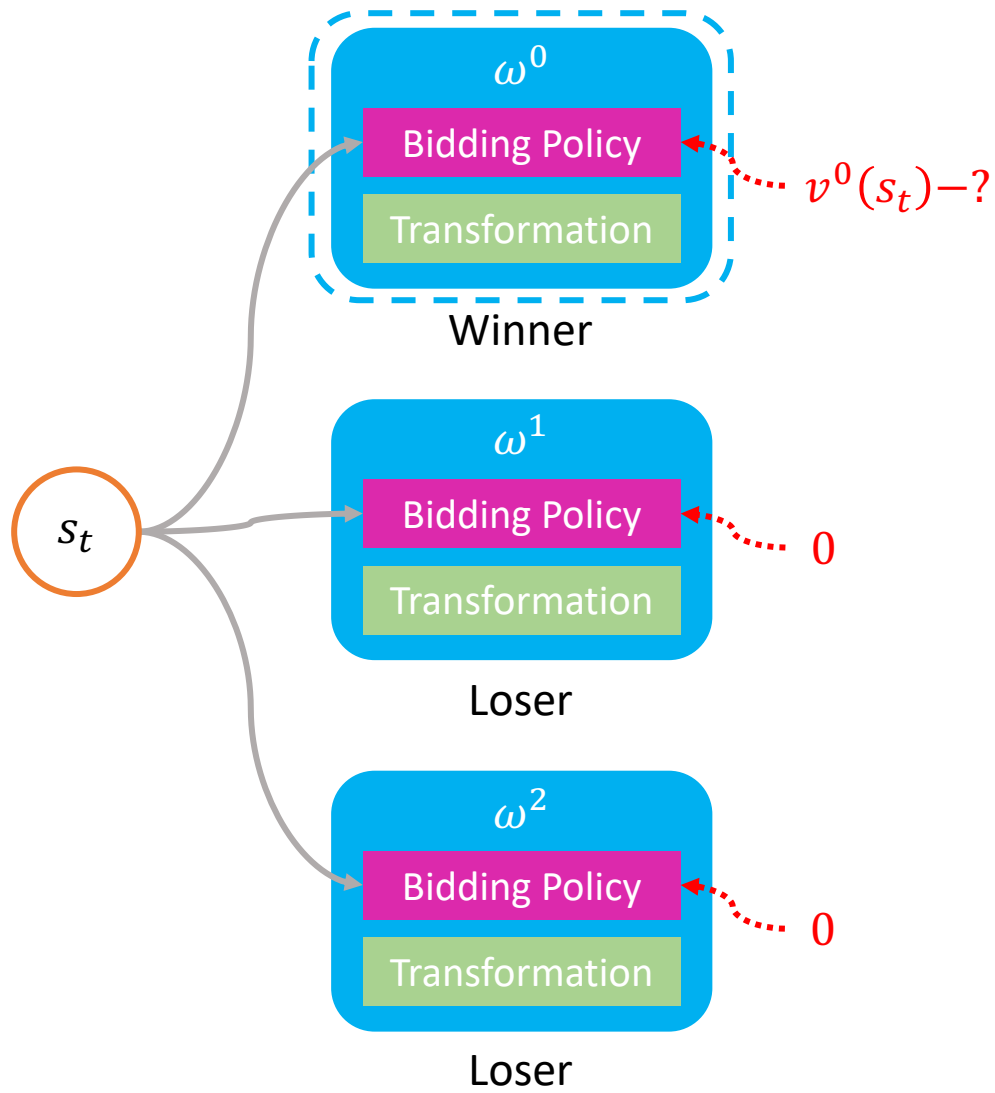
Question

What should the agents' utilities be?

Utilities?

Losers: $u^i(b) = 0$

What Should the Auction Mechanism be?



Assume

Each agent ω^k has a valuation $v^k(s_t)$ for state s_t

Question

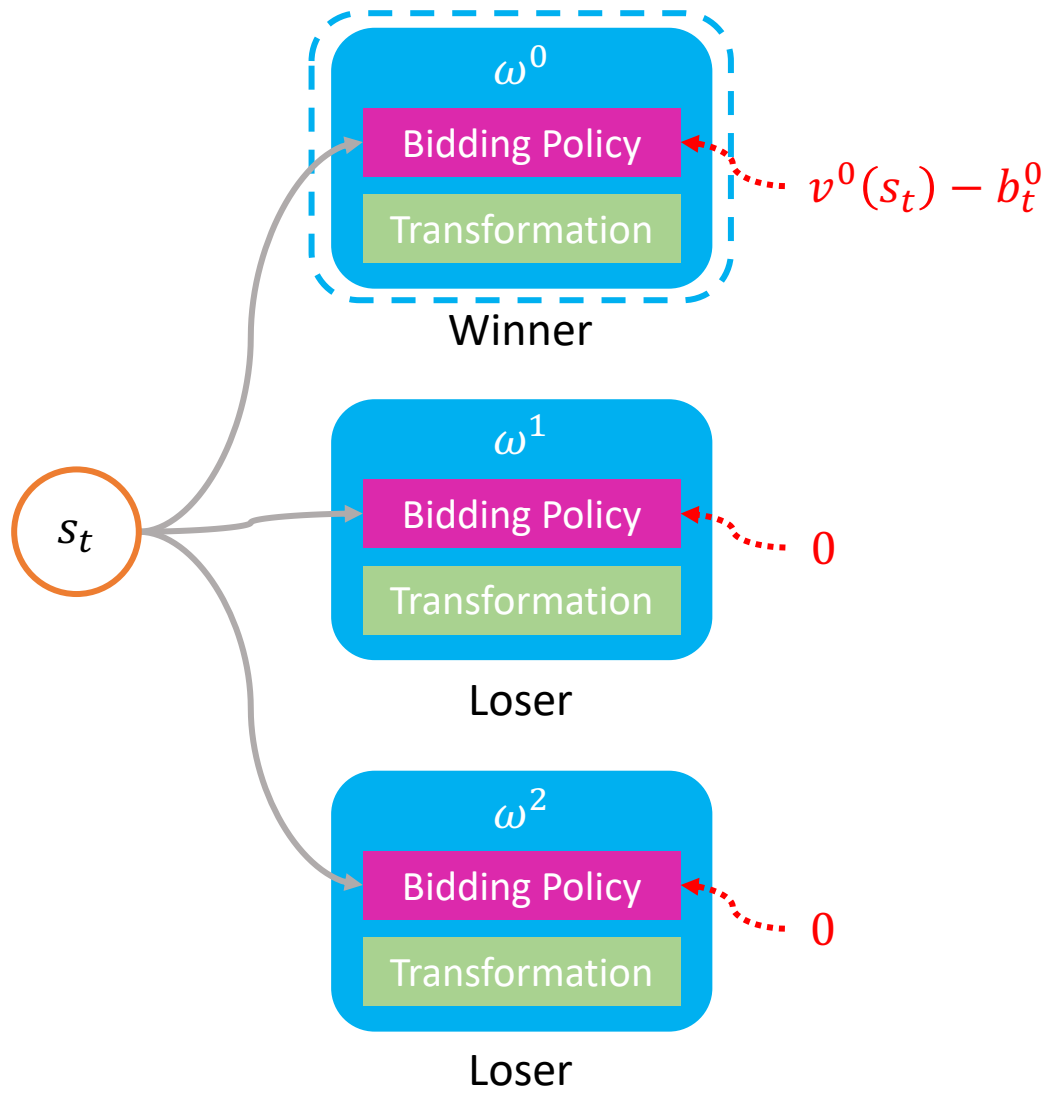
What should the agents' utilities be?

Utilities?

Losers: $u^i(b) = 0$

Winner: $u^i(b) = v^i - ?$

What Should the Auction Mechanism be?



Assume

Each agent ω^k has a valuation $v^k(s_t)$ for state s_t

Question

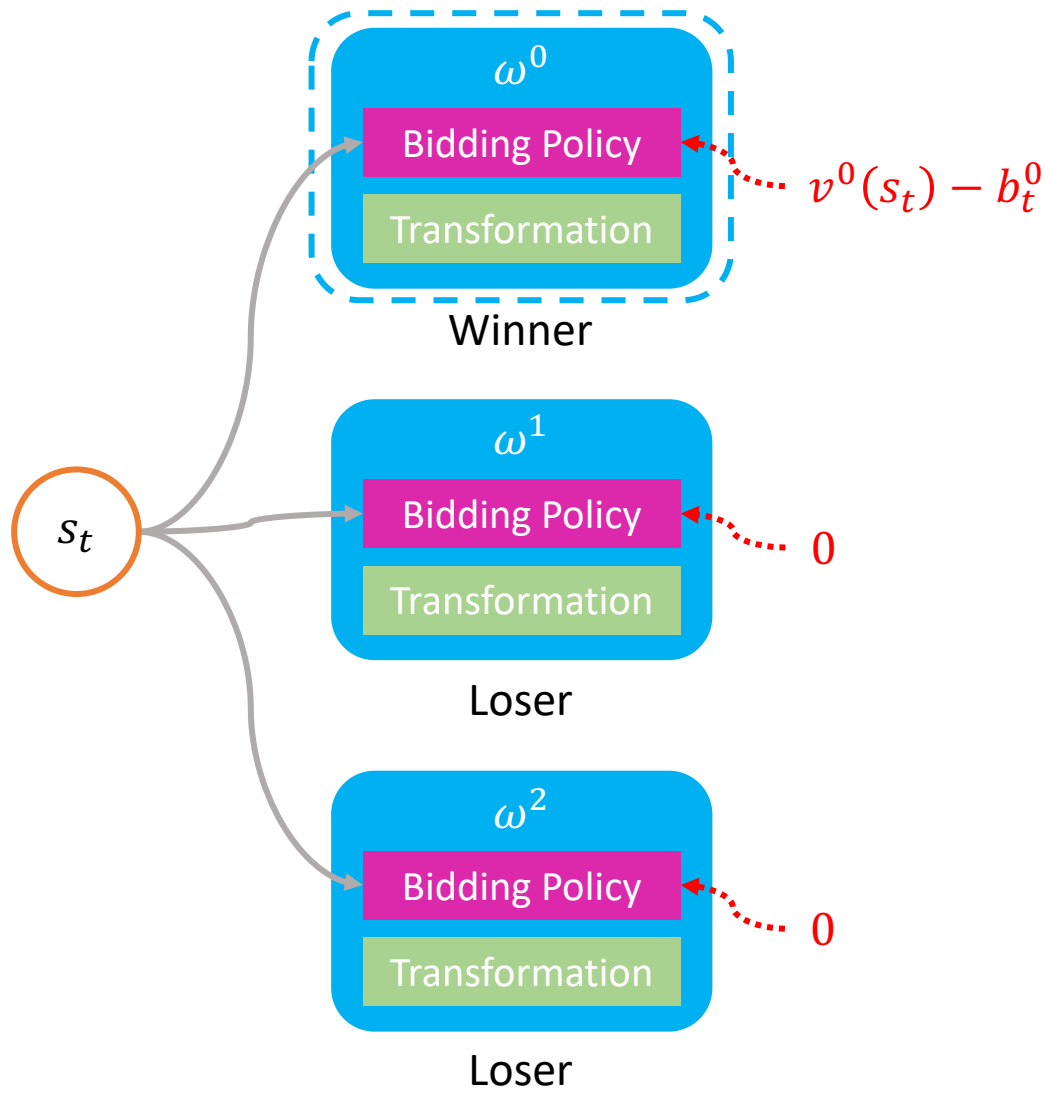
What should the agents' utilities be?

First Price Sealed-Bid Auction Utilities?

Losers: $u^i(b) = 0$

Winner: $u^i(b) = v^i - b$

What Should the Auction Mechanism be?



Assume

Each agent ω^k has a valuation $v^k(s_t)$ for state s_t

Question

What should the agents' utilities be?

First Price Sealed-Bid Auction Utilities?

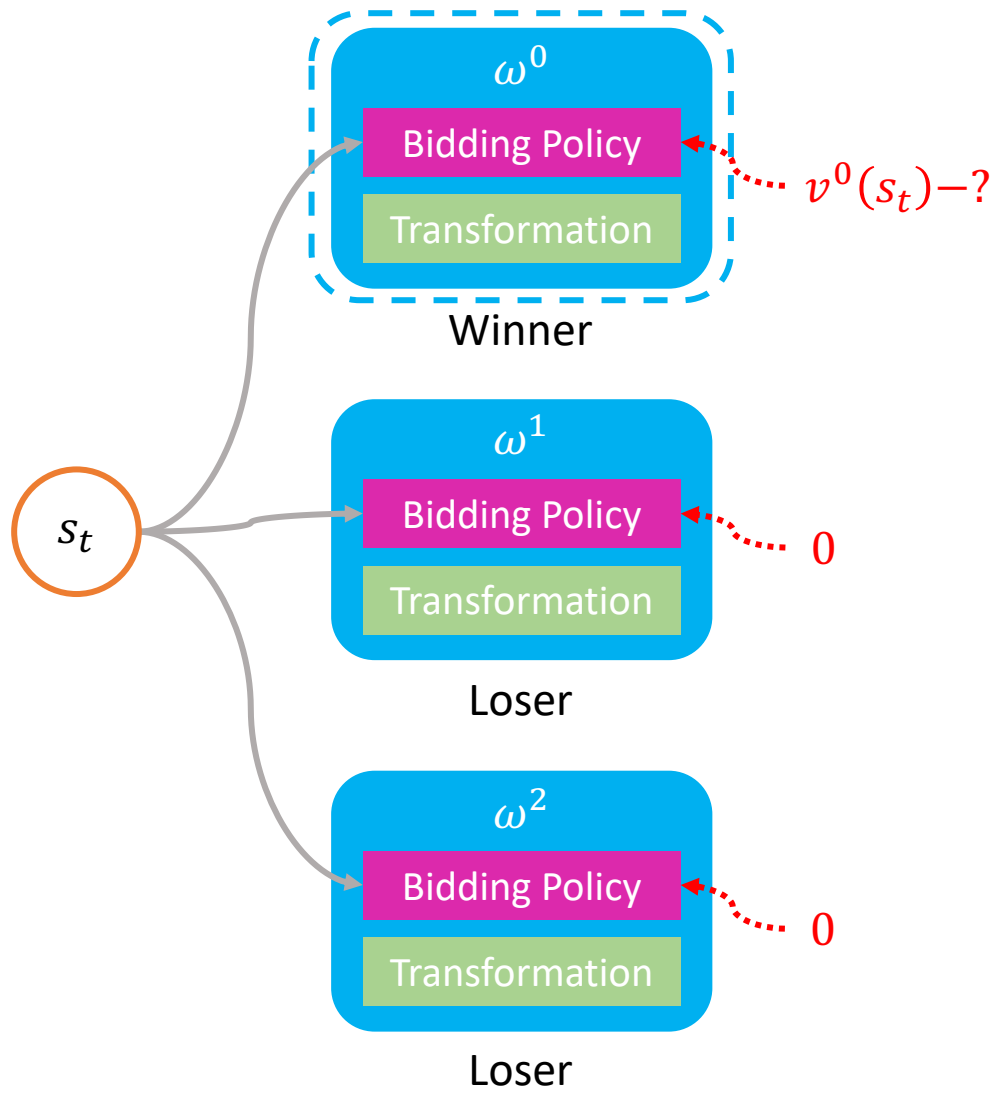
Losers: $u^i(b) = 0$

Winner: $u^i(b) = v^i - b$

Problem with First Price Sealed-Bid Auctions

There is no dominant strategy – the bid that optimizes an agent's utility depends on what other agents bid

What Should the Auction Mechanism be?



Assume

Each agent ω^k has a valuation $v^k(s_t)$ for state s_t

Question

What should the agents' utilities be?

Utilities?

Losers: $u^i(b) = 0$

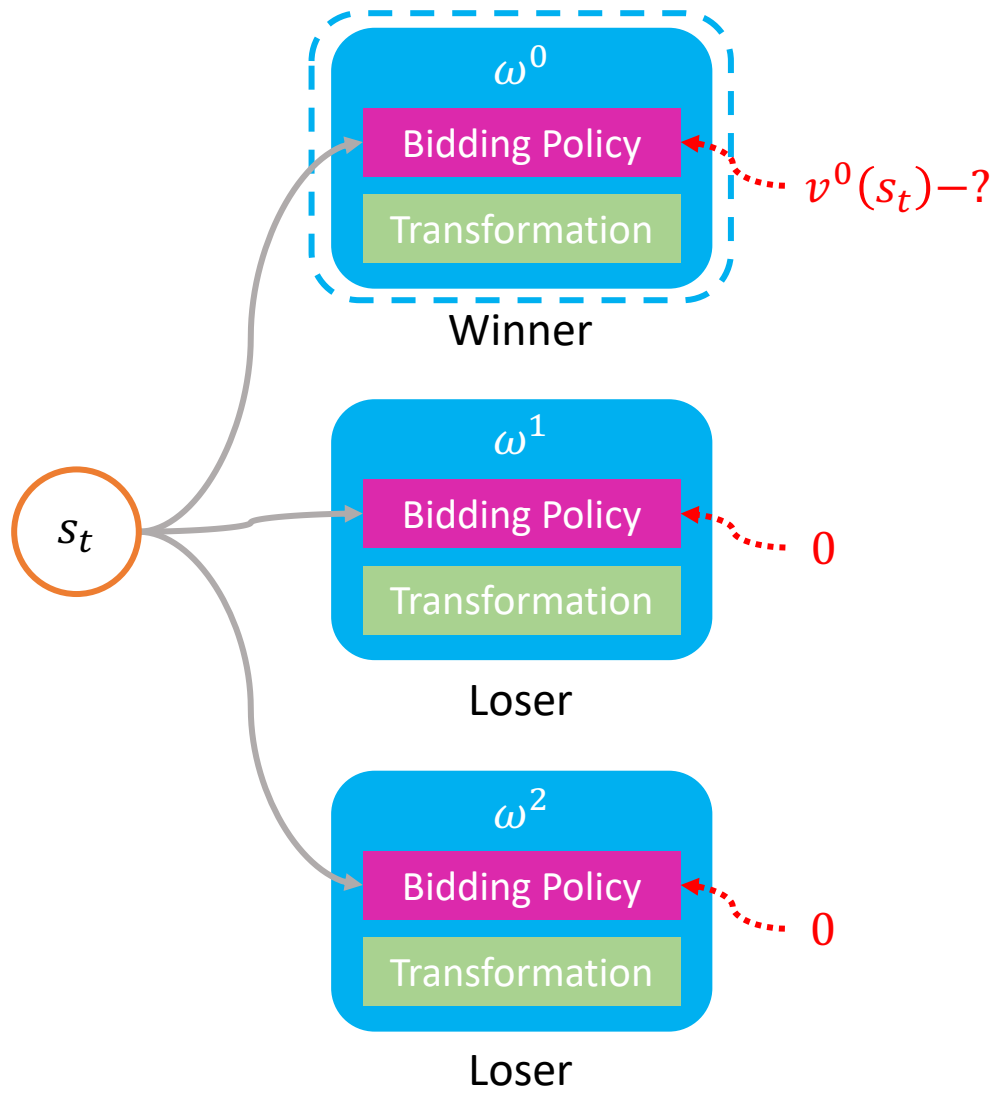
Winner: $u^i(b) = v^i - ?$

Want: Dominant Strategy Incentive Compatibility

The optimal strategy is to truthfully bid its own valuation:

$$b^i \leftarrow v^i$$

What Should the Auction Mechanism be?



Assume

Each agent ω^k has a valuation $v^k(s_t)$ for state s_t

Question

What should the agents' utilities be?

Utilities?

Losers: $u^i(b) = 0$

Winner: $u^i(b) = v^i - ?$

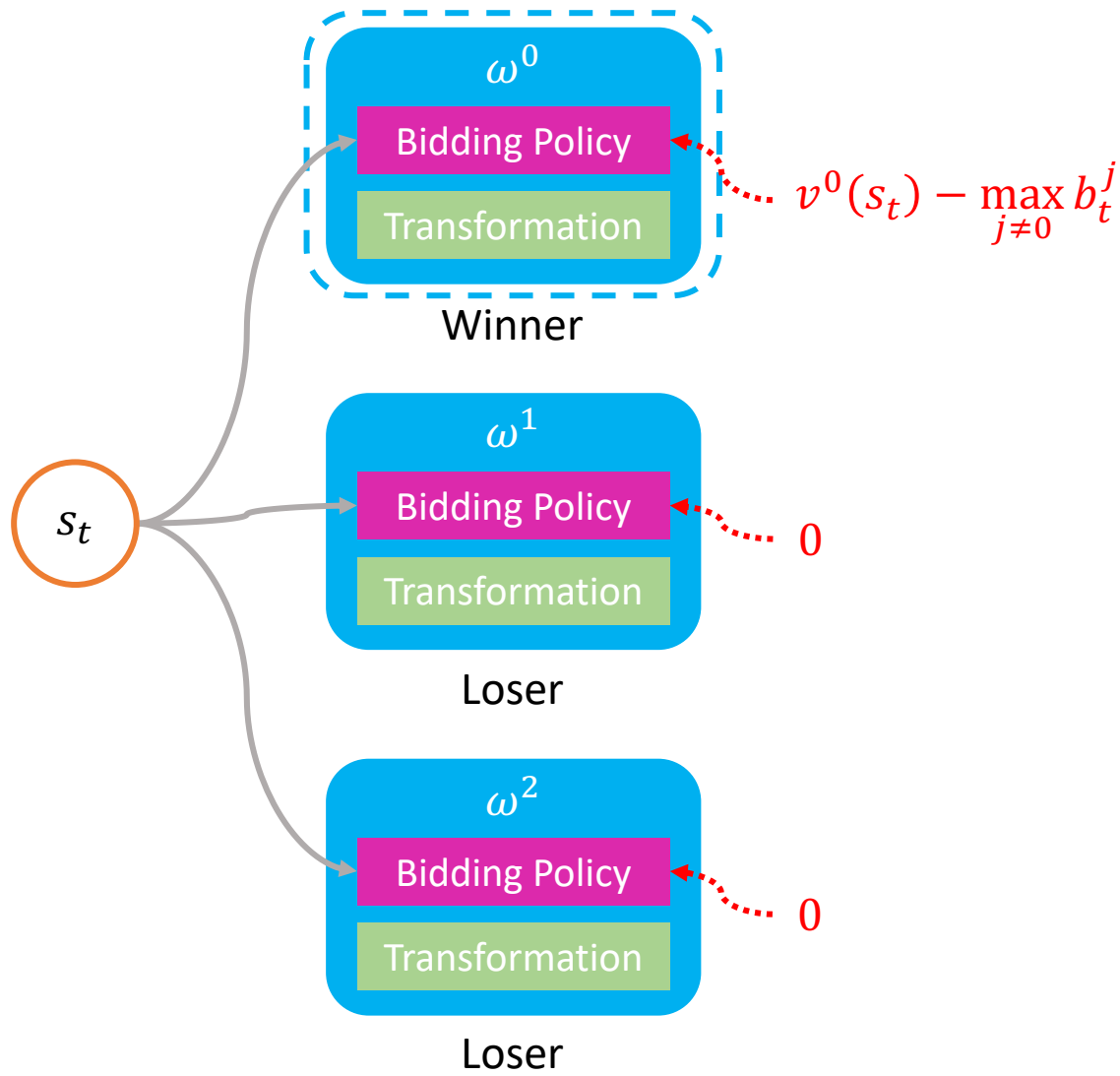
Want: Dominant Strategy Incentive Compatibility

The optimal strategy is to truthfully bid its own valuation:

$$b^i \leftarrow v^i$$

Implication: Set $v^k(s_t) = Q^*(s_t, \omega^k)$!

Vickrey Auction



Assume

Each agent ω^k has a valuation $v^k(s_t)$ for state s_t

Question

What should the agents' utilities be?

Vickrey Auction Utilities!

Losers: $u^i(b) = 0$

Winner: $u^i(b) = v^i - \max_{j \neq i} b^j$

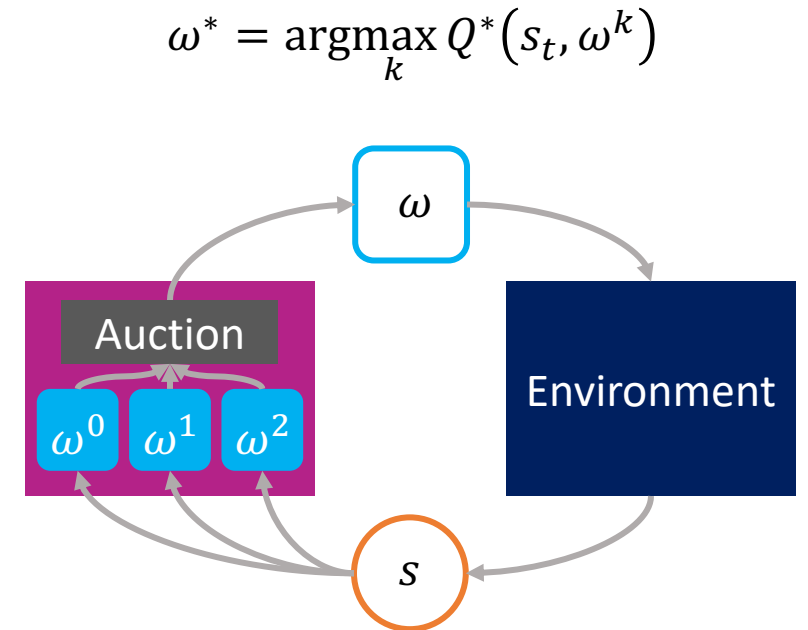
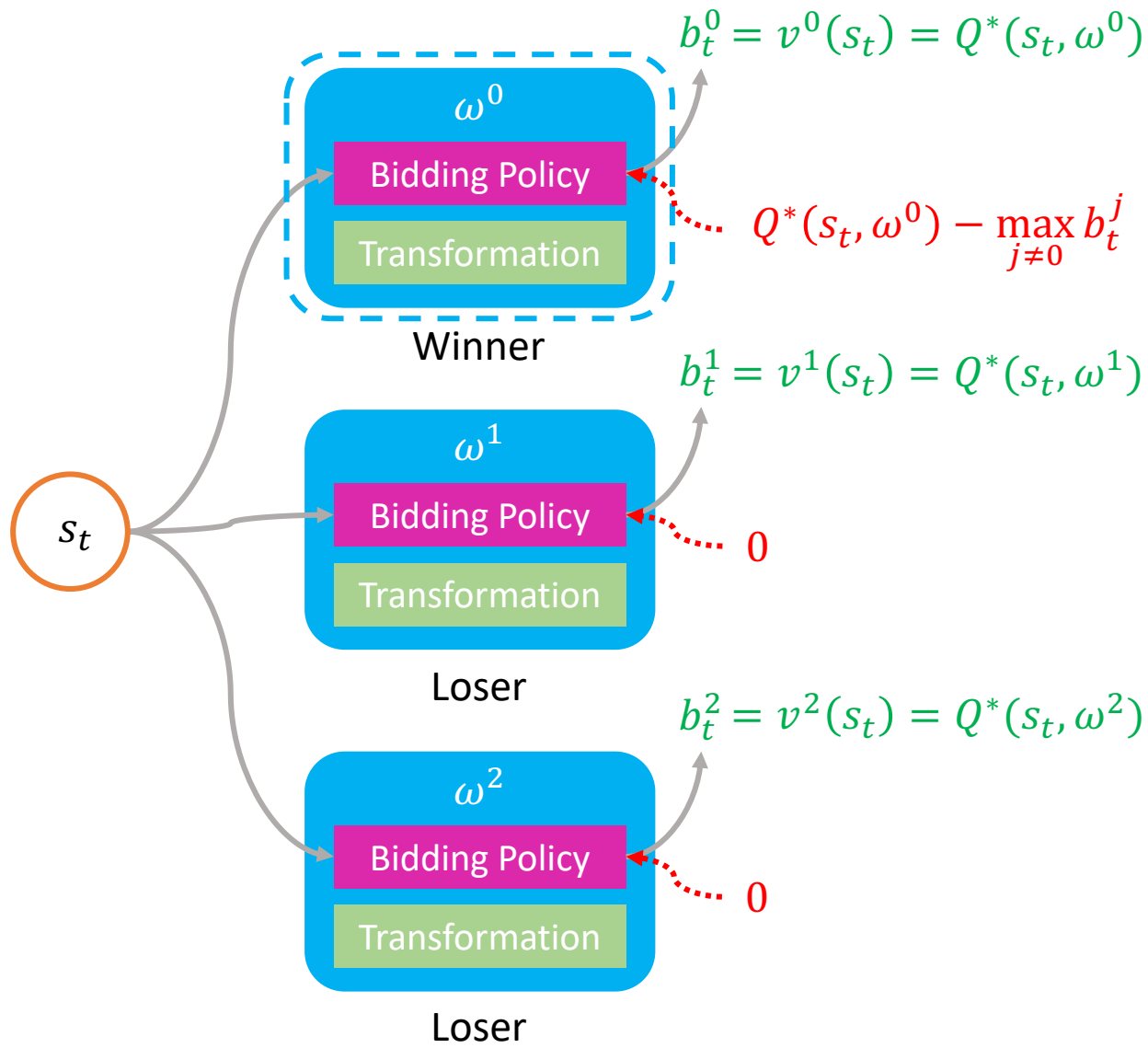
Want: Dominant Strategy Incentive Compatibility

The optimal strategy is to truthfully bid its own valuation:

$$b^i \leftarrow v^i$$

Implication: Set $v^k(s_t) = Q^*(s_t, \omega^k)$!

A Recipe for Relating Local and Global Objectives



$$\omega^* = \operatorname{argmax}_k Q^*(s_t, \omega^k)$$

Implication: Set $v^k(s_t) = Q^*(s_t, \omega^k)$!

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

For what auction mechanism would these optimal bids be an equilibrium strategy?

By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a **dominant strategy** to truthfully bid $Q^*(s, \omega^i)$.

How can we adapt this auction mechanism for discrete-action MDPs?

How can we avoid suboptimal equilibria?

How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?

But wait...

Optimal Q values are usually unknown!

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

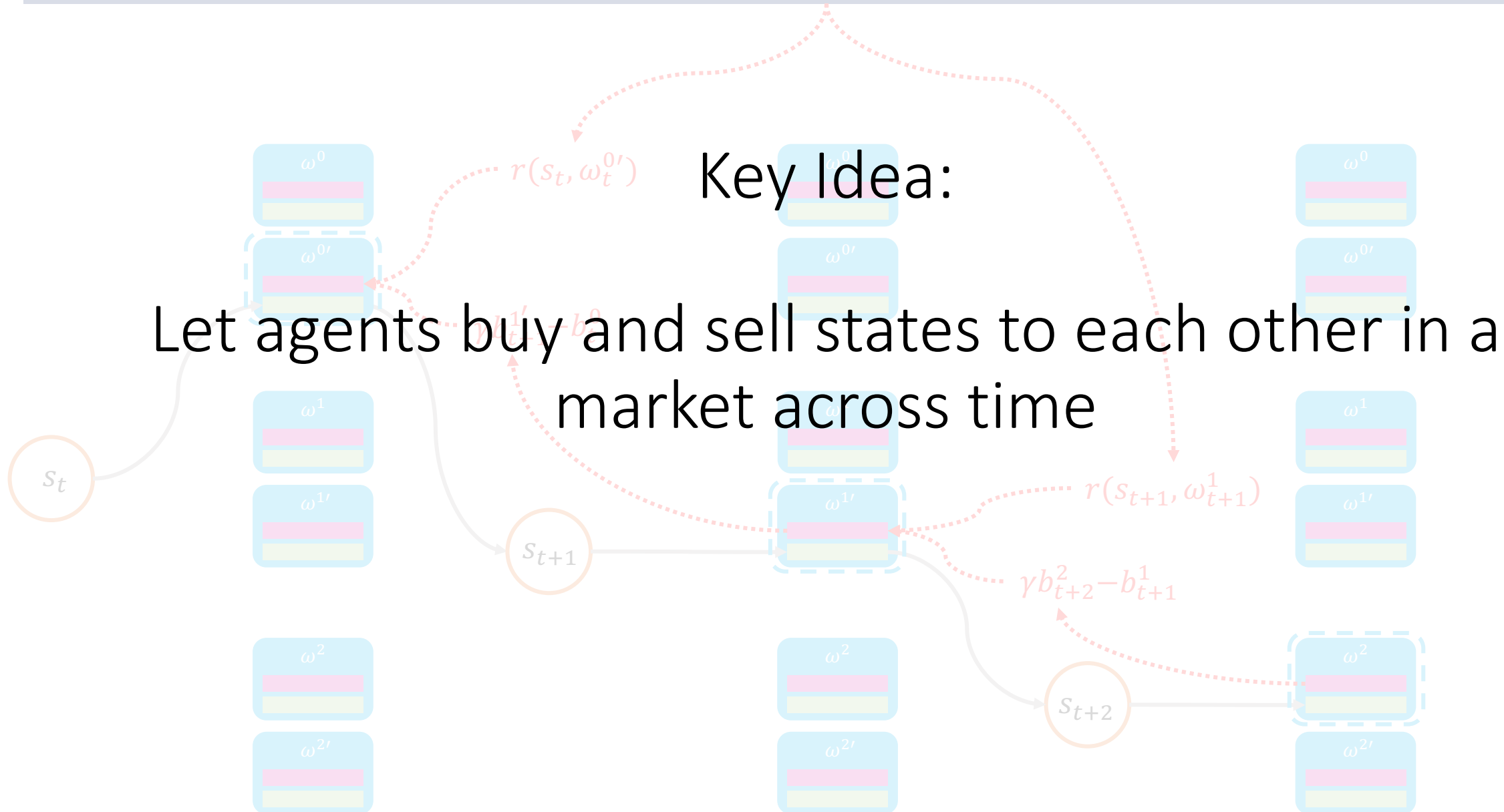
For what auction mechanism would these optimal bids be an equilibrium strategy?

By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a **dominant strategy** to truthfully bid $Q^*(s, \omega^i)$.

How can we adapt this auction mechanism for discrete-action MDPs?

How can we avoid suboptimal equilibria?

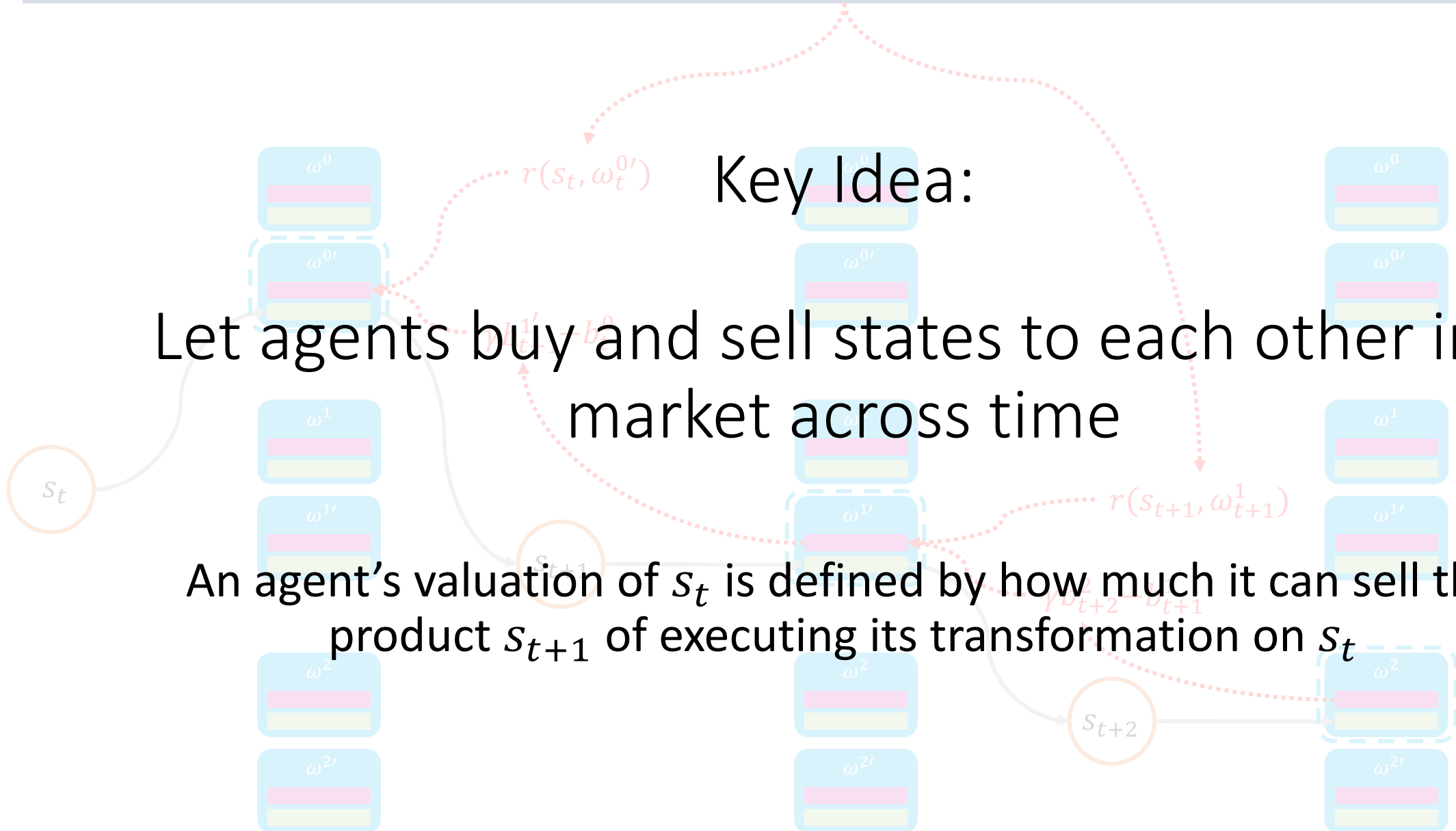
How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?

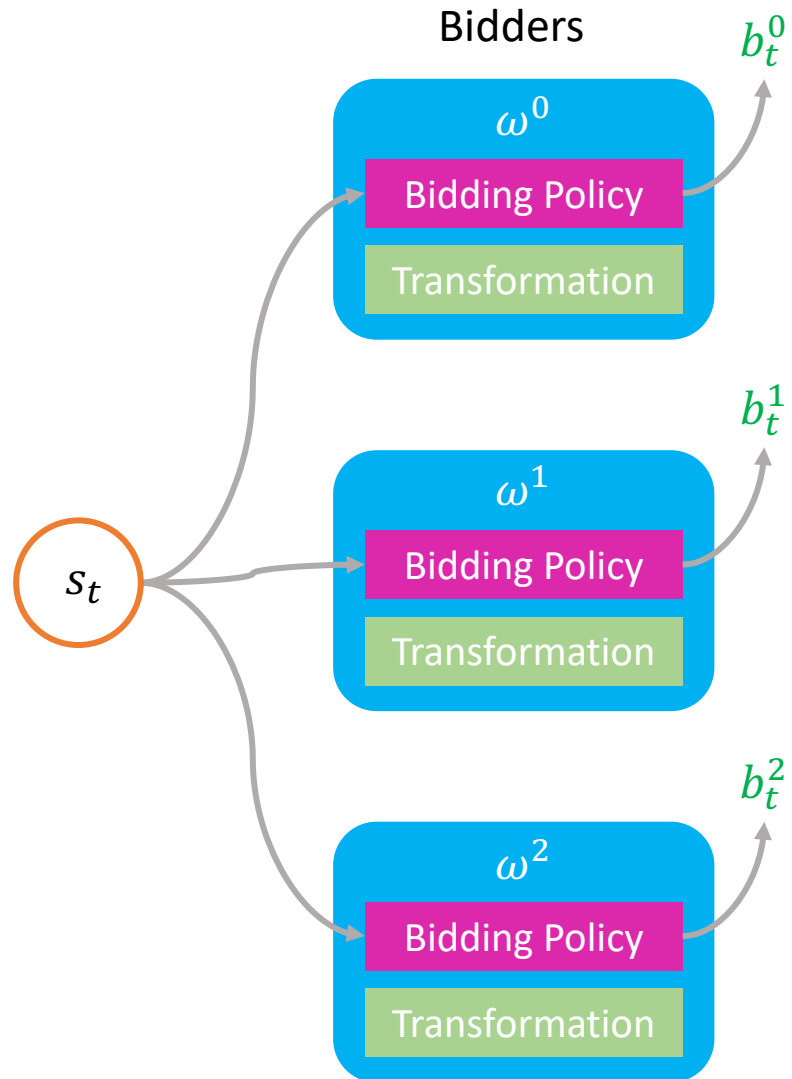


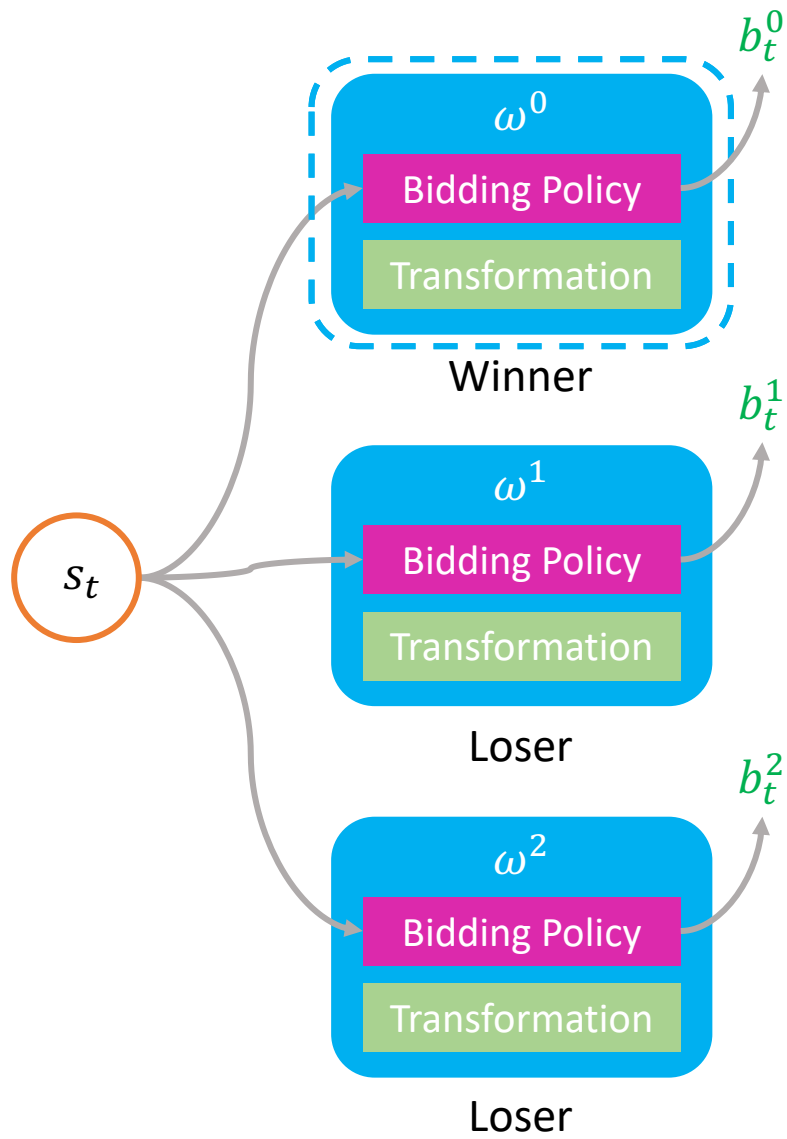
Key Idea:

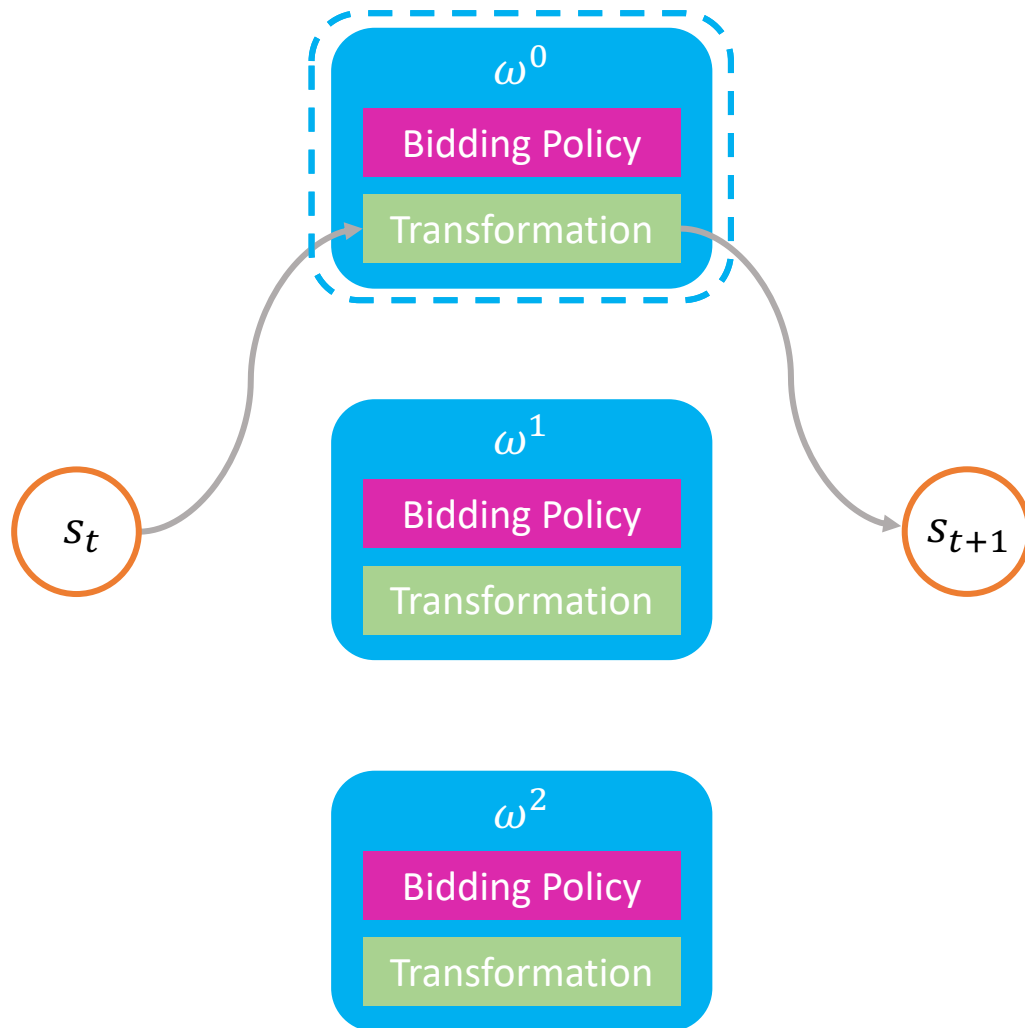
Let agents buy and sell states to each other in a market across time

An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t

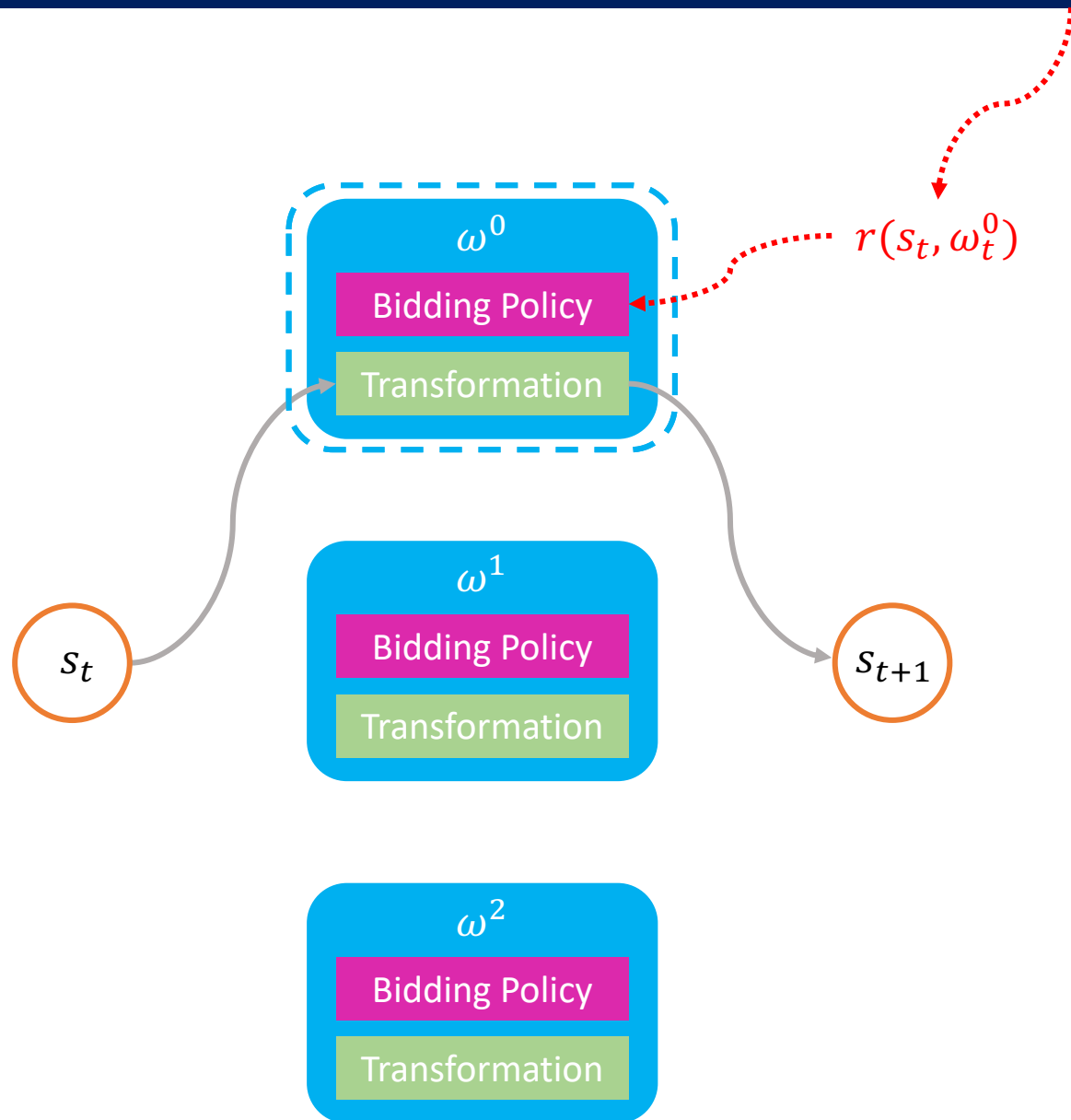




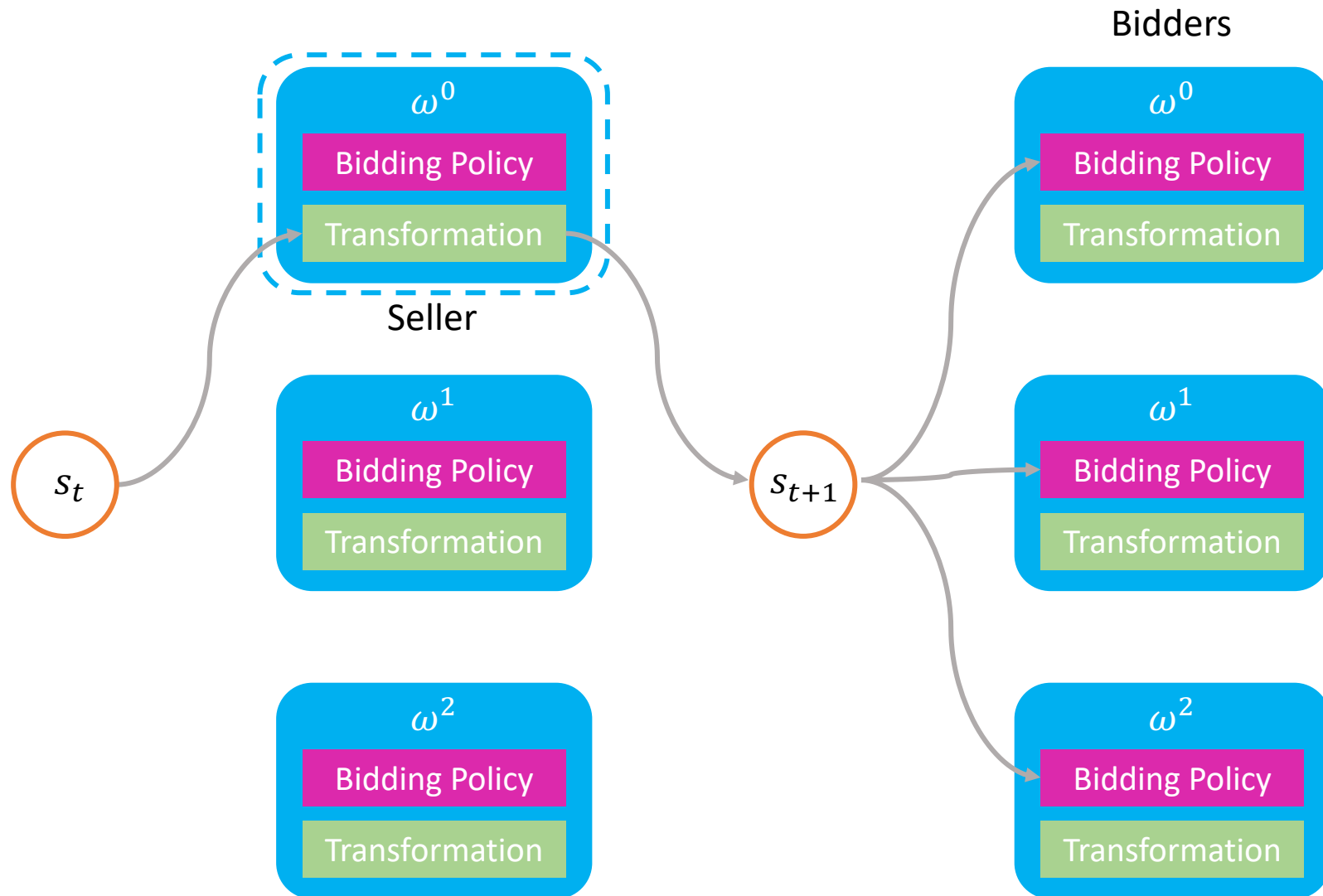




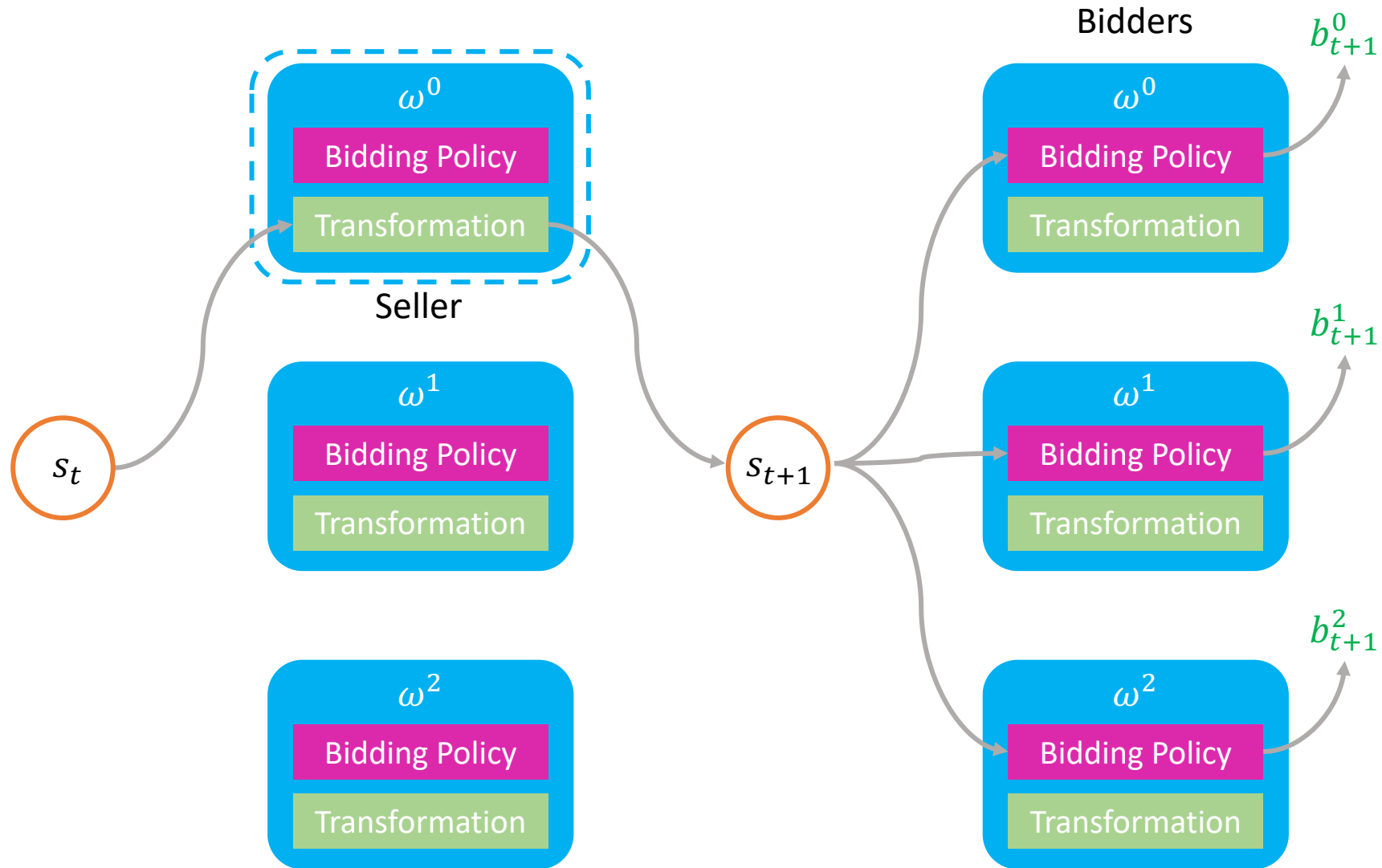
Environment



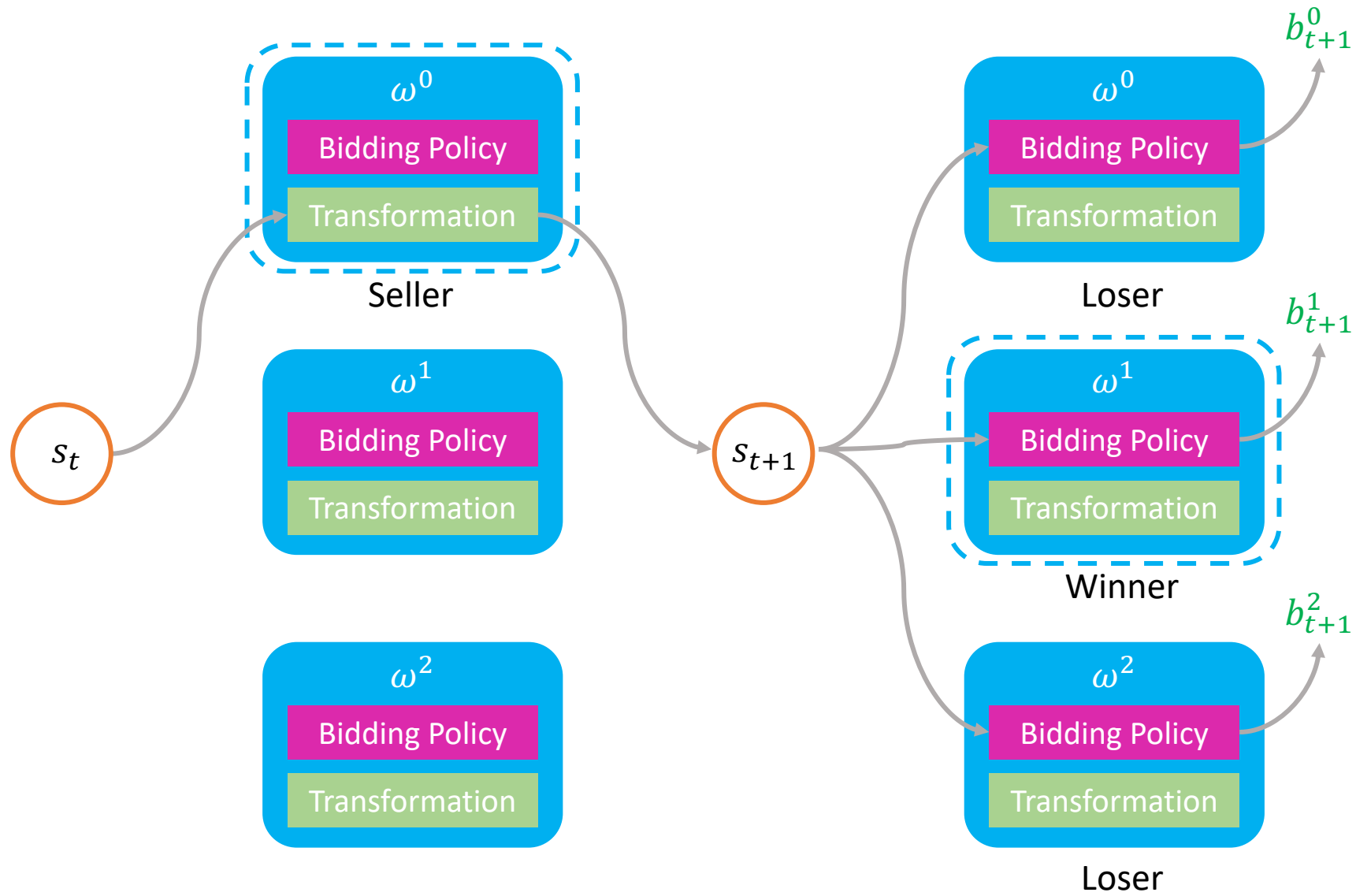
Environment



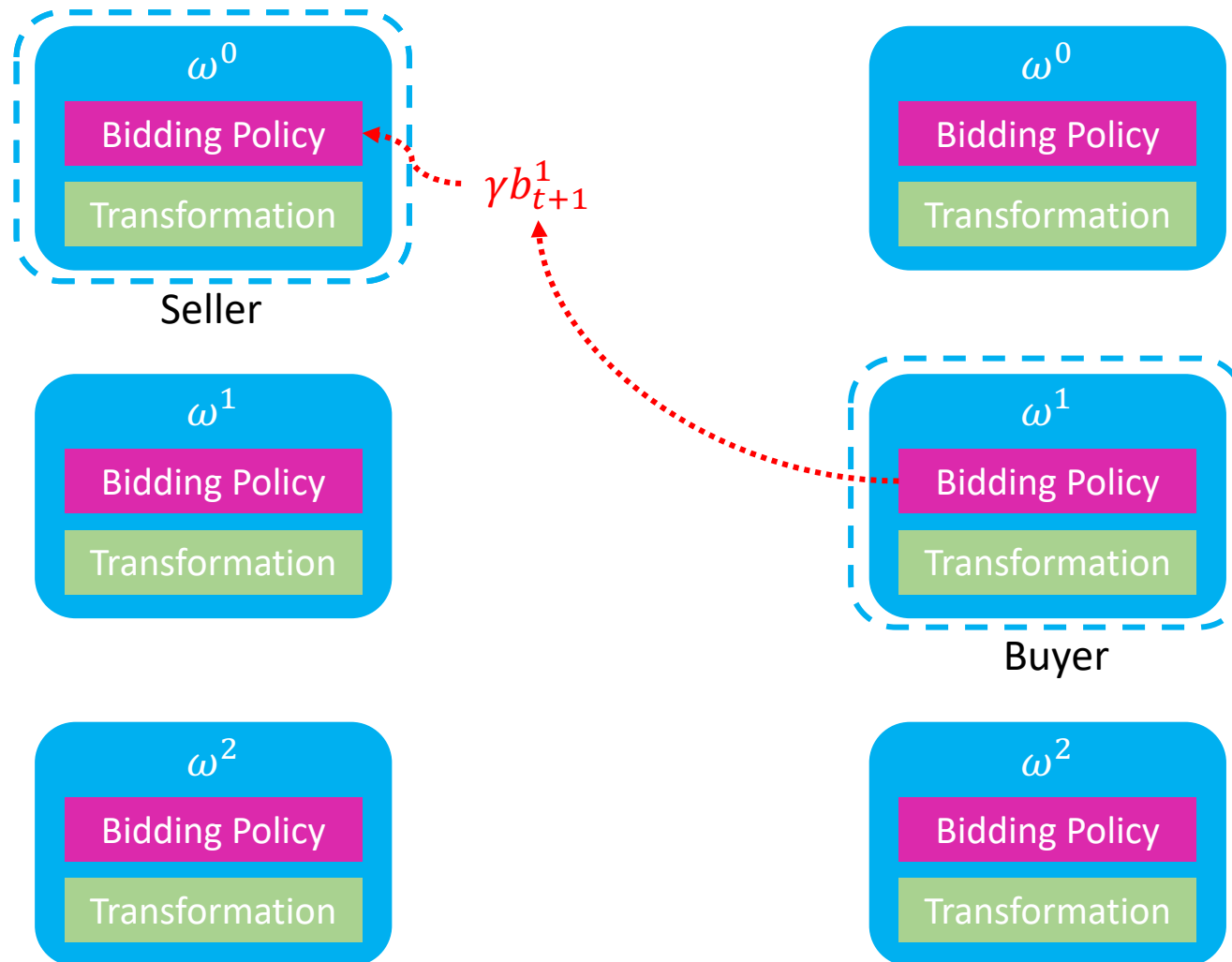
Environment



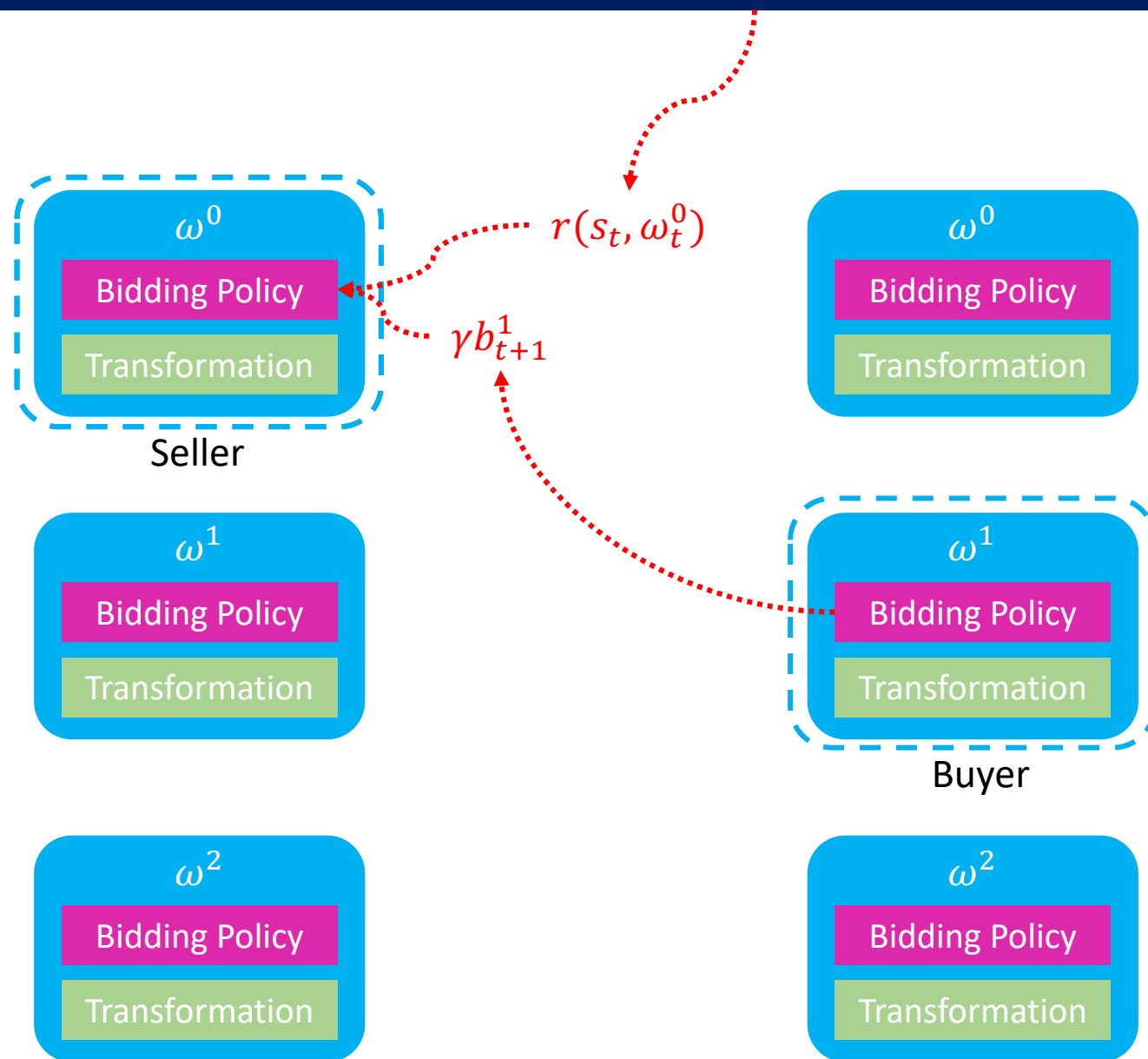
Environment



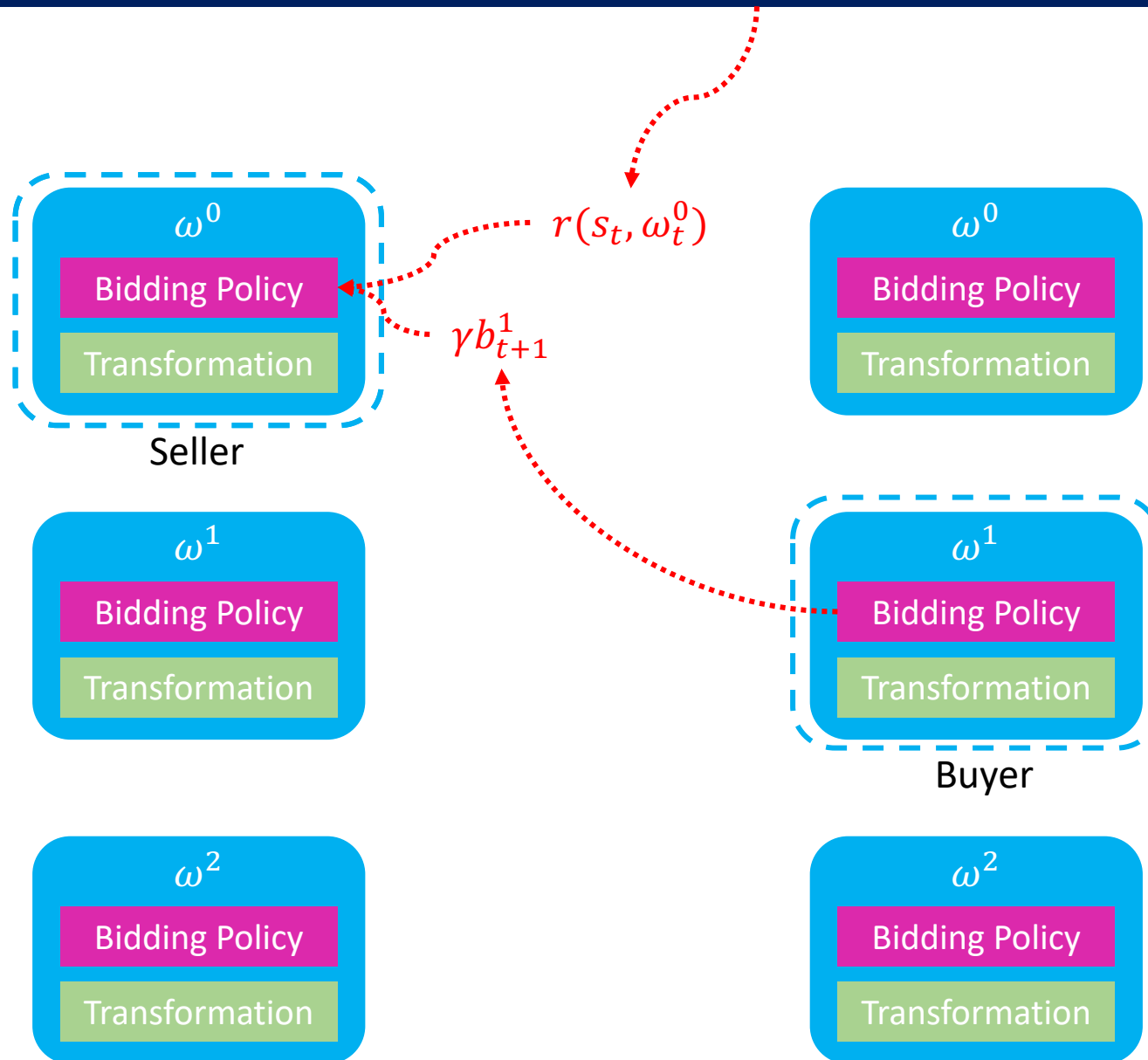
Environment



Environment



Environment



Valuations

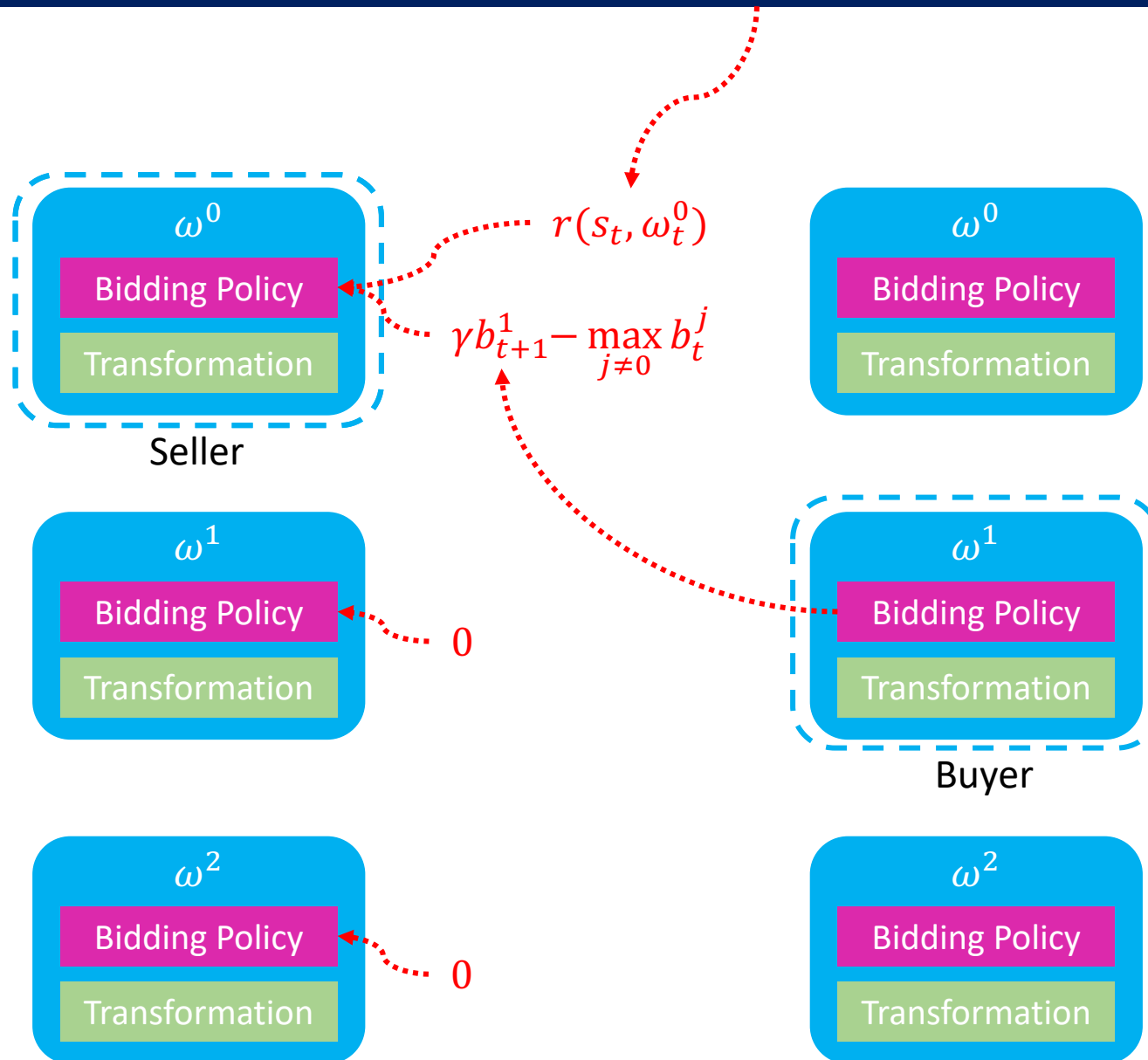
Before:

$$v^i(s_t) = Q^*(s_t \omega_t^i)$$

Now:

$$v^i(s_t) = r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k$$

Environment



Valuations

Before:

$$v^i(s_t) = Q^*(s_t \omega_t^i)$$

Now:

$$v^i(s_t) = r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k$$

Utilities

Winner's utility

$$u^i(b) = v^i - \max_{j \neq i} b^j$$

Loser's utility

$$u^i(b) = 0$$

Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

For what auction mechanism would these optimal bids be an equilibrium strategy?

By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a **dominant strategy** to truthfully bid $Q^*(s, \omega^i)$.

How can we adapt this auction mechanism for discrete-action MDPs?

Temporally couple the agents in a market: An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t .

How can we avoid suboptimal equilibria?

How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?

Proposition: If the utilities are defined as below, it is a Nash equilibrium for every primitive to bid their optimal Q value in the Global MDP.

Valuations

Utilities

Before:

$$v^i(s_t) = Q^*(s_t, \omega_t^i)$$

Winners:

$$u^i(b) = \left[r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k \right] - \max_{j \neq i} b^j$$

Now:

$$v^i(s_t) = r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k$$

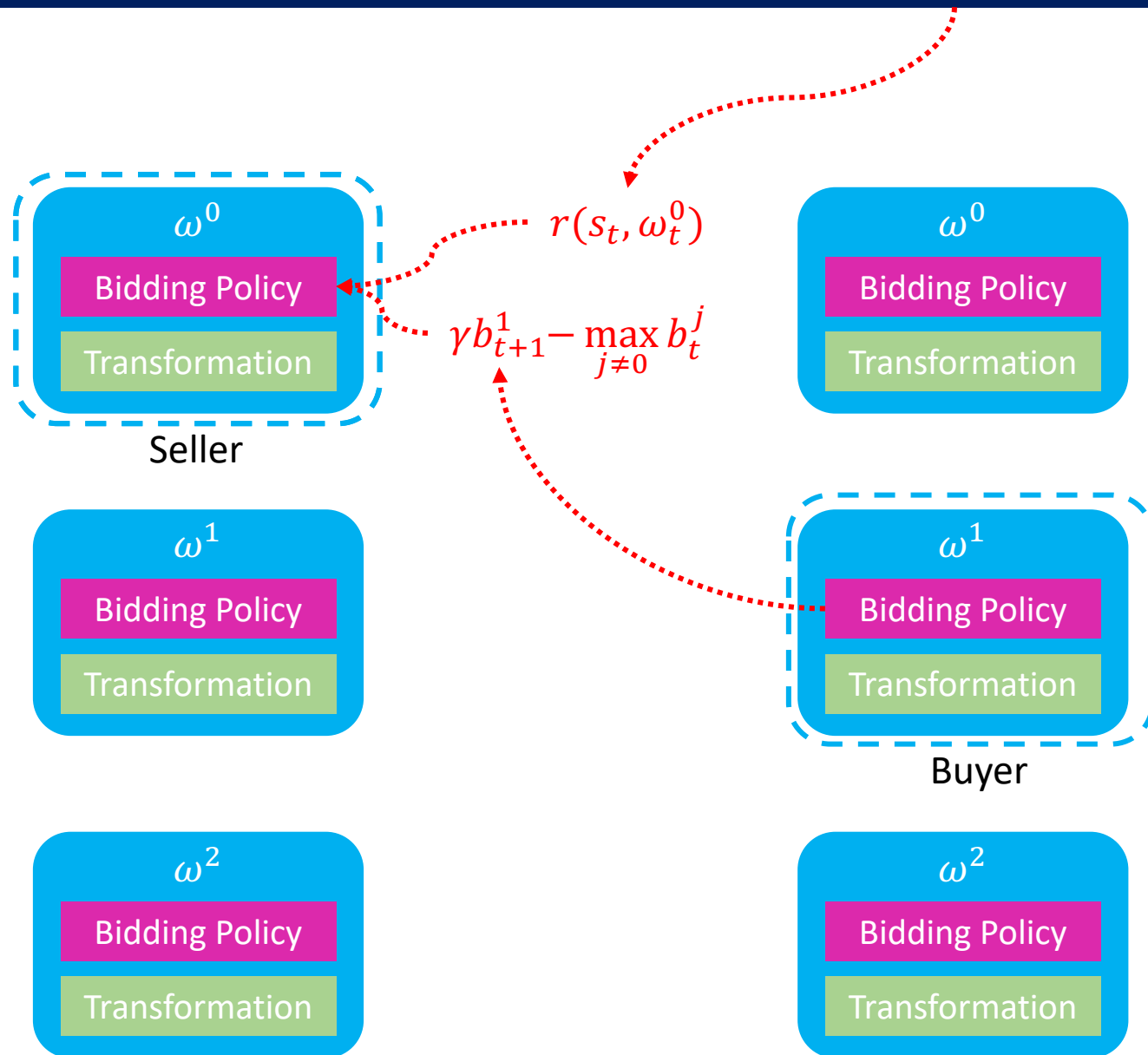
Losers:

$$u^i(b) = 0$$

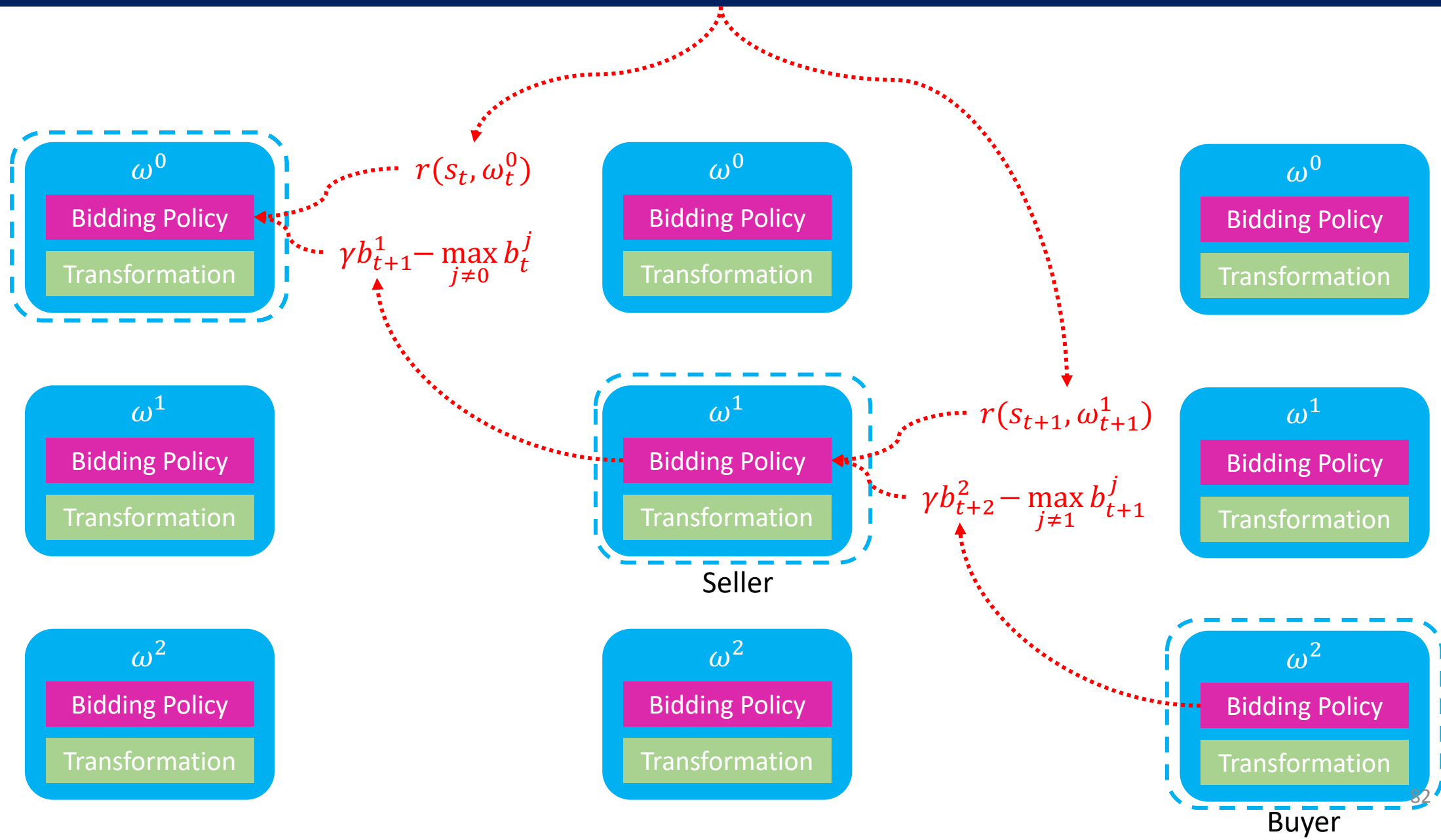
But wait...

Utility is not conserved!

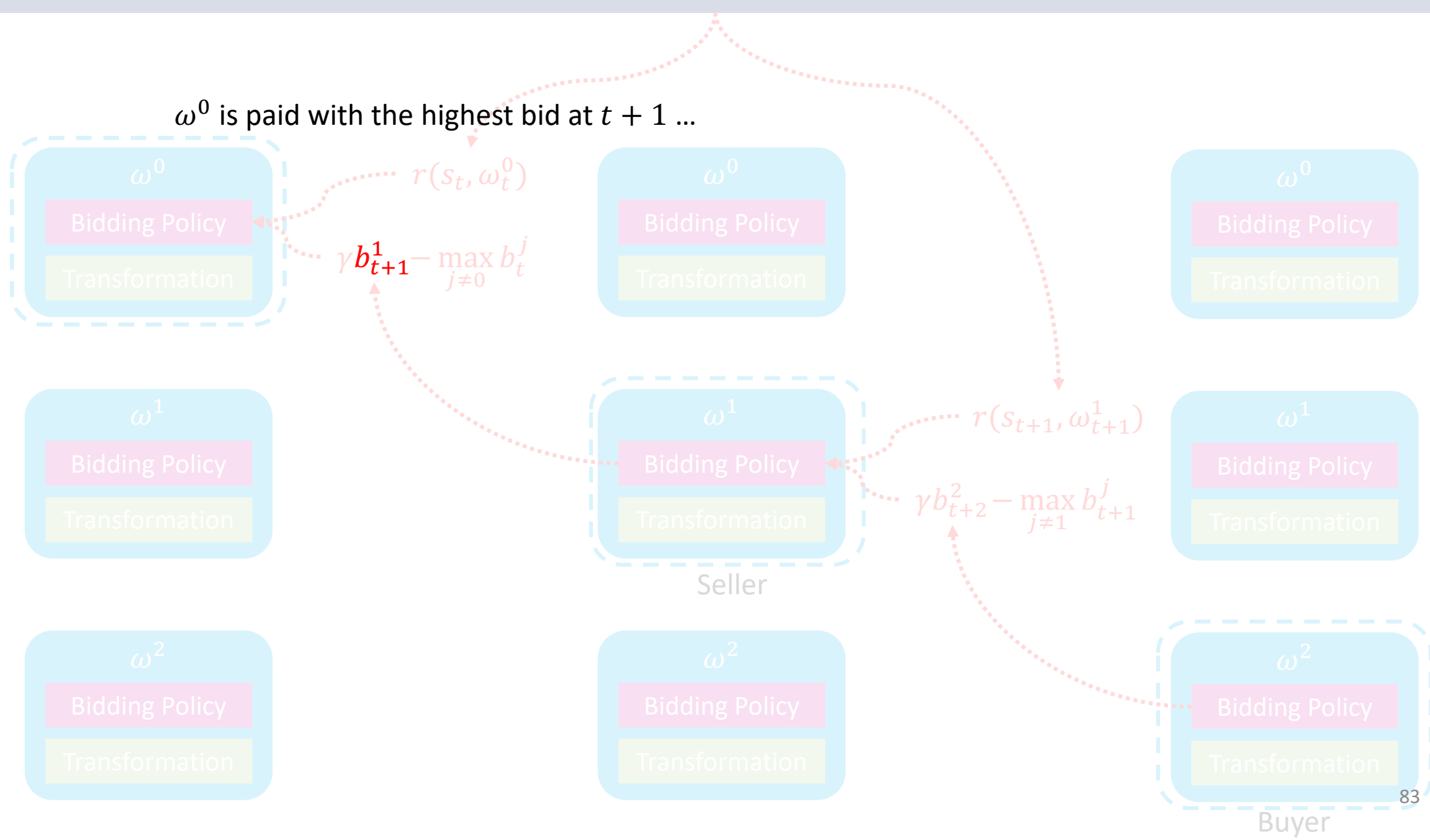
Environment



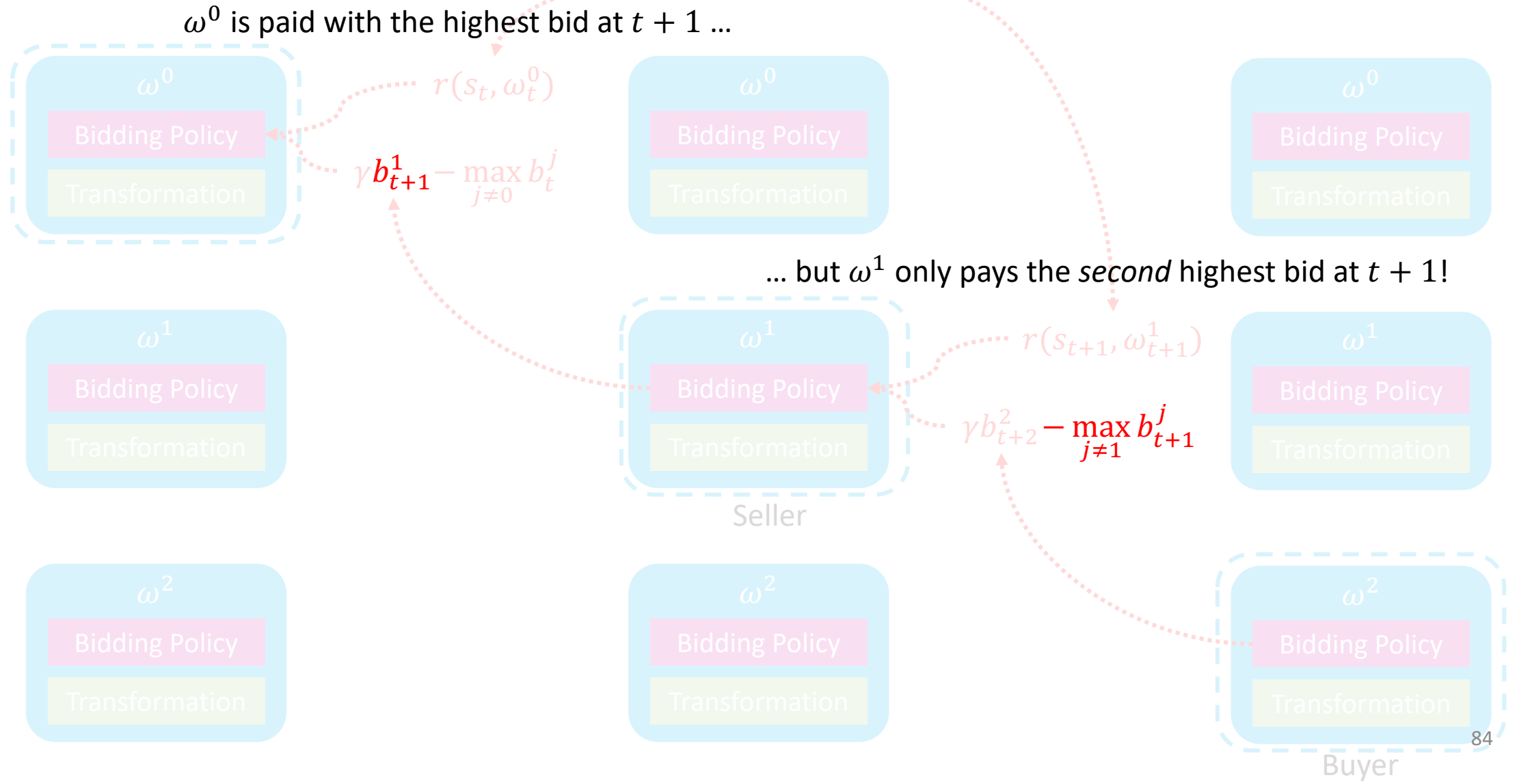
Environment



Environment



Environment



Roadmap

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

For what auction mechanism would these optimal bids be an equilibrium strategy?

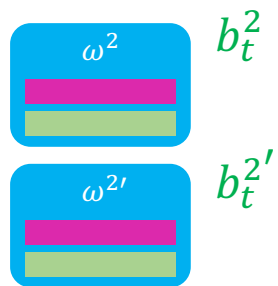
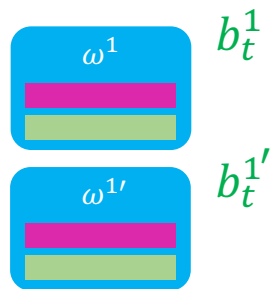
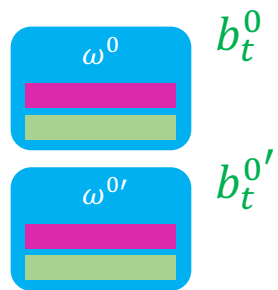
By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a **dominant strategy** to truthfully bid $Q^*(s, \omega^i)$.

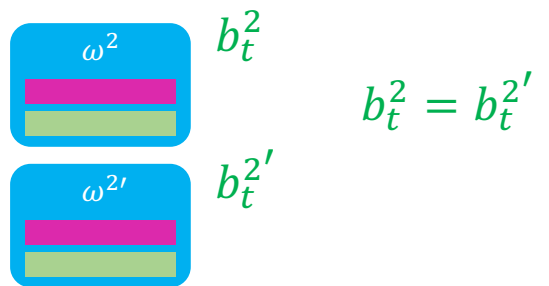
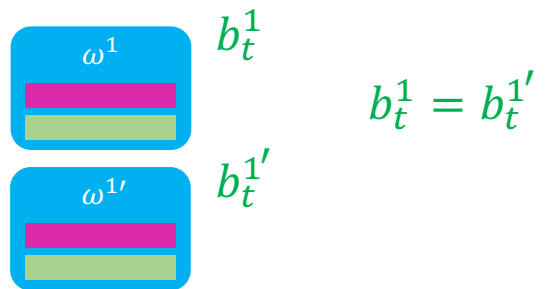
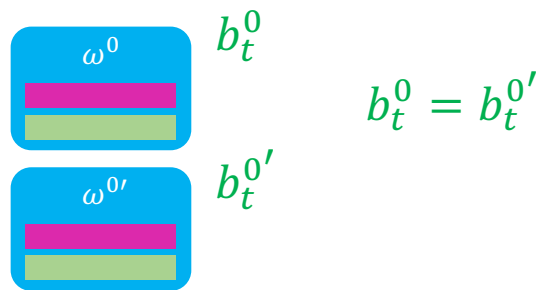
How can we adapt this auction mechanism for discrete-action MDPs?

Temporally couple the agents in a market: An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t .

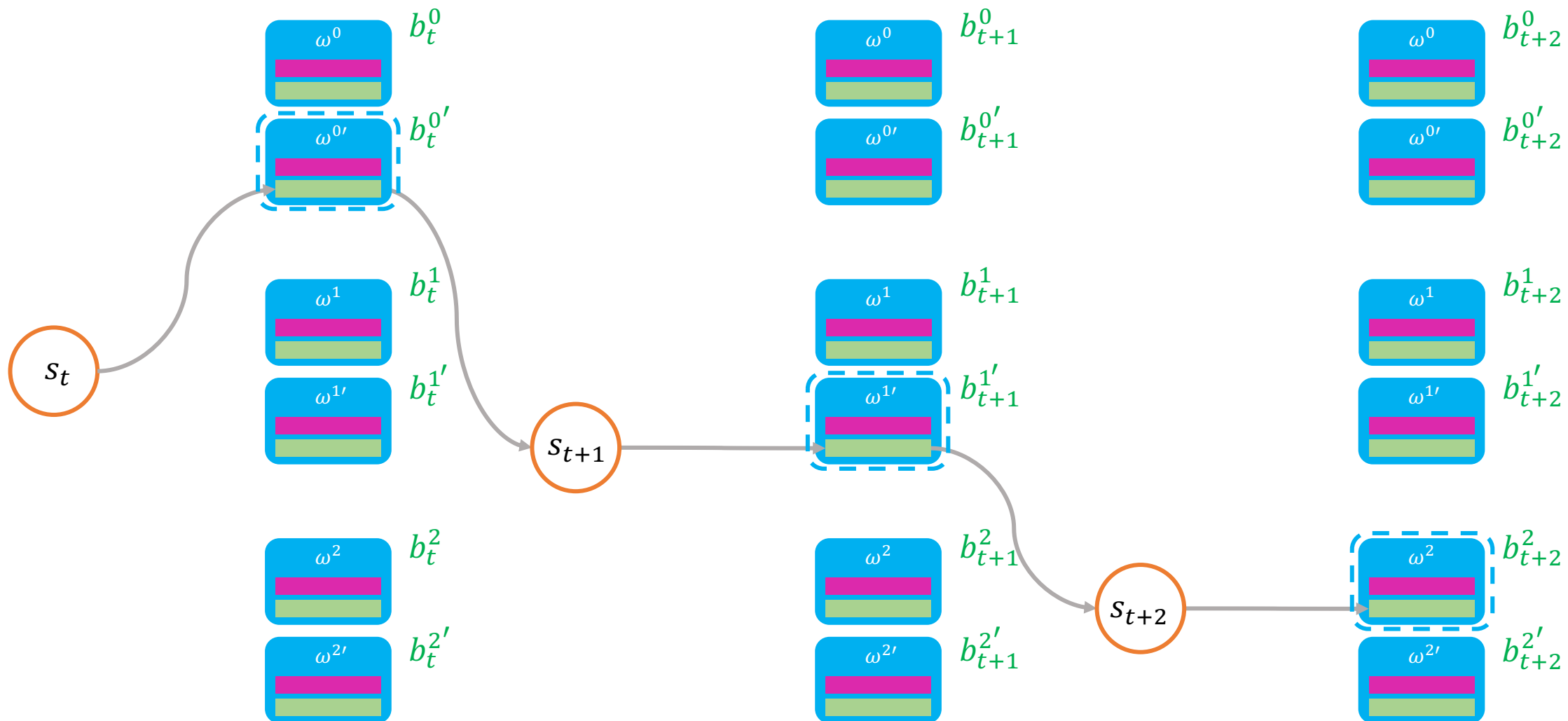
How can we avoid suboptimal equilibria?

How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?

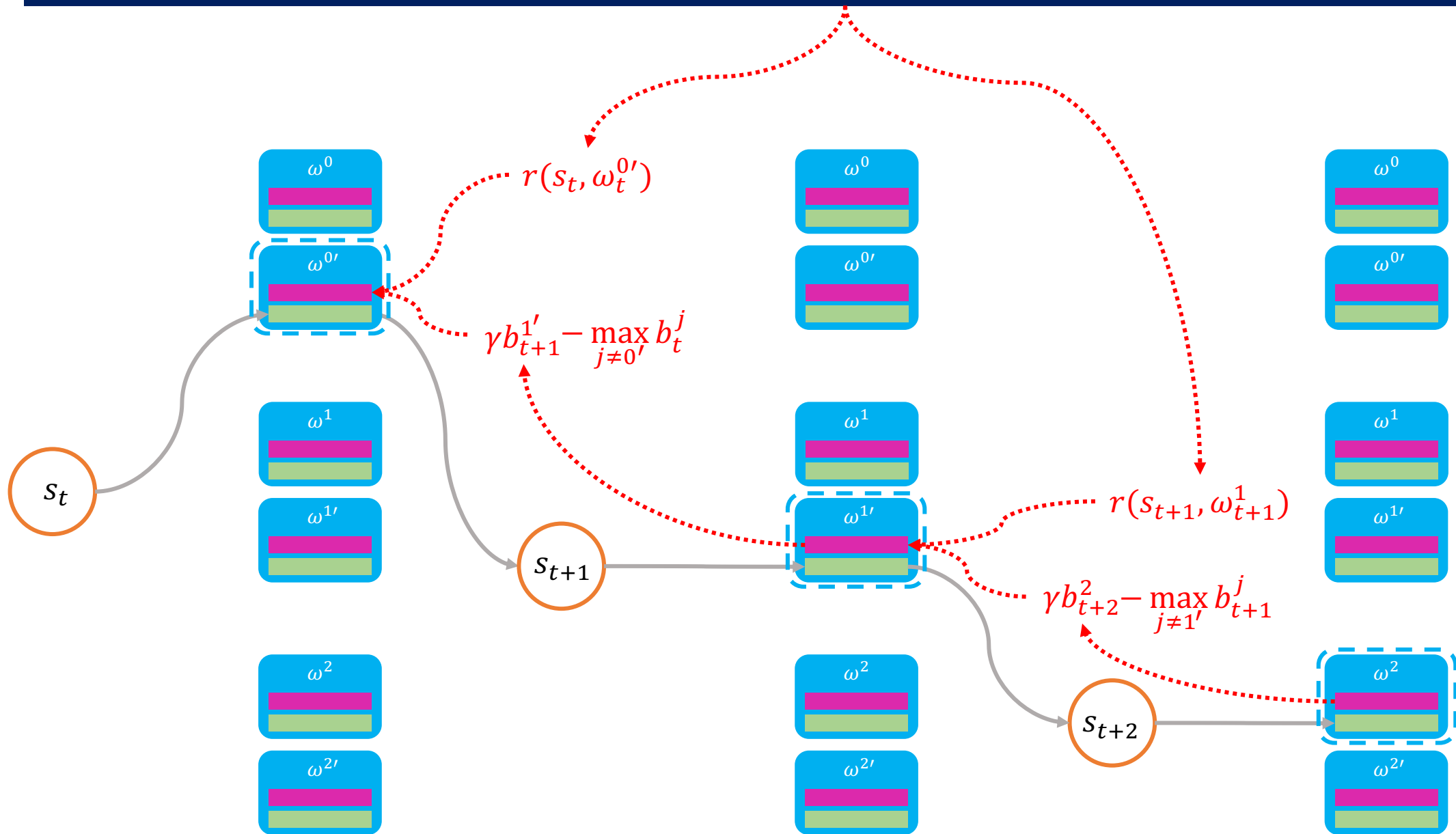




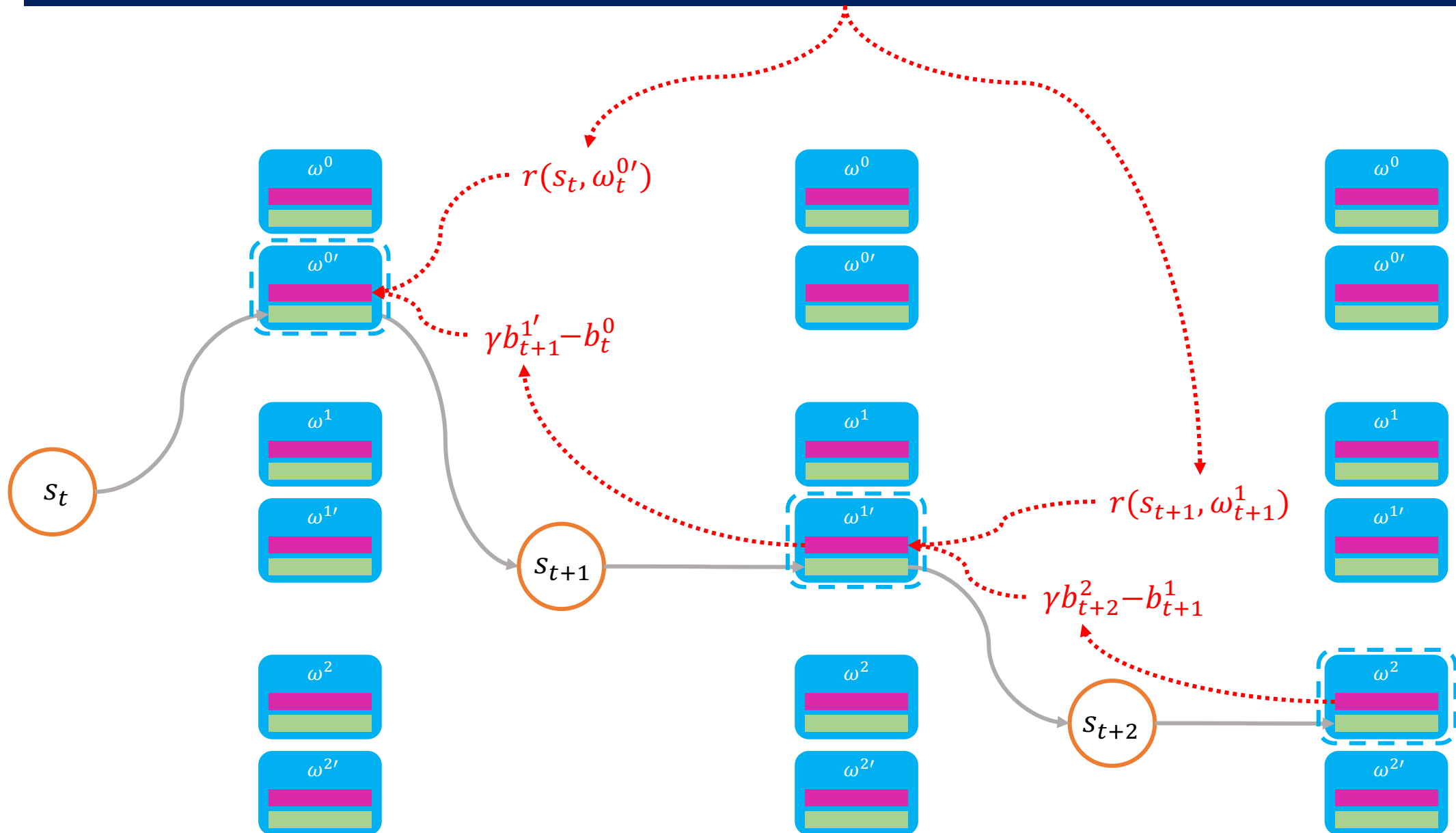
Environment



Environment

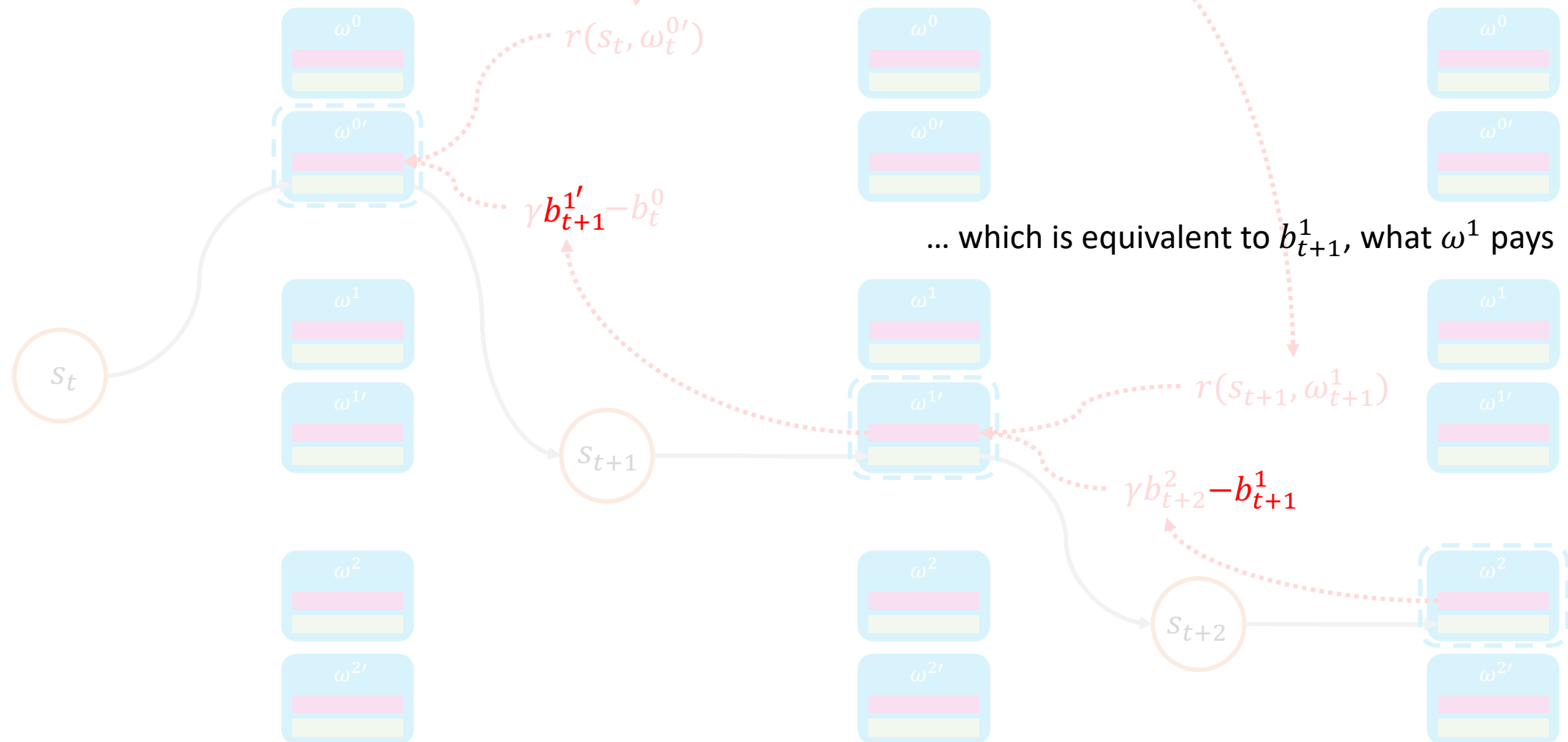


Environment

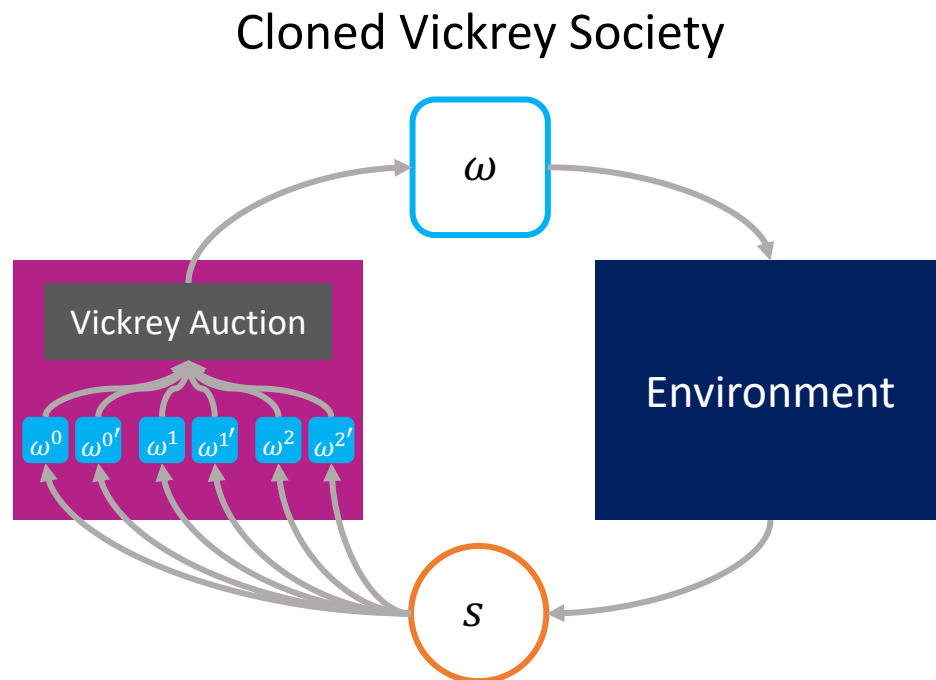


Environment

ω^0 is paid with the $b_{t+1}^{1'}$...



Main Result: Cloned Vickrey Society



Utilities

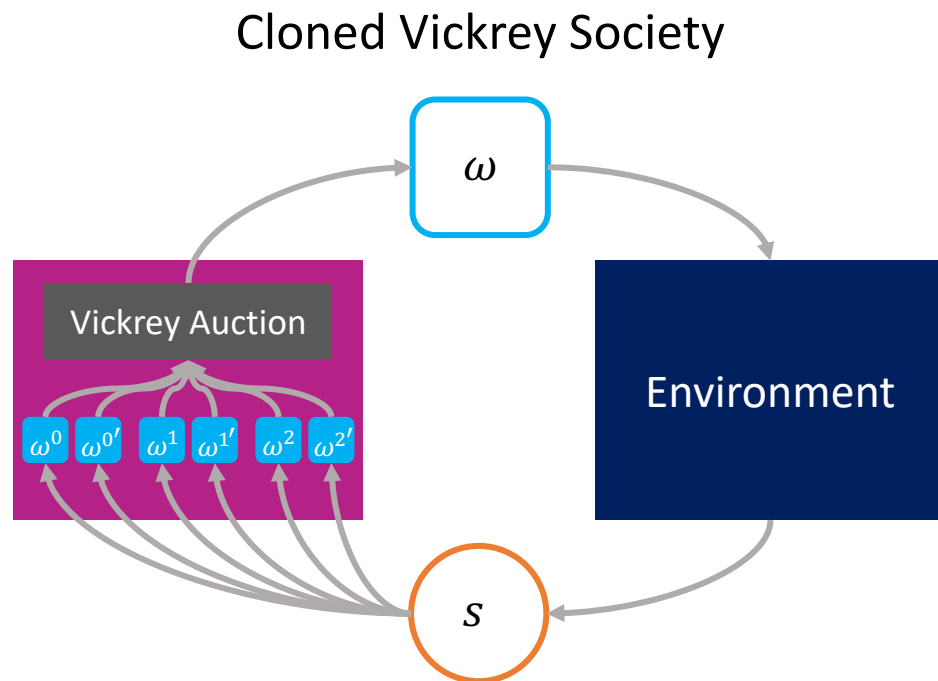
Winners:

$$u^i(b) = \left[r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k \right] - \max_{j \neq i} b^j$$

Losers:

$$u^i(b) = 0$$

Theorem: In a Cloned Vickrey Society, it is a Nash equilibrium for every primitive to bid their optimal Q value in the Global MDP and utility is conserved.



Utilities

Winners:

$$u^i(b) = \left[r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k \right] - \max_{j \neq i} b^j$$

Losers:

$$u^i(b) = 0$$

Roadmap

Question	Key Idea
What should the optimal bids be for the solution of the Global MDP to emerge?	Define the optimal bid as the optimal Q value $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .
For what auction mechanism would these optimal bids be an equilibrium strategy?	By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a dominant strategy to truthfully bid $Q^*(s, \omega^i)$.
How can we adapt this auction mechanism for discrete-action MDPs?	Temporally couple the agents in a market: An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t .
How can we avoid suboptimal equilibria?	Redundancy enforces credit conservation that helps avoid suboptimal equilibria.
How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?	

From Equilibria to Learning Objectives

Each agent learns a bidding policy by optimizes their utility as reward:

Winners:

$$u^i(b) = \left[r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k \right] - \max_{j \neq i} b^j$$

Losers:

$$u^i(b) = 0$$

Train bidding policies using standard reinforcement learning algorithms

Decentralized Reinforcement Learning

Each agent learns a bidding policy by optimizes their utility as reward:

Winners:

$$u^i(b) = \left[r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k \right] - \max_{j \neq i} b^j$$

Losers:

$$u^i(b) = 0$$

Train bidding policies using standard reinforcement learning algorithms

Society: an emergent solution that is **global** in space and time

Agent: learns via credit assignment **local** in space and time

Contributions

Question	Key Idea
What should the optimal bids be for the solution of the Global MDP to emerge?	Define the optimal bid as the optimal Q value $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .
For what auction mechanism would these optimal bids be an equilibrium strategy?	By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a dominant strategy to truthfully bid $Q^*(s, \omega^i)$.
How can we adapt this auction mechanism for discrete-action MDPs?	Temporally couple the agents in a market: An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t .
How can we avoid suboptimal equilibria?	Redundancy enforces credit conservation that helps avoid suboptimal equilibria.
How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?	Define the auction utility as the agents' reinforcement learning objective, yielding a decentralized reinforcement learning algorithm for the Global MDP.

Contributions

Assumptions

Assume the agents ω^i know their valuations as
$$v^i(s_t) = Q^*(s_t, \omega_t^i)$$

Dominant strategy equilibrium in auction = solution to Global MDP

Pro: provable dominant strategy equilibrium

Con: assumes optimal Q-values are known

Key Idea

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a **dominant strategy** to truthfully bid $Q^*(s, \omega^i)$.

Temporally couple the agents in a market: An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t .

Redundancy enforces credit conservation that helps avoid suboptimal equilibria.

Define the **auction utility** as the agents' reinforcement learning objective, yielding a **decentralized reinforcement learning algorithm** for the Global MDP.

Contributions

Assumptions

Assume the agents ω^i know their valuations as
$$v^i(s_t) = Q^*(s_t, \omega_t^i)$$

Dominant strategy equilibrium in auction = solution to Global MDP

Pro: provable dominant strategy equilibrium

Con: assumes optimal Q-values are known

Assume the agents ω^i know their valuations as
$$v^i(s_t) = r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k$$

Nash equilibrium in auction = solution to Global MDP

Pro: does not assume optimal Q-value is known

Con: assumes valuations are known

Key Idea

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a **dominant strategy** to truthfully bid $Q^*(s, \omega^i)$.

Temporally couple the agents in a market: An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t .

Redundancy enforces credit conservation that helps avoid suboptimal equilibria.

Define the **auction utility** as the agents' reinforcement learning objective, yielding a **decentralized reinforcement learning algorithm** for the Global MDP.

Contributions

Assumptions

Assume the agents ω^i know their valuations as
$$v^i(s_t) = Q^*(s_t, \omega_t^i)$$

Dominant strategy equilibrium in auction = solution to Global MDP

Pro: provable dominant strategy equilibrium

Con: assumes optimal Q-values are known

Assume the agents ω^i know their valuations as
$$v^i(s_t) = r(s_t, \omega_t^i) + \gamma \max_k b_{t+1}^k$$

Nash equilibrium in auction = solution to Global MDP

Pro: does not assume optimal Q-value is known

Con: assumes valuations are known

Assume the agents ω^i learn their valuations through interaction.

Nash equilibrium in auction = solution to Global MDP

Pro: does not assume valuations are known

Con: difficult to prove convergence to equilibrium

Key Idea

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a **dominant strategy** to truthfully bid $Q^*(s, \omega^i)$.

Temporally couple the agents in a market: An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t .

Redundancy enforces credit conservation that helps avoid suboptimal equilibria.

Define the **auction utility** as the agents' reinforcement learning objective, yielding a **decentralized reinforcement learning algorithm** for the Global MDP.

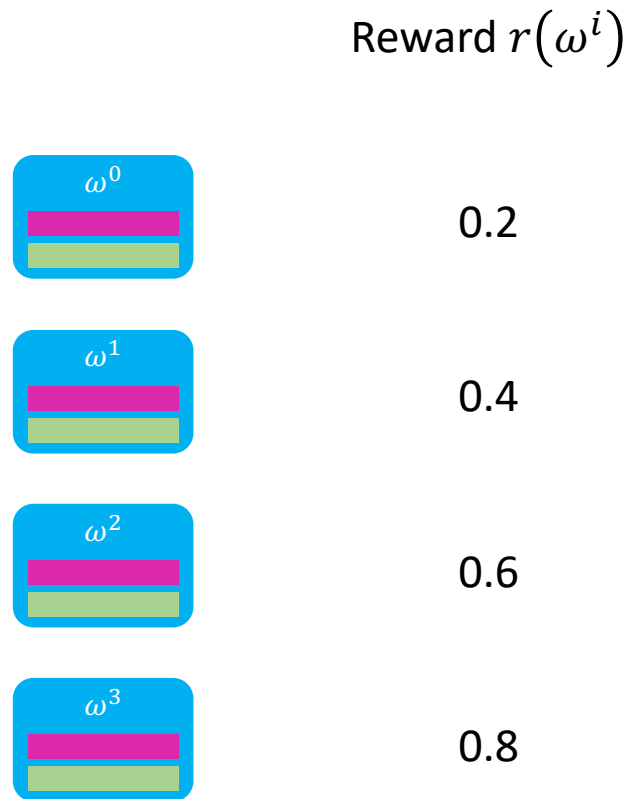
Numerical Simulations

1. How closely do the bids the agents learn match their optimal Q-values?
2. Does the solution to the global objective emerge from the competition among the agents?
3. How does redundancy affect the solutions the agents converge to?
4. Does the modularity of such a decentralized system offer benefit in transferring to new tasks?

Warm-Up: Bandit




Warm-Up: Bandit



Global Objective for the Society
Maximize reward

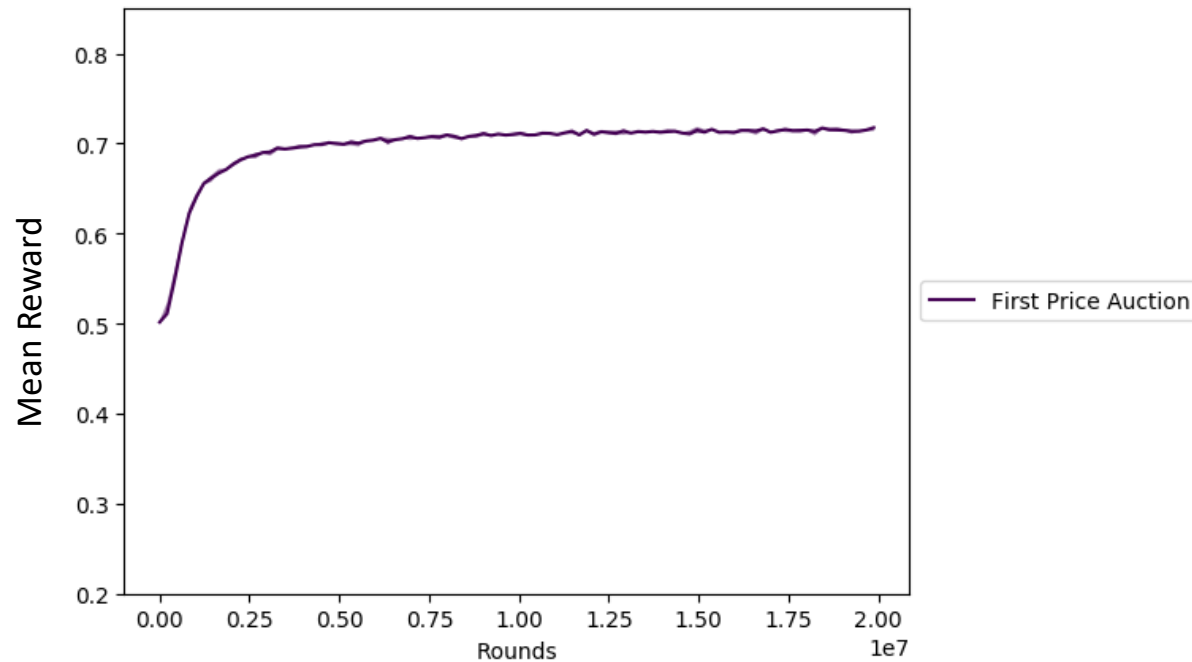
Local Objectives for the Agents
Maximize utility in the auction

Warm-Up: Bandit

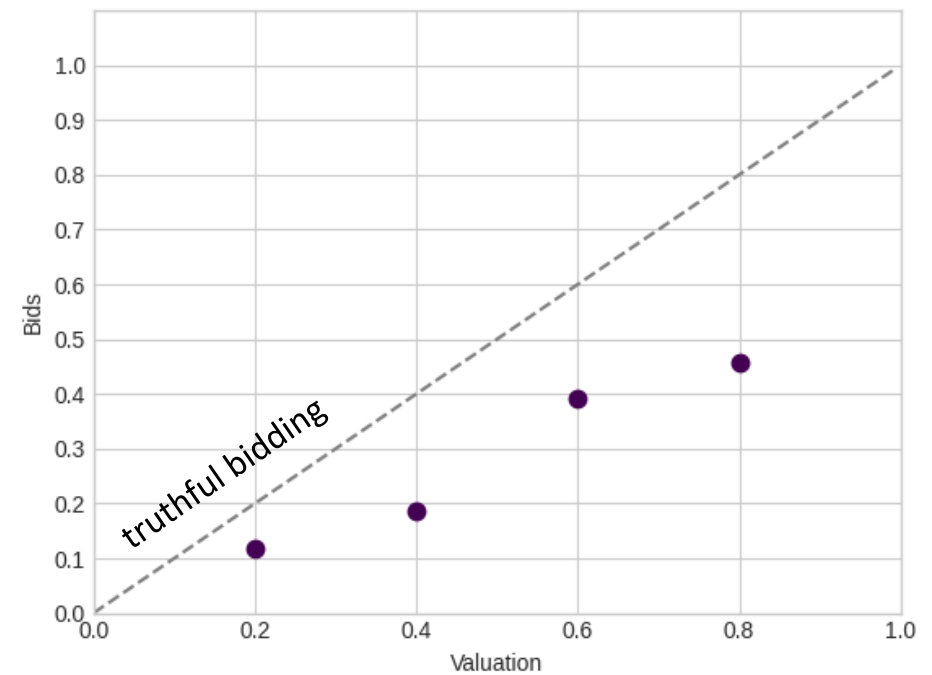
	Reward $r(\omega^i)$	Truthful Bid b^i	
	0.2	0.2	Global Objective for the Society Maximize reward
	0.4	0.4	Local Objectives for the Agents Maximize utility in the auction
	0.6	0.6	
	0.8	0.8	

Warm-Up: Bandit

Does the solution to the global objective emerge from the competition among the agents?

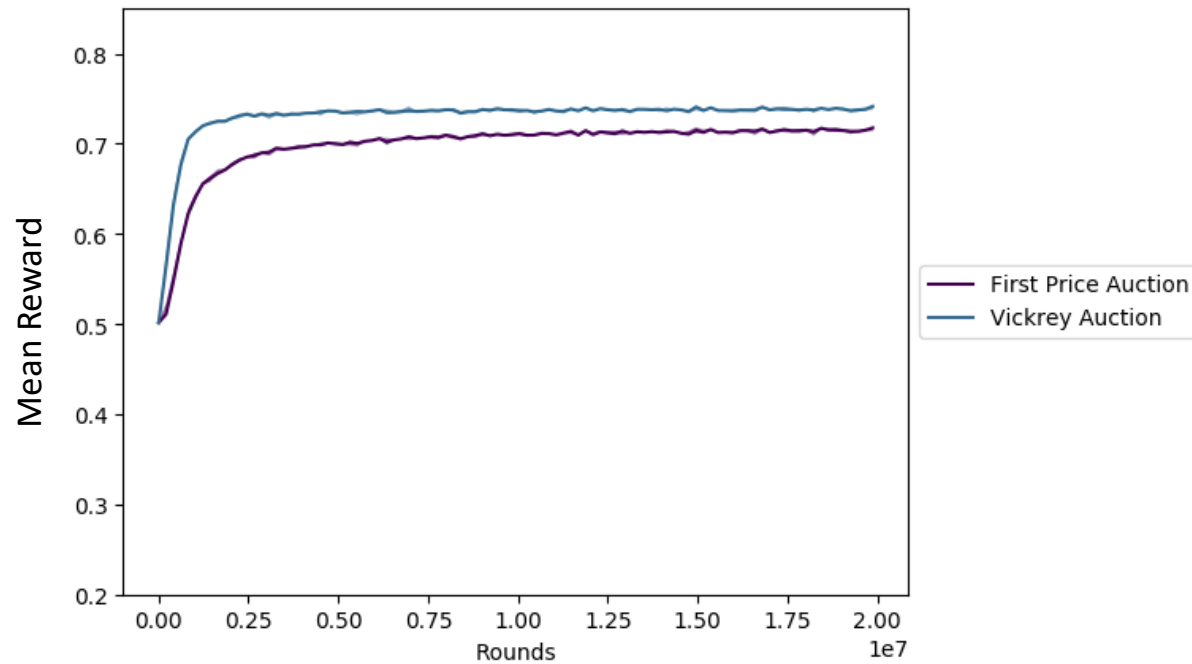


How closely do the bids the agents learn match their optimal Q-values?

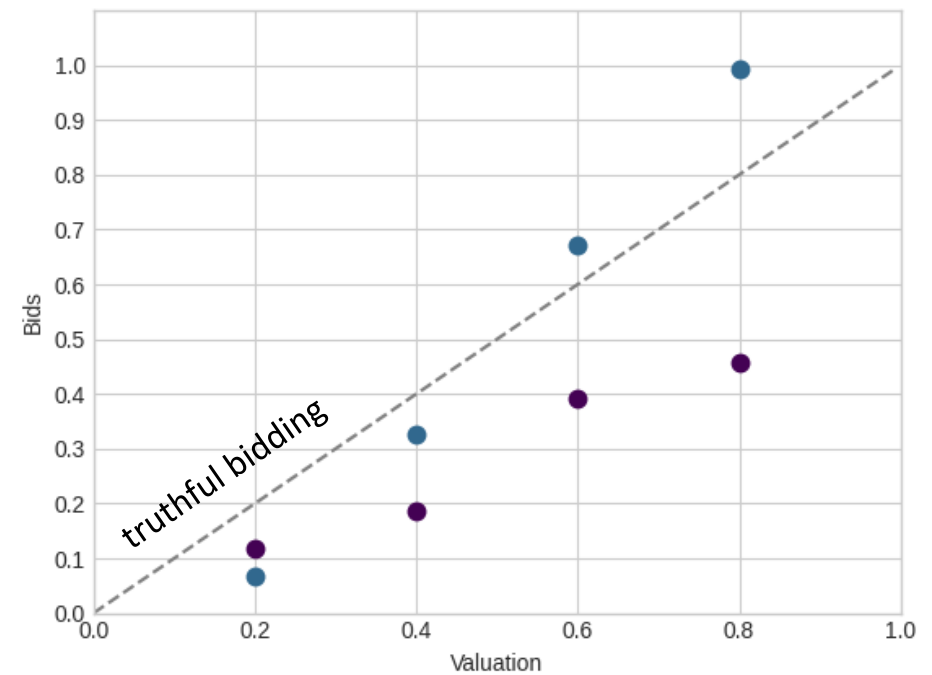


Warm-Up: Bandit

Does the solution to the global objective emerge from the competition among the agents?

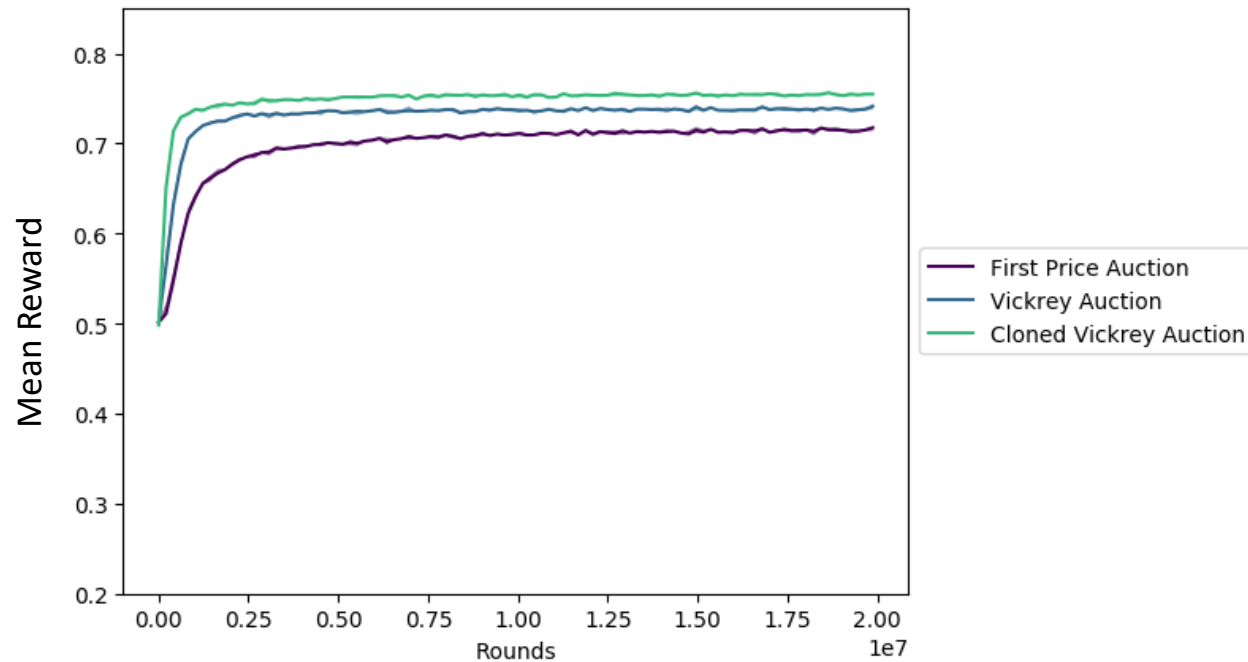


How closely do the bids the agents learn match their optimal Q-values?

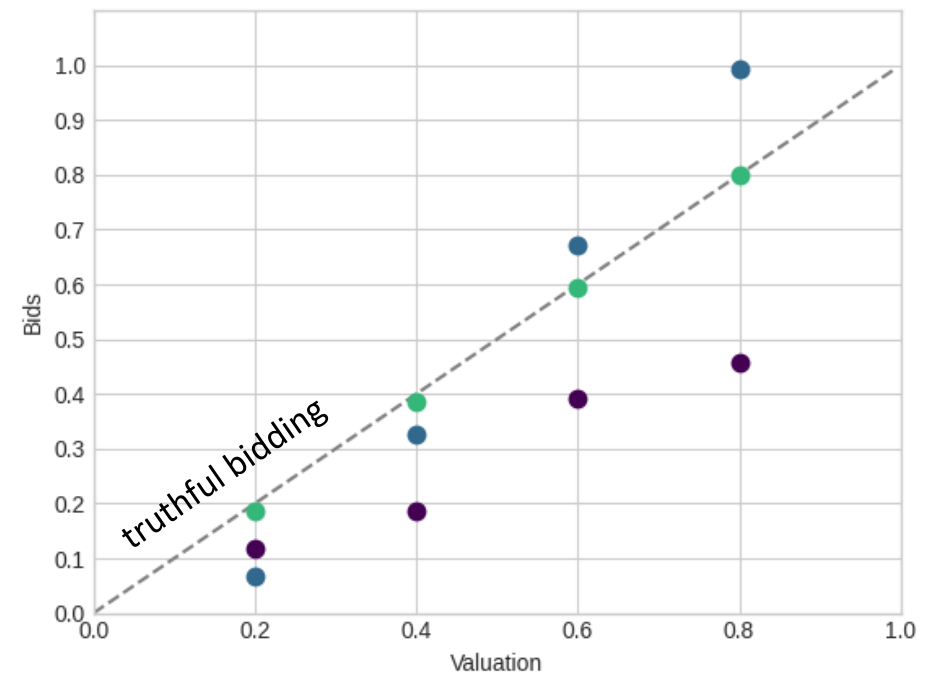


Warm-Up: Bandit

Does the solution to the global objective emerge from the competition among the agents?

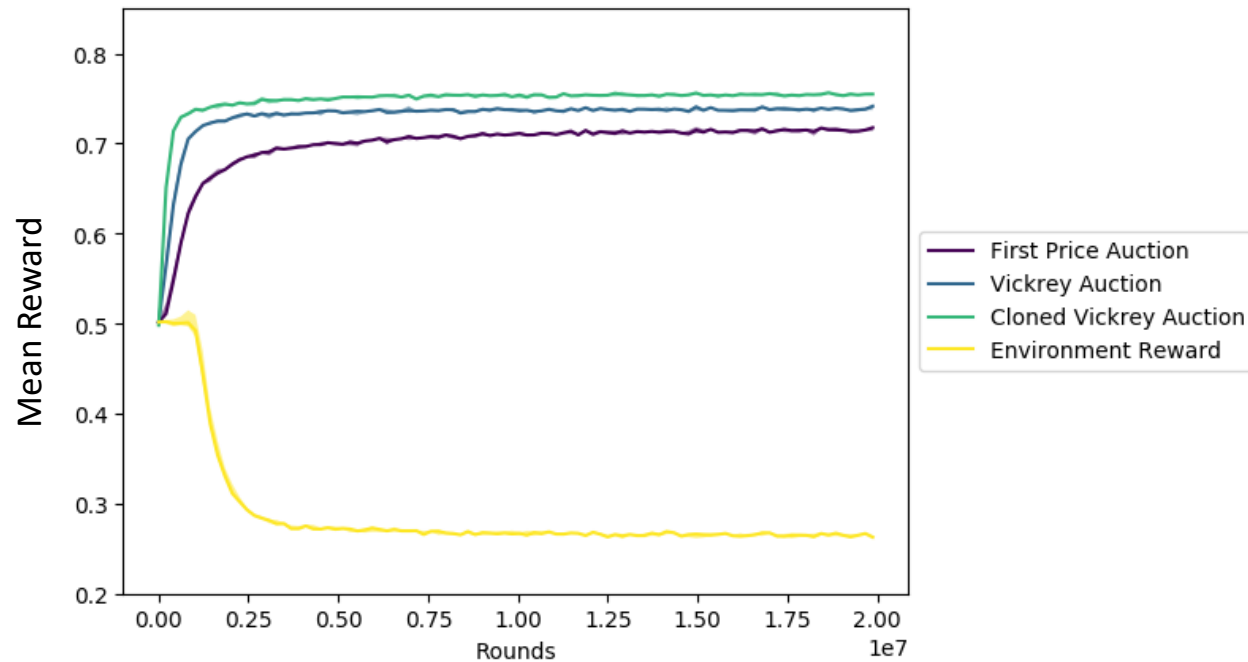


How closely do the bids the agents learn match their optimal Q-values?

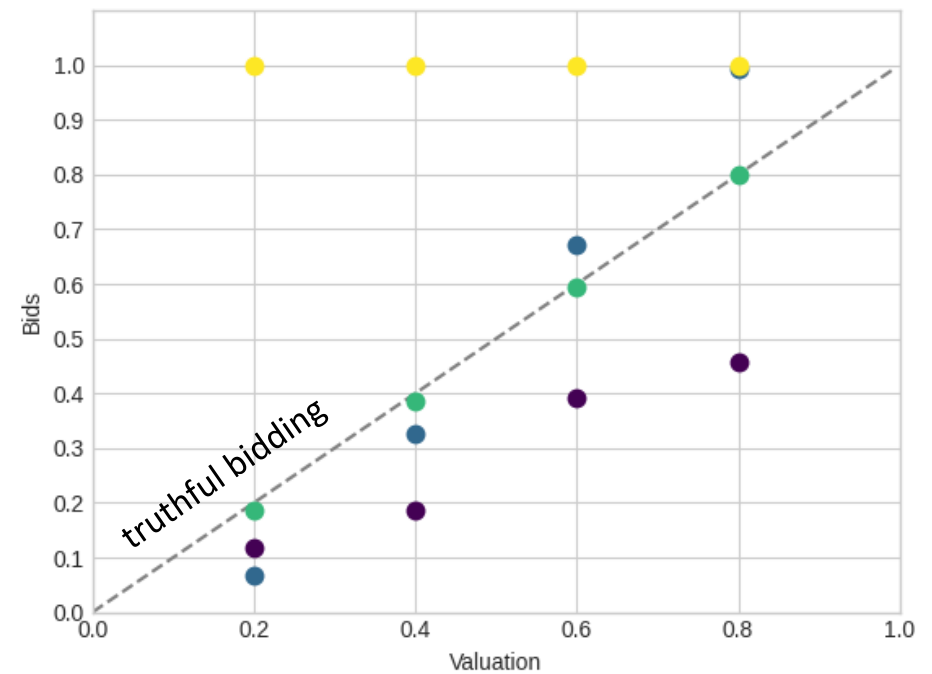


Warm-Up: Bandit

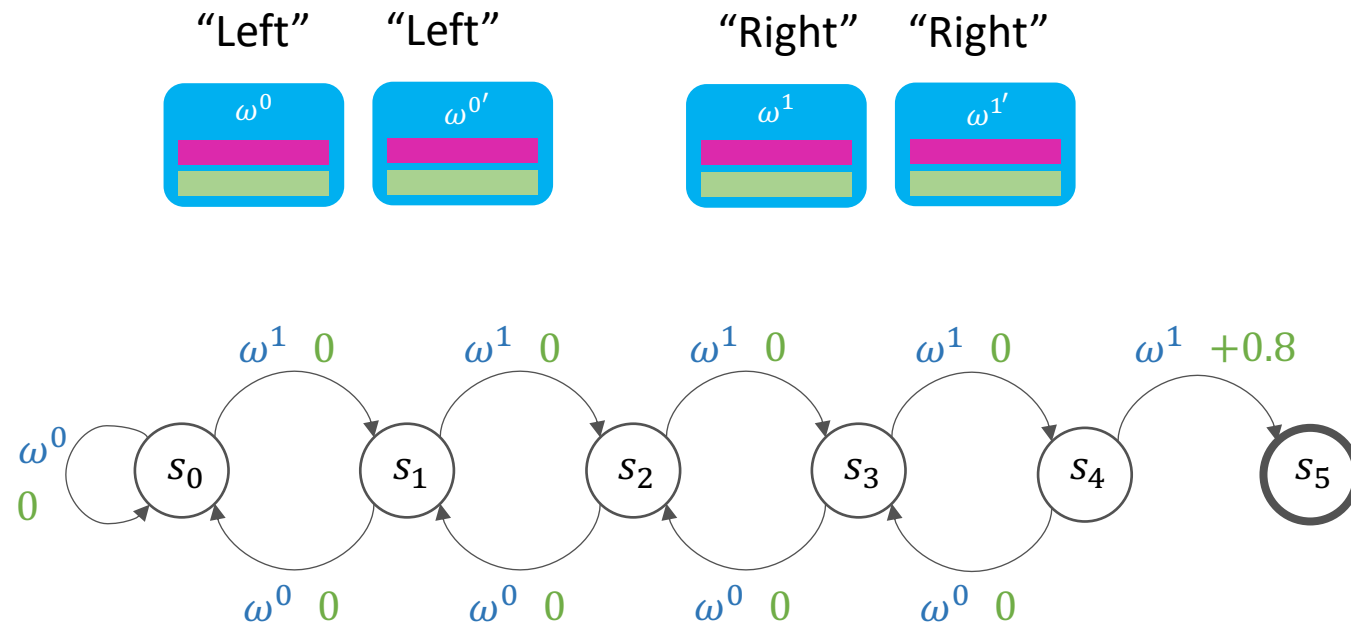
Does the solution to the global objective emerge from the competition among the agents?



How closely do the bids the agents learn match their optimal Q-values?



Multi-Step MDP

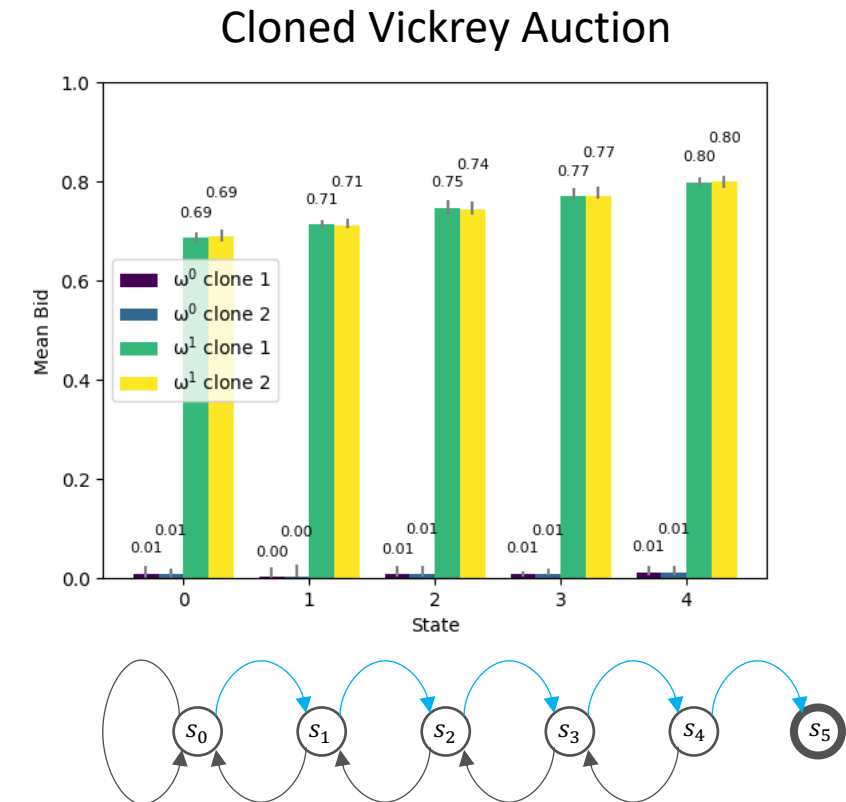


Global Objective for the Society
Maximize return

Local Objectives for the Agents
Maximize utility in the auction

Multi-Step MDP

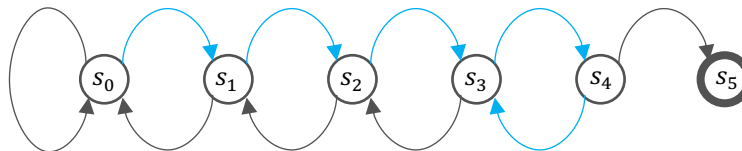
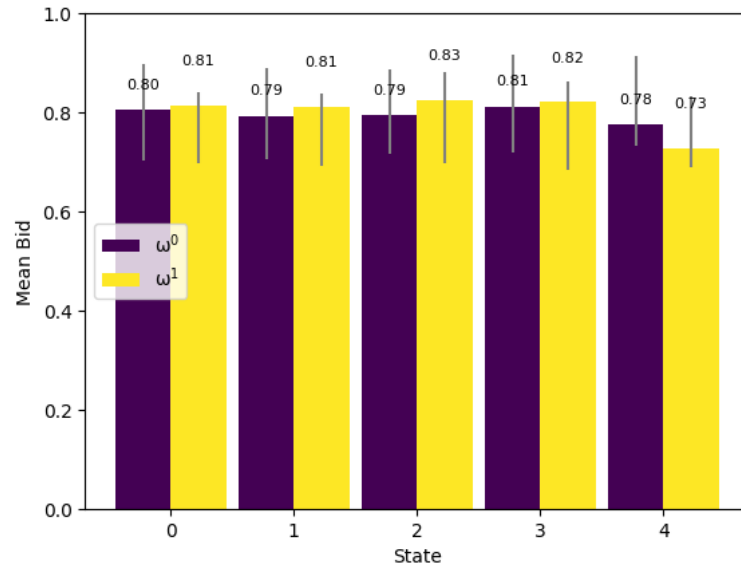
How closely do the bids the agents learn match their optimal Q-values?



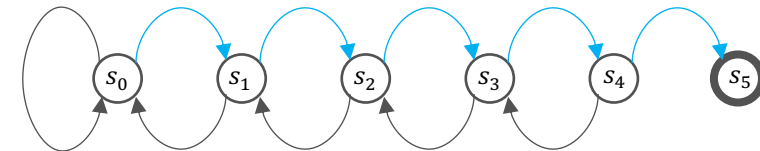
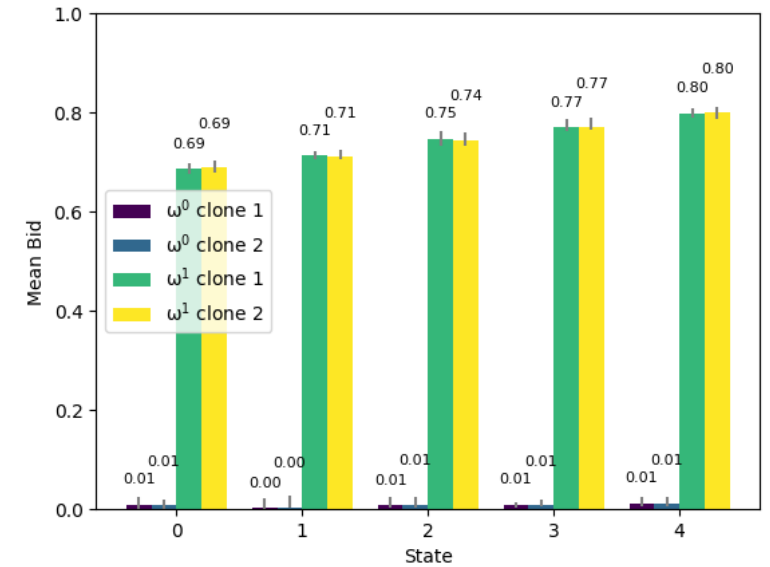
Multi-Step MDP

How closely do the bids the agents learn match their optimal Q-values?

Vickrey Auction



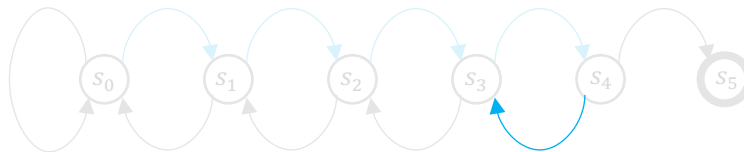
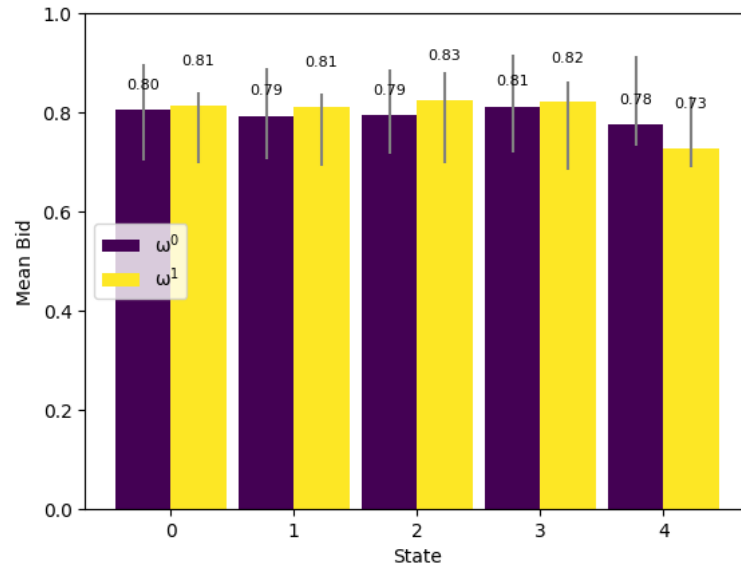
Cloned Vickrey Auction



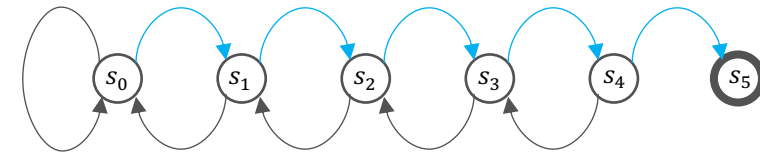
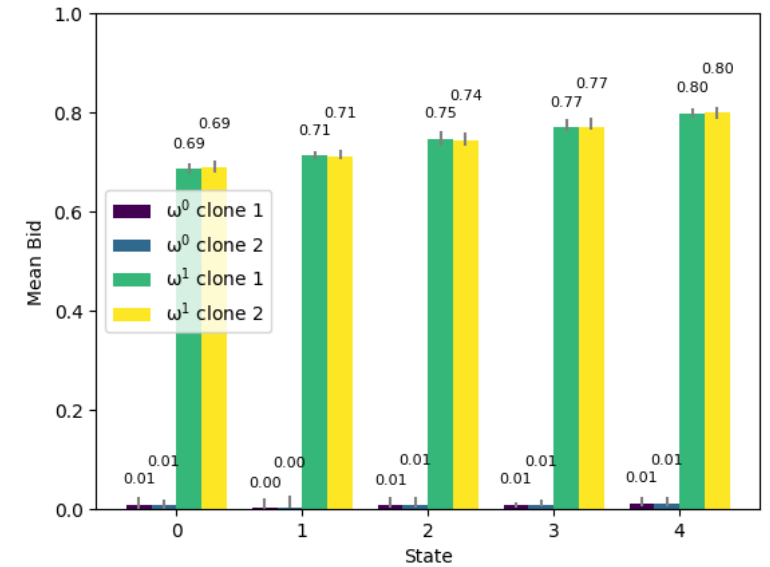
Multi-Step MDP

How closely do the bids the agents learn match their optimal Q-values?

Vickrey Auction



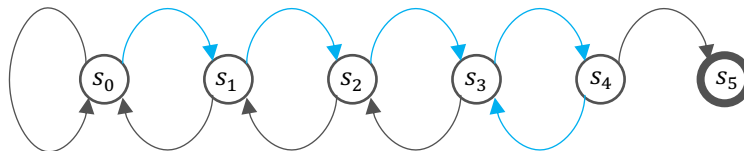
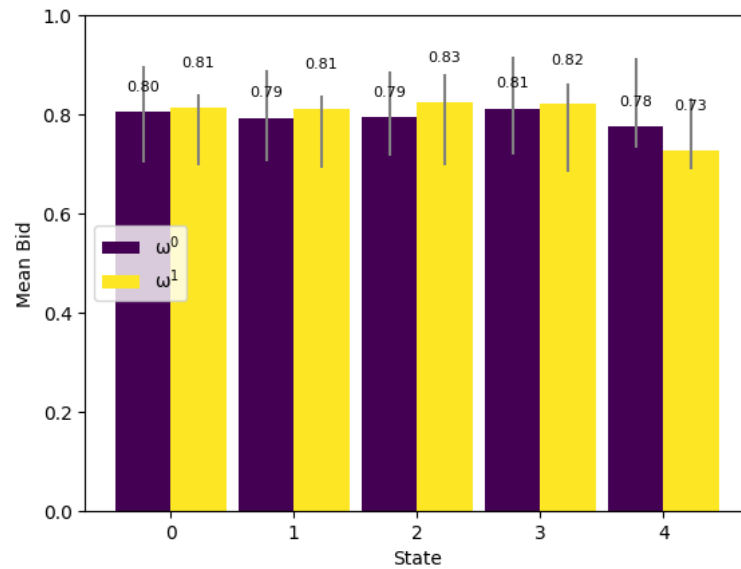
Cloned Vickrey Auction



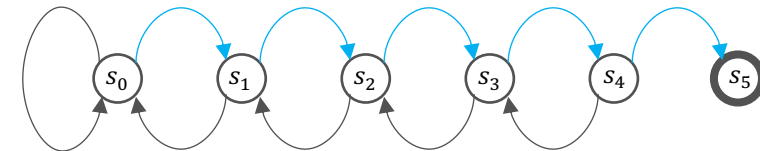
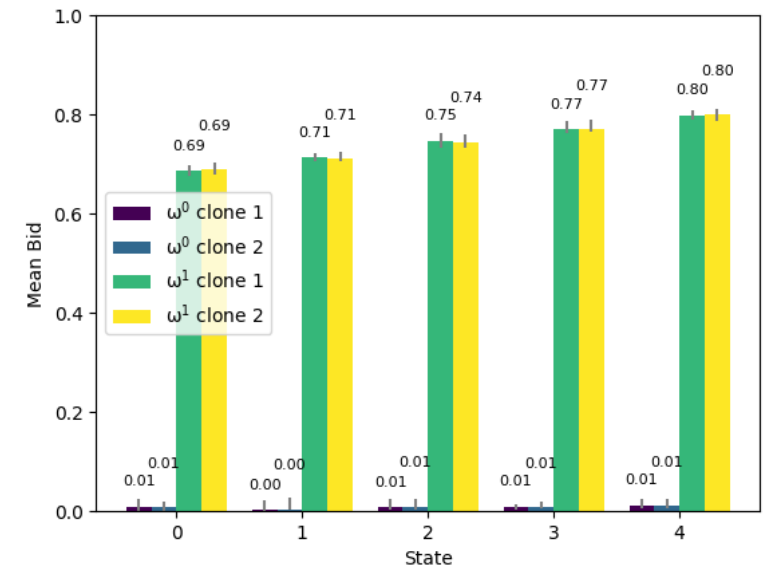
Multi-Step MDP

How closely do the bids the agents learn match their optimal Q-values?

Vickrey Auction



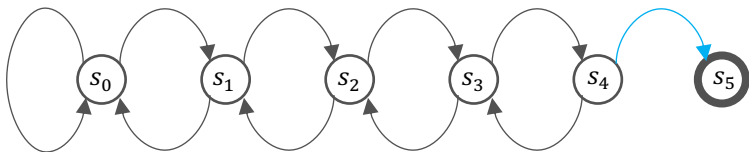
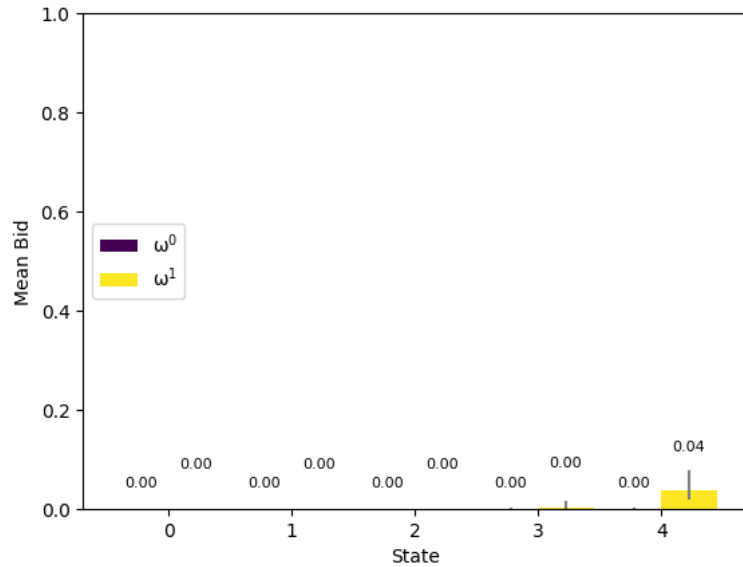
Cloned Vickrey Auction



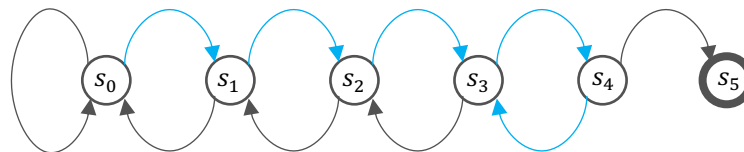
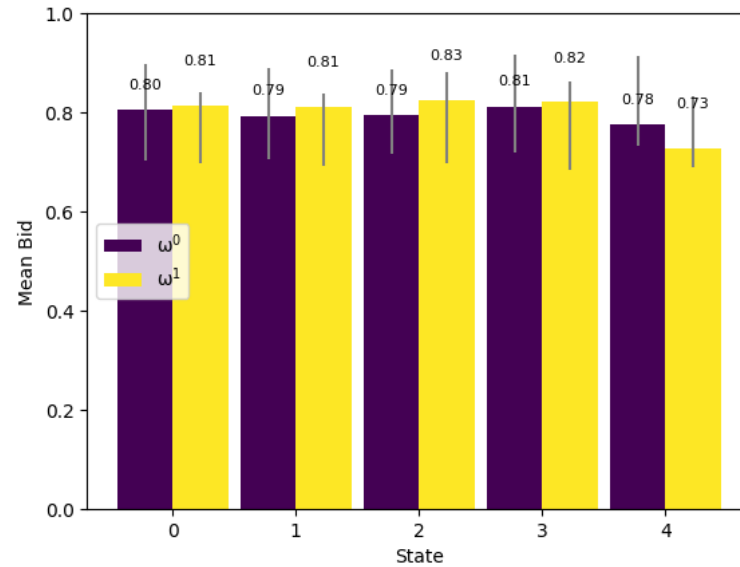
Multi-Step MDP

How closely do the bids the agents learn match their optimal Q-values?

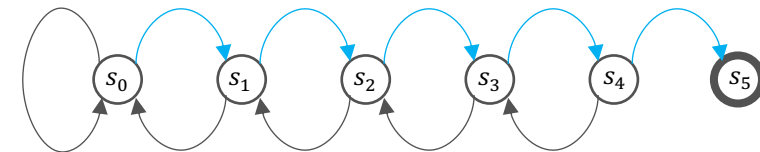
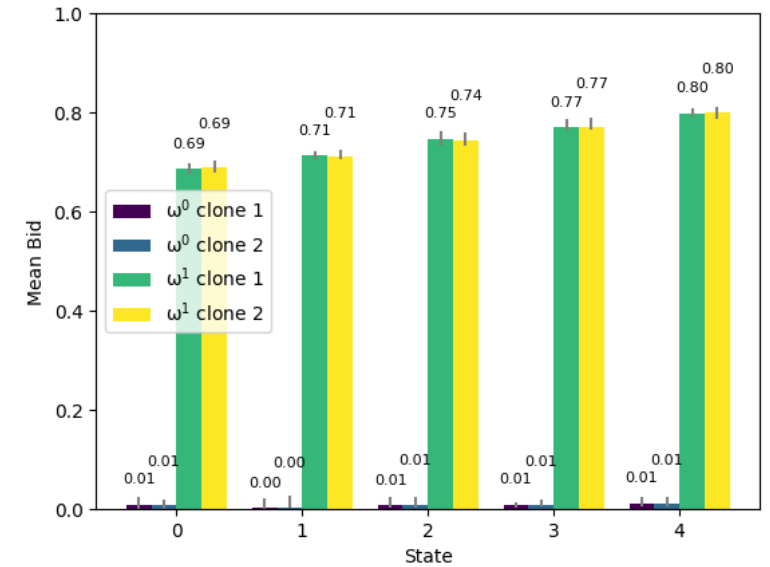
First Price Auction



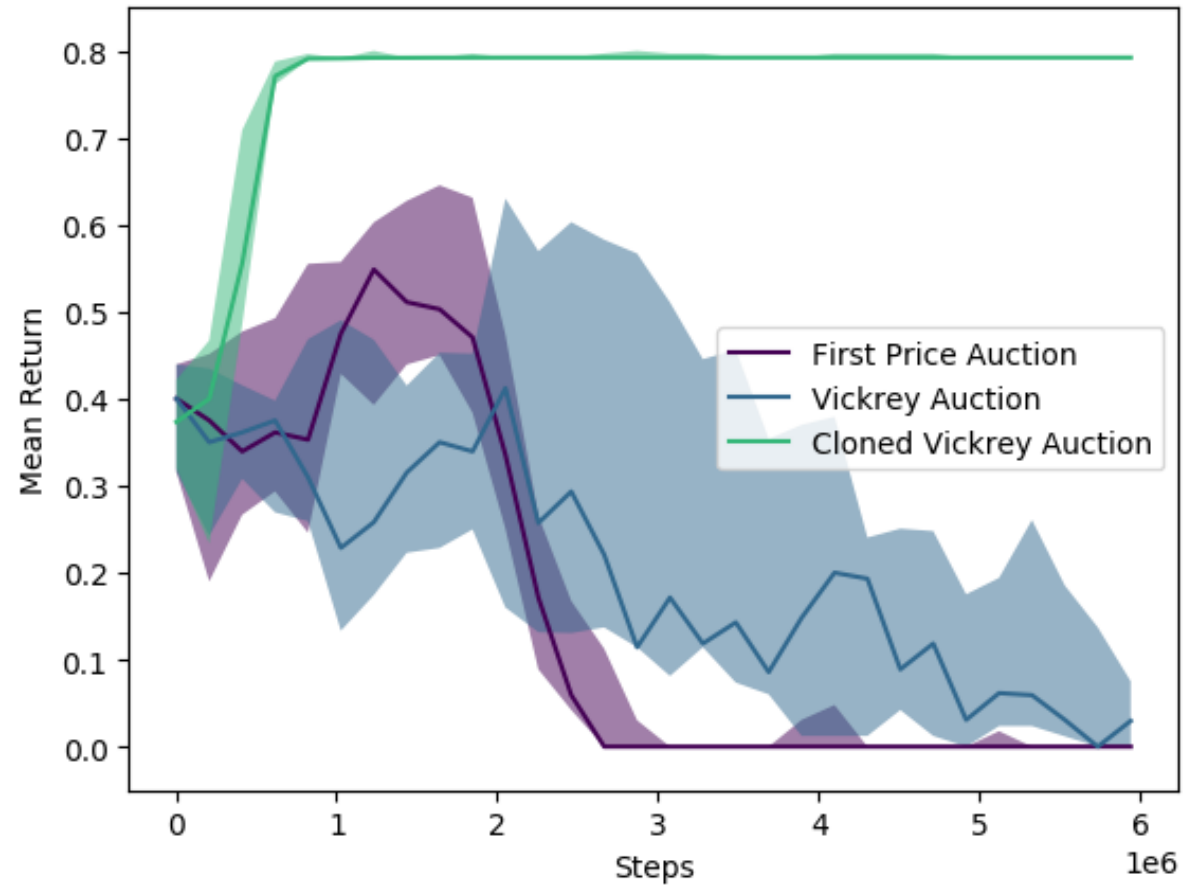
Vickrey Auction



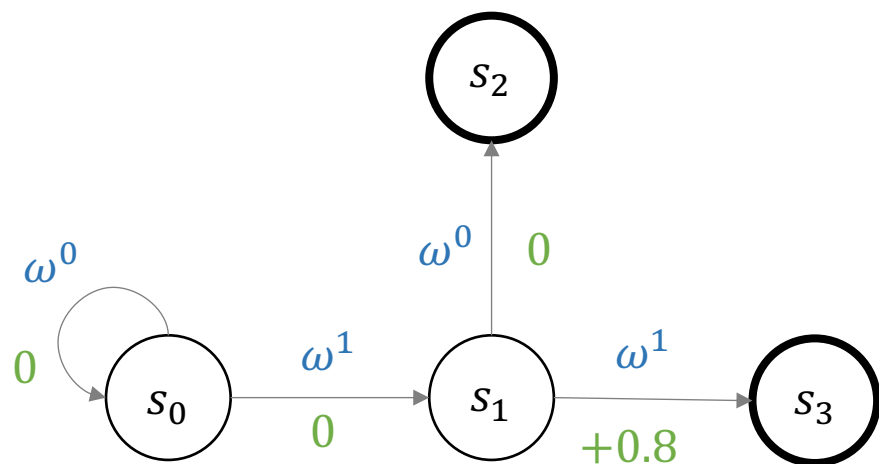
Cloned Vickrey Auction



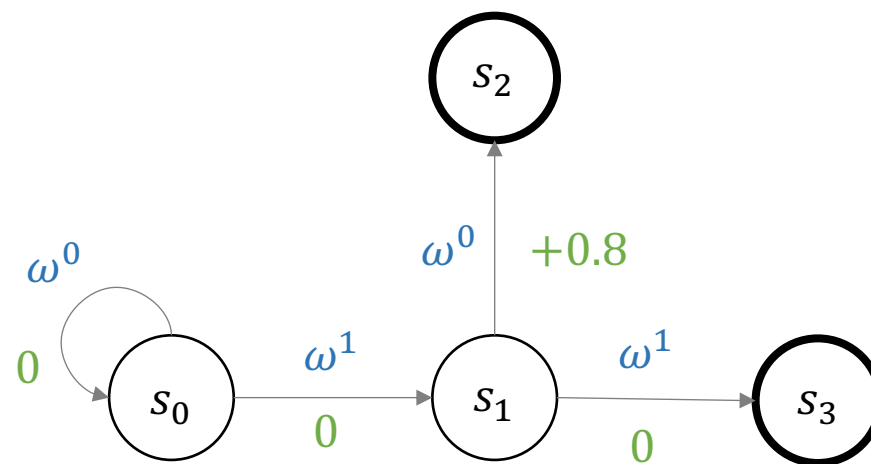
Multi-Step MDP



Transfer



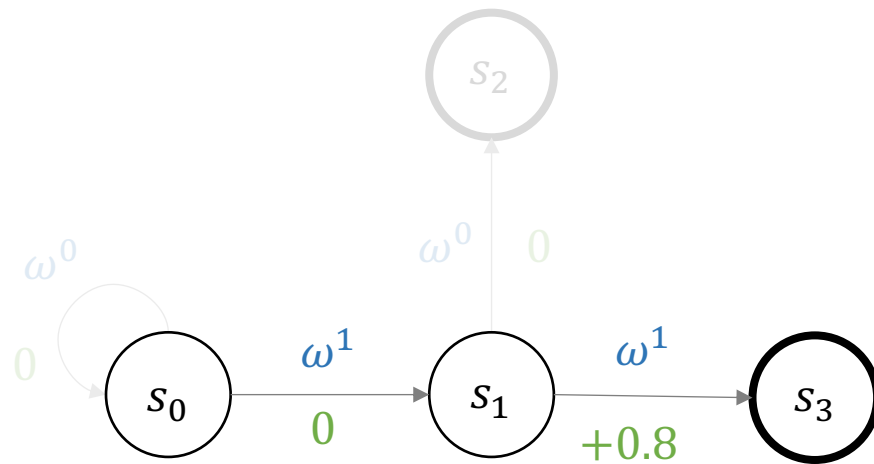
Pre-training Task



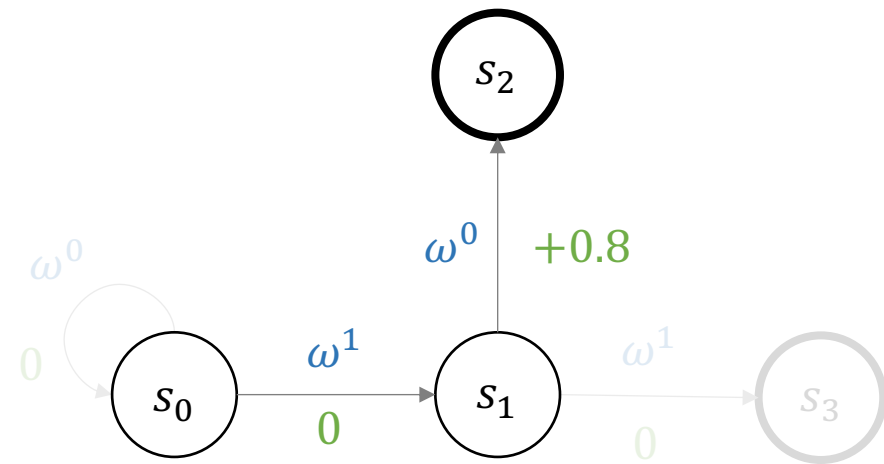
Transfer Task

Transfer

Optimal Policy for the Society



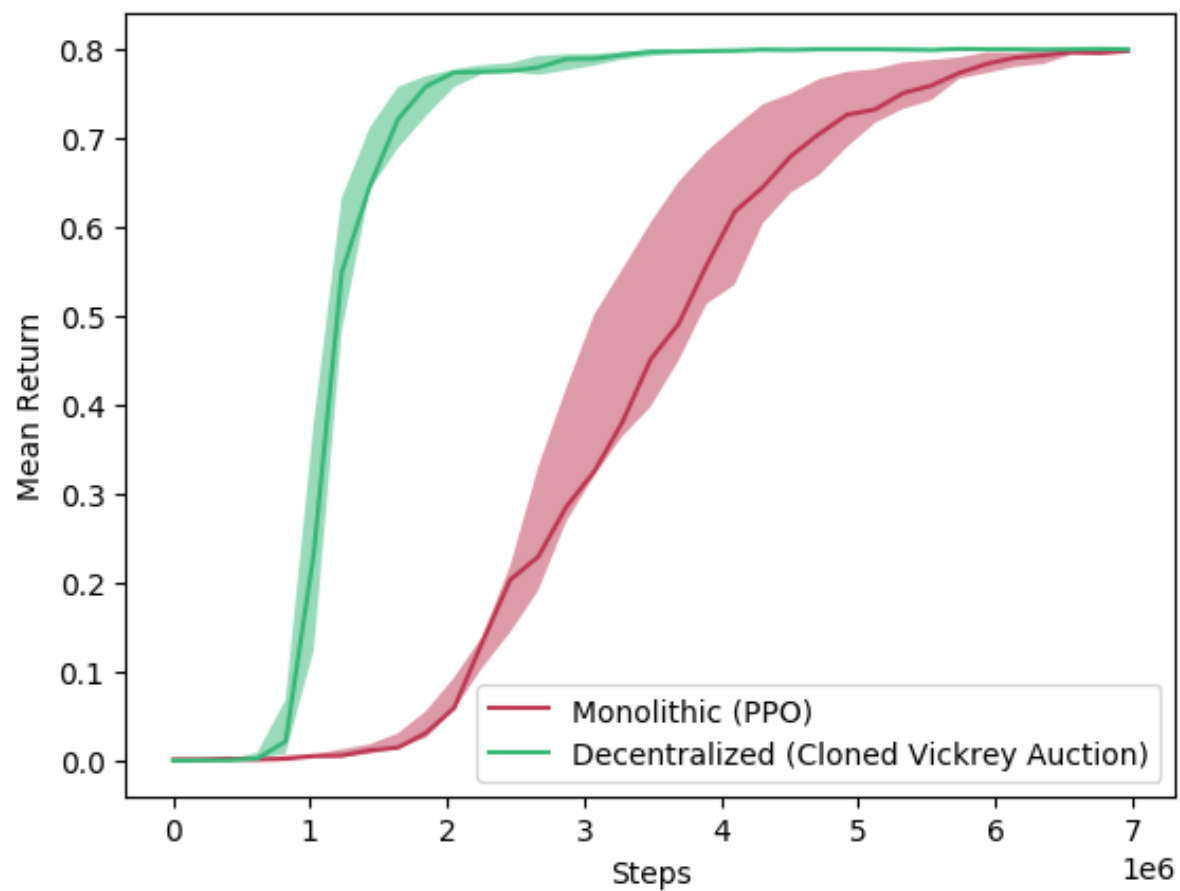
Pre-training Task



Transfer Task

Transfer

Continuing to Train on the Transfer Task



Contributions

Question	Key Idea
What should the optimal bids be for the solution of the Global MDP to emerge?	Define the optimal bid as the optimal Q value $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .
For what auction mechanism would these optimal bids be an equilibrium strategy?	By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a dominant strategy to truthfully bid $Q^*(s, \omega^i)$.
How can we adapt this auction mechanism for discrete-action MDPs?	Temporally couple the agents in a market: An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t .
How can we avoid suboptimal equilibria?	Redundancy enforces credit conservation that helps avoid suboptimal equilibria.
How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?	Define the auction utility as the agents' reinforcement learning objective, yielding a decentralized reinforcement learning algorithm for the Global MDP.

<https://sites.google.com/view/clonedvickreysociety>

Contributions

Cloned Vickrey Society

A society of agents that implements global decision making via local economic transactions.

Question

Key Idea

What should the optimal bids be for the solution of the Global MDP to emerge?

Define the optimal bid as the **optimal Q value** $Q^*(s_t, \omega^i)$ for activating agent ω^i at state s_t .

For what auction mechanism would these optimal bids be an equilibrium strategy?

By defining the agents' valuations $v^i(s)$ as $Q^*(s, \omega^i)$, under the Vickrey auction it is a **dominant strategy** to truthfully bid $Q^*(s, \omega^i)$.

How can we adapt this auction mechanism for discrete-action MDPs?

Temporally couple the agents in a market: An agent's valuation of s_t is defined by how much it can sell the product s_{t+1} of executing its transformation on s_t .

How can we avoid suboptimal equilibria?

Redundancy enforces credit conservation that helps avoid suboptimal equilibria.

How can we translate the auction mechanism into a decentralized reinforcement learning algorithm?

Define the **auction utility** as the agents' reinforcement learning objective, yielding a **decentralized reinforcement learning algorithm** for the Global MDP.

<https://sites.google.com/view/clonedvickreysociety>