T10. A Concise Introduction to Cooperative Multi-Agent Reinforcement Learning

AAMAS'2025 Tutorial

Part 1

Christopher Amato & Frans A. Oliehoek





Outline

Part 1:

- Multiagent decision problems
- Toward multiagent RL
- Some fundaments from multiagent planning
- Multiagent RL through the lens of Dec-POMDPs



available from my website





Part 1a: Multiagent decision problems







Why Multiagent Systems (MASs)?

why do intelligent agents (humans) interact the way they do?

subjective perspective:

 focus on protagonist agent

Why MASs?

Delft

• will need to interact with other agents!





Engineer

(John Nash – foto by Elke Wetzig) Amato&Oliehoek - Cooperative MARL I have a bunch of robots... how do I get them to do something useful?



objective perspective:

focus on team

Why MASs?

- flexible, robust, cost-effective.
- reduces complexity: agent can take care of some aspects itself.



Why Multiagent Reinforcement Learning?

protagonist needs to learn about other agents design tool...! (MAS are hard to program!)







(John Nash – foto by Elke Wetzig) Amato&Oliehoek - Cooperative MARL



MASs under Uncertainty













Uncertainty galore: ►Outcome Uncertainty ►Partial Observability ►about others

″UDelft



And will this work...?

- Well... MASs and MARL is complex...!
- but some encouraging developments...
 - Alpha Go
 - Deepstack
 - Capture the flag:







Single agent problems







Single-Agent Decision Making

- An MDP
 - *S* set of states
 - A set of actions
 - P_T transition function: P(s'|s,a)
 - *R* reward function: *R*(*s*,*a*) or *R*(*s*,*a*,*s'*)
 - *T* horizon (finite)
- A POMDP
 - O set of observations
 - *P*₀ observation function: *P*(*o* | *a*,*s'*)
- Initial state distribution: b₀







Example: Predator-Prey Domain

- Predator-Prey domain
 - 1 agent: predator
 - prey: part of environment
 - on a torus
- Formalization:
 - states (-3,4)
 - actions N,W,S,E
 - transitions failing to move, prey moves
 - rewards reward for capturing





Partial Observability

- Now: partial observability
 - E.g., limited range of sight
- MDP + observations
 - explicit observations
 - observation probabilities
 - noisy observations (detection probability)



o = 'nothing'





Partial Observability

- Now: partial observability
 - E.g., limited range of sight
- MDP + observations
 - explicit observations
 - observation probabilities
 - noisy observations (detection probability)



o = (-1, 1)





"Solving" single agent problems

Fully observable

F.O. planning

- generalized policy iteration
- online planning: e.g. MCTS

Partially observable

P.O. planning

- belief MDP, POMDP VI
- online planning e.g. POMCP

Model / simulator available





"Solving" single agent problems

Fully observable

F.O. planning

- generalized policy iteration
- online planning: e.g. MCTS

FORL

- (deep) Q-learning, SARSA, etc.
- (deep) model-based RL

In practice: RL methods is (almost only) used when a simulator is available Partially observable

P.O. planning

- belief MDP, POMDP VI
- online planning e.g. POMCP

Model / simulator available

PORL

- policy gradient with FSC
- deep RL with recurrent layers

No model / simulator available





[Corneil et al. 2018 ICML]



Multiple Agents







Multiple Agents (Fully observable)

- Now: multiple agents
 - fully observable

- Formalization:
 - states ((3,-4), (1,1), (-2,0))
 - actions {N,W,S,E}
 - **joint** actions {(N,N,N), (N,N,W),...,(E,E,E)}
 - transitions
 - rewards

probability of failing to move, prey moves

reward for capturing jointly





Multiple Agents (Fully observable)

- Multiagent MDP [Boutilier 1996]
 - n agents
 - joint actions *a* = <*a*₁,...*a*_n>
 - transitions and rewards depend on these joint actions



- Solution:
 - Treat as normal MDP with 1 'puppeteer agent' \rightarrow optimal policy $\pi^* \rightarrow$ specifies *a* for each *s*
 - Every agent simple executes its part: a_i



Multiple Agents & Partial Observability

- Now both...
 - partial observability
 - multiple agents
- "Decentralized POMDP" (Dec-POMDP) framework*
- both
 - joint actions and
 - joint observations

*See, e.g., "A Concise Introduction to Decentralized POMDPs" http://www.fransoliehoek.net/publications/htmlfiles/b2hd-OliehoekAmato16book.html Delft Amato&Oliehoek - Cooperative MARL





The Formal Dec-POMDP Model

- A Dec-POMDP:
 - *n* agents
 - S set of states
 - *A* set of **joint** actions *a* = <*a*₁,...,*a*_n>
 - *P*₇ transition function: *P(s'*|*s*,*a*)
 - *O* set of **joint** observations **o** = <*o*₁,...,*o*_n>
 - P_o observation function: P(o | a, s')
 - *R* reward function: *R*(*s*,*a*)
 - *T* horizon (finite)



►agents indices are generally subscript



Goal:

- Find the optimal **joint** policy $\pi^* = \langle \pi_{1,} \pi_2 \rangle$
- What is the optimal one?
 - Define **value** as the expected (discounted) sum of rewards: $\begin{bmatrix} T \\ T \end{bmatrix}$

$$V(\pi) = \boldsymbol{E} \left[\sum_{t=0}^{T-1} \boldsymbol{R}(\boldsymbol{s}_t, \boldsymbol{a}_t) \mid \pi, \boldsymbol{b}^0 \right]$$

• an optimal joint policy is one with maximal value





Acting in Dec-POMDPs...

- No clear reductions to single agent case...
- Idea: use joint belief, *b(s)* [Pynadath and Tambe 2002]
 - compute *b(s)* using joint actions and observations
 - Problem: agents do not know those during execution***
- So... now what?
 - How do we plan, when we have the model?
 - How do we learn, otherwise?

******* "puppeteer reduction" requires broadcasting observations!

Under instantaneous, cost-free, noise-free communication this is optimal [Pynadath and Tambe 2002]



"Solving" multiagent problems



*** Again MARL methods are (almost only) used when a simulator is available



Part 1b: Towards multiagent RL (Two first ideas)







Idea 1: Individual learning

• E.g. just use Q-learning per agent:

$$Q(s,a_i) := (1-\alpha) Q(s,a_i) + \alpha [r+\gamma V(s')]$$



- Can work well...but the environment is changing
 - "non-stationary learning problem"
- Convergence?
 - State observable \rightarrow converges to a **local** optimum
 - Otherwise \rightarrow no guarantees



Idea 2: Coupled Learners

• Multiagent MDP: MDP with joint actions...



- E.g. "joint Q-learning": $Q(s, a) = (1-\alpha) Q(s, a) + \alpha [r + \gamma max_{a'} Q(s', a')]$
 - A.k.a.: "Joint action learners"
 - Note: can be implemented decentrally in FORL (each agent executes its component)
 - Guarantees:
 - Same guarantees as normal Q-learning
 - will converge in the limit
 - Scalability: **exponential** in number of agent...



Limitations so far....

- OK. These were reasonable ideas... but:
 - individual learners \rightarrow no guarantees (non-stationarity)
 - coupled learners \rightarrow not scalable
- Plus... limited applicability
 - agents need to observe full state
 - in most MASs this is not possible

We need formal frameworks that:
► can deal with partial observability
► and represents knowledge: who learned what?







*** Again MARL methods are (almost only) used when a simulator is available



Part 1c: Some fundaments from multiagent planning







Dec-POMDPs: Off-line / On-line phases

When is planning and/or learning taking place?



- True online learning
- Planning / Simulation based planning
 - plan offline given model
 - do 'learning' offline using simulator
 - limits applicability... but common assumption in MARL (and in fact in nearly all RL)





Running Example

• 2 generals problem

Delft



Running Example

- 2 generals problem
- S { s_L, s_s } A_i – { (O)bserve, (A)ttack } O_i – { (L)arge, (S)mall }

Transitions

► Both Observe → no state change

► At least 1 Attack \rightarrow reset (50% probability s_L, s_S)

Observations

- ► Probability of correct observation: 0.85
- ► E.g., P(<L, L> | s_L) = 0.85 * 0.85 = 0.7225
- ►(reset is not observed!)

suppose T=3, what do you think is optimal in this problem?

Rewards

- ►1 general attacks: he loses the battle
 - R(*,<A,O>) = -10
- Both generals Observe: small cost
 - R(*,<0,O>) = -1
- Both Attack: depends on state
 - R(s_L,<A,A>) = -20
 - R(s_s,<A,A>) = +5

Policy Domain

- What do policies look like?
 - In general (action-observation) histories \rightarrow actions
 - in MDP/POMDP: compact representations: states/beliefs
- For Dec-POMDPs: no such representation known!
 → If we want optimal policies: we are stuck with histories
- Of course, can try and compress...
 - \rightarrow approximate internal states I_i
 - cf. RNNs

more general picture





The general picture



33



Observation histories and policy trees

- What type of histories?
 - observation histories

In cooperative case **deterministic** policies only need OHs:

$$\vec{o}_i = (o_i^1, \dots, o_i^t)$$



Policies for Two Generals...

Optimal policy for 2 generals, h=3 value=-2.86743

General 1: () --> observe (o_small) --> observe (o_large) --> observe (o_small,o_small) --> attack (o_small,o_large) --> attack (o_large,o_small) --> attack (o_large,o_large) --> observe General 2: () --> observe (o_small) --> observe (o_large) --> observe (o_small,o_small) --> attack (o_small,o_large) --> attack (o_large,o_small) --> attack (o_large,o_large) --> observe





Policies for Two Generals...

Optimal policy for 2 generals, h=3 value=-2.86743

General 1: () --> observe (o_small) --> observe (o_large) --> observe (o_small,o_small) --> attack (o_small,o_large) --> attack (o_large,o_small) --> attack (o_large,o_large) --> observe Anything that seems strange...?

General 2: () --> observe (o_small) --> observe (o_large) --> observe (o_small,o_small) --> attack (o_small,o_large) --> attack (o_large,o_small) --> attack (o_large,o_large) --> observe




Policies for Two Generals...

Anything that seems strange...? Optimal policy for 2 generals, h=3 value=-2.86743 General 1: General 2: () --> obs<u>erve</u> ()--> observe (o_small) --> observe (o_small) --> observe (o_large) --> observe (o_large) --> observe (o_small_o_small) --> attack (o_small,o_small) --> attack (o_small,o_large) --> attack (o_small,o_large) --> attack (o_large,o_small) --> attack (o_large,o_small) --> attack (o_large,o_large) --> observe (o_large,o_large) --> observe





Coordination vs. Exploitation

• Inherent trade-off:

coordination vs. exploitation of local information

- Ignore own observations → 'open loop plan'
 - E.g., "ATTACK on 2nd time step"
 - + maximally predictable
 - low quality
- Ignore coordination \rightarrow 'individual POMDP plan'
 - E.g., try to form an 'individual belief' $b_i(s)$
 - (e.g., assume other agents act random...)
 - + uses local information
 - likely to result in mis-coordination

• Optimal policy should balance between these!





Dec-POMDP Planning Techniques & Optimal value functions







Value of a Joint Policy - history perspective

• Iterative:

$$V(\boldsymbol{\pi}) = \sum_{t=0}^{T-1} \sum_{\boldsymbol{h}_t \in \boldsymbol{\mathcal{H}}_t} \Pr(\boldsymbol{h}_t | \boldsymbol{\pi}, b_0) \sum_{\boldsymbol{a}_t \in \boldsymbol{\mathcal{A}}} R(\boldsymbol{h}_t, \boldsymbol{a}_t) \boldsymbol{\pi}(\boldsymbol{a}_t | \boldsymbol{h}_t),$$

$$R(\boldsymbol{h}_t, \boldsymbol{a}_t) = \sum_{s_t \in \mathcal{S}} R(s_t, \boldsymbol{a}_t) \Pr(s_t | \boldsymbol{h}_t, b_0)$$





Value of a Joint Policy - history perspective

• Iterative: $V(\boldsymbol{\pi}) = \sum_{t=0}^{T-1} \sum_{\boldsymbol{h}_t \in \boldsymbol{\mathcal{H}}_t} \Pr(\boldsymbol{h}_t | \boldsymbol{\pi}, b_0) \sum_{\boldsymbol{a}_t \in \boldsymbol{\mathcal{A}}} R(\boldsymbol{h}_t, \boldsymbol{a}_t) \boldsymbol{\pi}(\boldsymbol{a}_t | \boldsymbol{h}_t),$

$$R(\boldsymbol{h}_t, \boldsymbol{a}_t) = \sum_{s_t \in \mathcal{S}} R(s_t, \boldsymbol{a}_t) \Pr(s_t | \boldsymbol{h}_t, b_0)$$

 Recursive: Bellman equation on joint AO-histories

$$V^{\boldsymbol{\pi}}(\boldsymbol{h}_t) = \sum_{\boldsymbol{a}_t \in \boldsymbol{\mathcal{A}}} \boldsymbol{\pi}(\boldsymbol{a}_t | \boldsymbol{h}_t) Q^{\boldsymbol{\pi}}(\boldsymbol{h}_t, \boldsymbol{a}_t)$$
$$Q^{\boldsymbol{\pi}}(\boldsymbol{h}_t, \boldsymbol{a}_t) = R(\boldsymbol{h}_t, \boldsymbol{a}_t) + \sum_{\boldsymbol{o}_{t+1} \in \boldsymbol{\mathcal{O}}} \Pr(\boldsymbol{o}_{t+1} | \boldsymbol{h}_t, \boldsymbol{a}_t) V^{\boldsymbol{\pi}}(\boldsymbol{h}_{t+1})$$



value propagation on a tree



Value of a Joint Policy - subtree policy perspective

• Sub-tree policies:



• Different formulation of value:

$$V(s, q^{\tau=k}) = R(s, a) + \Sigma_o P(o|s, a) V(s', q^{\tau=k-1})$$

• no need to explicitly remember AOHs... but number of q's is huge!





Brute Force Search

- We can compute the value of a joint policy $V(\pi)$
- So the **stupidest algorithm** is:
 - compute $V(\pi)$, for all π
 - select a π with maximum value
- Number of joint policies is huge! (doubly exponential in horizon *h*)
- Clearly intractable...

h	num. joint policies		
1	4		
2	64		
3	16384		
4	1.0737e+09		
5	4.6117e+18		
6	8.5071e+37		
7	2.8948e+76		
8	3.3520e+153		



Why Dec-POMDP planning?

- Finding an optimal plan is (very!) intractable: **NEXP-complete**
- So... should we not just give up?
- Perhaps, but many reasons to care:
 - Interesting..!
 - Understand the problem better
 - Problems do not get easier by ignoring their complexity (and we want collaborating agents... right?)
 - Theory of MDPs (e.g., value functions) are the foundation of RL
 → for effective MARL, we need Dec-POMDP theory.









Bottom-up vs. Top-down

- DP constructs bottom-up... alternative: top-down
 - → heuristic search [Szer et al. 2005, Oliehoek et al. 2008]





Amato&Oliehoek - Cooperative MARL



- Incrementally construct all (joint) policies
 - 'forward in time'

elft





- Incrementally construct all (joint) policies
 - Iforward in time

1 partial joint policy

50



- Incrementally construct all (joint) policies
 - 'forward in time'

1 partial joint policy

51



- Incrementally construct all (joint) policies
 - Iforward in time

1 partial joint policy

52



- Incrementally construct all (joint) policies
 - 'forward in time'

elft





• Creating **ALL** joint policies \rightarrow tree structure!



Root node: unspecified joint policy





• Creating **ALL** joint policies \rightarrow tree structure!



• Creating **ALL** joint policies \rightarrow tree structure!



• Creating **ALL** joint policies \rightarrow tree structure!











Amato&Oliehoek - Cooperative MARL





- Too big to create completely...
- Idea: use **heuristics**
 - avoid going down non-promising branches!



- Apply A* → **Multiagent A*** [Szer et al. 2005]
- Techniques for further scaling [Oliehoek et al. 2013 JAIR]:
 - lossless clustering of histories
 - Incremental expansion







How About Optimal Value Functions?

- We saw a value function **for a given joint (sub-tree) policy** ...how about **optimal** value functions...?
- E.g., how about something like

 $Q^{*}(h,a) = R(h,a) + \Sigma_{o} P(o \mid h,a) V^{*}(h' = <h,a,o>)$ $V^{*}(h') = max_{a'} Q^{*}(h',a')$

- ?
- A bit tricky...
 - would work for "multiagent POMDP"
 - but for not Dec-POMDPs: policies need to be decentralized!



not possible...! requires agent 1 to select different actions, while it gets the same observation ('S')



Reinterpreting the GMAA search tree

• Can view this as the decision making of the "planner"







Optimal value function of "Plan-time MDP"

- Leads to "Plan-time MDP"
 - states $\varphi_t = \langle \delta_0, ..., \delta_{t-1} \rangle$
 - actions δ_t
- That has Bellman optimality equations:

 $Q^*(\varphi_t, \delta_t) = R(\varphi_t, \delta_t) + V^*(\varphi_{t+1} = \langle \varphi_t, \delta_t \rangle)$

 $V^*(\varphi_t) = \max_{\delta} Q^*(\varphi_t, \delta)$

• With $R(\varphi_t, \delta_t)$ the expected reward at stage t: $R(\varphi_t, \delta_t) = \Sigma_s \Sigma_h P(s_t, h_t | \varphi_t) R(s_t, \delta_t(h_t))$

(no bold to increase readability, but all entities are 'joint') **TUDelft**Amato&Oliehoek - Cooperative MARE





The view with Plan-time statistics

• The σ_t allow for reuse [Oliehoek 2013 IJCAI]



 And it turns out that value function is PWLC on σ_t ... [Dibangoye et al. 2013 IJCAI]
 Dolft



... so yes it is a (special case of) POMDP!

- A plan-time non-observable MDP (NOMDP)! [Oliehoek and Amato '14]
 - States: $\langle s_t, h_t \rangle$ (or $\langle s_t, o_t \rangle$)
 - Actions: δ_t
 - Observations: NULL
 - In this NOMDP, the planner's belief is the "plan-time statistics" $\sigma(s_t, h_t)$
- Extend to deal with common (i.e. shared) information [Nayyar et al 2013]:
 → Plan-time POMDP
 - States: $\langle s_t, h_t \rangle \leftarrow h_t$ are joint histories of private information
 - Actions: δ_t
 - Observations: {o_{common}}
 - For each history of common observations h_{common} we get a different $\sigma(s_t, h_t)$
 - Select actions based on $<h_{common}, \sigma(s_t, h_t) >$
 - E.g., used in MARL for Hanabi in "Bayesian action decoder" [Foerster et al. 2019]

Terminology in decentralized control:

plan-time ... =
"the designer's approach"





Of course... histories still don't scale

- E.g., for infinite horizon...? \rightarrow do some compression on memory
- One option: finite-state controllers
 - each agent has information state *I_i*
 - and updates in some way: $I_i' = \iota(I_i, a_i, o)$
- Can incorporate in definition of Dec-POMDP problem...:

Definition 7 (ISA-Dec-POMDP). A Dec-POMDP with information state abstraction (ISA-Dec-POMDP) is a Dec-POMDP framework together with the specification of the sets $\{\mathcal{I}_i\}$ of information states (ISs).

For an ISA-Dec-POMDP, using the notation of the agent components defined above, there are two optimizations that need to be performed jointly:

- 1. the optimization of the joint action selection policy $\boldsymbol{\pi} = \langle \pi_1, \ldots, \pi_n \rangle$, and
- 2. the optimization of the joint information state function $\iota = \langle \iota_1, \ldots, \iota_n \rangle$.





Plan-time sufficient statistics & NOMDP formulation

• Can extend these concepts...

Definition 8 (Plan-time ISA sufficient statistic). The plan-time sufficient statistic for an ISA-Dec-POMDP is $\sigma_t(s,I) \triangleq \Pr(s,I|\boldsymbol{\delta}_0,...,\boldsymbol{\delta}_{t-1}).$

> **Definition 9** (Plan-time ISA-NOMDP). Given the internal state update functions, an ISA-Dec-POMDP can be converted to a *plan-time ISA-NOMDP* $\mathcal{M}_{PT-ISA-NOMDP}$ for a Dec-POMDP \mathcal{M} is a tuple $\mathcal{M}_{PT-ISA-NOMDP}(\mathcal{M}) = \langle \check{S}, \check{A}, \check{T}, \check{R}, \check{\mathcal{O}}, \check{O}, \check{h}, \check{b_0} \rangle$, where:

- \check{S} is the set of augmented states, each $\check{s} = \langle s, I \rangle$.
- \check{A} is the set of actions, each \check{a} corresponds to a joint decision rule δ (which is a joint (stationary) action selection policy) in the ISA-Dec-POMDP.
- \check{T} is the transition function:

$$\begin{split} \check{T}(\langle s, I' \rangle \,|\, \langle s, I \rangle, \boldsymbol{\delta}) &= & \Pr(s', I' | s, \boldsymbol{a} = \boldsymbol{\delta}(I)) \\ &= & \Pr(s' | s, \boldsymbol{a}) \sum_{\boldsymbol{o}} \iota(I' | I, \boldsymbol{a}, \boldsymbol{o}) \Pr(\boldsymbol{o} | \boldsymbol{a}, s') \end{split}$$

- \check{R} is the reward function: $\check{R}(\langle s, I \rangle, \delta) = R(s, \delta(I)).$
- $\tilde{\mathcal{O}} = \{NULL\}$ is the observation set which only contains the NULL observation.
- \check{O} is the observation function that specifies that NULL is received with probability 1 (irrespective of the state and action).

l i s | delft ₆₇

Amato&Oliehoek - Cooperative MARL

For more info:

Oliehoek & Amato. "Dec-POMDPs as non-observable MDPs." (2014)

TUDelft

Try out some Dec-POMDPs?

• Try it out!

Product ~ Solution	ns	ing	Q Sign in Sign up	
MADPToolbox / MAI Code ① Issues 2	Public [↑] Pull requests ⊙ Actions ⊞ Projects □ Wiki ①) Security 🗠 Insig	Notifications 😲 Fork 22 🏠 Star 77	
약 master ▾ १ 1 Branch	♥ 2 Tags Q Go to file	<> Code 🔸	About	
🕼 oliehoek Merge pull request #15 from laurimi/libdai-compile-fix 🚥 a6c1bb7 · 5 years ago 🕚 95 Commits			The Multiagent Decision Process (MADP) Toolbox - planning and 2 #include "JESPExhaustivePlanner.h"	
🖿 config	Import of MADP 0.3.1	10 years ago	<pre>int main() www.fransoliehoek.net/madp int main() www.fransoliehoek.net/madp if results and the set of the set o</pre>	
doc	exclude file update	5 years ago		
🖿 git	exclude file update	5 years ago		
m4	exclude file update	5 years ago		
p roblems	problems added for 'make check'	8 years ago		
src src	explicitly reference std::	5 years ago		
🗋 .travis.yml	trying brewsci	5 years ago		

https://github.com/MADPToolbox/MADP





Part 1d: Multiagent RL through the lens of Dec-POMDPs





Non-stationarity revisited...

- So why are we talking about Dec-POMDPs...?
 - \rightarrow gives us the tools to formalize 'non-stationarity'
 - before... "other agent changes"
 - but could not specify *how*...
 - the other agent changes due to individual observations
 - but individual learners do not represent those...!





Non-staticBut my other agents are Q-learners...?• So why are
 \rightarrow gives us• they can
be seen
as policies
 $\pi_i(a_i | h_i) !$ In particular, let's think of π_i as an individual Q-learning agent... it receives its
own observations and remembers its own actions in \vec{h}_i^t , and uses this for updates
 a_i add action selection:• So why are
 \rightarrow gives us the seen
 $a_i (a_i | h_i) !In particular, let's think of <math>\pi_i$ as an individual Q-learning agent... it receives its
own observations and remembers its own actions in \vec{h}_i^t , and uses this for updates
 a_i add action selection:• $f_i(a_i | h_i) !<math>Q_i(o_i^{t-1}, a_i^{t-1}) \leftarrow (1 - \alpha)Q_i(o_i^{t-1}, a_i^{t-1}) + \alpha(r + \gamma \max_{a_i} Q(o_i^t, a_i))$
 $\pi_i(a_i | \vec{h}_i^t) = epsilon-greedy(Q(o_i^t, \cdot), a_i, \epsilon)$

So, even though it does not need to remember \vec{h}_i^t , each \vec{h}_i^t deterministically induces an updated Q_i and thus action selection probabilities.

Example: predicting rewards

• before... "o

Let's predict the reward for agent *i* in a Dec-POMDP, given π_{i}

$$R_{i}(\vec{h}_{i}^{t}, a_{i}) = \sum_{s^{t}, \vec{h}_{-i}^{t}} \Pr(s^{t}, \vec{h}_{-i}^{t} | \vec{h}_{i}^{t}) \sum_{\boldsymbol{a}_{-i}} \pi_{-i}(\boldsymbol{a}_{-i} | \vec{h}_{-i}^{t}) R_{i}(s^{t}, a_{i}, \boldsymbol{a}_{-i})$$

with

- $\pi_{-i}(\boldsymbol{a}_{-i}|\vec{\boldsymbol{h}}_{-i}^t) = \pi_1(a_1^t|\vec{h}_1^t) \times \cdots \times \pi_n(a_n^t|\vec{h}_n^t)$ the application of the policies of the other agents
- $\Pr(s^t, \vec{h}_{-i}^t | \vec{h}_i^t)$ the belief over both states and AOHs of other agents, given our own AOH \vec{h}_i^t . This can be computed, given π_{-i} , in a similar way as the normal belief update in a POMDP.

DELFT 71

S

MARL: an objective perspective









Objective perspective

- Objective perspective:
 reason for the entire team
- Two approaches to decision making:
 - centralized: "puppeteer"

• decentralized







Centralized Objective Approach

• Centralize all information (communication or fully observable)

centralized problem is...



fully observable: multiagent MDP

(standard) RL methods

Example: joint Q-learning!

Overcoming the scalability hurdle is topic of much research.

Amato&Oliehoek - Cooperative MARL



partially observable: multiagent POMDP




Decentralized Objective Approach ("Dec-POMDP RL")

- Learn for all agents in the team
- Truly decentralized **execution**
- optionally: offline training phase
 - CTDE (with centralized components)



- Based on all kinds of RL techniques:
 - value-based (e.g., Q-learning)
 - policy search
 - actor-critic

Policy search methods:

- no dependence on Markov property
- ► policy gradient
 - can be decentralized [Peshkin et al. 2000]
 - has convergence guarantees
 - many recent deep versions





Policy gradient for partially observable RL

- Let's look at policy gradient ("REINFORCE") in P.O. settings (e.g., POMDPs)
- history: $h_t = (o_0, a_0, \dots, o_{t-1}, a_{t-1}, o_t)$
- value:

$$V(\pi_{\theta}) = \sum_{t=0}^{T-1} \sum_{h_t \in \mathcal{H}_t} \Pr(h_t | \pi_{\theta}) \sum_{a_t} \pi_{\theta}(a_t | h_t) R(h_t, a_t)$$

Cf. the value of a joint policy we saw before:

$$V(\boldsymbol{\pi}) = \sum_{t=0}^{T-1} \sum_{\boldsymbol{h}_t \in \boldsymbol{\mathcal{H}}_t} \Pr(\boldsymbol{h}_t | \boldsymbol{\pi}, b_0) \sum_{\boldsymbol{a}_t \in \boldsymbol{\mathcal{A}}} R(\boldsymbol{h}_t, \boldsymbol{a}_t) \boldsymbol{\pi}(\boldsymbol{a}_t | \boldsymbol{h}_t),$$
Amato&Oliehoek - Cooperative MARL



Derive gradient...

• Gradient:

$$\nabla_{\theta} V(\pi_{\theta}) = \nabla_{\theta} \left[\sum_{t=0}^{T-1} \sum_{h_t \in \mathcal{H}_t} \Pr(h_t | \pi_{\theta}) \sum_{a_t} \pi_{\theta}(a_t | h_t) R(h_t, a_t) \right]$$

{next slide} = $\sum_{t=0}^{T-1} \sum_{h_t, a_t} \Pr(h_t, a_t | \pi_{\theta}) R(h_t, a_t) \nabla_{\theta} \log \Pr(h_t, a_t | \pi_{\theta})$





Derive gradient...

• Gradient:

$$\nabla_{\theta} V(\pi_{\theta}) = \nabla_{\theta} \left[\sum_{t=0}^{T-1} \sum_{h_{t} \in \mathcal{H}_{t}} \Pr(h_{t}|\pi_{\theta}) \sum_{a_{t}} \pi_{\theta}(a_{t}|h_{t}) R(h_{t},a_{t}) \right]$$

$$\left\{ \text{next slide} \right\} = \sum_{t=0}^{T-1} \sum_{h_{t},a_{t}} \Pr(h_{t},a_{t}|\pi_{\theta}) R(h_{t},a_{t}) \nabla_{\theta} \log \Pr(h_{t},a_{t}|\pi_{\theta}) \right]$$

$$\left(\text{with } \nabla_{\theta} \log \Pr(h_{t},a_{t}|\pi_{\theta}) = \nabla_{\theta} \log \left[\Pr(a_{t}|h_{t};\pi_{\theta}) \left[\prod_{k=0}^{t-1} \Pr(a_{k}|h_{k};\pi_{\theta}) \Pr(o_{k+1}|h_{k},a_{k}) \right] P(o_{0}) \right] \right]$$

$$= \nabla_{\theta} \left[\log P(o_{0}) + \sum_{k=0}^{t} \log \Pr(a_{k}|h_{k};\pi_{\theta}) + \sum_{k=0}^{t-1} \log \Pr(o_{k+1}|h_{k},a_{k}) \right]$$
For ly middle term depends on θ

$$= \sum_{k=0}^{t} \nabla_{\theta} \log \Pr(a_{k}|h_{k};\pi_{\theta})$$

$$Delft$$

$$A mato&Oliehoek - Cooperative MARL$$

(Full deriv.)

TUDelft

• for completeness:

$$\nabla_{\theta} V(\pi_{\theta}) = \nabla_{\theta} \left[\sum_{t=0}^{T-1} \sum_{h_{t} \in \mathcal{H}_{t}} \Pr(h_{t} | \pi_{\theta}) \sum_{a_{t}} \pi_{\theta}(a_{t} | h_{t}) R(h_{t}, a_{t}) \right]$$

$$= \sum_{t=0}^{T-1} \nabla_{\theta} \underbrace{\sum_{h_{t} \in \mathcal{H}_{t}} \Pr(h_{t} | \pi_{\theta}) \sum_{a_{t}} \pi_{\theta}(a_{t} | h_{t}) R(h_{t}, a_{t})}_{R_{t}(\theta)}$$

$$= \sum_{t=0}^{T-1} \sum_{h_{t} \in \mathcal{H}_{t}} \sum_{a_{t}} \nabla_{\theta} \left[\Pr(h_{t} | \pi_{\theta}) \pi_{\theta}(a_{t} | h_{t}) R(h_{t}, a_{t}) \right]$$
{prod. rule:}
$$= \sum_{t=0}^{T-1} \sum_{h_{t} \in \mathcal{H}_{t}} \sum_{a_{t}} \left[\nabla_{\theta} \left(\Pr(h_{t} | \pi_{\theta}) \pi_{\theta}(a_{t} | h_{t}) \right) R(h_{t}, a_{t}) + \Pr(h_{t} | \pi_{\theta}) \pi_{\theta}(a_{t} | h_{t}) \nabla_{\theta} R(h_{t}, a_{t}) \right]$$

$$= \sum_{t=0}^{T-1} \sum_{h_{t}, a_{t}} \left[\left(\nabla_{\theta} \Pr(h_{t}, a_{t} | \pi_{\theta}) \right) R(h_{t}, a_{t}) + \Pr(h_{t} | \pi_{\theta}) \pi_{\theta}(a_{t} | h_{t}) \cdot 0 \right]$$

$$= \sum_{t=0}^{T-1} \sum_{h_{t}, a_{t}} \left[\left(\nabla_{\theta} \Pr(h_{t}, a_{t} | \pi_{\theta}) \right) R(h_{t}, a_{t})$$

$$= \sum_{t=0}^{T-1} \sum_{h_{t}, a_{t}} \Pr(h_{t}, a_{t} | \pi_{\theta}) \left(\frac{\nabla_{\theta} \Pr(h_{t}, a_{t} | \pi_{\theta})}{\Pr(h_{t}, a_{t} | \pi_{\theta})} \right) R(h_{t}, a_{t})$$
{log trick:}
$$= \sum_{t=0}^{T-1} \sum_{h_{t}, a_{t}} \Pr(h_{t}, a_{t} | \pi_{\theta}) R(h_{t}, a_{t} | \pi_{\theta})$$
Amato&Oliehoek - Cooperative MARL

1

Multiagent PG

Delft

• Single agent:
$$\nabla_{\theta} V(\pi_{\theta}) = \sum_{t=0}^{T-1} \sum_{h_t, a_t} \Pr(h_t, a_t | \pi_{\theta}) R(h_t, a_t) \sum_{k=0}^t \nabla_{\theta} \log \Pr(a_k | h_k; \pi_{\theta})$$

• Multiagent:

$$\nabla_{\theta} V(\boldsymbol{\pi}_{\theta}) = \sum_{t=0}^{T-1} \sum_{\boldsymbol{h}_{t},\boldsymbol{a}_{t}} \Pr(\boldsymbol{h}_{t},\boldsymbol{a}_{t} | \boldsymbol{\pi}_{\theta}) R(\boldsymbol{h}_{t},\boldsymbol{a}_{t}) \sum_{k=0}^{t} \nabla_{\theta} \log \Pr(\boldsymbol{a}_{k} | \boldsymbol{h}_{k}; \boldsymbol{\theta})$$

$$\{\text{policy is decentralized}\} = \sum_{t=0}^{T-1} \sum_{\boldsymbol{h}_{t},\boldsymbol{a}_{t}} \Pr(\boldsymbol{h}_{t},\boldsymbol{a}_{t} | \boldsymbol{\pi}_{\theta}) R(\boldsymbol{h}_{t},\boldsymbol{a}_{t}) \sum_{k=0}^{t} \nabla_{\theta} \log \prod_{i=1}^{n} \Pr(\boldsymbol{a}_{i,k} | \boldsymbol{h}_{i,k}; \boldsymbol{\theta}_{i})$$

$$\{\text{gradient is per agent}\} = \sum_{t=0}^{T-1} \sum_{\boldsymbol{h}_{t},\boldsymbol{a}_{t}} \Pr(\boldsymbol{h}_{t},\boldsymbol{a}_{t} | \boldsymbol{\pi}_{\theta}) R(\boldsymbol{h}_{t},\boldsymbol{a}_{t}) \sum_{i=1}^{n} \sum_{k=0}^{t} \nabla_{\theta_{i}} \log \Pr(\boldsymbol{a}_{i,k} | \boldsymbol{h}_{i,k}; \boldsymbol{\theta}_{i})$$

Peshkin L, Kim KE, Meuleau N, Kaelbling LP. Learning to cooperate via policy search. In Proceedings of the Sixteenth conference on Uncertainty in artificial intelligence 2000 Jun 30 (pp. 489-496). **IFAAMAS Influential paper award 2024**

Upshot:

convergent coop. MARL that can be decentralized! (only return needs to be observed)



MARL: a subjective perspective







Subjective Perspective

- So now let's further formalize...
 - Other agents: part of the environment
 - Best response model = multiagent environment + models of other agents



Frans A. Oliehoek and Christopher Amato. Best Response Bayesian Reinforcement Learning for Multiagent Systems with State Uncertainty. In Proceedings of the Ninth AAMAS Workshop on Multi-Agent Sequential Decision Making in Uncertain Domains (MSDM), 2014.



Models & Environment

- Multiagent environment
 - $MEA_i = \langle S, \{A_j\}, \{O_j\}, T, O, R_i \rangle$

- Models of other agents
 - think: finite state controllers
 - $m_j = \langle A_j, O_j, IS_j, \pi_j, \beta_j, I_j \rangle$
 - but... very general!
 - no real restrictions on internal states or functions
 - i.e., policy and belief update can be computational procedures
 - so includes MDPs, POMDPs, etc.







Amato&Oliehoek - Cooperative MARL

Best-response model (BRM)

 a_i

- Formalized as a (non-standard) POMDP
- States: $\overline{s}_i = \langle s, I_{-i} \rangle$
- Actions:
- Observations: o_i
- Dynamics function:

$$\bar{D}_{i}(\bar{s}_{i}', o_{i}|\bar{s}_{i}, a_{i}) = \sum_{a_{-i}} \sum_{o_{-i}} T(s'|s, a) O(o|a, s') \beta_{-i}(I_{-i}'|I_{-i}, a_{-i}, o_{-i}) \pi_{-i}(a_{-i}|I_{-i})$$

• Rewards:

elft

$$\overline{R}_i(\overline{s}_i, a_i) = \sum_{a_{-i}} R_i(s, a) \pi_{-i}(a_{-i}|I_{-i})$$





Planning / Learning for BRMs

- Since a BRM is a POMDP...
 - If you have all the models \rightarrow use POMDP solver
 - otherwise \rightarrow POMDP RL
 - yes... difficult... (but all methods apply)
- Open question: is RL for a BRM easier (or more difficult) than other POMDP RL?
 - easier to have priors of 'sane' behavior for agents (rather than other environmental aspects) ?
 - but if other agents are learning... ...continual extrapolation problem?



https://worldmodels.github.io/







Summary Part 1

- MARL is complex...! (and interesting!)
- Two initial ideas...
 - fully centralized \rightarrow issues with scalability (action spaces, communication)
 - fully decentralized \rightarrow issues with convergence
- Dec-POMDP: framework that allows reasoning over private observations
- Fundaments of multiagent planning
 - 2 generals: need to also factor in 'predictability'
 - Heuristic search: past joint policy matters
 → because it determines the distribution over states and knowledge
 - Enables formulation of "plan-time models" or "designer's approach"
- Dec-POMDP perspective to MARL:
 - helps to understand 'non-stationarity'
 - objective approach: policy gradient still works (since value of a given joint policy is analogue to a POMDP)
 - subjective approach: formalize a "best-response model" → RL in difficult POMDPs







- Most references are in:
 - Frans A. Oliehoek and Christopher Amato. **A Concise Introduction to Decentralized POMDPs**, SpringerBriefs in Intelligent Systems, Springer, May 2016.
- Some more details:
 - Frans A. Oliehoek. Decentralized POMDPs. In Wiering, Marco and van Otterlo, Martijn, editors, *Reinforcement Learning: State of the Art*, Adaptation, Learning, and Optimization, pp. 471–503, Springer Berlin Heidelberg, Berlin, Germany, 2012.
 - Frans A. Oliehoek and Christopher Amato. **Dec-POMDPs as Non-Observable MDPs**. IAS technical report IAS-UVA-14-01, Intelligent Systems Lab, University of Amsterdam, 2014.





Bonus: Further MARL Topics





Further Topics in MARL

- Multiagent (reinforcement) learning is active topic of research
 - scaling up
 - deep learning
 - learning to communicate
 - multiagent approaches to ML (e.g., GANs)
 - ad-hoc teamwork: coordination without training











Scaling multiagent (PO)MDPs

- E.g., even in stateless setting: Q(a) too large...
- Coordination graphs [Guestrin et al. NIPS 2001]
 - address by factorizing...
 - E.g., $Q(\mathbf{a}) \approx Q_{1,2}(a_1,a_2) + Q_{2,3}(a_2,a_3)$



- Benefits:
 - compact representation
 - can optimize:

 $\max_{a} [Q_{1,2}(a_{1},a_{2}) + Q_{2,3}(a_{2},a_{3})]$

• with inference or COP techniques (variable elimination, max-plus, etc.)



Scaling multiagent (PO)MDPs

• back to sequential setting...







Factored Value Functions [Guestrin NIPS'01]

• Approximate with factored (Q)-value function



'e' denotes subsets of agents and state variables





Amato&Oliehoek - Cooperative MARL



Large state spaces: deep MARL

- Huge state spaces
 - even 2 intersection source problems are intractable
 - (even 1 intersection is...)
- Deep Q-learning
 - encode state as matrix





(b) Traffic situation

(c) Simplified example of state representation in 8×8 matrix.

95



Elise Van der Pol and Frans A. Oliehoek. Coordinated Deep Reinforcement Learners for Traffic Light Control. In NIPS'16 Workshop on Learning, Inference and Control of Multi-Agent Systems 2016. 1 1 i s Amato&Oliehoek - Cooperative MARL

Scaling via "Transfer Planning"

• Define **source problem** for each Q-component



Amato&Oliehoek - Cooperative MARL

96

Scaling via "Transfer Planning"

- 'Solve' source problems independently
 - use Q-function as components











Scaling via "Transfer Planning"

- 'Solve' source problems independently
 - use Q-function as components



Empirical Results



99

Communication

- instantaneous, cost-free, and noise-free:
 - Dec-MDP \rightarrow multiagent MDP (MMDP)
 - Dec-POMDP → multiagent POMDP (MPOMDP)
- but in practice:
 - probability of failure
 - delays
 - costs
- Also: implicit communication! (via observations and actions)



Explicit Communication

- perform a particular information update (e.g., sync) as in the MPOMDP:
 - each agent broadcasts its information, and
 - each agent uses that to perform joint belief update
- Other approaches:
 - Communication cost [Becker et al. 2005]
 - Delayed communication [Hsu et al. 1982, Spaan et al. 2008, Oliehoek & Spaan 2012]
 - Communicate every k stages [Goldman & Zilberstein 2008]





Implicit Communication

• Encode communications by actions and observations



• Embed the **optimal meaning** of messages by finding the optimal plan [Goldman and Zilberstein 2003, Spaan et al. 2006]





Implicit Communication

• Encode communications by actions and observations



• Embed the **optimal meaning** of messages by finding the optimal plan [Goldman and Zilberstein 2003, Spaan et al. 2006]





Implicit Communication

• Encode communications by actions and observations





Me in 2014:

• Embed the **optimal meaning** (plan [Goldman and Zilberstein 2(

"E.g. communication bit doubles the #actions and observations!

Clearly, useful... but intractable for general settings (perhaps for analysis of very small communication systems)"

but then...



DELF

unit

Deep learning of communication

• ...these scalability issues can be overcome (at least to some extend) by deep learning...

Amato&Oliehoek - Cooperative MARL





CommNet [Sukhbaatar, et al. 2016]

DELFT

unit

105



Deep MARL: where is the field...?

- MARL for traffic:
 - http://www.fransoliehoek.net/trafficvideo
- Learning to communicate:
 - https://youtu.be/KhtdEvJ1F6Q?t=41
 - https://arxiv.org/abs/1810.11187
- Starcraft II via 'value factorization'
 - https://arxiv.org/pdf/1803.11485.pdf
 - https://www.youtube.com/watch?v=4WBi8xI_8YA
- Capture the flag
 - https://deepmind.com/blog/capture-the-flag-science/





Deep MARL: where is the field...?

- "From Motor Control to Team Play in Simulated Humanoid Football"
 - https://youtu.be/KHMwq9pv7mg?t=249
- Learning a Decentralized Multi-arm Motion Planner
 - https://multiarm.cs.columbia.edu/
- Learning to fly
 - https://github.com/utiasDSL/gym-pybullet-dr ones
- Playing Stratego
 - https://www.science.org/doi/10.1126 /science.add4679







More Dec-POMDP solver improvements





MAA* via Bayesian Games

- Each node \rightarrow a ϕ^t
- decision problem for stage t

	$\vec{\theta_{2}^{t=0}}$)	()			
			a_2	\bar{a}_2			
		a_1	+2.	75 -4.	1		
	()	\bar{a}_1	-0.	9 + 0.	3		
			•				
	$\vec{ heta}_2^{t=1}$		(a_2, o_2)		(a_2, \bar{o}_2)		
	$ec{ heta}_1^{t=1}$		a_2	\bar{a}_2	a_2	\bar{a}_2	
	(α, α)	a_1	-0.3	+0.6	-0.6	+4.0	
	(a_1, o_1)	\bar{a}_1	-0.6	+2.0	-1.3	+3.6	
		a_1	+3.1	+4.4	-1.9	+1.0	
	(a_1, o_1)	\bar{a}_1	+1.1	-2.9	+2.0	-0.4	
	(\bar{a}_1, o_1)	a_1	-0.4	-0.9	-0.5	-1.0	
		\bar{a}_1	-0.9	-4.5	-1.0	+3.5	
	(\bar{a}_1, \bar{o}_1)						
Delft Amato							0&O



MAA* via Bayesian Games - 2

MAA* perspective



- node $\rightarrow \phi^t$
- joint decision rule δ maps OHs to actions
- Expansion: appending all next-stage decision rules: φ^{t+1}=(φ^t,δ^t)

BG perspective



- node \rightarrow a BG
- joint BG policy β maps 'types' to actions
- Expansion: enumeration of all joint BG policies φ^{t+1}=(φ^t, β^t)

direct correspondence: δ 🔄 β


MAA* via Bayesian Camoo 2

direct co

MAA* perspective



- node $\rightarrow \phi^t$
- joint decision rule δ maps OHs to actions
- Expansion: appending all next decision rules: φ^{t+1}=(φ^t,δ^t)

What is the point?

Generalized MAA* [Oliehoek & Vlassis '07]
 Unified perspective of MAA* and 'BAGA' approximation [Emery-Montemerlo et al. '04]
 No direct improvements...

However...

- BGs provide abstraction layer
- Facilitated two improvements that lead to state-of-the-art performance [Oliehoek et al. '13]

116

- Clustering of histories
- Incremental expansion

MAA* Limitations

- Number of children grows doubly exponential with nodes depth
- For a node last stage, number of children is

$$O(|A_*|^{n|O_*|^{h-1}})$$

• Total number of joint policies $O(|A_*|^{(n|O_*|^h-1)/(|O_*|-1)})$

→ MAA* can only solve 1 horizon longer than brute force search... [Seuken & Zilberstein '08]

- Techniques to overcome this problem [Oliehoek et al. 2013 JAIR]:
 - lossless clustering of histories
 - Incremental expansion

Amato&Oliehoek - Cooperative MARL

DELFT

Lossless Clustering

 Two types (=action-observation histories) in a BG are probabilistically equivalent iff

$P(\vec{\theta}_{-i} \vec{\theta}_{i,a}) = P(\vec{\theta}_{-i} \vec{\theta}_{i,b})$		_		_	
			ō	\vec{s}_{2}^{2}	
$P(s \vec{\theta}_{i},\vec{\theta}_{i},\vec{q}) = P(s \vec{\theta}_{i},\vec{\theta}_{i},\vec{h})$	\vec{o}_{1}^{2}	$(o_{ m HL},\!o_{ m HL})$	$(o_{ m HL},\!o_{ m HR})$	$(o_{ m HR}, o_{ m HL})$	$(o_{ m HR},\!o_{ m HR})$
(1 - i) i, u $(1 - i) i, v$	$(o_{\mathrm{HL}}, o_{\mathrm{HL}})$	0.261	0.047	0.047	0.016
	$(o_{ m HL},\!o_{ m HR})$	0.047	0.016	0.016	0.047
	$(o_{ m HR},\!o_{ m HL})$	0.047	0.016	0.016	0.047
	$(o_{ m HR}, o_{ m HR})$	0.016	0.047	0.047	0.261
		(a) The jo	oint type prob	pabilities.	
	I		\vec{o}	$\frac{2}{2}$	
	$ec{o}_1^2$	$(o_{ m HL},\!o_{ m HL})$	$(o_{ m HL}, o_{ m HR})$	$(o_{ m HR},\!o_{ m HL})$	$(o_{ m HR}, o_{ m HR})$
	$(o_{ m HL}, o_{ m HL})$	0.999	0.970	0.970	0.5
	$(o_{ m HL}, o_{ m HR})$	0.970	0.5	0.5	0.030
	$(o_{ m HR}, o_{ m HL})$	0.970	0.5	0.5	0.030
	$(o_{ m HR}, o_{ m HR})$	0.5	0.030	0.030	0.001
	(b) The induced	joint beliefs.	Listed is the	e probability	$\Pr(s_l \vec{\theta}^2, b^0)$ o

(b) The induced joint beliefs. Listed is the probability $\Pr(s_l | \vec{\theta}^2, b^0)$ of the tiger being behind the left door.





Lossless Clustering

• Two types (=action-observation histories) in a BG are **probabilistically equivalent** iff

$P\left(\vec{\Theta}_{-i} \vec{\Theta}_{i,a}\right) = P\left(\vec{\Theta}_{-i} \vec{\Theta}_{i,b}\right)$					
			ō	\vec{b}_{2}^{2}	
$P(s \dot{\theta}_{-i},\dot{\theta}_{i,a}) = P(s \dot{\theta}_{-i},\dot{\theta}_{i,b})$	\vec{o}_1^2	$(o_{ m HL}, o_{ m HL})$	$(o_{ m HL},\!o_{ m HR})$	$(o_{ m HR}, o_{ m HL})$	$(o_{ m HR}, o_{ m HR})$
	$(o_{ m HL},\!o_{ m HL})$	0.261	0.047	0.047	0.016
	$(o_{ m HL}, o_{ m HR})$	0.047	0.016	0.016	0.047
	$(o_{ m HR}, o_{ m HL})$	0.047	0.016	0.016	0.047
	$(o_{ m HR}, o_{ m HR})$	0.016	0.047	0.047	0.261
		(a) The jo	oint type prob	pabilities.	
	I		\vec{o}	$\frac{2}{2}$	
	\vec{o}_1^2	$(o_{ m HL},\!o_{ m HL})$	$(o_{ m HL}, o_{ m HR})$	$(o_{ m HR},\! o_{ m HL})$	$(o_{ m HR}, o_{ m HR})$
	$(o_{\mathrm{HL}}, o_{\mathrm{HL}})$	0.999	0.970	0.970	0.5
	$(o_{ m HL}, o_{ m HR})$	0.970	0.5	0.5	0.030
	$(O_{\rm HR}, O_{\rm HL})$	0.970	0.5	0.5	0.030
_	$(o_{ m HR}, o_{ m HR})$	0.5	0.030	0.030	0.001
	(b) The induced	joint beliefs.	Listed is the	e probability	$\Pr(s_l \vec{\theta}^2, b^0)$ of

(b) The induced joint beliefs. Listed is the probability $\Pr(s_l | \vec{\theta}^2, b^0)$ of the tiger being behind the left door.





- Key idea: even though nodes can have many children, only few are useful.
 - i.e., only few will be selected for further expansion
 - others will have too low heuristic value



- if we can generate the nodes in increasing heuristic order
 - \rightarrow can avoid expansion of redundant nodes











Select for expansion \rightarrow







































TUDelft



Incremental Expansion: How?

• How do we generate the next-best child?

- Node ↔ BG, so...
 - find the solutions of the BG (in decreasing order of value)
 - i.e., 'incremental BG solver'
 - Modification of BaGaBaB [Oliehoek et al. 2010]
 - stop searching when next solution found
 - save search tree for next time visited.





Some Results

ŤU

Delft

	problem primitives			
	n	$ \mathcal{S} $	$ \mathcal{A}_i $	$ \mathcal{O}_i $
DEC-TIGER	2	2	3	2
BroadcastChannel		4	2	2
$\operatorname{Grid}\operatorname{Small}$	2	16	5	2
Cooperative Box Pushing Recycling Robots		100	4	5
		4	3	2
Hotel 1	2	16	3	4
FIREFIGHTING	2	432	3	2

		h	MILP	DP-LPC	DP-IPG	$GMAA - Q_{BG}$		$Q_{ m BG}$			
									IC	ICE	heur
ocult					Broa	dcastCh	ANNEL, ICE	solvable to h	= 900		
esuit	5				2	0.38	≤ 0.01	0.09	≤ 0.01	≤ 0.01	≤ 0.01
					3	1.83	0.50	56.66	≤ 0.01	≤ 0.01	≤ 0.01
					4	34.06	*	*	≤ 0.01	≤ 0.01	≤ 0.01
					5	48.94			≤ 0.01	≤ 0.01	≤ 0.01
					DEC-	FIGER, IC	E solvable to	o $h = 6$			
					2	0.69	0.05	0.32	≤ 0.01	≤ 0.01	≤ 0.01
					3	23.99	60.73	55.46	≤ 0.01	≤ 0.01	≤ 0.01
		roblom	. primi	ives	4	*	—	2286.38	0.27	≤ 0.01	0.03
	ł	broblen	r prinn	lives	5			_	21.03	0.02	0.09
	n	$ \mathcal{S} $	$ \mathcal{A}_i $	$ {\cal O}_i $	FireF	IGHTING	(2 agents, 3)	houses, 3 fire	levels), IC	CE solvab	le to $h \gg 1000$
					2	4.45	8.13	10.34	≤ 0.01	≤ 0.01	≤ 0.01
DEC TICED	2	2	9	0	3	—	—	569.27	0.11	0.10	0.07
DEC-1IGER	2	2	3	2	4			—	950.51	1.00	0.65
STCHANNEL	2	4	2	2	GRID	Small, IC	E solvable t	o $h = 6$			
GRIDSMALL	2	16	5	2	2	6.64	11.58	0.18	0.01	≤ 0.01	≤ 0.01
GRIDOMALL	2	10	0	2	3	*	—	4.09	0.10	≤ 0.01	0.42
ox Pushing	2	100	4	5	4			77.44	1.77	≤ 0.01	67.39
ING ROBOTS	2	4	3	2	Recy	cling Ro	BOTS, ICE s	solvable to h =	= 70		
	_	-	0	-	2	1.18	0.05	0.30	≤ 0.01	≤ 0.01	≤ 0.01
Hotel 1	2	16	3	4	3	*	2.79	1.07	≤ 0.01	≤ 0.01	≤ 0.01
deFiguring	2	132	3	2	4		2136.16	42.02	≤ 0.01	≤ 0.01	0.02
KET IGHTING	2	432	3	2	5		—	1812.15	≤ 0.01	≤ 0.01	0.02
					Hote	L 1, ICE s	solvable to h	a = 9			
					2	1.92	6.14	0.22	≤ 0.01	≤ 0.01	0.03
					3	315.16	2913.42	0.54	≤ 0.01	≤ 0.01	1.51
					4	—	—	0.73	≤ 0.01	≤ 0.01	3.74
					5			1.11	≤ 0.01	≤ 0.01	4.54
					9			8.43	0.02	≤ 0.01	20.26
				10			17.40	#	#		
'—' mem	lory	r limi	it vio	lations	15			283.76			
'*' time	lim	it our	orrun	q	Coop	ERATIVE I	Box Pushin	${}^{\rm AG}({\rm Q}_{\rm POMDP})$, ICE solv	x to h	4 = 4
≁ ume.	11111	10 000	errun	15	2	3.56	15.51	1.07	≤ 0.01	≤ 0.01	≤ 0.01
'#' houri	isti	e hot	tlong	ek	3	2534.08	—	6.43	0.91	0.02	0.15
# neur	19010		Am:	tox()liphopk -	4			1138.61	*	328.97	0.6313
			AIII	atodonenoek -	Cooper			Min Plan		un	it 15

Some Results

TUDelft



Scalability w.r.t. #agents

	h	V^*	$T_{GMAA*}(s)$	$T_{IC}(s)$	$T_{ICE}(s)$			
Recycling Robots								
	3	10.660125	≤ 0.01	≤ 0.01	≤ 0.01			
	4	13.380000	713.41	≤ 0.01	≤ 0.01			
	5	16.486000	_	≤ 0.01	≤ 0.01			
	6	19.554200		≤ 0.01	≤ 0.01			
	10	31.863889		≤ 0.01	≤ 0.01			
	15	47.248521		≤ 0.01	≤ 0.01			
	20	62.633136		≤ 0.01	≤ 0.01			
	30	93.402367		0.08	0.05			
	40	124.171598		0.42	0.25			
	50	154.940828		2.02	1.27			
	70	216.479290		-	28.66			
	80			_	—			
		Broad	DCASTCHAN	NEL				
	4	3.890000	≤ 0.01	≤ 0.01	≤ 0.01			
	5	4.790000	1.27	≤ 0.01	≤ 0.01			
	6	5.690000	_	≤ 0.01	≤ 0.01			
	7	6.590000		≤ 0.01	≤ 0.01			
	10	9.290000		≤ 0.01	≤ 0.01			
	25	22.881523		≤ 0.01	≤ 0.01			
	50	45.501604		≤ 0.01	≤ 0.01			
	100	90.760423		≤ 0.01	≤ 0.01			
	250	226.500545		0.06	0.07			
	500	452.738119		0.81	0.94			
	700	633.724279		0.52	0.63			
	800			-	_			
	900	814.709393		9.57	11.11			
	1000				_			

Cases that compress well

* excluding heuristic



131

- Generate all policies in a special way:
 - rightarrow from 1 stage-to-go policies $Q^{\tau=1}$
 - construct all 2-stages-to-go policies $Q^{\tau=2}$, etc. \triangleright



- Generate all policies in a special way:
 - rightarrow from 1 stage-to-go policies $Q^{\tau=1}$



- Generate all policies in a special way:
 - rightarrow from 1 stage-to-go policies $Q^{\tau=1}$



- Generate all policies in a special way:
 - rightarrow from 1 stage-to-go policies $Q^{\tau=1}$



- Generate all policies in a special way:
 - ▷ from 1 stage-to-go policies $Q^{\tau=1}$



- Generate all policies in a special way:
- rightarrow from 1 stage-to-go policies $Q^{\tau=1}$ a new **Exhaustive backup operation** \boldsymbol{a}_{i} t =t q_2^{τ}

UDelft

- Generate all policies in a special way:
 - ▷ from 1 stage-to-go policies $Q^{\tau=1}$



- All actions
- All assignments of q^{τ} to observations

• (obviously) this scales very poorly...



• (obviously) this scales very poorly...





140

unit

• (obviously) this scales very poorly...

 $Q_1^{ au=3}$

ቆእ ቆእ *ቆ*እ ቆእ

 $Q_2^{\tau=3}$

 \mathbf{A} *ቆ*እ ቆእ



• (obviously) this scales very poorly.	h	num. indiv. policies
$Q_1^{ au=3}$	1	2
శీశి శీశి శీశి శీశి శీశి శీశి శీశి శీశి	శీశి శీశి శీశి శీశి శి	8
\$& \$& \$& \$& \$& \$& \$& \$& \$& \$& \$& \$& \$& \$	&& & & & & & & & & & & & & & & & & & &	128
This does not get us anywhere!	\$& & & & 4	32768
but		2.1475e+09
DUL	is a s a 6	9.2234e+18
ϕ	۲ ۲ ۲	1.7014e+38
శీశి శీశి శీశి శీశి శీశి శీశి శీశి శీశి	శీశి శీశి శీశి శీశి శి	5.7896e+76
శీసి శీసి శీసి శీసి శీసి శీసి శీసి శీసి	శిశి శిశి శిశి శిశి శిశి శిశి శిశి	ቆቆ ቆቆ ቆ
Delft Amato&Oliehoek -	Cooperative MARL	ellis delft ₁₄₂

- Perhaps not all those Q_i^{τ} are useful!
 - Perform **pruning** of 'dominated policies'!
- Algorithm [Hansen et al. 2004] $Q_i^{\tau=1} = A_i$

```
Initialize Q1(1), Q2(1)
for tau=2 to h
  Q1(tau) = ExhaustiveBackup(Q1(tau-1))
  Q2(tau) = ExhaustiveBackup(
                                     Note: cannot prune independently!
  Prune(Q1,Q2,tau)
                                      ▶ usefulness of a q_1 depends on Q_2
end
                                      ► and vice versa
                                       → Iterated elimination of policies
                                      ► how? linear programming.
                        Amato&Oliehoek - Cooperative MARL
                                                                                143
```

Initialization



Amato&Oliehoek - Cooperative MARL

DELFT

unit

144

• Exhaustive Backups gives





Amato&Oliehoek - Cooperative MARL

• Pruning agent 1...

Delft



Hypothetical Pruning (not the result of actual pruning)



• Pruning agent 2...



S

unit

DELFT

147

• Pruning agent 1...



• Etc...



149

• Etc...



 $Q_1^{\tau=3}$

• Exhaustive backups:

ቆኤ ፊኤ

We avoid generation of many policies!

ፈዬ

 $Q_2^{\tau=3}$


$Q_1^{\tau=3}$

Exhaustive backups:

£\$ £\$ £\$ £\$ £\$ £\$ £\$ £\$

 $Q_{2}^{\tau=3}$

Amato&Oliehoek - Cooperative MARL

• Pruning agent 1...

 $Q_1^{\tau=3}$ $Q_{2}^{\tau=3}$ ፚ፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟ ፚ፟፟፟፟፟፟፟፟፟፟ ፚ፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟ ቆ፟፟፟፟፟፟፟ ቆ፟፟፟፟፟ ቆ፟፟፟፟፟፟፟፟ ቆ፟፟፟፟፟፟፟ ቆ፟፟፟፟፟ አቆ፝ አቆ፝ አቆ፝ ፚ፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟ ቆ፟፟፟፟፟ ቆ፟፝፝፝፝፟ ቆ፝፟፝ ቆ፟፝፝፟ ቆ፟፝፝፟ ቆ፟፝፝፟ ቆ፟፝ ቆි እ ቆ እ ቆ እ ቆ እ ቆ እ ፈዬ ፈዬ ፈዬ ፈዬ ፈዬ ፈዬ ፈዬ ፚ፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟ ቆි ቆቆ ቆቆ **ቆ**፟፟፟፟፟፟፟፟፟፟ አቆ፝፟ አቆ፝ አቆ፝ አቆ አቆ አቆ አቆ እ ፚ፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟፟ 1 S DELFT Amato&Oliehoek - Cooperative MARL

153

• Pruning agent 2...



Amato&Oliehoek - Cooperative MARL

• Etc...



Amato&Oliehoek - Cooperative MARL

unit

• Etc...



Incremental Policy Generation - 1

Bottleneck: exhaustive backup



 $Q_i^{\tau} = \bigcup_a Q_i^{\tau,a}$ $Q_{i}^{\tau,a} = (+)_{o} Q_{i}^{\tau,a,o}$ $Q_i^{\tau,a,o} = BackProject(Q_i^{\tau-1})$



Amato&Oliehoek - Cooperative MARL

Incremental Policy Generation - 1



Incremental Policy Generation - 2

 IPG [Amato et al 2009]: some states may be unreachable (for specific a,o)
→ prune only over reachable sub-space

