# The Size of Message Set Needed for the Optimal Communication Policy

Tatsuya Kasai, Hayato Kobayashi, and Ayumi Shinohara

Graduate School of Information Sciences, Tohoku University, Japan

The 7th European Workshop on Multi-Agent Systems (EUMAS 2009)
Ayia Napa, Cyprus
Dec17-18, 2009

# Background

- Multi-agent coordination with communication.

> Main objective :
> To find the optimal *action policy* $\delta^A$ and *communication policy* $\delta^M$
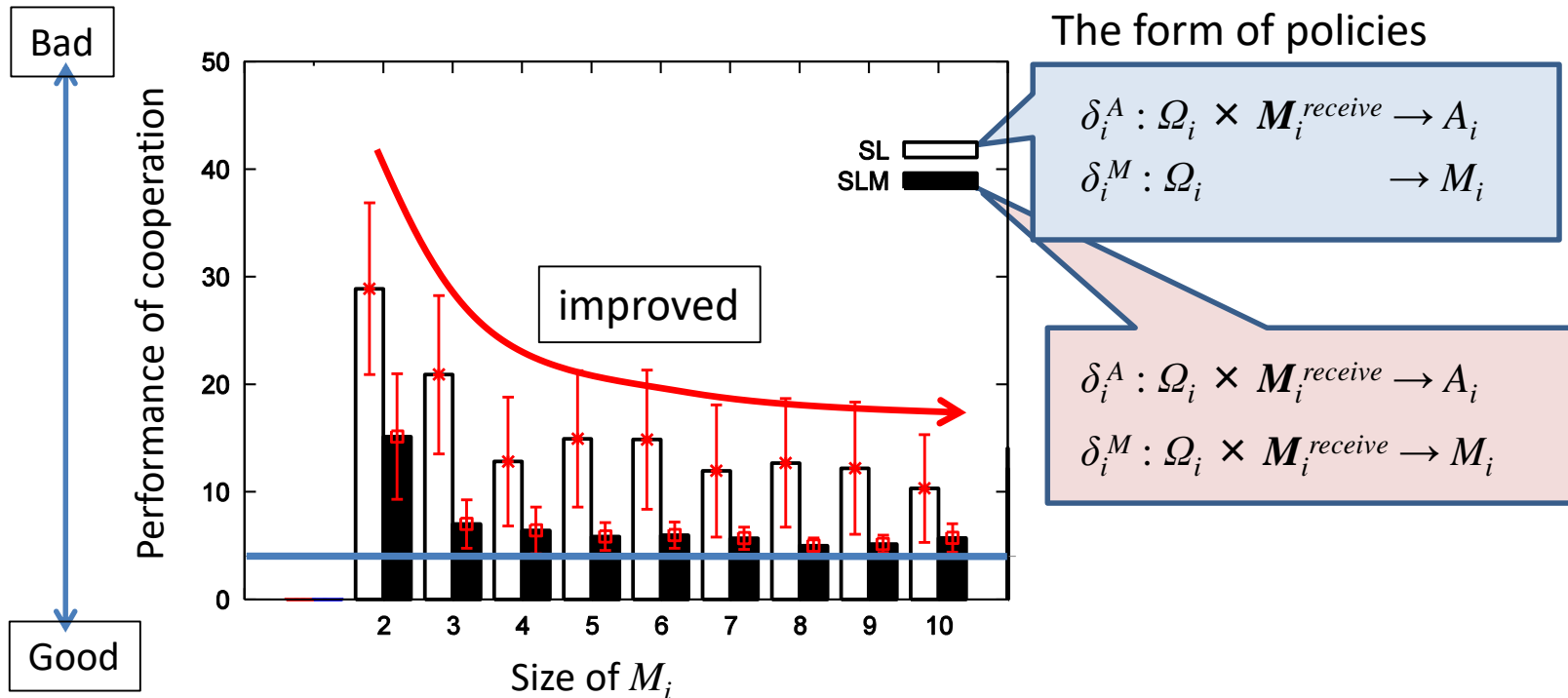
- We are interested in an approach based on autonomous learning.

- Definition of policies for agent $i$ in our proposed methods

| | Signal Learning (SL) [Kasai+ 08] | Signal Learning with Messages (SLM) [Kasai+ AAMAS09] |
|---|---|---|
| *Action Policy* | $\delta_i^A : \Omega_i \times M_i^{receive} \rightarrow A_i$ | |
| *Communication Policy* | $\delta_i^M : \Omega_i \rightarrow M_i$ | $\delta_i^M : \Omega_i \times M_i^{receive} \rightarrow M_i$ |

A set of observations

A set of received messages

A set of actions

A set of messages to send other agents

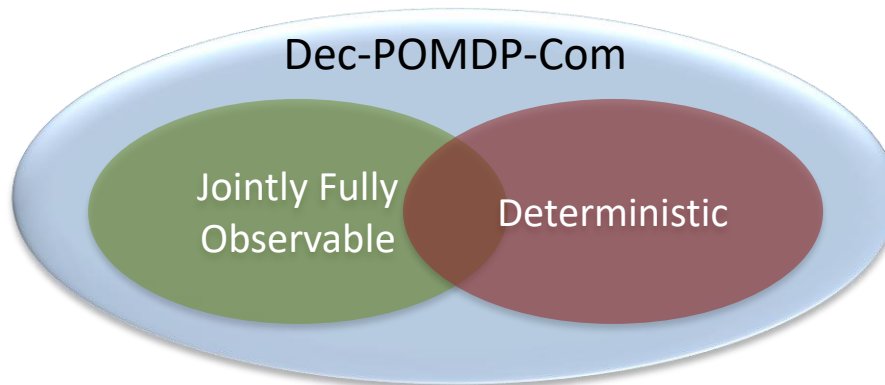(SL and SLM are based on Multi-Agent Reinforcement Learning framework)

# Motivation

- Actual learning results of SL and SLM [Kasai+ AAMAS09]
  - The performance of cooperation when the size of $M_i$ is increases.



The form of policies

$$\delta_i^A : \Omega_i \times \boldsymbol{M}_i^{receive} \rightarrow A_i$$
$$\delta_i^M : \Omega_i \qquad\qquad \rightarrow M_i$$

$$\delta_i^A : \Omega_i \times \boldsymbol{M}_i^{receive} \rightarrow A_i$$
$$\delta_i^M : \Omega_i \times \boldsymbol{M}_i^{receive} \rightarrow M_i$$

- We have an interest about how much size of $M_i$ for constructing the optimal policy ?

# Scheme of talk

- We show *minimum required sizes $|M_i|$* for achieving the optimal policy for
  - Signal Learning on *Jointly Fully Observable Dec-POMDP-Com*
  - Signal Learning with Messages on *Deterministic Dec-POMDP-Com*
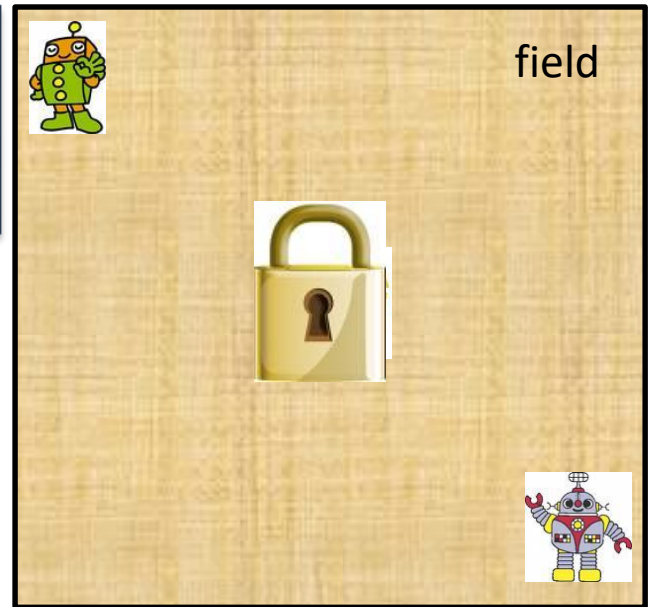
# Outline

☑Background

☑Scheme of talk

☐Review : Dec-POMDP-Com [Goldman+ 04]

☐Constrained model

   🌐 Jointly Fully Observable Dec-POMDP-Com [Goldman+ 04]

   🌐 Deterministic Dec-POMDP-Com (we define)

☐Theoretical analysis

☐Conclusion

# Dec-POMDP-Com [Goldman+ 04]

(**De**centralized **P**artially **O**bservable **M**arkov **D**ecision **P**rocess with **Com**munication)

- A decentralized multi-agent system, where agents can communicate with each other and only observe the restricted information.

Example of model

- Two agents get a treasure cooperatively.
- The treasure is locked.
- Both agents must reach the treasure at the same time to open the lock.



field

# Dec-POMDP-Com [Goldman+ 04]

(Decentralized Partially Observable Markov Decision Process with Communication)

- A decentralized multi-agent system, where agents can communicate with each other and only observe the restricted information.
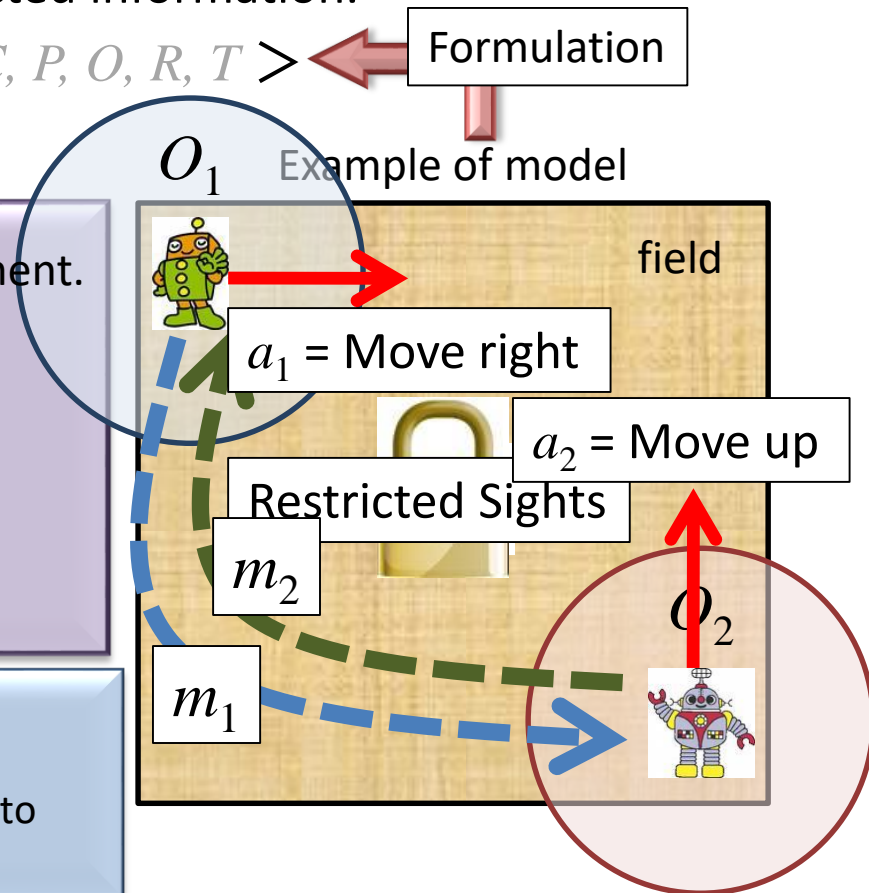
- Dec-POMDP-Com := $< I, S, \Omega, A, M, C, P, O, R, T >$ ⟵ Formulation

1step for agent $i$ on Dec-POMDP-Com

$O_1$   Example of model

1. Receive an observation $O_i$ from the environment.

2. Send a message $m_i$ to the other agents.

3. Perform an action $a_i$ in the environment.

Repeat until both agent arrive at the treasure.

- Two agents get a treasure cooperatively.
- The treasure is locked.
- Both agents must reach the treasure at the same time to open the lock.

field

$a_1$ = Move right

$a_2$ = Move up

Restricted Sights

$m_2$

$O_2$

$m_1$

# Dec-POMDP-Com [Goldman+ 04]

(Decentralized Partially Observable Markov Decision Process with Communication)

- A decentralized multi-agent system, where agents can communicate with each other and only observe the restricted information.
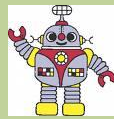
- Dec-POMDP-Com := $\langle I, S, \Omega, A, M, C, P, O, R, T \rangle$ ← Formulation

Example of model

- A set of agents' indices

- e.g., $I = \{1, 2\}$

= 1        = 2

$O_1$

$a_1$        field

$a_2$

$m_2$        $O_2$

$m_1$
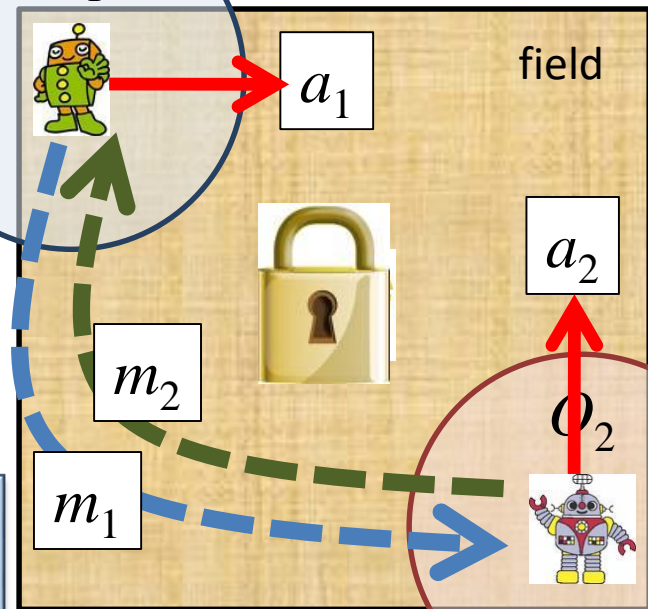
- Two agents get a treasure cooperatively.
- The treasure is locked.
- Both agents must reach the treasure at the same time to open the lock.

# Dec-POMDP-Com [Goldman+ 04]

(Decentralized Partially Observable Markov Decision Process with Communication)

- A decentralized multi-agent system, where agents can communicate with each other and only observe the restricted information.
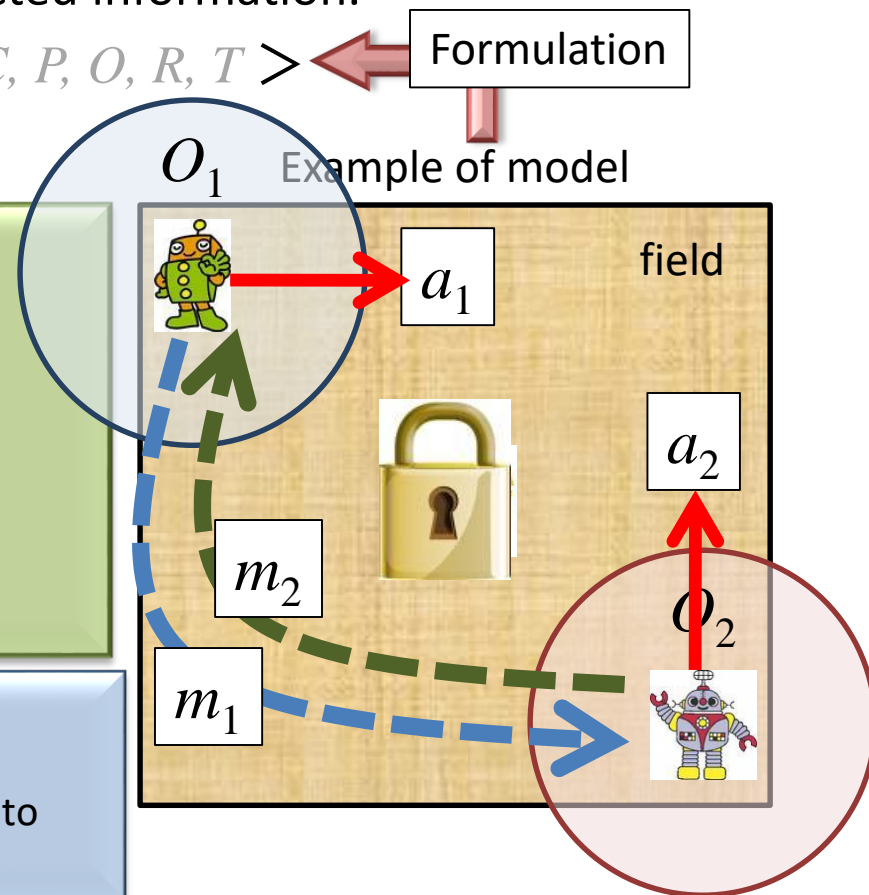
- Dec-POMDP-Com := $< I, S, \Omega, A, M, C, P, O, R, T >$ ← Formulation

$O_1$ Example of model

- A set of global states

- e.g., $s = ($position of agent 1,
          position of agent 2,
          position of treasure $)$
          $, s \in S$

$a_1$

$a_2$

$O_2$

field

$m_2$

$m_1$

- Two agents get a treasure cooperatively.
- The treasure is locked.
- Both agents must reach the treasure at the same time to open the lock.

# Dec-POMDP-Com [Goldman+ 04]

(Decentralized Partially Observable Markov Decision Process with Communication)

- A decentralized multi-agent system, where agents can communicate with each other and only observe the restricted information.
- Dec-POMDP-Com := $< I, S, \Omega, A, M, C, P, O, R, T >$

Formulation

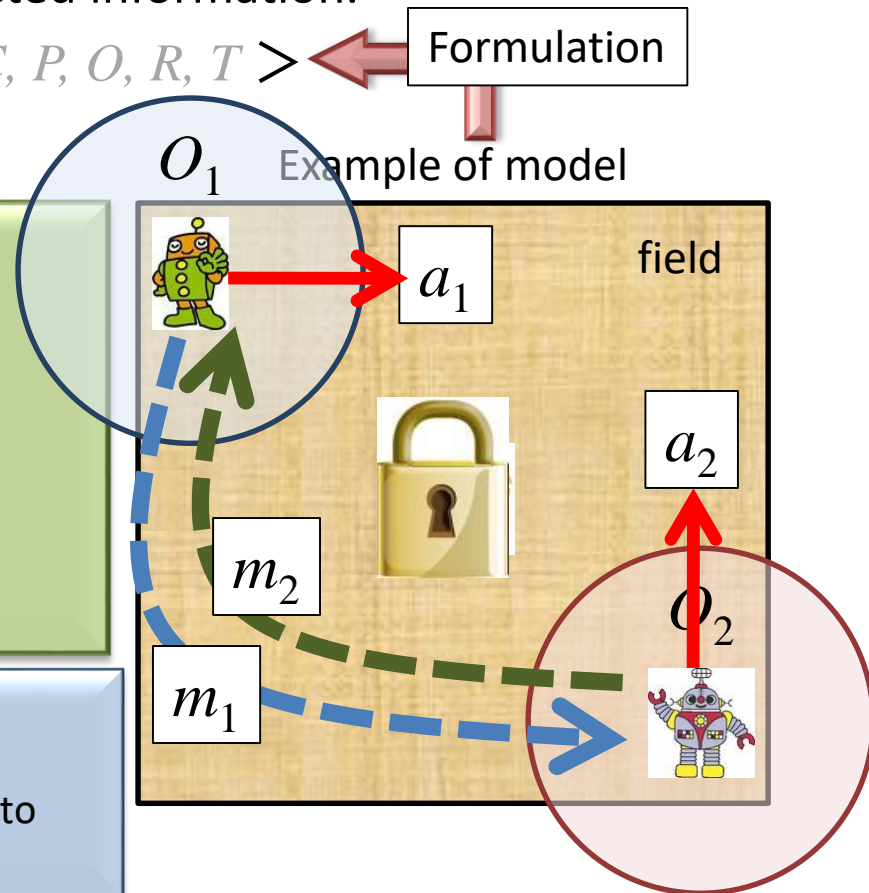$O_1$  Example of model

field

- $\Omega$ : a set of joint observations
- $\Omega = \Omega_1 \times \Omega_2$, where $\Omega_i$ is a set of observations for agent $i$

- $A$ : a set of joint actions
- $A = A_1 \times A_2$

$a_1$

$a_2$

$O_2$

$m_2$

$m_1$

- Two agents get a treasure cooperatively.
- The treasure is locked.
- Both agents must reach the treasure at the same time to open the lock.

# Dec-POMDP-Com [Goldman+ 04]

(**De**centralized **P**artially **O**bservable **M**arkov **D**ecision **P**rocess with **Com**munication)

- A decentralized multi-agent system, where agents can communicate with each other and only observe the restricted information.

- Dec-POMDP-Com := $< I, S, \pmb{\Omega}, A, \pmb{M}, \pmb{C}, P, O, R, T >$ ← Formulation

$O_1$  Example of model

- $M$ : a set of joint messages
- $M = M_1 \times M_2$

- $C : M \to \mathcal{R}$ is a cost function
- $C(m)$ represent the total cost of transmitting the messages sent by all agents.

$a_1$   field

$a_2$

$m_2$

$O_2$

$m_1$

- Two agents get a treasure cooperatively.
- The treasure is locked.
- Both agents must reach the treasure at the same time to open the lock.

# Dec-POMDP-Com [Goldman+ 04]

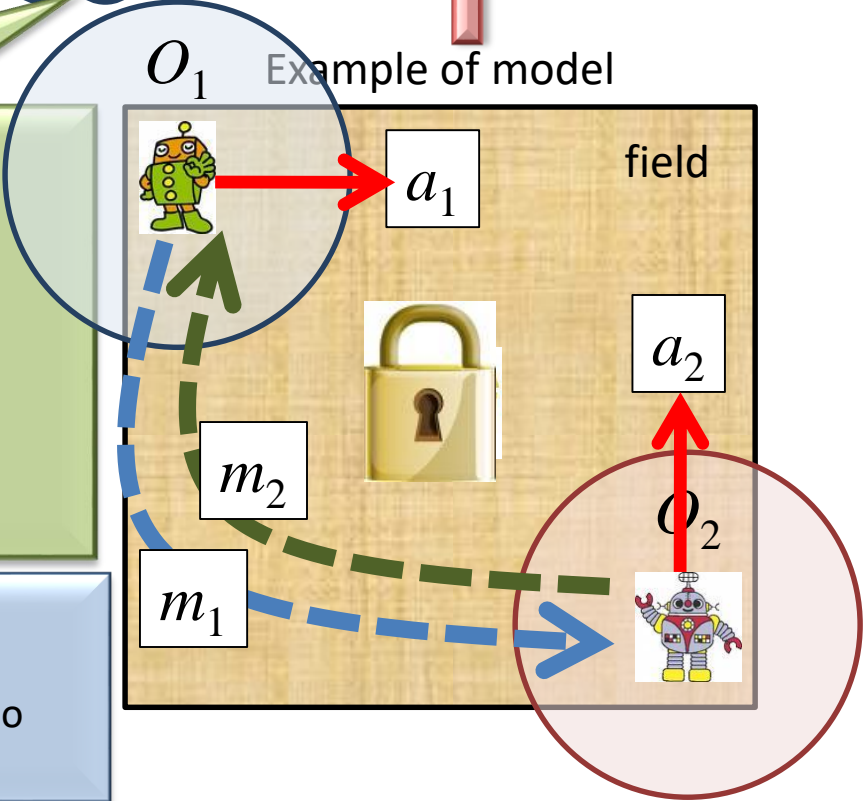(Decentralized Partially Observable Markov Decision Process with Communication)

- A decentralized multi-agent system, where agents can communicate with each other and only observe the restricted information.
- Dec-POMDP-Com := $< I, S, \Omega, A, M, C, P, O, R, T >$

Formulation

Example of model

$O_1$

- $P$ : a transition probability function

- $O$ : an observation probability function

field

$a_1$

$a_2$

$m_2$

$O_2$

$m_1$

- Two agents get a treasure cooperatively.
- The treasure is locked.
- Both agents must reach the treasure at the same time to open the lock.

# Dec-POMDP-Com [Goldman+ 04]

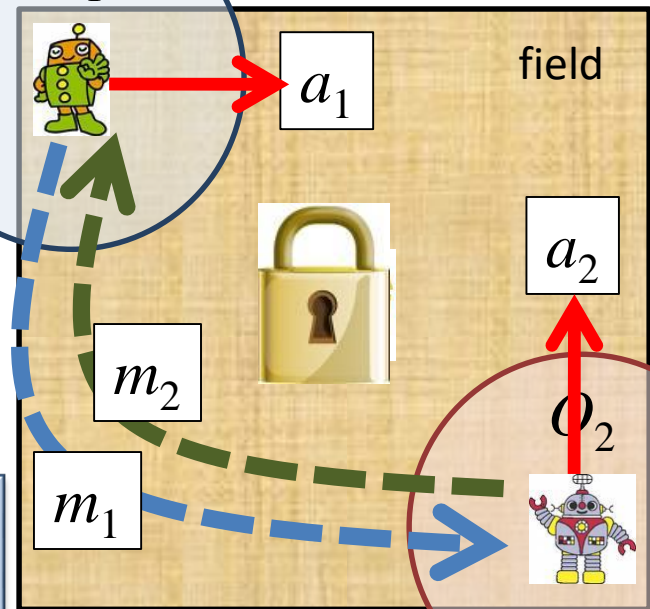(Decentralized Partially Observable Markov Decision Process with Communication)

- A decentralized multi-agent system, where agents can communicate with each other and only observe the restricted information.

- Dec-POMDP-Com := $< I, S, \Omega, A, M, C, P, O, R, T >$

Formulation

$O_1$    Example of model

- $R$ : a reward function
- e.g., the treasure obtained by agents

- $T$ : a time horizon

$a_1$    field

$a_2$

$m_2$

$O_2$

$m_1$

- Two agents get a treasure cooperatively.
- The treasure is locked.
- Both agents must reach the treasure at the same time to open the lock.
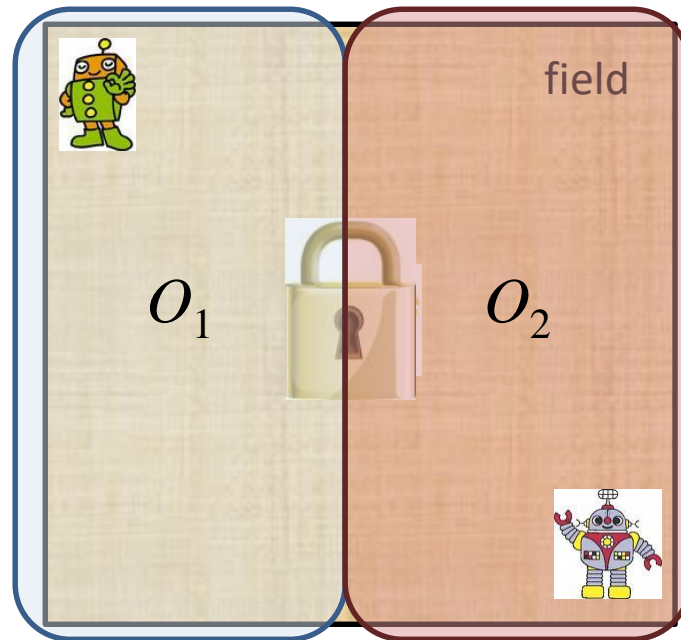
# Outline

☑Background

☑Scheme of talk

☑Review : Dec-POMDP-Com [Goldman+ 04]

☐Constrained model

- 🌐 Jointly Fully Observable Dec-POMDP-Com [Goldman+ 04]
- 🌐 Deterministic Dec-POMDP-Com (we define)

☐Theoretical analysis

☐Conclusion

# Jointly Fully Observable Dec-POMDP-Com

- The Dec-POMDP-Com such that the combination of the agents' observations leads to the global state.

Jointly fully Observable Dec-POMDP-Com



field

$O_1$

$O_2$

$o_1 + o_2 =$ global state
(That is Jointly fully observable)
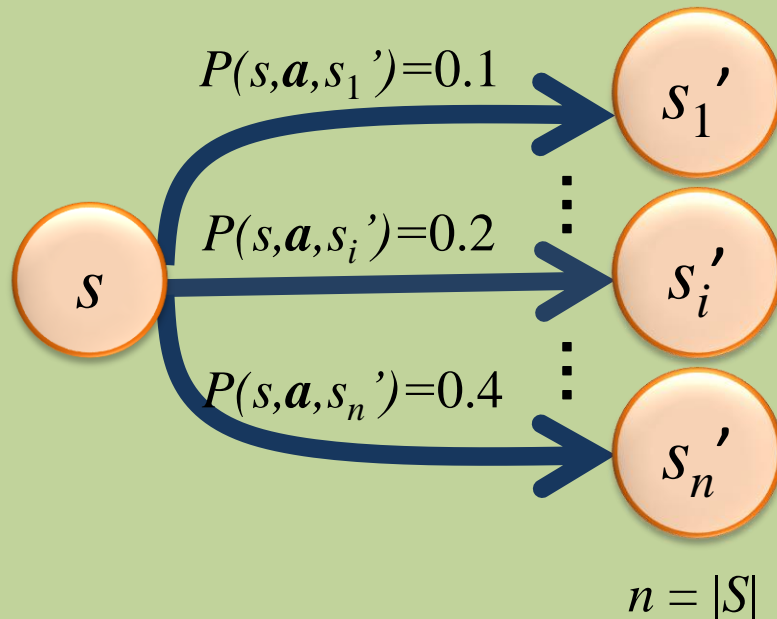
# Deterministic Dec-POMDP-Com

- The model where $P$ and $O$ on the definition are constrained.
- Dec-POMDP-Com := $\langle I, S, \Omega, A, M, C, P, O, R, T \rangle$
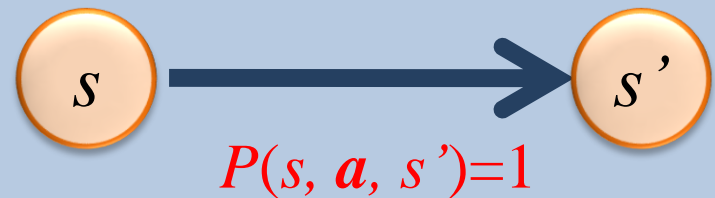
# Deterministic Dec-POMDP-Com

- The model where $P$ and $O$ on the definition are constrained.
- Dec-POMDP-Com := $\langle I, S, \boldsymbol{\Omega}, A, \boldsymbol{M}, C, P, O, R, T \rangle$

**Restriction 1 : Deterministic transitions**

- $P$ is a transition probability function

$P(s,\boldsymbol{a},s_1')=0.1$

$P(s,\boldsymbol{a},s_i')=0.2$

$P(s,\boldsymbol{a},s_n')=0.4$

$s_1'$

$s_i'$

$s_n'$

$s$

$n = |S|$

- For any state $s \in S$ and any joint action $\boldsymbol{a} \in A$, there exists a state $s' \in S$ such that $P(s, \boldsymbol{a}, s') = 1$.

$s \longrightarrow s'$

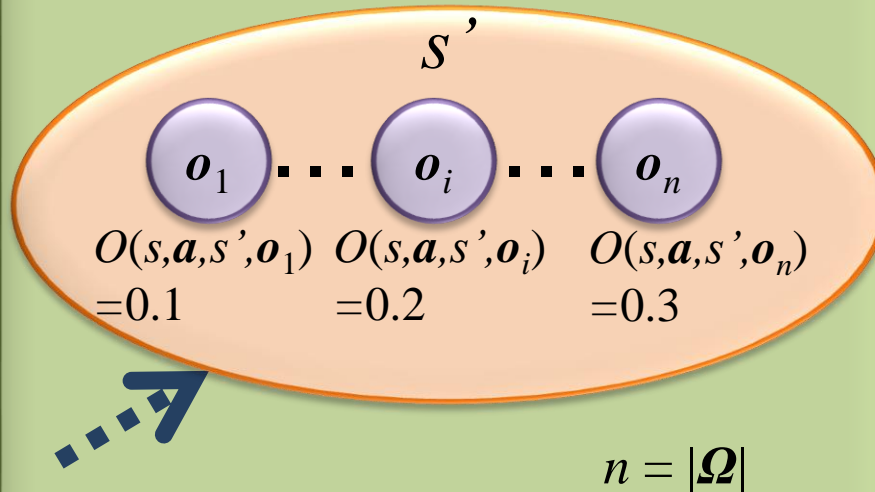$P(s, \boldsymbol{a}, s')=1$

The next global state is decided uniquely.

# Deterministic Dec-POMDP-Com

- The model where $P$ and $O$ on the definition are constrained.
- Dec-POMDP-Com $:= \langle I, S, \boldsymbol{\Omega}, A, \boldsymbol{M}, C, P, O, R, T \rangle$

**Restriction 2 : Deterministic observable**

- $O$ is a observation probability function

$$S'$$

$o_1 \cdots o_i \cdots o_n$

$O(s,\boldsymbol{a},s',\boldsymbol{o}_1)$   $O(s,\boldsymbol{a},s',\boldsymbol{o}_i)$   $O(s,\boldsymbol{a},s',\boldsymbol{o}_n)$
$=0.1$            $=0.2$            $=0.3$

$n = |\boldsymbol{\Omega}|$

- For any state $s$, $s' \in S$ and any joint action $\boldsymbol{a} \in A$, there exists a joint observation $\boldsymbol{o} \in \boldsymbol{\Omega}$ such that $O(s, \boldsymbol{a}, s', \boldsymbol{o}) = 1$.

$$S'$$

$o$

$O(s,\boldsymbol{a},s',\boldsymbol{o})=1$

The current observation is decided uniquely.

# Deterministic Dec-POMDP-Com

- The model  where $P$ and $O$ on the definition are constrained.

- Dec-POMDP-Com := $\langle\, I,\, S,\, \mathbf{\Omega},\, A,\, M,\, C,\, P,\, O,\, R,\, T\, \rangle$

| Restriction 1 : Deterministic transitions | The next global state is decided uniquely. |
| Restriction 2 : Deterministic observable | The current observation  is decided uniquely. |

- When Dec-POMDP-Com  has Restriction 1 and 2,

  it is called Deterministic Dec-POMDP-Com

# Outline

☑Background

☑Scheme of talk

☑Review : Dec-POMDP-Com [Goldman+ 04]

☑Constrained model

- 🌐 Jointly Fully Observable Dec-POMDP-Com [Goldman+ 04]
- 🌐 Deterministic Dec-POMDP-Com (we define)
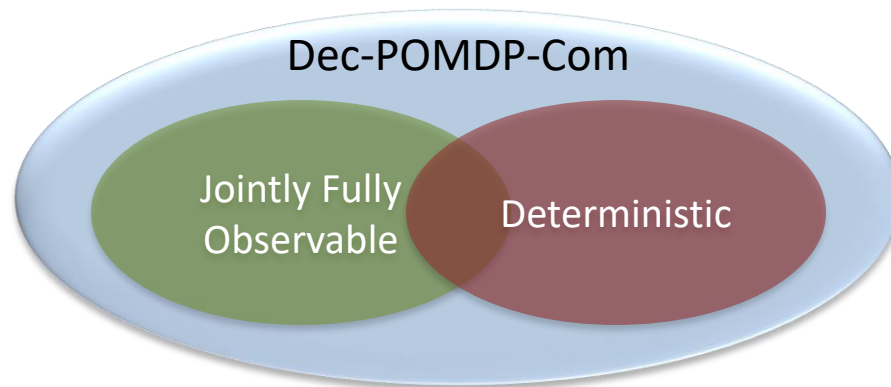
☐Theoretical analysis

☐Conclusion

# Main results

- Corollary 1 :

  *Minimum required sizes* $|M_i|$ for Signal Learning on *Jointly Fully Observable* Dec-POMDP-Com

- Theorem 2 :

  *Minimum required sizes* $|M_i|$ for Signal Learning with Messages on *Deterministic* Dec-POMDP-Com
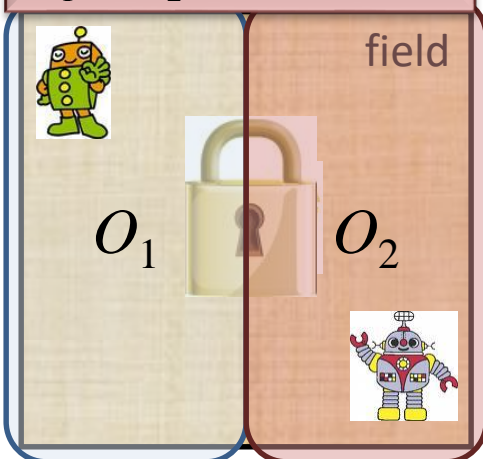
# Theorem 1 [Goldman 04]

**Theorem 1**

For any *jointly fully observable Dec-POMDP-Com*, the following equation holds.

$$\forall M, \quad \max_{\delta \in D} V_{\delta, M}^{T}(s_0) \leqq \max_{\delta \in D} V_{\delta, M'}^{T}(s_0)$$

the value of the optimal joint policy with respect to any joint message set *M*

the value of the optimal joint policy with respect to the joint message set *M' := Ω*.

$O_1 + O_2$ = global state

field

$O_1$   $O_2$

This theorem means that the optimal communication policy of each agent is to send its own observation in jointly fully observable Dec-POMDP-Com (i.e. for agent $i$, $m_i := o_i$).

Each agent can always know the current global state by own observation and received message.
(e.g., for agent 1, $o_1 + m_2 = o_1 + o_2$ = global state)
Therefore, each agent always perform the optimal actions.

# Corollary 1

**Corollary 1**

*For any jointly fully observable Dec-POMDP-Com, if the size $|M_i|$ of the message set of each agent i satisfies the condition,*

$$|M_i| \geqq |\Omega_i|$$

*then the following equation holds:*

$$\max_{\delta \in D^{SL}} V_\delta^T(s_0) = \max_{\delta' \in D} V_{\delta'}^T(s_0)$$

*the value of the optimal joint policy on SL*

*the value of the optimal joint policy with history*

From theorem 1 by Goldman,
the optimal communication policy of each agent is
to send its own observation in jointly fully observable Dec-POMDP-Com.

Therefore, If each agent has $|M_i|$ that is larger than $|\Omega_i|$,
it is possible to constructing the optimal policy such that each agent can
send its own observation.

# Theorem 2

**Theorem 2**

*For any deterministic Dec-POMDP-Com, if the size |M_i| of the message set of each agent i satisfies the condition,*

$$|M_i| \geqq \max_{j \in I} \max_{o \in \Omega_j} |S_j^{obs}(o)|$$

*then the following equation holds:*

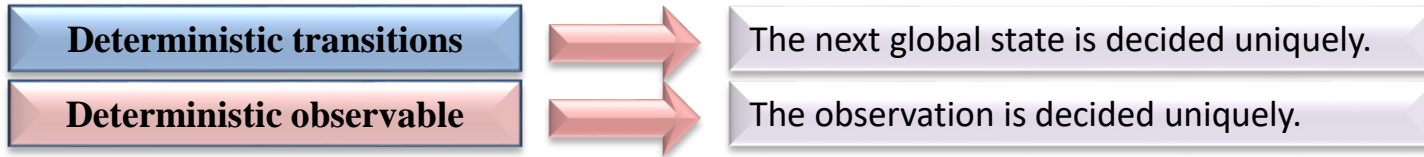$$\max_{\delta \in D^{SLM}} V_{\delta}^T(s_0) = \max_{\delta' \in D} V_{\delta'}^T(s_0)$$

*the value of the optimal joint policy on SLM*

*the value of the optimal joint policy with history*

First, I explain a function $S_j^{obs}$ .

$$S_j^{obs}$$

- Deterministic Dec-POMDP-Com has the following properties.

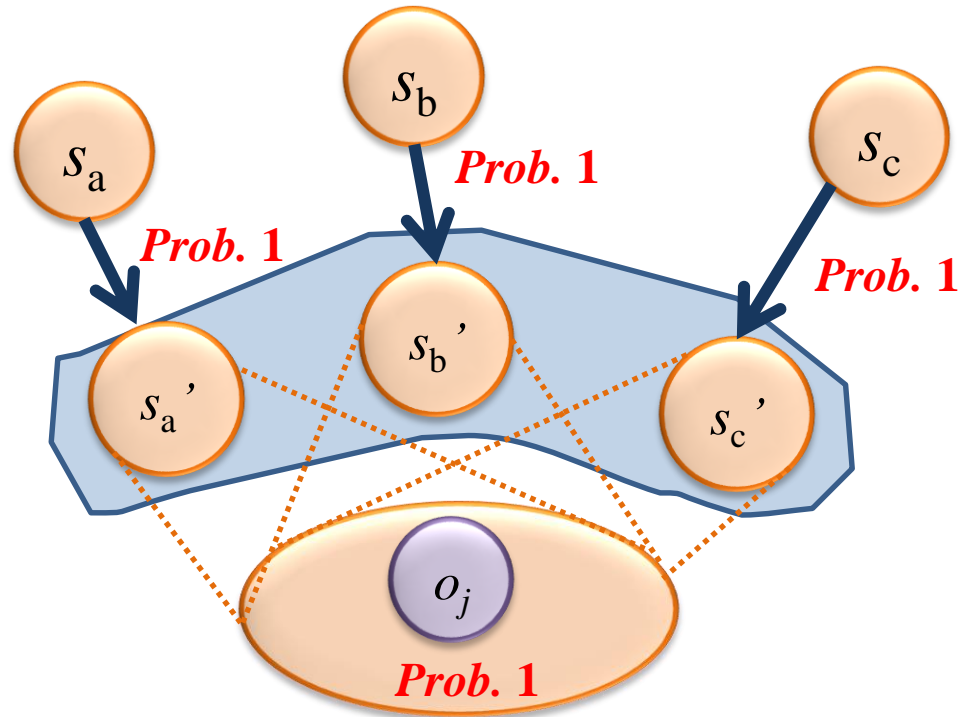| Deterministic transitions | ⇒ | The next global state is decided uniquely. |
| Deterministic observable | ⇒ | The observation is decided uniquely. |

There exists some transitions such that agent $j$ observes the same observation.

From the properties on deterministic Dec-POMDP-Com, we can compute the following function.

$$S_j^{obs}(o_j) = \{ s_a{}', s_b{}', s_c{}' \}$$

$S_j^{obs}$ returns the set of all states where agent $j$ observes $o_j$.
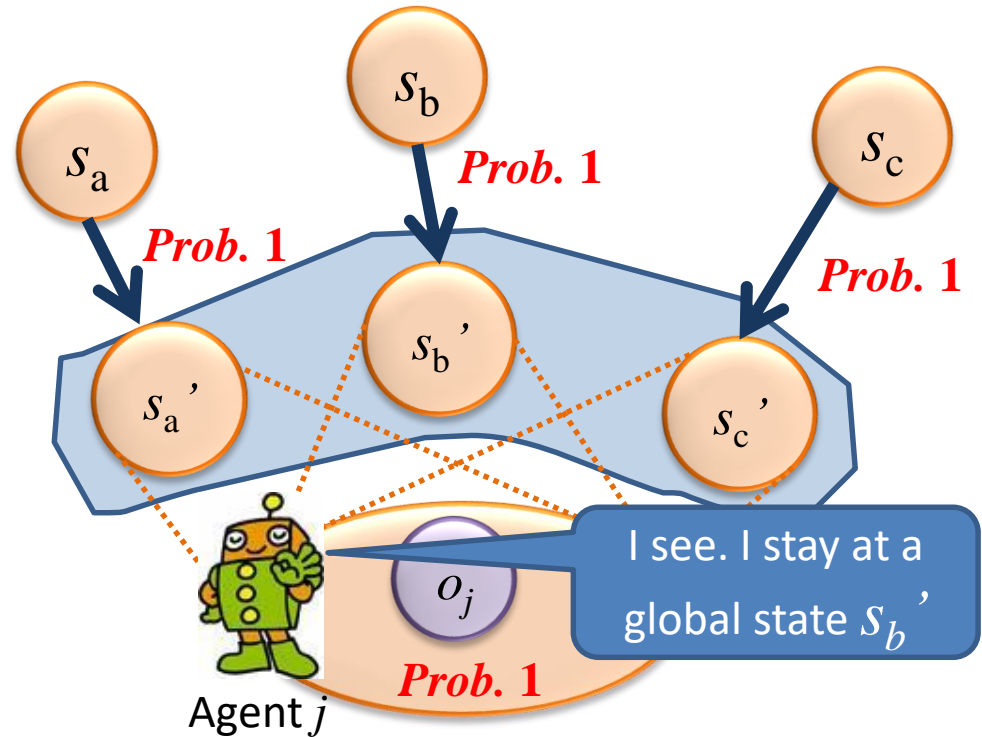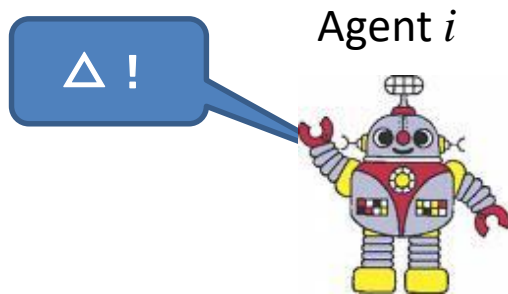
# Proof sketch of Theorem2

The condition of theorem 2 is $|M_i| \geqq \max_{j \in I} \max_{o \in \Omega_j} |S_j^{obs}(o)|$
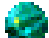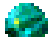
The condition shows that
agent $j$ can know the global state based on the message received from
agent $i$ by setting the set of message which have the maximum size
of $S_j^{obs}(o_j)$ .

$S_j^{obs}(o_j) = \{s_a', s_b', s_c'\}$

$\downarrow \quad \downarrow \quad \downarrow$

$M_i = \{\square, \triangle, \bigcirc\}$



△ !

Agent $i$

$s_a$

$s_b$

$s_c$

Prob. 1

Prob. 1

Prob. 1

$s_a'$

$s_b'$

$s_c'$

$o_j$

Prob. 1

Agent $j$

I see. I stay at a global state $s_b'$

# Outline

☑Background

☑Scheme of talk

☑Review : Dec-POMDP-Com [Goldman+ 04]

☑Constrained model

- 🌐 Jointly Fully Observable Dec-POMDP-Com [Goldman+ 04]

- 🌐 Deterministic Dec-POMDP-Com (we define)

☑Theoretical analysis

☐Conclusion

# Conclusion

- We defined deterministic Dec-POMDP-Com for theoretical analysis

| **Restriction 1 : Deterministic transitions** | The next global state $s'$ is decided uniquely. |
| **Restriction 2 : Deterministic observable** | The observation $o$ is decided uniquely. |

- We showed *Minimum required sizes $|M_i|$* for

  - Signal Learning on *Jointly Fully Observable Dec-POMDP-Com* at corollary 1

  - Signal Learning with Messages on *Deterministic Dec-POMDP-Com* at Theorem 2