**NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE**

# SC3000/CZ3005
# Artificial Intelligence

## Markov Decision Process

Prof Bo AN

*Research area*: artificial intelligence, computational game theory, reinforcement learning, optimization
www.ntu.edu.sg/home/boan
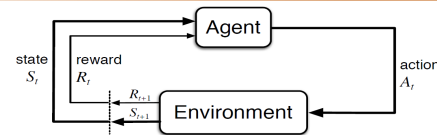*Email*: boan@ntu.edu.sg
*Office*: N4-02b-55

---

# Lesson Outline

- Introduction
- Markov Decision Process
- Two methods for solving MDP
  - Value iteration
  - Policy iteration

---

# Introduction

- We consider a framework for decision making under uncertainty
- Markov decision processes (MDPs) and their extensions provide an extremely general way to think about how we can act optimally under uncertainty
- For many medium-sized problems, we can use the techniques from this lecture to compute an optimal decision policy
- For large-scale problems, approximate techniques are often needed (more on these in later lectures), but the paradigm often forms the basis for these approximate methods

---
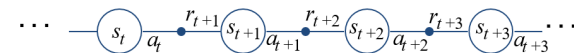
# The Agent-Environment Interface



Agent and environment interact at discrete time steps: $t = 0,1,2,\dots$
Agent:
1. observes state at step $t$: $s_t \in S$
2. Produces action at step $t$: $a_t \in A(s_t)$
3. Gets resulting reward: $r_{t+1}$ and the next state: $s_{t+1} \in S$

*(handwritten top)* ↑ model (knows tran function)
MDP vs RL ↓ We don't know the model

# Making Complex Decisions

- Make a sequence of decisions
  - Agent's utility depends on a sequence of decisions
  - Sequential Decision Making

- Markov Property
  - Transition properties depend only on the current state, not on previous history (how that state was reached)
  - Markov Decision Processes

*(handwritten)* deterministic ·Go. w̅ some assumptions· stochastic

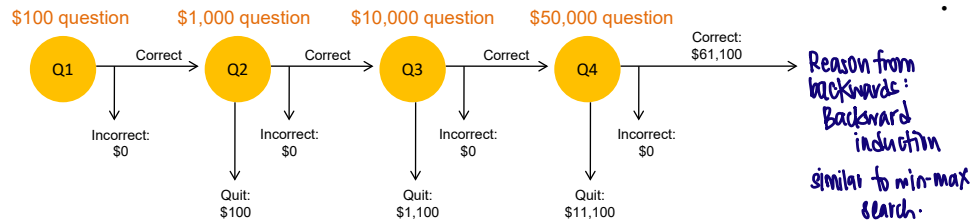# Markov Decision Processes

- Formulate the agent-environment interaction as an MDP
- Components:
  - *Markov* **States** $s$, beginning with initial state $s_0$
  - **Actions** $a$
    - Each state $s$ has actions $A(s)$ available from it
  - **Transition model** $P(s' \mid s, a)$ *(handwritten)* ↙ take action a at state s ⇒ prob that reach state s'
    - *assumption*: the probability of going to $s'$ from $s$ depends only on $s$ and $a$ and not on any other past actions or states
  - **Reward function** $R(s)$, or r(s)
- **Policy** $\pi(s)$: the action that an agent takes in any given state
  - The "solution" to an MDP *(handwritten)* ↓ Goal: Find policy ⇒ How to make the decision.
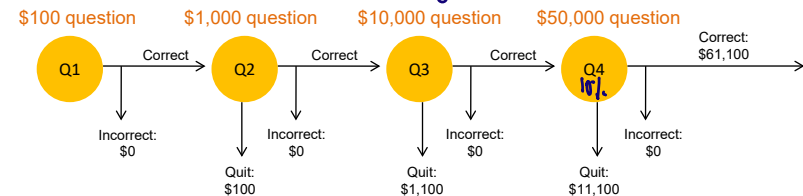
# Game Show

- A series of questions with increasing level of difficulty and increasing payoff
- Decision: at each step, take your earnings and quit, or go for the next question
  - If you answer wrong, you lose everything

$100 question    $1,000 question    $10,000 question    $50,000 question

Q1 —Correct→ Q2 —Correct→ Q3 —Correct→ Q4 —Correct: $61,100→

Q1 Incorrect: $0
Q2 Incorrect: $0 / Quit: $100
Q3 Incorrect: $0 / Quit: $1,100
Q4 Incorrect: $0 / Quit: $11,100

*(handwritten)* Reason from backwards: Backward induction similar to min-max search.

# Game Show

- Consider $50,000 question
  - Probability of guessing correctly: 1/10
  - Quit or go for the question?
- What is the expected payoff for continuing?
  0.1 * 61,100 + 0.9 * 0 = 6,110   *(handwritten)* 6110 < 11,100 ∴ After 4 you should not try.
- What is the optimal decision?

$100 question    $1,000 question    $10,000 question    $50,000 question

Q1 —Correct→ Q2 —Correct→ Q3 —Correct→ Q4 *(handwritten 10%)* —Correct: $61,100→

Q1 Incorrect: $0
Q2 Incorrect: $0 / Quit: $100
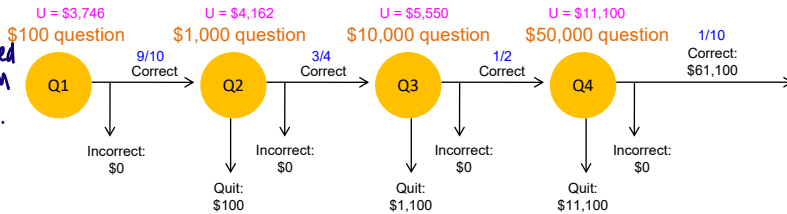Q3 Incorrect: $0 / Quit: $1,100
Q4 Incorrect: $0 / Quit: $11,100

- What should we do in Q3?
  - Payoff for quitting: $1,100
  - Payoff for continuing: 0.5 * $11,100 = $5,550
- What about Q2?
  - $100 for quitting vs. $4,162 for continuing
- What about Q1?

*5550 > 1100*
*∴ you should try.*
*after Q3*

*↳75% x 5550*

*Obviously need to try. if you don't is $0.*

| U = $3,746 | U = $4,162 | U = $5,550 | U = $11,100 |
|---|---|---|---|
| $100 question | $1,000 question | $10,000 question | $50,000 question |



R(s) = -0.04 for every non-terminal state

---

Transition model:

*if you want to go up 80% it will go up, 10% goes left, 10% goes right*

0.1    0.8    0.1

*execution is not perfect*

R(s) = -0.04 for every non-terminal state

---

---

Transition model:

0.8
0.1    0.1

R(s) = -0.04 for every non-terminal state

# Grid World



Optimal policy when R(s) = -0.04 for every non-terminal state

# Atari Video Games



Execute actions to get the maximum accumulated rewards

# Solving MDPs

- MDP components:
  - **States** $s$
  - **Actions** $a$
  - **Transition model** $P(s' \mid s, a)$ ⇒ *after taking some action ⇒ prob of reaching another state*
  - **Reward function** $R(s)$
- The solution:
  - **Policy** $\pi(s)$ mapping from states to actions
  - How to find the optimal policy?
    *max total acc reward*

# Maximizing Accumulated Rewards

→ *continue until terminal state.*

- The optimal policy should maximise the accumulated *rewards* over given a trajectories like $\tau = \langle S_1, A_1, R_1, \ldots, S_T, A_T, R_T \rangle$ under some policies:

*if we treat γ=1, we are just summing up total reward.*

$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \cdots \gamma^K R_{t+K} = \sum_{k=0}^{K} \gamma^k R_{t+k}$$

→ *what we care about is the total reward.*

- How to define the accumulated rewards of a state sequence?

*-γ ~ | usually i.e 0.99, 0.98*

  - Discounted sum of rewards of individual states
  - Problem: infinite state sequences
  - If finite, LP can be applied

# Accumulated Rewards

*Why we use* (handwritten)
*discounted value* (handwritten)
*-money total is worth more compared to tmr.* (handwritten)

- Normally, we would define the *accumulated rewards* of trajectories as the discounted sum of the rewards
- **Problem:** infinite time horizon
- **Solution:** *discount* the individual rewards by a factor $\gamma$ between 0 and 1:

$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \cdots$$
$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k} \leq \frac{R_{max}}{1-\gamma} \quad (0 < \gamma < 1)$$

- Sooner rewards count more than later rewards
- Makes sure the total *accumulated rewards* stays bounded
- Helps algorithms converge

*If there's no discount goes to ∞ computer cannot handle.* (handwritten)

# Value Function

*$V^\pi(s)$ = expected utility starting in s & acting according to $\pi$* (handwritten)
*$V^{\pi*}(s)$ = expected utility starting in s & acting optimally* (handwritten)

- The "true" value of a state, denoted $V(s)$, is the expected sum of discounted rewards if the agent executes an *optimal* policy starting in state $s$

$$V^\pi(s) = \mathbb{E}_\pi[G_t \mid S_t = s]$$

- Similarly, we define the action-value of a state-action pair as

$$Q^\pi(s,a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$$

- The relationship between Q and V

$$V^\pi(s) = \sum_{a \in A} Q^\pi(s,a)\pi(a|s)$$

*policy* (handwritten)
*prob we will take the action a in state s* (handwritten)

*Action value function* (handwritten)
*$Q^\pi(s,a)$ = expected utility taking a in s & the following $\pi$* (handwritten)

*$Q(s,a_1)=1$  $Q(s,a_2)=10$  $Q(s,a_3)=8$* (handwritten)
*Returns* (handwritten)
*$\pi'(s)$ ⇒ Returns max $Q(s,a)$ ⇒ $a_2$* (handwritten)
*$V(s)$ ⇒ Returns highest value ⇒ 10.* (handwritten)

# Finding the Value Function of States

*Bellman eqn.* (handwritten)

Max node — S

*$a'$   $a''$   $a'''$* (handwritten)
*If you take action a & state s, this is what you get* (handwritten)

Chance node — s, a

*$P(s' | s, a)$   $r(s,a,s')$* (handwritten)
*When you take can action, due to uncertainty you might reach other unexpected states* (handwritten)

S' — *$V(s')$* (handwritten)

- What is the expected value of taking action $a$ in state $s$?

*prob that you will reach diff $s'$* (handwritten)
*prob that you reach diff $s'$* (handwritten)
*after you reach $s'$, you get immediate reward* (handwritten)

$$Q(s,a) = \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma * V(s')]$$

*plus future.* (handwritten)

- How do we choose the optimal action?

$$\pi^*(s) = \operatorname{argmax}_{a \in A} \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V(s')]$$

*$Q(s,a)$* (handwritten)

- What is the recursive expression for $V(s)$ in terms of the utilities of its successor states?

*$Q(s,a)$* (handwritten)

$$V(s) = \max_{a \in A} \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V(s')]$$

*optimal value at this state s* (handwritten)
*you choose the best action that max a function* (handwritten)

# The Bellman Equation

- Recursive relationship between the accumulated rewards of successive states:

$$V(s) = \max_{a \in A} \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V(s')]$$

- For *N* states, we get *N* equations in *N* unknowns
  - Solving them solves the MDP
  - We could try to solve them through expectimax search, but that would run into trouble with infinite sequences
  - Instead, we solve them algebraically
  - Two methods: **value iteration** and **policy iteration**

*you can take many action* (handwritten)
*∴ you take action that max* (handwritten)

S

s, a — Choose optimal action $a$

Receive reward $r(s,a,s')$

s'

End up here with $P(s' | s, a)$
Get value $V(s')$
(discounted by $\gamma$)

# Method 1: Value Iteration

- Start out with every $V_0(s) = 0$
- Iterate until convergence
  - During the *i*th iteration, update the value of each state according to this rule:

$$V_{i+1}(s) = \max_a \sum_{s' \in S} P(s'|s,a)\left[R(s,a,s') + \gamma V_i(s')\right]$$

- In the limit of infinitely many iterations, guaranteed to find the correct values
  - In practice, don't need an infinite number of iterations…

# Value Iteration: Example

- A simple example to show how VI works
  - State, action, reward (non-zero) and transition probability are shown in the figure
  - We use this example as an MDP and solve it using VI and PI



These structure can be easily implemented by `dict` of Python or `HashMap` of Java

# Value Iteration: Example (cont'd)

- Given the states and actions are finite, we can use matrices to represent the value function $V(s)$
- Pseudo-code of VI

  1. Initialize $V_0(s)$, for all $s$
  2. For $i = 0, 1, 2 …$
  3. $V_{i+1}(S) = \max_a \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V_i(s')]$ for all state $s$

# Value Iteration: Example (cont'd)

- How VI works in an iteration? Given iteration $i = 0$

For all state find

$$V_{i+1}(S) = \max_a \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V_i(s')]$$

We can instead calculate Q(s, a) values for each $s$ and $a$ and get the best V(s)

$$Q_{i+1}(s,a) = \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V_i(s')]$$

Finally

$$V_{i+1}(S) = \max_a \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V_i(s')] = \max_a Q_{i+1}(s,a)$$

# Value Iteration: Example (cont'd)

- We use the previous example and solve it using VI
- Construct a Q table, Q(s, a)
  - 3 rows (3 states) and 2 columns (2 actions)
  - $V_0 = [0, 0, 0]$  # 3 states

$Q_0$

| | $a_0$ | $a_1$ |
|---|---|---|
| $s_0$ | 0 | 0 |
| $s_1$ | 0 | 0 |
| $s_2$ | 0 | 0 |

$V_0$

| |
|---|
| 0 |
| 0 |
| 0 |

$Q_1$

| | $a_0$ | $a_1$ |
|---|---|---|
| $s_0$ | 0 | 0 |
| $s_1$ | 3.5 | 0 |
| $s_2$ | 0 | -0.3 |

$V_1$

| |
|---|
| 0 |
| 3.5 |
| 0 |

The initial Q and V table

$$Q_{i+1}(s,a) = \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V_i(s')]$$

$$Q_1(s_1, a_0) = P(s_0|s_1, a_0)[r(s_1, a_0, s_0) + \gamma V_o(s_0)] + P(s_1|s_1, a_0)[r(s_1, a_0, s_1) + \gamma V_o(s_1)] + P(s_2|s_1, a_0)[r(s_1, a_0, s_2) + \gamma V_o(s_2)]$$

$$Q_1(s_1, a_0) = 0.7 * [5 + 0] + 0.1 * [0 + 0] + 0.2 * [0 + 0] = 3.5$$

Repeat this process until it converges

---

# Value Iteration: Example (cont'd)

- Iterating the process (code available on later slides)
  - Discount factor 0.9

```
iter   0  |   V(s0) = 0.000   V(s1) = 0.000   V(s2) = 0.000

iter   1  |   V(s0) = 0.000   V(s1) = 3.500   V(s2) = 0.000

iter   2  |   V(s0) = 0.000   V(s1) = 3.815   V(s2) = 1.890

iter   3  |   V(s0) = 1.701   V(s1) = 4.184   V(s2) = 2.060

. . . . . .

iter  63  |   V(s0) = 8.020   V(s1) = 11.160   V(s2) = 8.912

iter  64  |   V(s0) = 8.021   V(s1) = 11.161   V(s2) = 8.913

iter  65  |   V(s0) = 8.022   V(s1) = 11.162   V(s2) = 8.915
```

Very good, the values converge now

---

# Value Iteration: Example (cont'd)

- Put the optimal values on the graph

---

# Value Iteration: Example (cont'd)

during this process, we don't calc the process
- just apply Bellman eqn
- until conv & the use argmax to find the policy

- Use $V^*$ to find optimal policy
  - aka optimal actions in each state

$$\pi^*(s) = \arg\max_a \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V_i(s')] = \arg\max_a Q_i(s,a)$$

$\pi^*$

| | $a^*$ |
|---|---|
| $s_0$ | $a_1$ |
| $s_1$ | $a_0$ |
| $s_2$ | $a_0$ |

Done! We get the optimal policy of the example

Python Code on Colab: https://colab.research.google.com/drive/1DnYIr3QJxpfs_rR_jAAUrqHZMjvsjGSx?usp=sharing

# Method 2: Policy Iteration

- Start with some initial policy $\pi_0$ and alternate between the following steps:
  - **Policy evaluation:** calculate $V^{\pi_i}(s)$ for every state $s$, *like VI*
  - **Policy improvement:** calculate a new policy $\pi_{i+1}$ based on the updated utilities

$$\pi_{i+1}(s) = \text{argmax}_{a \in A} \sum_{s'} P(s'|s,a)[r(s,a,s') + \gamma V^{\pi_i}(s')]$$



starting
$V$ $\pi$

$V = V\pi$

$V^*$
$\pi^*$

$\pi = greedy(V)$

evaluation
$V \rightarrow V^\pi$

$\pi$

$V$

$\pi \rightarrow greedy(V)$

improvement

Keep repeating
· this until
value converges.

Policy evaluation Estimate $v_\pi$
  Any policy evaluation algorithm
Policy improvement Generate $\pi' \geq \pi$
  Any policy improvement algorithm

$\pi^* \rightleftarrows V^*$

---

# Policy Iteration: Main Steps

- Unlike VI, policy iteration has to maintain a policy - chosen actions from all states - and estimate $V^{\pi_i}$ based on this policy.
  - Iterate this process until convergence (like VI)
- Steps of PI
  1. Initialization
  2. Policy Evaluation (calculating the $V$)
  3. Policy Improvement (calculate the policy $\pi$)

---

# Policy Iteration: Detailed Steps

1. Initialization
   a. Initialize $V(s)$ and $\pi(s)$ for all state $s$

2. Policy Evaluation (calculate the $V$)
   a. Repeat
   b. $\Delta \leftarrow 0$
   c. For each state $s$:
   d. $v \leftarrow V(s)$
   e. $V(s) \leftarrow \sum_{s'} p(s'|s, \pi(s))[R(s, \pi(s), s') + \gamma V(s')]$
   f. $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
   g. until $\Delta < \omega$ (a small positive number)

3. Policy Improvement (calculate the new policy $\pi$)
   a. *policystable* $\leftarrow$ *True*
   b. For each state $s$:
   c. $b \leftarrow \pi(s)$
   d. $\pi(s) \leftarrow \text{argmax}_a \sum_{s'} p(s'|s, a)[R(s, a, s') + \gamma V(s')]$
   e. If $b \neq \pi(s)$, then *policystable* $\leftarrow$ *False*
   f. If *policystable* $\leftarrow$ *True*, then stop; else go to step 2;

---

# Policy Iteration: Policy Evaluation

- Use the example used in VI
- Start iteration $i = 0$, initialize random $\pi$, $\gamma = 0.99$, and $V(s) = 0$ for all state $s$

$$V^{\pi_i}(s) = \sum_{s'} P(s, \pi_i(s), s')[R(s, \pi_i(s), s') + \gamma V^{\pi_i}(s')]$$

$$V(s_0) \leftarrow P(s_0, \pi(s_0), s_0)[R(s_0, \pi(s_0), s_0) + \gamma V(s_0)] +$$
$$P(s_0, \pi(s_0), s_1)[R(s_0, \pi(s_0), s_1) + \gamma V(s_1)] +$$
$$P(s_0, \pi(s_0), s_2)[R(s_0, \pi(s_0), s_2) + \gamma V(s_2)]$$

$$V(s_1) \leftarrow P(s_1, \pi(s_1), s_0)[R(s_1, \pi(s_1), s_0) + \gamma V(s_0)] +$$
$$P(s_1, \pi(s_1), s_1)[R(s_1, \pi(s_1), s_1) + \gamma V(s_1)] +$$
$$P(s_1, \pi(s_1), s_2)[R(s_1, \pi(s_1), s_2) + \gamma V(s_2)]$$

$$V(s_2) \leftarrow P(s_2, \pi(s_2), s_0)[R(s_2, \pi(s_2), s_0) + \gamma V(s_0)] +$$
$$P(s_2, \pi(s_2), s_1)[R(s_2, \pi(s_2), s_1) + \gamma V(s_1)] +$$
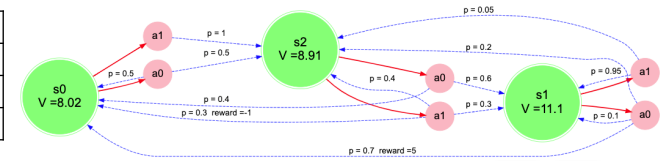$$P(s_2, \pi(s_2), s_2)[R(s_2, \pi(s_2), s_2) + \gamma V(s_2)]$$

## Policy Iteration: Policy Evaluation (cont'd)

- Use the example used in VI
- Start iteration $i = 0$, $\gamma = 0.99$, initialize random $\pi = [a_1, a_0, a_1]$ and $V(s) = 0$ for states $s_0$, $s_1$ and $s_2$

$V(s_0) \leftarrow 0.0[0.0 + \gamma V(s_0)] + 0.0[0.0 + \gamma V(s_1)] + 1.0[0.0 + \gamma V(s_2)]$

$V(s_1) \leftarrow 0.7[5.0 + \gamma V(s_0)] + 0.1[0.0 + \gamma V(s_1)] + 0.2[0.0 + \gamma V(s_2)]$

$V(s_2) \leftarrow 0.3[-1 + \gamma V(s_0)] + 0.3[0.0 + \gamma V(s_1)] + 0.4[0.0 + \gamma V(s_2)]$

*(handwritten annotations:)* Reward =0; transition to $s_0, s_1$ is 0 & $s_2 = 1$; 70% goes so with reward =5; 10% stays in $s_1$ According to initial policy; Reward of reaching state; prob of going to state by taking $a_1$

- Values are calculated asynchronously
- The converged values are used for policy improvement

$V(s_0) = 0$

$V(s_1) = 0.7 * 5.0 = 3.5$

$V(s_2) = 0.3 * (-1) + 0.3 * (0.99 * 3.5) = 0.7395$

... loop this process as instructed in step 2 ...

*(handwritten: takes $a_1$ goes to $s_2$ only; p = 0.05; p = 1; p = 0.2; p = 0.5; p = 0.5; p = 0.4; p = 0.6; p = 0.95; p = 0.3; p = 0.4; p = 0.3 reward =-1; p = 0.7 reward =5; takes $a_0$; p = 0.1)*

---

## Policy Iteration: Example (cont'd)

- The optimal Q values are
  - Nearly the same as that of VI (as shown in the figure below)

| | $V^*$ |
|---|---|
| $s_0$ | 8.03 |
| $s_1$ | 11.2 |
| $s_2$ | 8.9 |

*(diagram with states: s0 V =8.02, s2 V =8.91, s1 V =11.1; p = 0.05, p = 1, p = 0.5, p = 0.2, p = 0.95, p = 0.5, p = 0.4, p = 0.6, p = 0.3, p = 0.4 reward =-1, p = 0.7 reward =5, p = 0.1)*

  - We can easily calculate the optimal policy. Can you try it?

$\pi^*$

| | $a^*$ |
|---|---|
| $s_0$ | $a_1$ |
| $s_1$ | $a_0$ |
| $s_2$ | $a_0$ |

*(handwritten: argmax)*

---

## Further Reading [AAAI'18: http://www.ntu.edu.sg/home/boan/papers/AAAI18_Malmo.pdf]

We Won 2017 Microsoft Collaborative AI Challenge

- Collaborative AI
  - *How can AI agents learn to recognise someone's intent (that is, what they are trying to achieve)?*
  - *How can AI agents learn what behaviours are helpful when working toward a common goal?*
  - *How can they coordinate or communicate with another agent to agree on a shared strategy for problem-solving?*

---

## Further Reading [AAAI'18: http://www.ntu.edu.sg/home/boan/papers/AAAI18_Malmo.pdf]

- Microsoft Malmo Collaborative AI Challenge
  - *Collaborative mini-game, based on an extension "stag hunt"*
  - *Uncertainty of pig movement*
  - *Unknown type of the other agent*
  - *Detection noise (frequency 25%)*
- Our team HogRider won the challenge (out of more than 80 teams from 26 countries)
  - *learning + game theoretic reasoning + sequential decision making + optimisation*



Catch the pig by HogRider
AMI Research Group
Nanyang Technological University