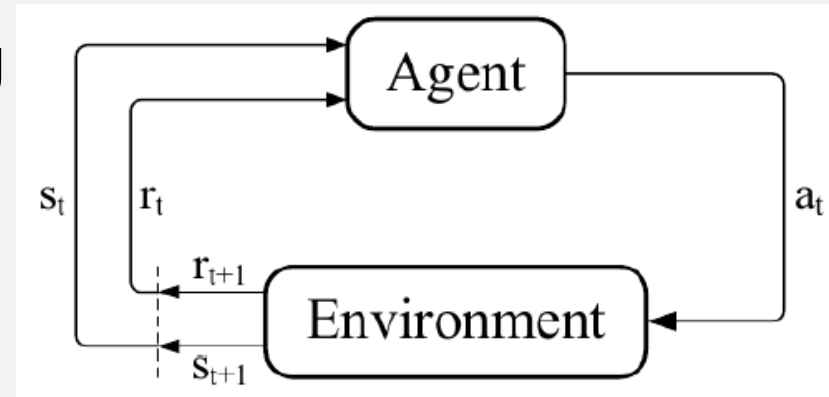


# Tanulás elosztott rendszerekben/3

# Multi Agent Reinforcement Learning

## Többágenses megerősítéses tanulás (MARL)



Kezdünk egy ágenssel.

Legyenek a környezeti **állapotai**  $s$ -ek, **cselekvései**  $a$ -k, az ágens cselekvéseit meghatározó **eljárás mód**  $\pi$ , ill. ágens **cselekvés-érték** függvénye  $Q(s,a)$ .

Az állapotok és a cselekvések közötti kapcsolatot az un. **Markov döntési folyamat (MDF)** írja le,  $T(s,a,s')$  **átmenet-valószínűségel**. Egyes állapotokban ágens  $r(s,a,s')$  közvetlen **megerősítést** kap.

Ágens célja megállapítani azt az optimális eljárásmódot, ami a **diszkont hátralévő jutalmat** (az  $s_k$  állapottól végtelen jövőbe) maximálja, ahol  $\gamma$  a diszkont faktor és  $r$  a megerősítés.

$$R_k = E \left\{ \sum_{j=0}^{\infty} \gamma^j r_{k+j+1} \right\}$$

Adott eljárásmód mellett az ágens **cselekvés-érték** függvényt tanul

$$Q^\pi(s,a) = E \left\{ \sum_{j=0}^{\infty} \gamma^j r_{k+j+1} \mid s_k = s, a_k = a, \pi \right\}$$

# Multi Agent Reinforcement Learning

## Többágenses megerősítéses tanulás (MARL)

A lehető legjobb eredmény az optimális cselekvés-érték függvény:

ami teljesíti az un. **Bellman egyenletet**:  $Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a)$

$$Q^*(s, a) = \sum_{s' \in \mathcal{S}} T(s, a, s') \left[ r(s, a, s') + \gamma \max_{a'} Q^*(s', a') \right]$$

Az ágens eljárásmodja **mohó**: ami optimális, ha Q is optimális.

$$\bar{\pi}(s) = \arg \max_a Q(s, a)$$

A Bellman-egyenlet ismeretlen  $r$  és  $T$  mellett az un. **Q-tanulással** oldható meg (jelen formában **időkülönbség Q-tanulással**):

$$\begin{aligned} \overset{\text{új}}{Q}(s_k, a_k) &= Q(s_k, a_k) + \alpha \left[ r + \gamma \max_{a'} Q(s_{k+1}, a') - Q(s_k, a_k) \right] \\ &= (1 - \alpha) Q(s_k, a_k) + \alpha \left[ r + \gamma \max_{a'} Q(s_{k+1}, a') \right] \end{aligned}$$

# Multi Agent Reinforcement Learning

## Többágenses megerősítéses tanulás (MARL)

Q-tanulás bizonyos feltételek mellett optimális Q-hoz konvergál.

A legfontosabb feltétel, hogy a **tanuló ágensnek nem nulla valószínűséggel ki kell próbálnia minden létező cselekvését.**

Nem tud tehát csak mohó lenni, a mohóságát **felfedezési** igénnyel kell vegyítenie.

A „mohóság + felfedezés” keverékviselkedést biztosítani tudjuk:

- **$\epsilon$ -mohósággal**: az ágens  $\epsilon$  valószínűséggel véletlen cselekvést választ, ill.  $1-\epsilon$  valószínűséggel mohó, vagy
- **Boltzmann-felfedezési modellel**, ahol egy  $a$  cselekvés megválasztásának valószínűsége egy  $s$  állapotban:

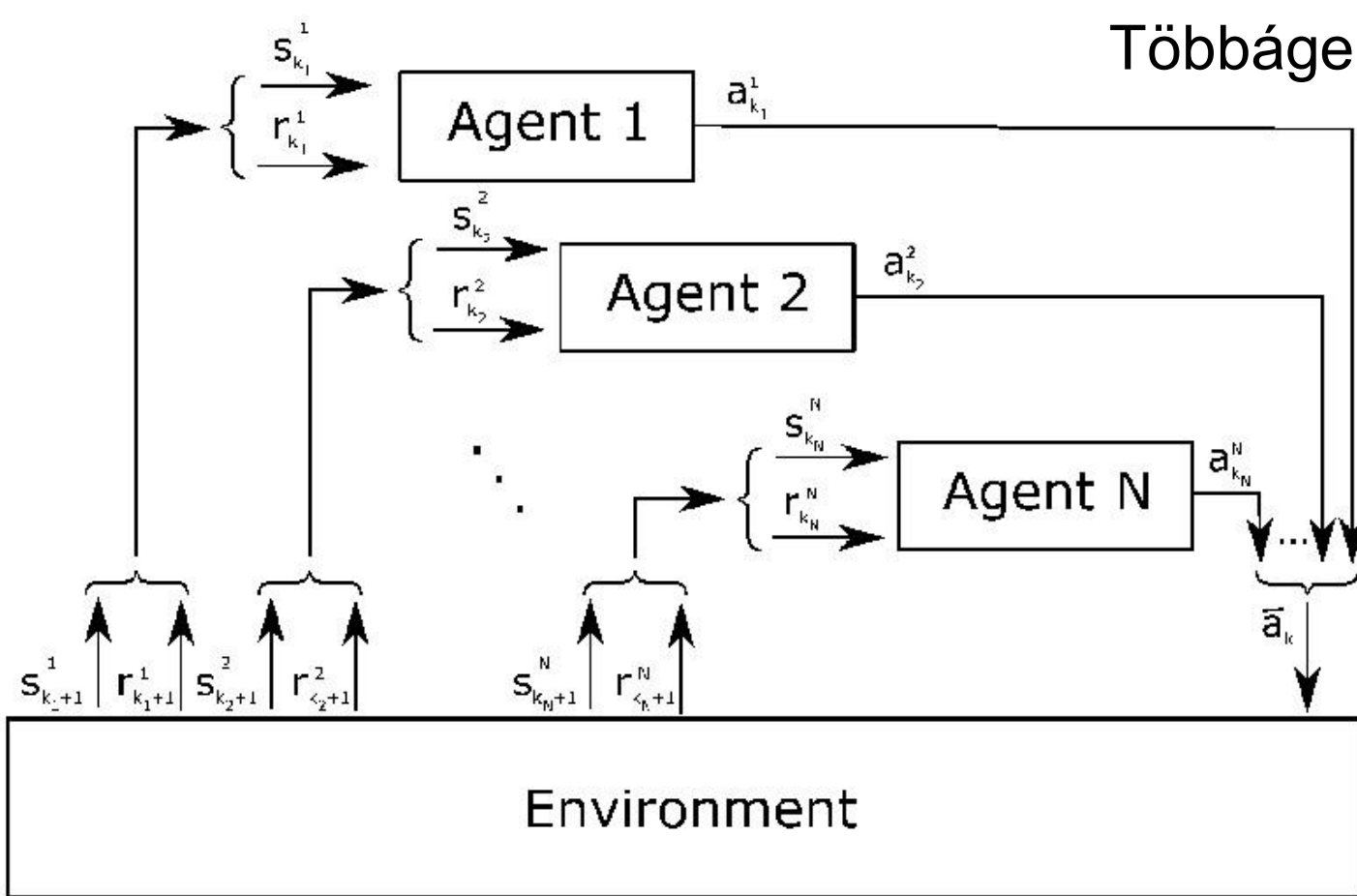
ahol a  $T$  „hőmérséklet” a két véglet között szabályoz.

ha  $T \rightarrow \infty$ , akkor a választás tisztán **(egyenletesen) véletlen**,

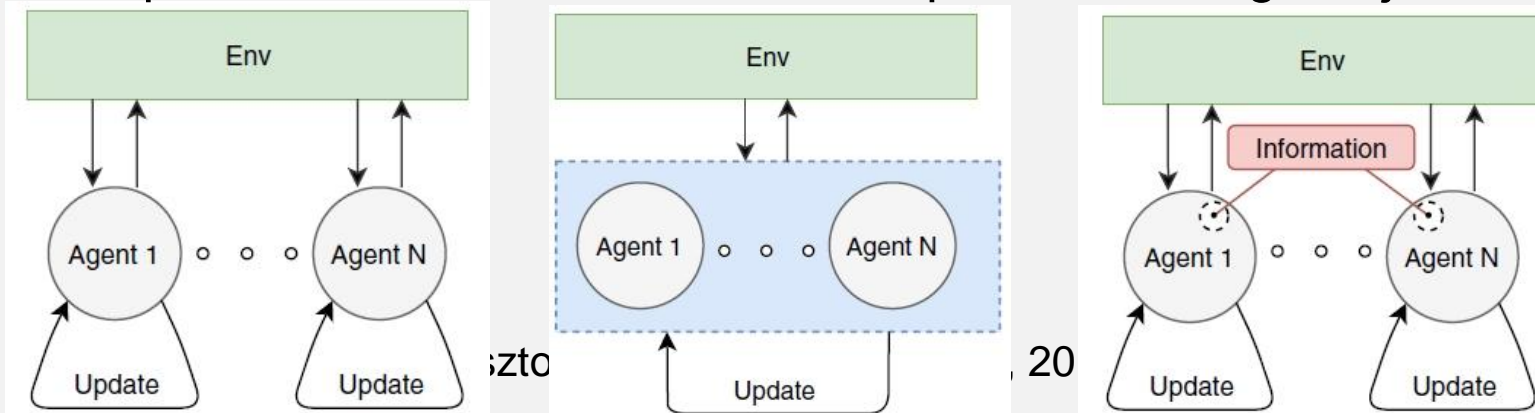
ha  $T \rightarrow 0$ , akkor a választás **mohó**.

$$\pi(s, a) = \frac{e^{Q(s, a)/T}}{\sum_{a'} e^{Q(s, a')/T}}$$

# Többágenses eset

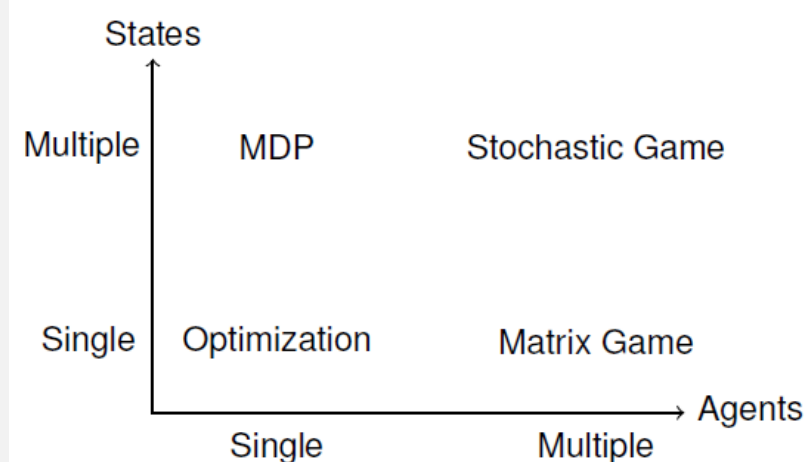


## Elosztott/központosított tanulás - Elosztott/központosított végrehajtás



## Többágenses eset:

- **(matrix)játék**, hasznossági mátrixsal definiált,
- **ismételt játék**, minden fordulóban ugyanazt a mátrixjátékot játsszák,
- **sztochasztikus játék**, a MDF többágenses kiterjesztése. Az állapotátmeneteket és a kapott megerősítést az összes ágens **együttes cselekvése** határozza meg, és az egyedi ágensek eljárás módjai mellett beszélünk az **együttes eljárásmódról** is.  
Mindegyik állapotban az ágensek új mátrix játékot játszanak, aminek mátrixát a tanult hasznosságok határozzák meg.



$$\left( N, S, A = \prod_{k \in N} A_k, T, \{R_k\} \right), \quad T : S \times A \times S \rightarrow [0, 1], \quad R_k : S \times A \rightarrow \mathfrak{R}$$

Megjegyzés:

- Egy mátrix játékban mindegyik ágens megerősítése/ hasznossága függ az állapottól és az összes ágens együttes cselekvésétől (joint action, joint learners)
- MDF a sztochasztikus játék egyágenses esete
- Ismételt játék a sztochasztikus játék egyetlenegy állapotú esete.

$$R_i(s, \mathbf{a})$$

Játék lehet **modell-alapú**. Akkor az ágens először megtanulja az ellenfél stratégiáját, majd talál rá a legjobb választ.

Lehet **model-nélküli** is, amikor az ágens az ellenfélre jó választ adó stratégiát tanulja meg anélkül, hogy az ellenfél stratégiáját explicite kitanulná.

Jelölje egy-egy ágens megerősítését generáló függvényt  $\rho_i$ . Beszélhetünk akkor

- **teljesen kooperatív** ágensrendszerekről  $\rho_1 = \dots = \rho_n$

- **teljesen versengő** ágensrendszerekről, ill.  $\rho_1 = -\rho_2$  (két ágens, zérus összegű)

-  $\rho_1 + \rho_2 + \dots + \rho_n = 0$ , több ágens esetén

- **vegyes** ágensrendszerekről (általános összegű, ahol semmilyen feltétel nem adható).

$$\rho_1 + \rho_2 + \dots + \rho_n > 0$$

Minden zérus-összegű mátrixjátéknak van NE-ja tiszta stratégiákban. (Neumann)

Minden általános összegű mátrixjátéknak van NE-ja. (Nash)

Minden teljesen versengő sztochasztikus játéknak van NE-ja (Shapley)

Minden általános (vegyes) sztochasztikus játéknak van NE-ja. (Fink)

# Többágenses megerősítéses tanulás problémái:

Alapvető problémák:

**nem stacionárius** – a szokásos (egy ágenses) bizonyítható konvergencia lehetetlen.  
**koordinálás igénye** (pl. több NE esetén)

Mi legyen a tanulás célja?

- (1) **Stabilitás** - Konvergencia stratégiában valamilyen egyensúlyhoz (pl. NE), ha a saját maga ellen játszik (**self-play**, minden ágens u.a-t az algoritmust használja)
- (2) **Adaptivitás** - Az ellenfél stratégiájának sikeres megtanulása.
- (3) Egy bizonyos hasznossági szintet **túlhaladó** nyereségek megszerzése.

Milyen tulajdonságokkal rendelkezzen egy tanulási algoritmus?

- (1) **Biztonságos** – garantálja legalább a minimax szintű nyereséget.
- (2) **Konzisztens** – legalább ilyen jó, mint az egyensúlyi esetre számított legjobb válasz (**best response**).
- (P1) **Konvergencia** - Konvergáljon egy stacionárius eljárás módhoz.
- (P2) **Racionalitás** - Ha az ellenfél egy stacionárius stratégiához konvergál, a tanulónak a legjobb válaszhoz kell konvergálnia.



Mechanizmus A legjobb válasz tanulása  
Egyensúly tanulása

## Egyensúlyok pro és kontra

### **Pro**

Egy egyensúly feltételeket fogalmaz meg, hogy a tanulásnak mikor le kell, le kellene állnia.

Egyensúlyban egyszerűbb a játék (kevesebb erőforrás, számítás)

### **Kontra**

NE nem egy előírás jellegű.

Több egyensúly is lehet (koordinálás).

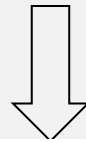
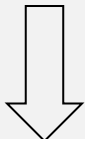
Valamilyen konkrét egyensúlyra vonatkozó ötletadás nem lehetséges (vagy csalás)

Az ellenfél nem biztos, hogy az egyensúlyra vágyik (és játszik).

NE egyensúly számításának komplexitása általános játékok esetén tiltó lehet.

# Egy ágenstől több ágensig

$$Q^{új}(s_k, a_k) = Q(s_k, a_k) + \alpha [r + \gamma \max_{a'} Q(s_{k+1}, a') - Q(s_k, a_k)]$$



Mások cselekvései is

$$Q^{új}(s_k, \mathbf{a}_k) = Q(s_k, \mathbf{a}_k) + \alpha [r + \gamma \max_{\mathbf{a}'} Q(s_{k+1}, \mathbf{a}') - Q(s_k, \mathbf{a}_k)]$$

Valami más, ami a cselekvések baráti,  
vagy adverz jellegére utal és eszerint  
számítja ki a várható jövőt.



$$Q^{új}(s_k, a_k) = Q(s_k, \mathbf{a}_k) + \alpha [r + \gamma [XYZ] - Q(s_k, \mathbf{a}_k)]$$

1 ágens  $\implies$  2 ágens  $\implies$  N ágens

# Teljes együttműködés

$$\rho_1 = \dots = \rho_n$$

Optimális együttes Q értékek parallel tanulása

$$\overset{új}{Q}(s_k, \mathbf{a}_k) = Q(s_k, \mathbf{a}_k) + \alpha [r + \gamma \max_{\mathbf{a}'} Q(s_{k+1}, \mathbf{a}') - Q(s_k, \mathbf{a}_k)]$$

és belőle egyenkénti optimális eljárasmód származtatása

$$\bar{\pi}_i^*(s) = \mathbf{a} = \arg \max_{a_i} \max_{a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_n} Q^*(s, \mathbf{a})$$

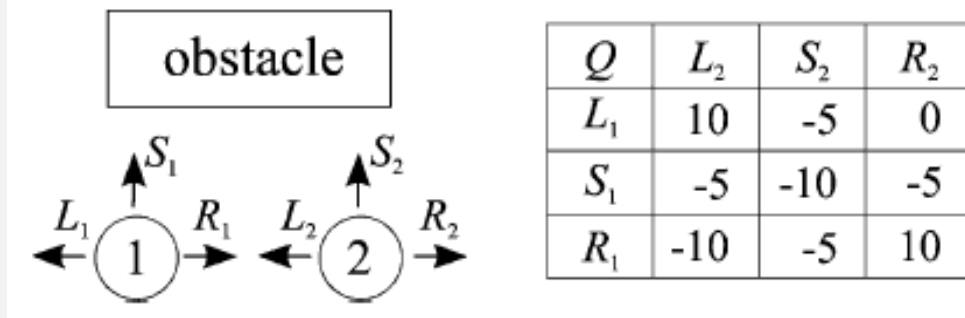
Együttműködés ellenére probléma a koordinálás.

Példa: formáció-mozgás

A két optimális helyzet ellenére, koordinálás hiányában ágensek szuboptimális helyzetben

$$Q(L_1, R_2)$$

végezhetnek. (ha a Q érték közös, mindkét optimális eset egy Nash egyensúly)



Q	L <sub>2</sub>	S <sub>2</sub>	R <sub>2</sub>
L <sub>1</sub>	10	-5	0
S <sub>1</sub>	-5	-10	-5
R <sub>1</sub>	-10	-5	10

$$Q(L_1, L_2) = Q(R_1, R_2) = 10$$

# Koordinálás kérdése

## Koordinálás-mentes

pl. **Team-Q**: optimális együttes cselekvés egyedi jellegét tételezi fel

**Distributed-Q-Learning**: lokális  $Q$  és  $\pi$  tanulás, de az egyedi  $Q$  frissítése csak akkor, ha növekszik (a közös optimális  $Q$  értéket is el fogja kapni).

A stratégia frissítése csak akkor, ha a  $Q$  érték növekszik.

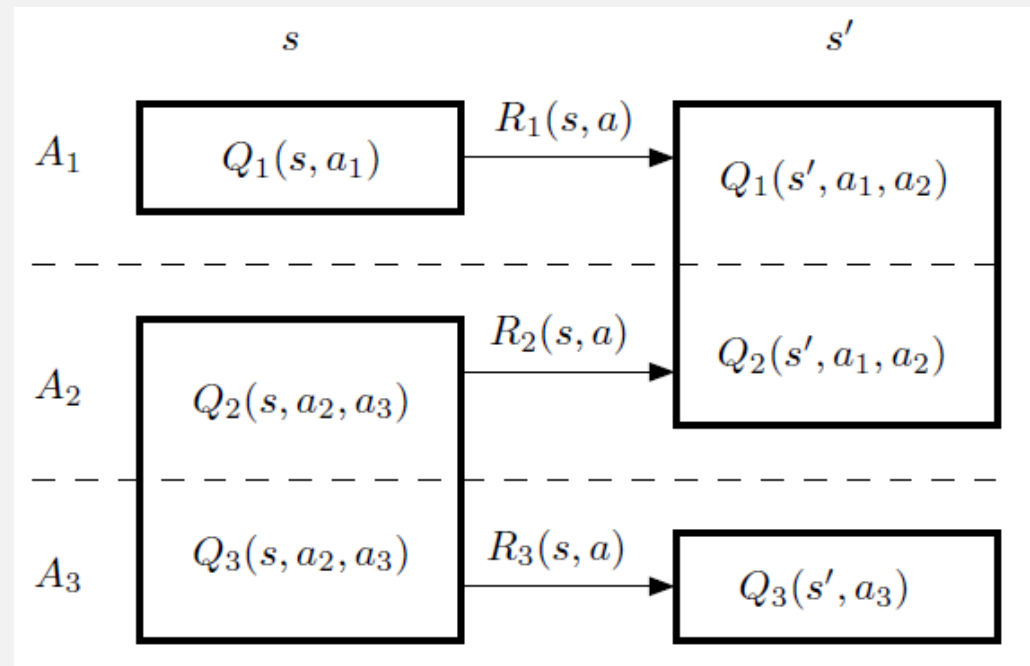
$$Q_i^{új} (s_k, a_{i,k}) = \max \left\{ Q_i (s_k, a_{i,k}), r + \gamma \max_{a'_i} Q_i (s_{k+1}, a'_i) \right\}$$

# Koordinálás kérdése

**Koordinálás-alapú** pl. együttes Q dekomponálása kisebb csoportosulások szerint.  
(koordinációs gráfok)

$$Q(s, \mathbf{a}) = Q_1(s, a_1) + Q_2(s, a_2, a_3) + Q_3(s, a_2, a_3)$$

$$Q(s', \mathbf{a}) = Q_1(s', a_1, a_2) + Q_2(s', a_1, a_2) + Q_3(s', a_3)$$



# Koordinálás kérdése

## Indirekt koordinálás

pl. tanulva mások (stacionárius probabilitisztikus) empirikus modelljét, hogy egyes cselekvéseit milyen gyakorisággal használják:

## JAL – Joint Action Learner

$C_j^i(a_j)$  – i-edik ágens hányszor látta, hogy az j-edik ágens egy  $a_j$  cselekvéséhez folyamodik, erre a legjobb válasz számítása

$$\hat{\sigma}_j^i(a_j) = \frac{C_j^i(a_j)}{\sum_{\tilde{a}_j \in A_j} C_j^i(\tilde{a}_j)}$$

**Frequency Maximum Q-value heurisztika:** mely cselekvések jók voltak a múltban?

$r_{\max}(a_i)$  - az  $a_i$ -re eddig kapott max megerősítés,  $C_{\max}^i(a_i)$  – ennek a gyakorisága és  $C^i(a_i)$  – az  $a_i$  cselekvés (önmagában) gyakorisága (az  $a_i$  -re több megerősítés is jöhetett, szórás csakis a mások cselekvései miatt, determinisztikus probléma).  
A számított Q értéket a Boltzmann-felfedezés képletében használja ( $\nu$  egy súly).

$$\tilde{Q}_i(a_i) = Q_i(a_i) + \nu \frac{C_{\max}^i(a_i)}{C^i(a_i)} r_{\max}(a_i)$$

## Explicit koordinálás

pl. társadalmi szabályok  
normatívák, törvények  
kommunikáció  
szerepek

pl. ágens 1 < ágens 2  
 $L < R < S$   
döntés ( $L_1, L_2$ )

# Teljes versengés

**Minimax Q-tanulás** (1. ágens esete,  $Q = Q_1 = -Q_2$ )

$$\begin{aligned} Q^{\text{új}}(s_k, a_{1,k}, a_{2,k}) = & Q(s_k, a_{1,k}, a_{2,k}) + \\ & \alpha \left[ r + \gamma \mathbf{m}_1(Q, s_{k+1}) - Q(s_k, a_{1,k}, a_{2,k}) \right] \end{aligned}$$

$$\mathbf{m}_1(Q, s) = \max_{\pi_1(s, \cdot)} \min_{a_2} \sum_{a_1} \pi_1(s, a_1) Q(s, a_1, a_2)$$

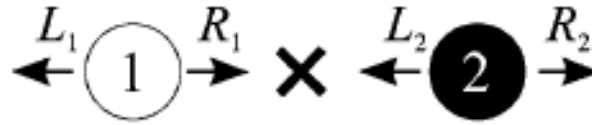
Az 1. ágens minimax haszna (Lineáris Programozással)

$$\pi_1(s_k, \cdot) = \arg \mathbf{m}_1(Q, s_k)$$

Garantált konvergencia (NE-hez) (self-play), de nem racionális.

# Teljes versengés

Minimax Q-tanulás, példa



$Q$	$L_2$	$R_2$
$L_1$	0	1
$R_1$	-10	10

Az 1. ágens szeretne elfoglalni a keresztet és elmenekülni.

A 2. ágens szeretne elkapni az 1. ágenst.

A Q táblázat az 1. ágens perspektíváját mutatja, a 2. ágens Q függvénye ennek (-1)-szerese.

A minimax megoldás az 1. ágensre:

Ha  $L_1$ -et lép, akkor a 2. minimalizálva  $L_2$ -et lép, eredményben 0.

Ha  $R_1$ -et lép, akkor a 2. minimalizálva szintén  $L_2$ -et lép, eredményben -10.

Az 1. ágensnek tehát  $L_1$ -et kell lépnie (minimumok maximuma), mert így legfeljebb 0-val megússza.

$(L_1, L_2)$  egyben egy NE is.



# Vegyes feladatok

Nincsenek feltételek megerősítésekre. Valamilyen egyensúly felé mégis húzni kell. Lehet pl. Nash-egyensúly, de mi van, ha több van ilyen?

**Egyedi** ágens Q-tanulása (a többi implicite a környezeti információban rejlik)

**Ágens-független** módszerek (egymástól függetlenül, de feltehetően egy közös egyensúly fel), pl. **Nash-Q-tanulás**.

$$Q_i^{új}(s_k, \mathbf{a}_k) = Q_i(s_k, \mathbf{a}_k) + \alpha \left[ r_i + \gamma \text{eval}_i \{Q_-(s_{k+1}, \cdot)\} - Q_i(s_k, \mathbf{a}_k) \right]$$

Q értékek által definiált kifizetési mátrix, m-edik pillanatban  $Q_{,m}(s_k, \cdot)$

$$\text{eval}_i \{Q_-(s, \cdot)\} = V_i(s, NE \{Q_-(s, \cdot)\})$$

$$\text{solve}_i \{Q_-(s, \cdot)\} = NE_i \{Q_-(s, \cdot)\}$$

$$\pi_i(s, \cdot) = \text{solve}_i \{Q_-(s_k, \cdot)\}$$

$NE$  az egyensúly kiszámítása,  $NE_i$  az  $i$ -ik ágens stratégiája az egyensúlyban és  $V_i$  az ágens várható haszna  $s$ -ben az egyensúlyban. Bizonyos feltételekkel  $NE$ -hoz konvergál. Mindegyik ágens számon tartja mások  $Q$  értékeit is!

# Vegyes feladatok

Nash-Q egyszerűbben: **FoF Friend-or-Foe tanulás**

Fel kell ismerni, hogy az ellenfél kooperál, vagy ellenünk dolgozik.

Kooperál esetében: JAL (Joint Action Learner) a megoldás, a Q közös maximumára való törekvés. Ha ellenünk van legyen a minimax-Q játék.

$$Nash_1(s, Q_1, Q_2) = \max_{a_1, a_2} Q_1(s, a_1, a_2)$$

$$Nash_1(s, Q_1, Q_2) = \max_{\pi_1(s, \cdot)} \min_{a_2} \sum_{a_1} \pi_1(s, a_1) Q_1(s, a_1, a_2)$$

**Ellenséges egyensúly (EE):** ha mások tőle eltérnek, az ágensünk helyzete változatlan lesz, vagy javul.

**Kooperatív egyensúly (KE):** ahol az ágensek maximális haszonhoz jutnak (a jóléti megoldás).

Ha a játéknak van EE egyensúlya, FOE-Q tanulás ezt megtanulja.

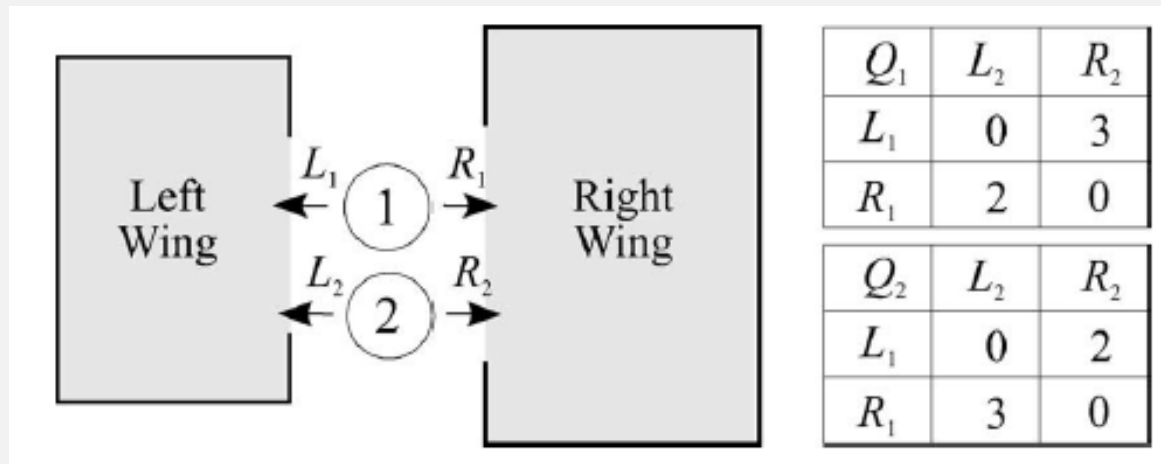
Ha a játéknak van KE egyensúlya, FRIEND-Q tanulás ezt megtanulja.

# Vegyes feladatok

Problémák egyensúlyi helyzetek koordinálásával itt is, pl.:

Két porszívóágens feladata a két szobából álló lakás kitakarítása.

Mindegyik jobban szeretne a bal szobát megkapni, mert ez kisebb.



Két Nash-egyensúly van:  $(L_1, R_2), (R_1, L_2)$

de ha a két ágens között nincs koordináció (nincs megegyezés melyik NE-re játszanak, akkor mindketten ugyanabban a szobában végeznek, kisebb hasznossággal.

$$Q_1(L_1, L_2) = Q_1(R_1, R_2) = Q_2(L_1, L_2) = Q_2(R_1, R_2) = 0$$

# Vegyes feladatok

Ágens-követő, ágens-tudatos módszerek (más ágensek modellezése, a modell felhasználása tanulásban: - érzékelés + stratégia-váltás)

## **AWESOME**

**(Adapt When Everyone is Stationary, Otherwise Move to Equilibrium)**

Induláskor  $\pi_i$  egyensúlyi stratégiát játszik (az  $i$ -ik ágens) és mások cselekvéseit követi. Minden  $N$ -ik kör végére a megfigyelt gyakoriságokból kiszámítja az  $s_j$ -t, az  $j$ -ik ágens (feltehetően kevert) stratégiabecslését.

Ha  $s_j$  minden  $j$ -ik játékos stratégiája egyensúlyi, akkor  $i$  folytatja az egyensúlyi stratégiájának bevetését.

Különben az  $i$ -ik ágens megjátssza az  $s_j$  stratégiára kiszámított legjobb válasz stratégiáját (best response).

Ha minden játékos AWESOME játékos, akkor a közös tanulás az egyensúlyi helyzethez konvergál és nem fog tőle eltérni.

# Vegyes feladatok

## IGA (Infinitesimal Gradient Ascent) –

ágensek cselekvéseinek a valószínűsége az,  
amit az ágensek tanulnak (2 ágens, 2-2 cselekvés),

$\alpha$  az 1. ágens 1. cselekvésének a valószínűsége és  
 $\beta$  a 2. ágens 1. cselekvésének a valószínűsége:

$$\alpha_{k+1} = \alpha_k + \delta_{1,k} \frac{\partial E\{r_1 | \alpha, \beta\}}{\partial \alpha}$$
$$\beta_{k+1} = \beta_k + \delta_{2,k} \frac{\partial E\{r_2 | \alpha, \beta\}}{\partial \beta}$$

$$\delta_{1,k} = \delta_{2,k} = \delta$$

# Vegyes feladatok

## WoLF-IGA (Win-or-Learn-Fast)

- győztes helyzetben ágens óvatos, kis  $\delta$ -val lassan tanul, nehogy az előnyös pozícióját elveszítse,
- vesztes esetben nagyobb  $\delta$ -val gyorsan igyekszik kikerülni a jelen helyzetből.

$$\alpha_{k+1} = \alpha_k + \delta_{1,k} \frac{\partial E\{r_1 | \alpha, \beta\}}{\partial \alpha}$$
$$\beta_{k+1} = \beta_k + \delta_{2,k} \frac{\partial E\{r_2 | \alpha, \beta\}}{\partial \beta}$$

$$\delta_{1,k}, \delta_{2,k} \in [\delta_{\min}, \delta_{\max}] > 0$$

Hogyan látszik meg, melyik a győzelmi, melyik a vesztes állapot?

### WoLF-elv:

Ha a várható hasznom, ahogy játszom és ahogy az ellenfél játszik, jobb, mintha a jelen játéka mellett én az egyensúlyi stratégiát játszanám, akkor győzelemre állok.

Ha rosszabbul állok az egyensúlyi stratégiámhoz képest (az ellenfél adott játéka mellett), akkor vesztesre állok.

Lucian Busoniu, Robert Babuska, and Bart De Schutter, A Comprehensive Survey of Multiagent Reinforcement Learning, IEEE Trans. on Systems, Man, and Cybernetics—Part C: Applications and Reviews, Vol. 38, No. 2, March 2008

# Többágenses mély megerősítéses tanulás

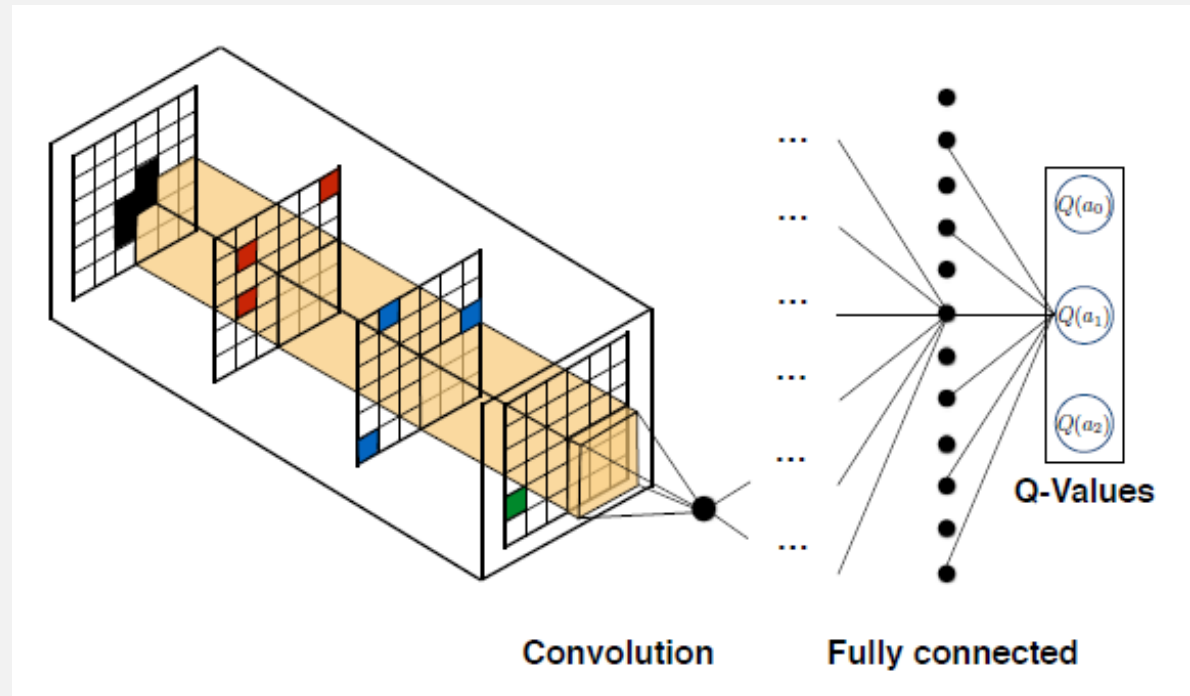
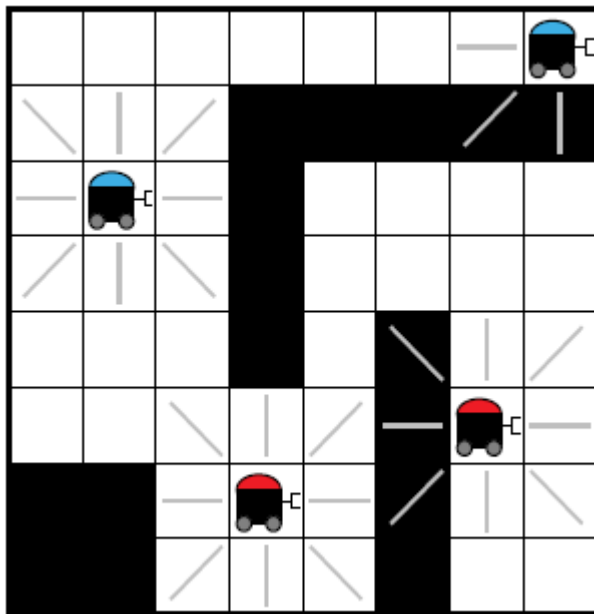
$$Q(s, \mathbf{a}) = Q(s, \mathbf{a}) + \alpha [r + \gamma \max_{\mathbf{a}'} Q(s', \mathbf{a}') - Q(s, \mathbf{a})], \quad Q(s, \mathbf{a} | \boldsymbol{\theta})$$

$$\min L(s, \mathbf{a} | \boldsymbol{\theta}_i) = (r + \gamma \max_{\mathbf{a}'} Q(s', \mathbf{a}' | \boldsymbol{\theta}_i) - Q(s, \mathbf{a} | \boldsymbol{\theta}_i))^2$$

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i + \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_i)$$

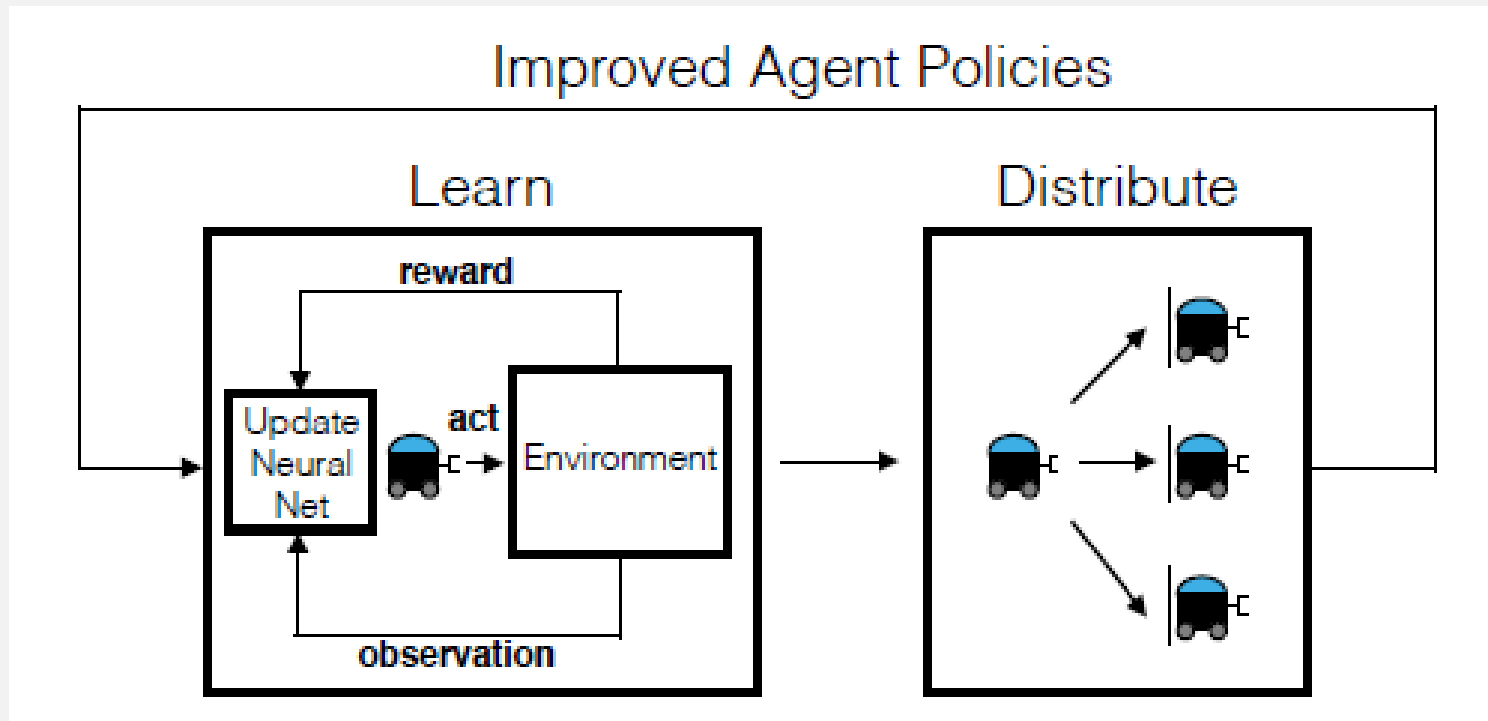
Környezet 2 dim, 2 ágens típus, diszkrét tér/idő.

Konvolúciós NN csatornák: háttér cs. (akadályok), ellenfelek cs., szövetségeselek cs., saját magam cs. :  $4 \times H \times H$ .



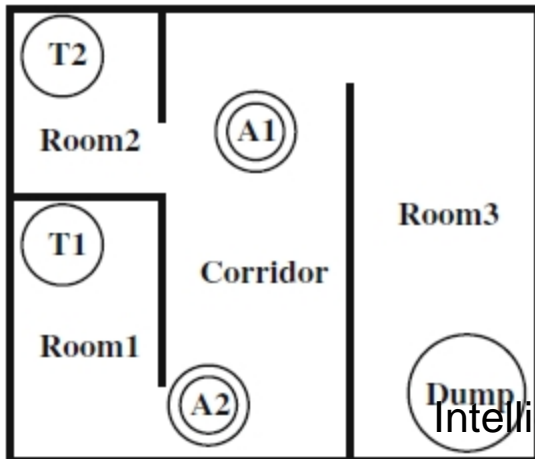
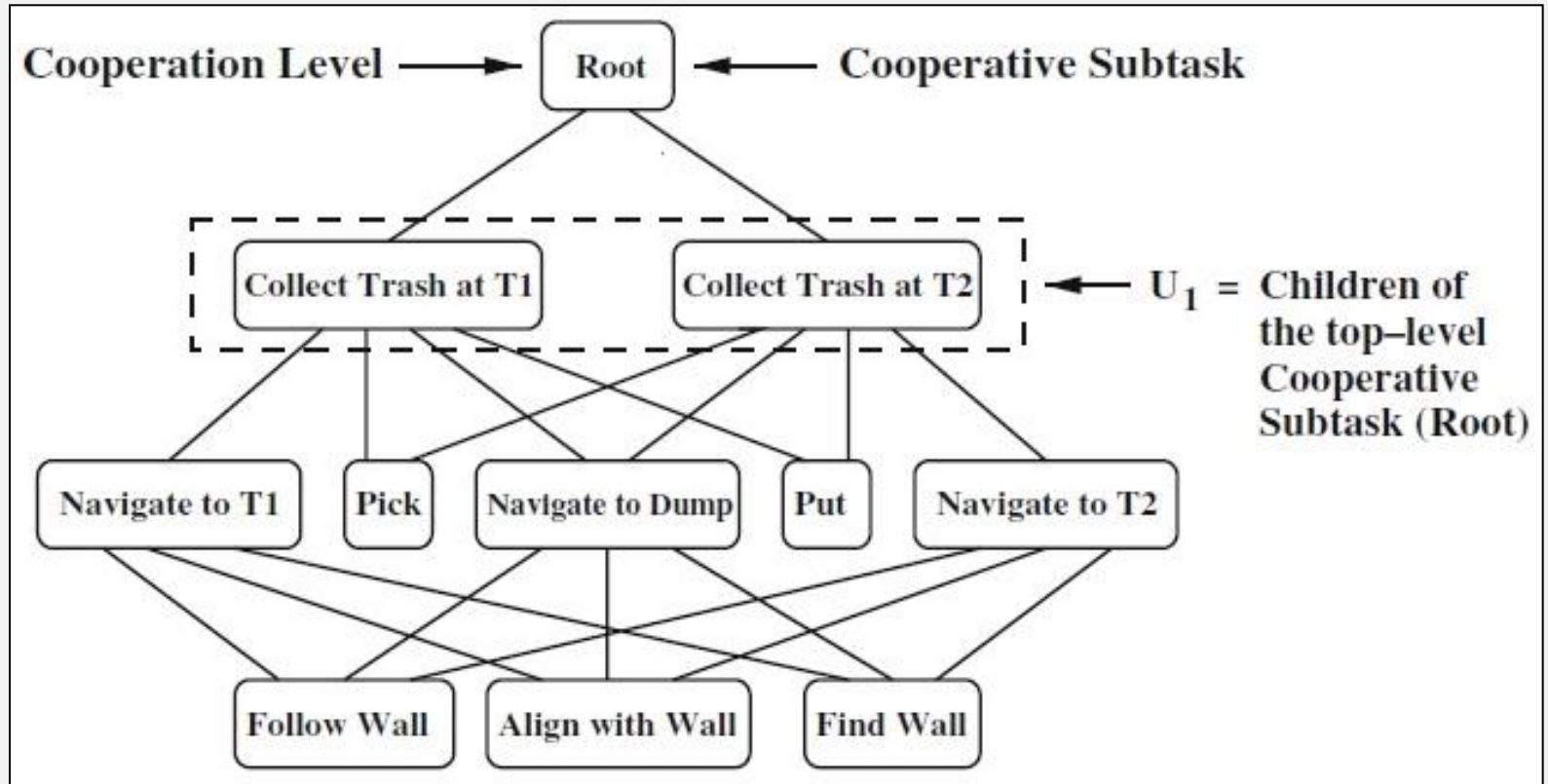
# Többágenses mély megerősítéses tanulás

Tanulás: egyszerre 1 ágens, a többinek stratégiája lefixált, a megtanult stratégia kiosztása a saját típusú ágensekre.

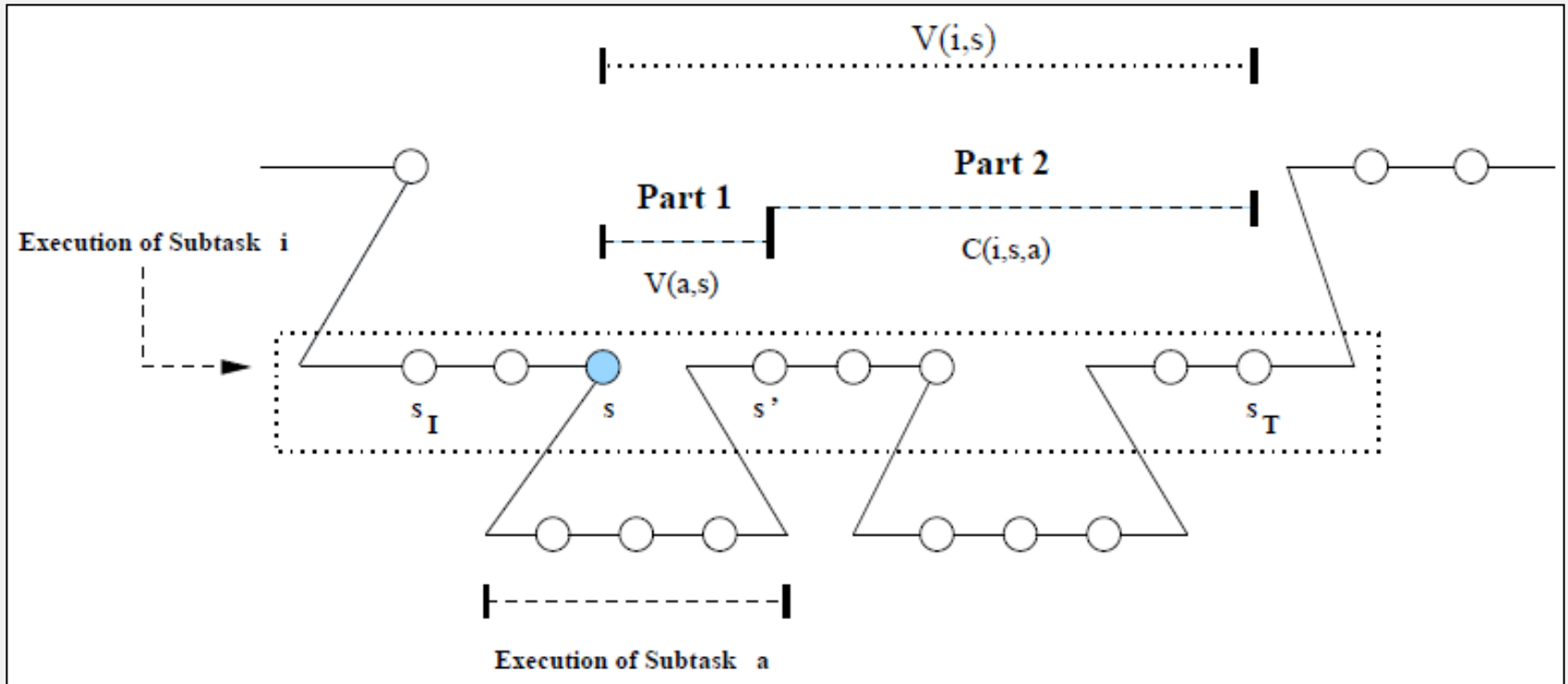




# Hierarchikus többágenses megerősítéses tanulás



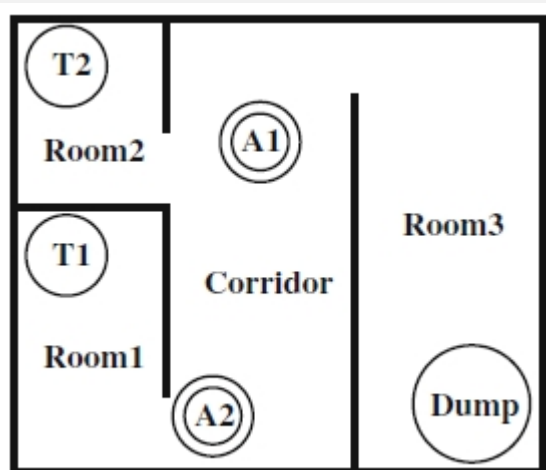
# Hierarchikus többágenses megerősítéses tanulás



$$Q^\pi(i, s, a) = V^\pi(a, s) + C^\pi(i, s, a)$$

$$V^\pi(i, s) = \begin{cases} Q^\pi(i, s, \pi_i(s)) \\ \sum_{s' \in S_i} P(s'|s, i) R(s'|s, i) \end{cases}$$

# Hierarchikus többágenses megerősítéses tanulás



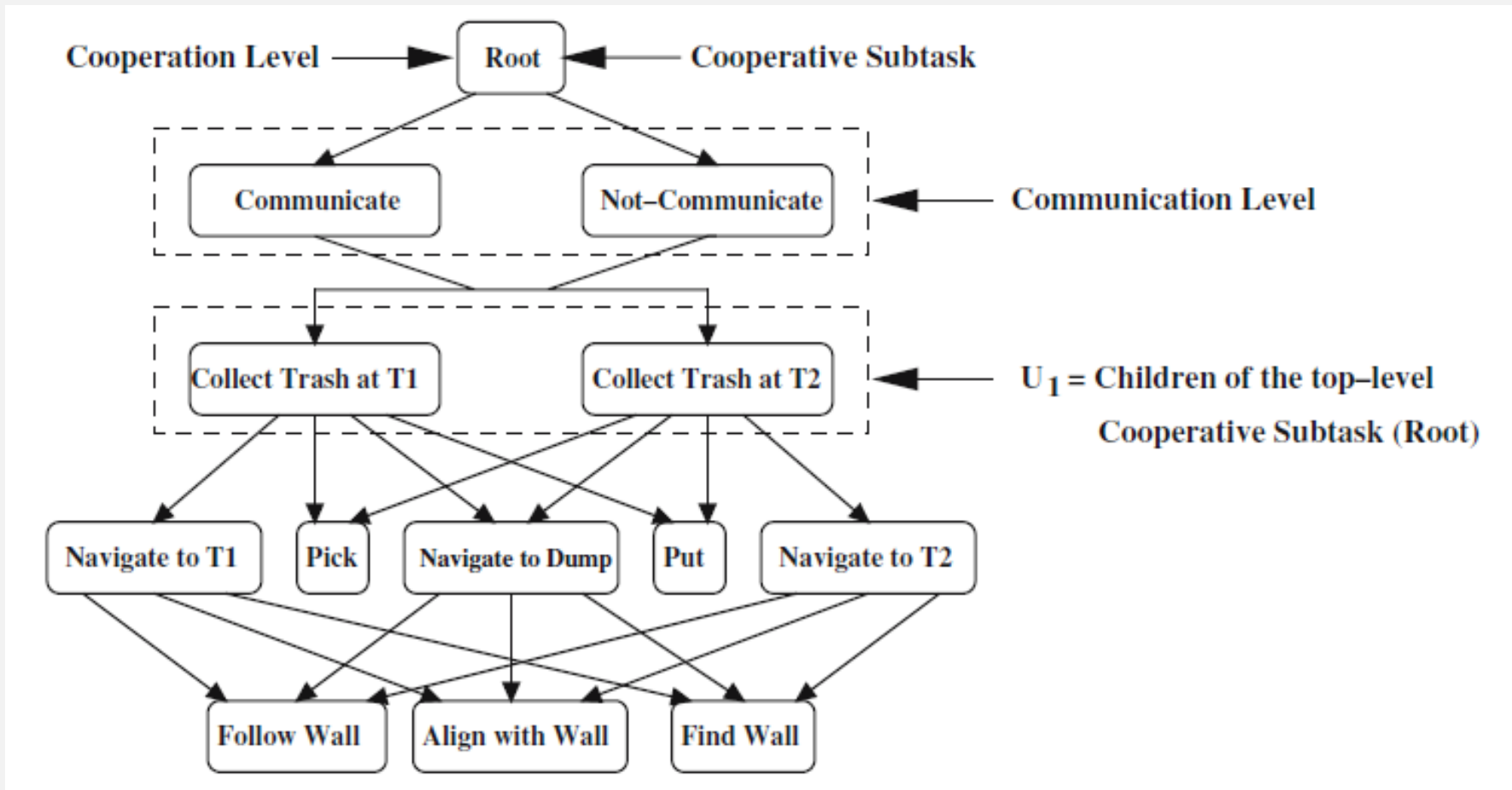
## Learned Policy for Agent 1

```
root
  navigate to T1
    go to location of T1 in room 1
  pick trash from T1
  navigate to Dump
    exit room 1
  enter room 3
    go to location of Dump in room 3
  put trash collected from T1 in Dump
end
```

## Learned Policy for Agent 2

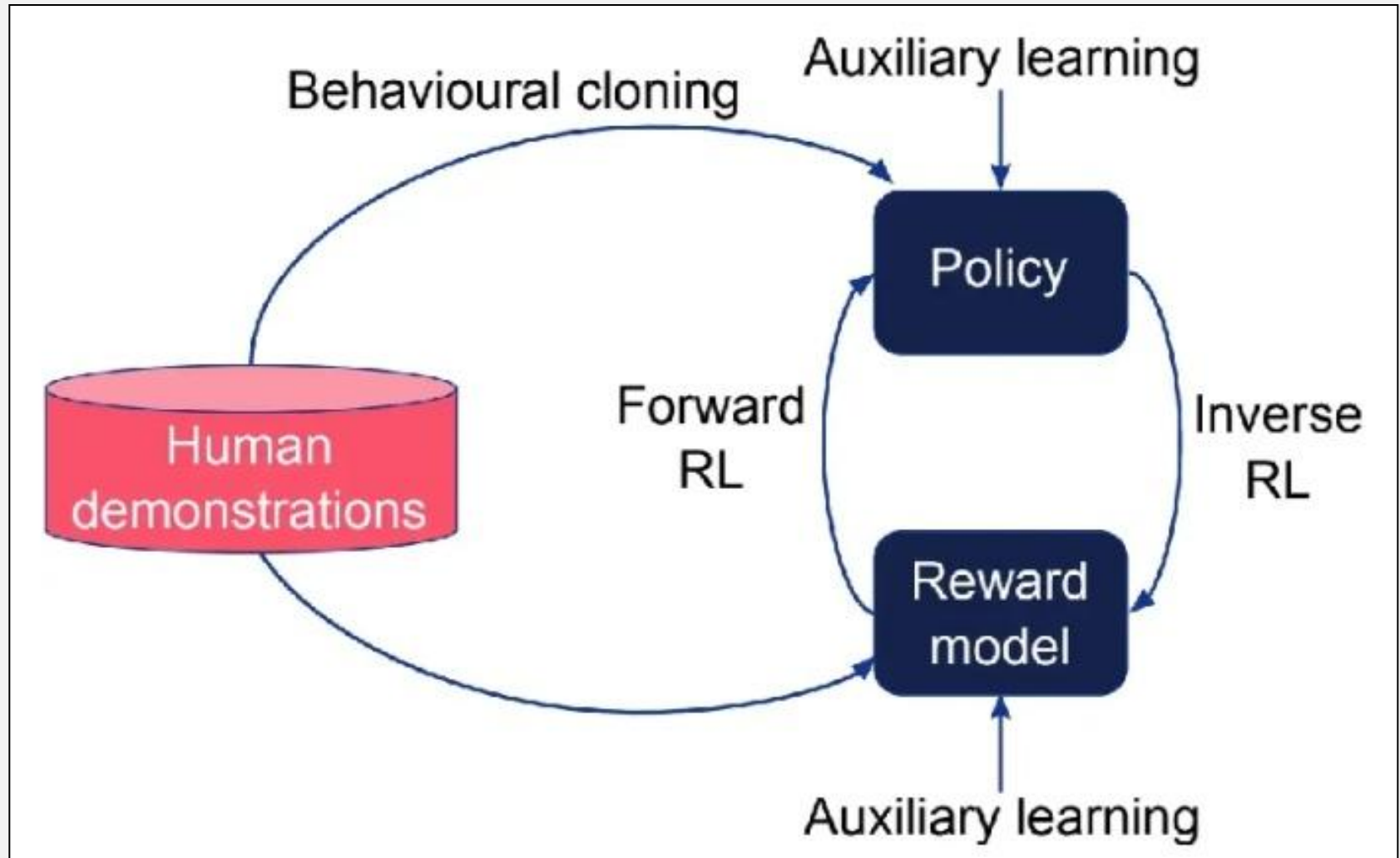
```
root
  navigate to T2
    go to location of T2 in room 2
  pick trash from T2
  navigate to Dump
    exit room 2
  enter room 3
    go to location of Dump in room 3
  put trash collected from T2 in Dump
end
```

# Hierarchikus többágenses megerősítéses tanulás



Taszkgráf kommunikációs cselekvésekkel

# Inverz megerősítéses tanulás



# Inverz megerősítéses tanulás

Ha adott a  $\pi^*$ , visszaállítható-e az R?

Ha adottak a végrehajtási pályák, visszaállítható-e az R?

## Bemenet:

Állapottér

Cselekvéshalmaz

Állapotátmeneti modell

**R megerősítés nincs!**

Tanári bemutató van!

$s_0, a_0, s_1, a_1, s_2, a_2, \dots$

(a tanári  $\pi^*$  eljárás mód nyoma)

**Kimenet: R?**

Tudományosan érdekes

**Inas tanulás/ Imitációs tanulás**

Kooperatív és versengő ágensek modellezése

Járművezetés, városi navigálás

Légi navigálás

Parkolás

Emberi útvonaltervezés

**Emberi célok azonosítása**

Négy lábú robotok mozgása

...

Azt az  $R^*$  megerősítésfüggvényt találjuk meg, ami a **szakértő viselkedését megmagyarázza!**

# Inverz megerősítéses tanulás

$$\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi^*] \geq \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi] \quad \forall \pi$$

(Számítási problémák)

Jellegvektor alapú megerősítés (regressziós számítás)

$$R(s) = w^\top \phi(s)$$

$$\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(s_t) | \pi] = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t w^\top \phi(s_t) | \pi] = w^\top \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) | \pi]$$

$$= w^\top \mu(\pi)$$

Várható kumulált diszkont jellegértékek

$$\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi^*] \geq \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi] \quad \forall \pi$$

$$w! \longrightarrow w^{*\top} \mu(\pi^*) \geq w^{*\top} \mu(\pi) \quad \forall \pi$$

# Inverz megerősítéses tanulás

Jellegektor alapú megerősítés (regressziós számítás)

Egyértelmű megoldás érdekében

$$w^{*\top} \mu(\pi^*) \geq w^{*\top} \mu(\pi) \quad \forall \pi$$

Standard max tartalék

$$\begin{aligned} \min_w \quad & \|w\|_2^2 \\ \text{s.t.} \quad & w^\top \mu(\pi^*) \geq w^\top \mu(\pi) + 1 \quad \forall \pi \end{aligned}$$

„Struktúrált predikciós” max tartalék

$$\begin{aligned} \min_w \quad & \|w\|_2^2 \\ \text{s.t.} \quad & w^\top \mu(\pi^*) \geq w^\top \mu(\pi) + m(\pi^*, \pi) \quad \forall \pi \end{aligned}$$

„Struktúrált predikciós” max tartalék gyengítő változókkal, stb.

$$\begin{aligned} \min_{w, \xi} \quad & \|w\|_2^2 + C\xi \\ \text{s.t.} \quad & w^\top \mu(\pi^*) \geq w^\top \mu(\pi) + m(\pi^*, \pi) - \xi \quad \forall \pi \end{aligned}$$



# Többágenses megerősítéses tanulás - Szimulációk

[http://busoniu.net/files/repository/2010-08-04\\_marl-1.3.zip](http://busoniu.net/files/repository/2010-08-04_marl-1.3.zip)

