



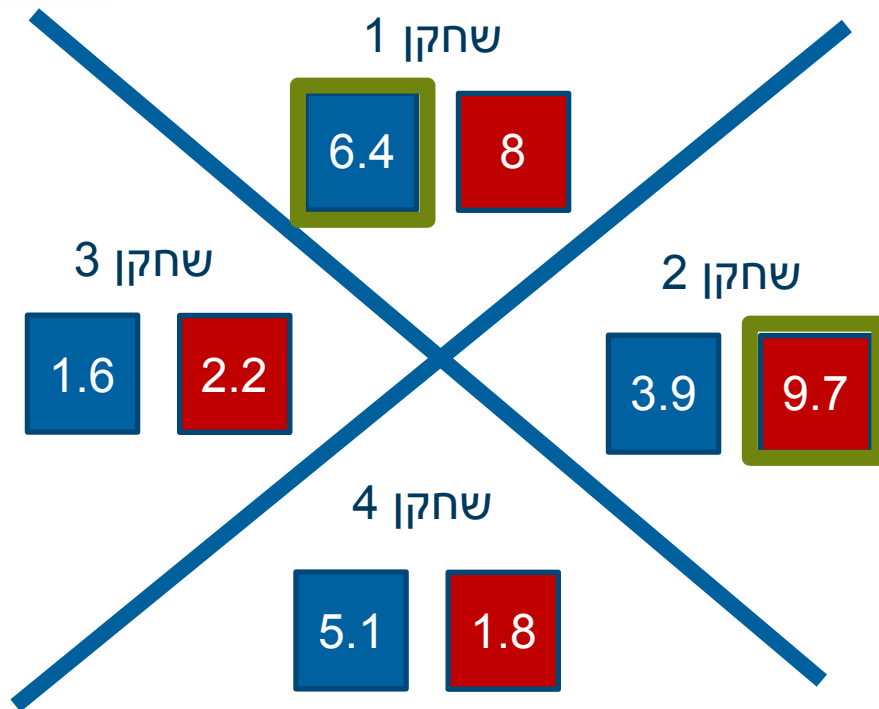
## למידת חיזוקים עמוקה מרובת משתמשים ברשתות תקשורת

ד"ר אושרי נפרסטק



## בעיית צמצום

# אתגר לימוד אסטרטגיה במשחק



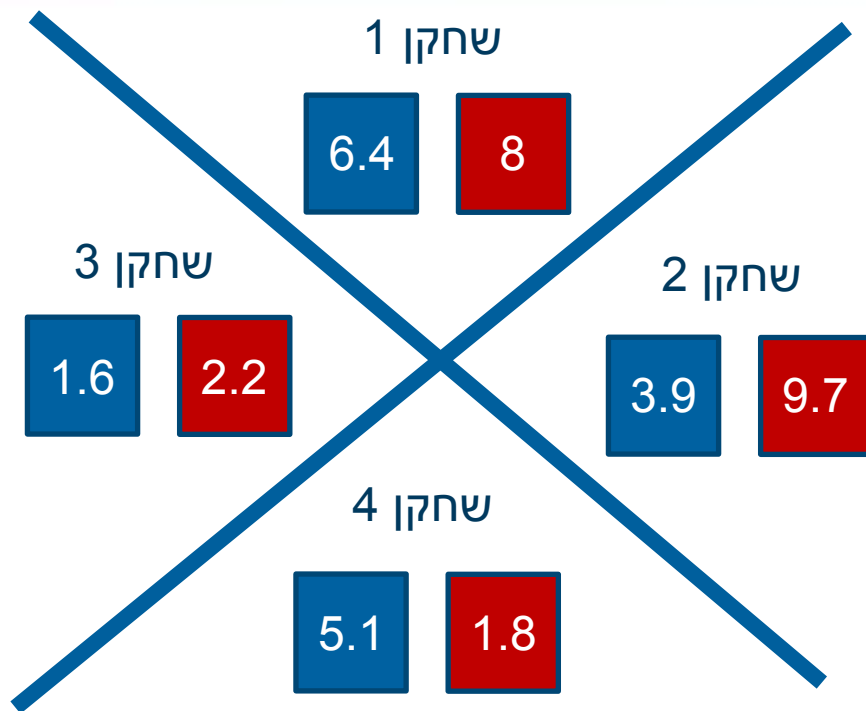
- ארבעה שחקנים ללא יכולת לתקשר ביניהם.
- לכל שחקן שני כפתורים, כחול ואדום.
- משחקים 20 פעמים את המשחק.
- לכל שחקן יש תגמול על כל כפתור שמפולג אחיד בין 0 ל 10.
- בכל שלב כל שחקן צריך להחליט האם ללחוץ או לא ללחוץ ועל איזה כפתור.
- אם רק שחקן אחד לוחץ על כפתור בצבע מסוים, הוא מקבל תגמול.
- אם יותר משחקן אחד לוחץ על צבע מסוים, אף אחד לא מקבל תגמול.
- בסוף המשחק התגמולים שהתקבלו מחולקים בצורה שווה.

# אתגר לימוד פרוטוקול תקשורת פשוט

המטרה:

ללמוד אסטרטגיה שממקסמת את הרווח.

רעיונות?



# איך נגשים לפתור את הבעיה?

אפשרות 1: לחשוב על אלגוריתם ולבדוק אותו או להוכיח שהוא טוב.

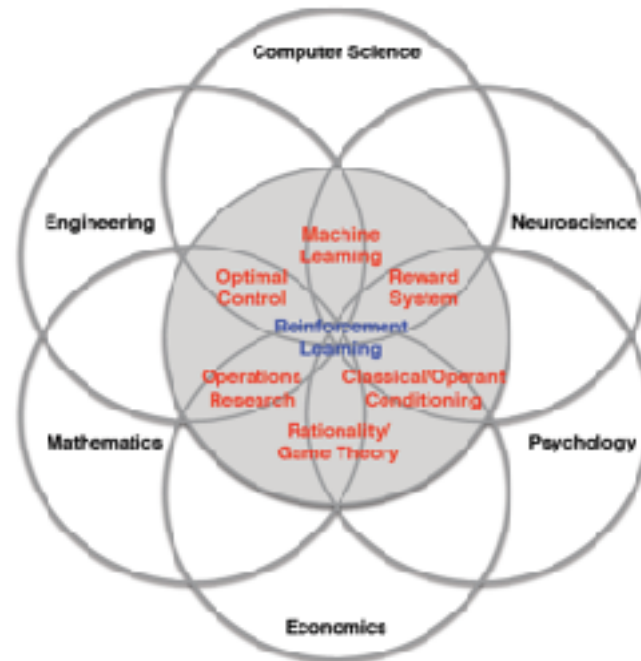
הבעיה: זה קשה ולא בטוח שנצליח.

אפשרות 2: לבנות מערכת לומדת למציאת אסטרטגיות טובות

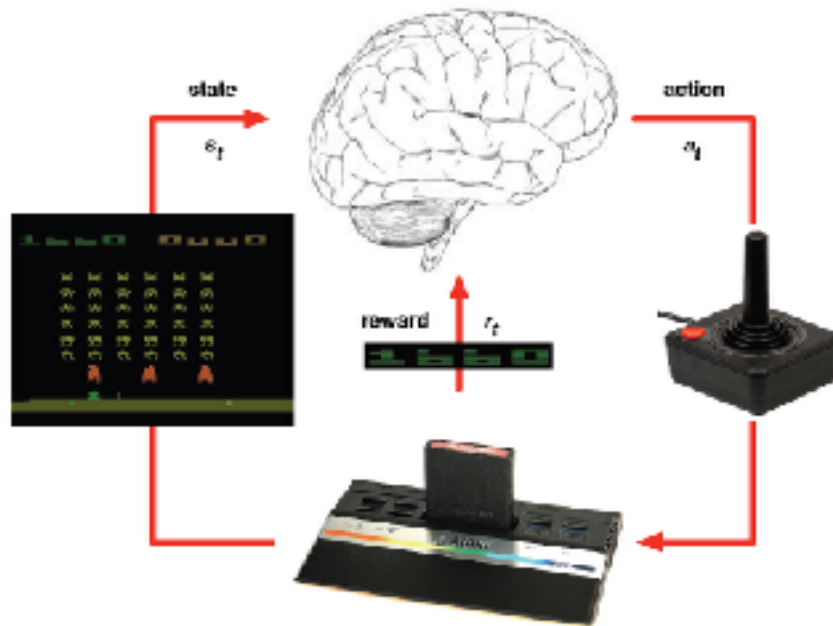
הבעיה: איך עושים את זה?

הפתרון: למידת חיזוקים עמוקה (Deep reinforcement learning)

## למידת חיזוקים



# למידת חיזוקים (David Silver)



- ▶ At each step  $t$  the agent:
  - ▶ Executes action  $a_t$
  - ▶ Receives observation  $o_t$
  - ▶ Receives scalar reward  $r_t$
- ▶ The environment:
  - ▶ Receives action  $a_t$
  - ▶ Emits observation  $o_{t+1}$
  - ▶ Emits scalar reward  $r_{t+1}$



## למידת חיזוקים (David Silver)

- ▶ Experience is a sequence of observations, actions, rewards

$$o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t$$

- ▶ The **state** is a summary of experience

$$s_t = f(o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t)$$

- ▶ In a fully observed environment

$$s_t = f(o_t)$$

# Major Components of an RL Agent

- ▶ An RL agent may include one or more of these components:
  - ▶ **Policy**: agent's behaviour function
  - ▶ **Value function**: how good is each state and/or action
  - ▶ **Model**: agent's representation of the environment

# Policy

- ▶ A **policy** is the agent's behaviour
- ▶ It is a map from state to action:
  - ▶ Deterministic policy:  $a = \pi(s)$
  - ▶ Stochastic policy:  $\pi(a|s) = \mathbb{P}[a|s]$

## Discounted Future Reward

Total Reward:

$$R = r_1 + r_2 + r_3 + \dots + r_N$$

Total Future Reward from point t:

$$R_t = r_t + r_{t+1} + r_{t+2} + \dots + r_N$$

ע"מ לתת עדיפות לטווח הקרוב והודאי יותר, נבטא את  $R_t$  בצורה הבאה:

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots + \gamma^N r_N = r_t + \gamma R_{t+1}$$

$Q_{\pi}(S_t, A_t) = E(R_{t+1})$  (or: How good is action A in state S?)

$$\pi(S, A) = \operatorname{argmax}_A Q(S, A)$$



# Discounted Future Reward

**Bellman's equation:**  $Q(S,A) = r + \gamma Q(S',A')$

**Implementation:**  $Q(S,A) = Q(S,A) + \underbrace{\underbrace{A \cdot A}_{\text{קצב למידה}} \underbrace{\{r + \gamma \max_{A'} Q(S',A') - Q(S,A)\}}_{\text{Prediction}}}_{\text{Target}}$

קצב למידה

Target

Prediction

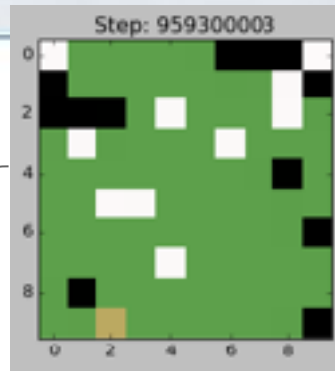
CONVERGENCE OF Q-LEARNING: A SIMPLE PROOF

Francisco S. Melo  
 Institute for Systems and Robotics  
 Instituto Superior Técnico,  
 Lisboa, PORTUGAL  
 fsmelo@isr.isr.tecnico.pt



$Q(S,A)$

$$Q(S,A) = Q(S,A) + \frac{1}{N} \{r + \gamma \max_{A'} Q(S',A') - Q(S,A)\}$$



	←	0	→
State #1	0	0	0
→ State #2	0	0	0
→ State #3	0	0	0
→ State #4	0	0	0
→ ⋮	0	0	0
→ ⋮	0	0	0
State #N		0	0

$Q(S,A)$

$\leftarrow$	0	$\rightarrow$
--------------	---	---------------

State #1

0 0 0

$\rightarrow$  State #2

0 0 0

$\rightarrow$  State #3

0 0 0

$\rightarrow$  State #4

0 0 0

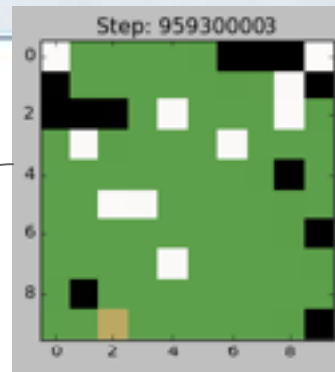
$\rightarrow$

0 0 0

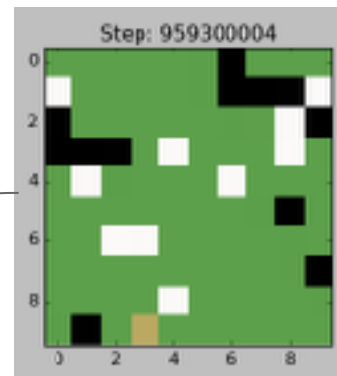
0 0 0

State #N

0 0



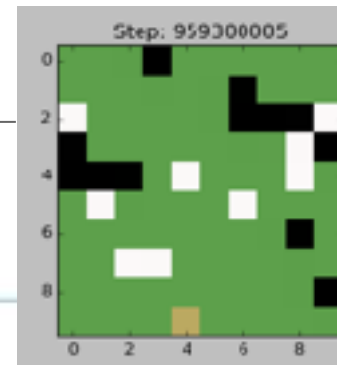
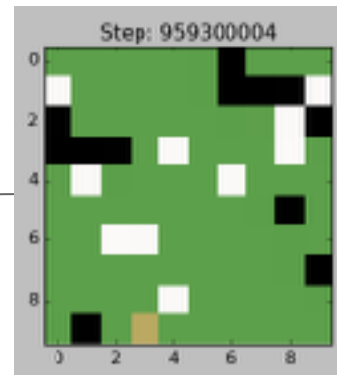
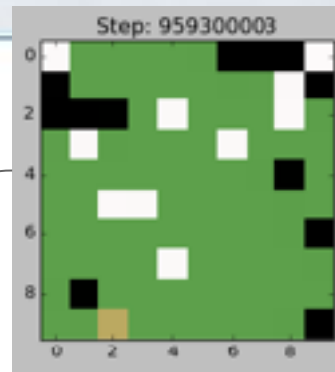
$$Q(S,A) = Q(S,A) + \frac{1}{A} \{r + \gamma \max_{A'} Q(S',A') - Q(S,A)\}$$



$Q(S,A)$

$$Q(S,A) = Q(S,A) + \alpha \{r + \gamma \max_{A'} Q(S',A') - Q(S,A)\}$$

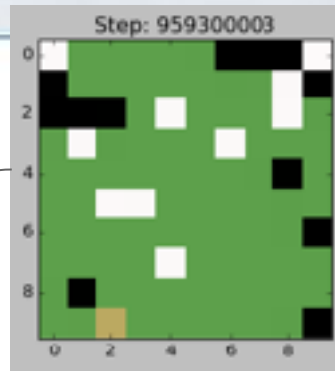
	←	0	→
State #1	0	0	0
→ State #2	0	0	0
→ State #3	0	0	0
→ State #4	0	0	0
→	0	0	$\alpha r$
	0	0	0
<del>State #N</del>	0	0	0



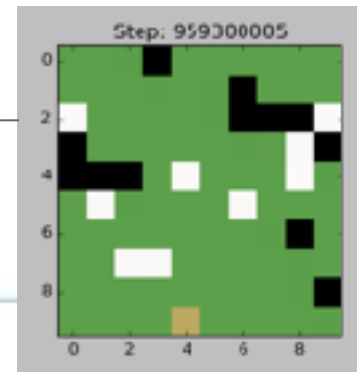
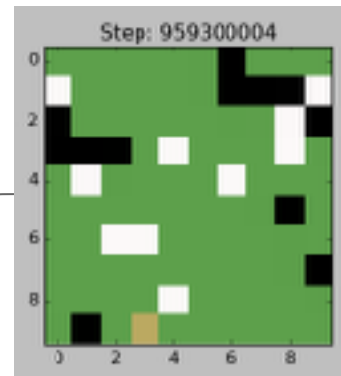


$Q(S,A)$

	←	0	⇒
State #1	0	0	0
→ State #2	0	0	0
→ State #3	0	0	0
→ State #4	0	0	0
→ ⋮	0	0	αr
→ ⋮	0	0	0
State #N	0	0	0



$$Q(S,A) = Q(S,A) + \frac{1}{N} \{r + \gamma \max_{A'} Q(S',A') - Q(S,A)\}$$



$\leftarrow$	0	$\Rightarrow$
--------------	---	---------------

State #1

0

0

$\alpha^2 \gamma r$

$\rightarrow$

State #2

0

0

0

$\rightarrow$

State #3

0

0

0

$\rightarrow$

State #4

0

0

0

$\rightarrow$

0

0

$\alpha r$

0

0

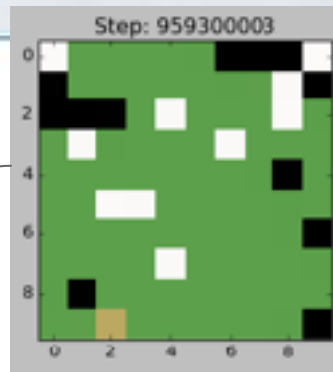
0

State #N

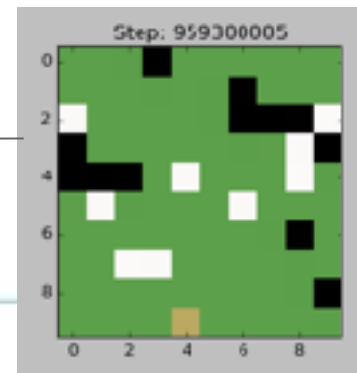
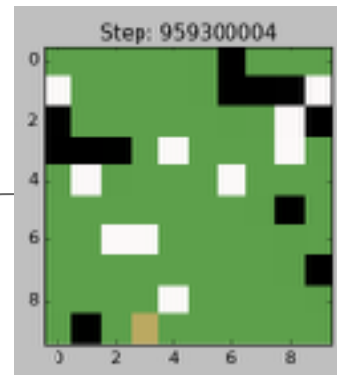
0

0

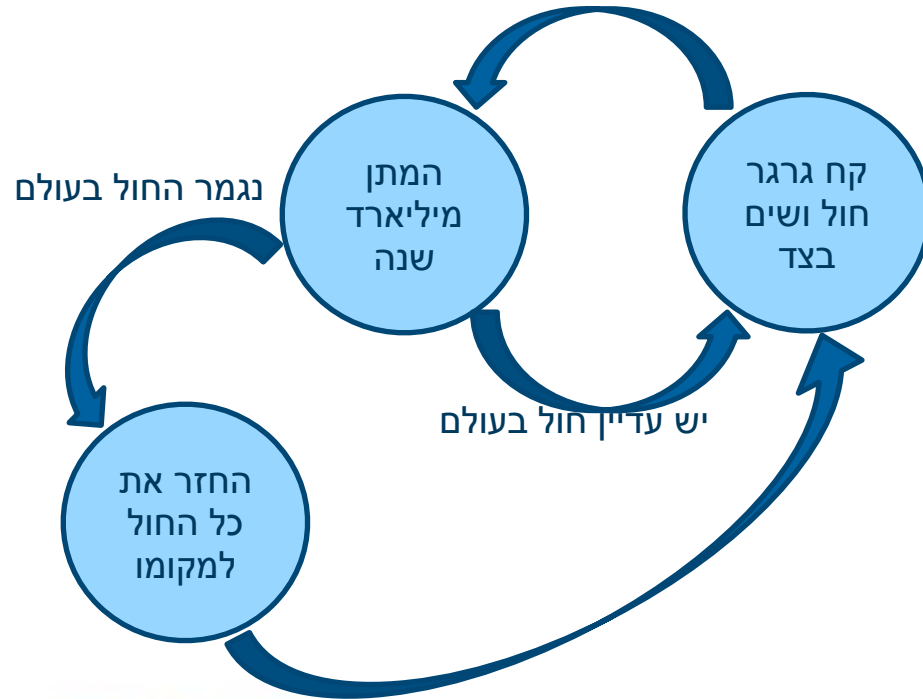
0



$$Q(S,A) = Q(S,A) + \frac{1}{4} \{ r + \gamma \max_{A'} Q(S',A') - Q(S,A) \}$$

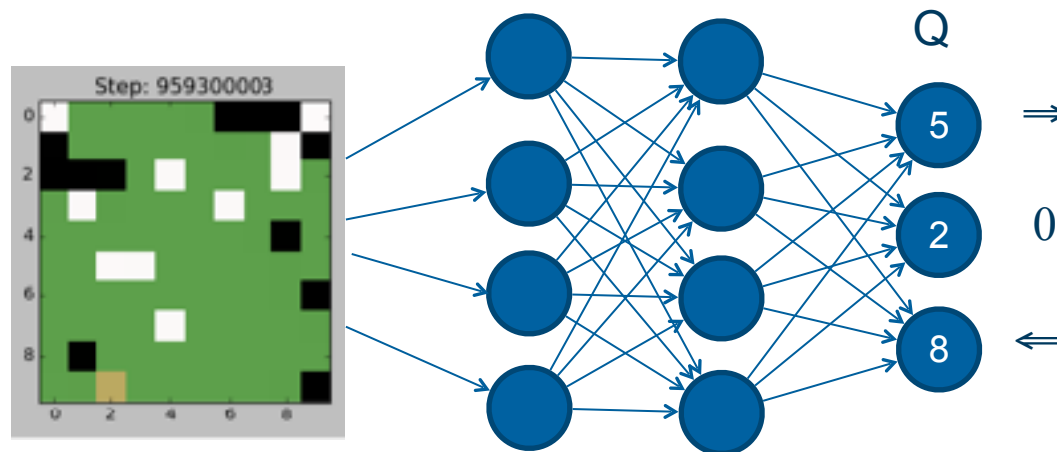


- מספר המצבים האפשריים בדוגמא הוא בערך  $10^{49}$ .
- כמה גדול המספר הזה? כמה זמן זה בשניות?



$\times 10000$

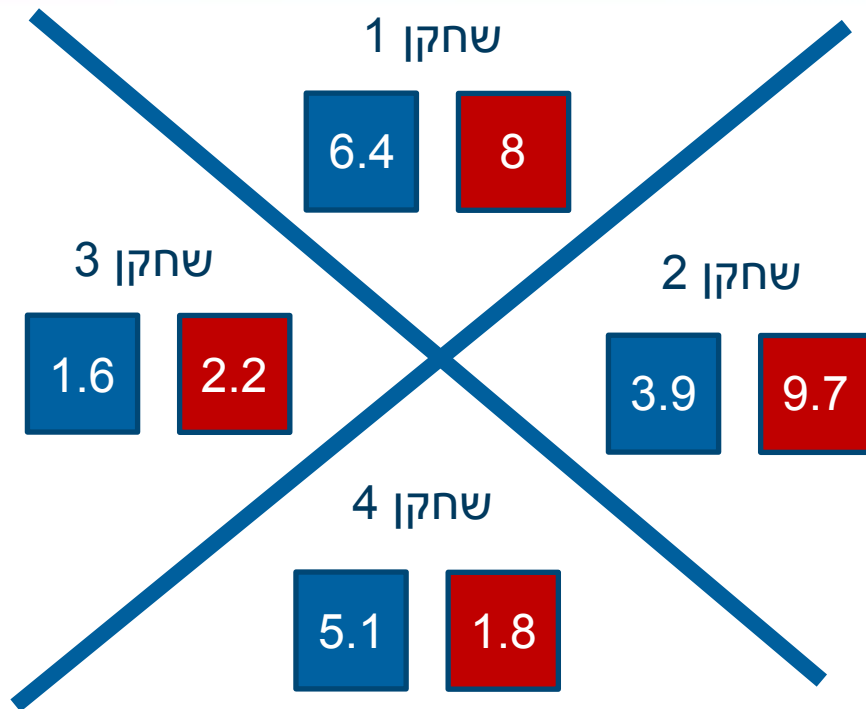
- הפתרון: נשתמש ברשת נוירונים בשביל ללמוד ייצוג של פעולות כתגובה למצבים.



## חזרה לבעיה

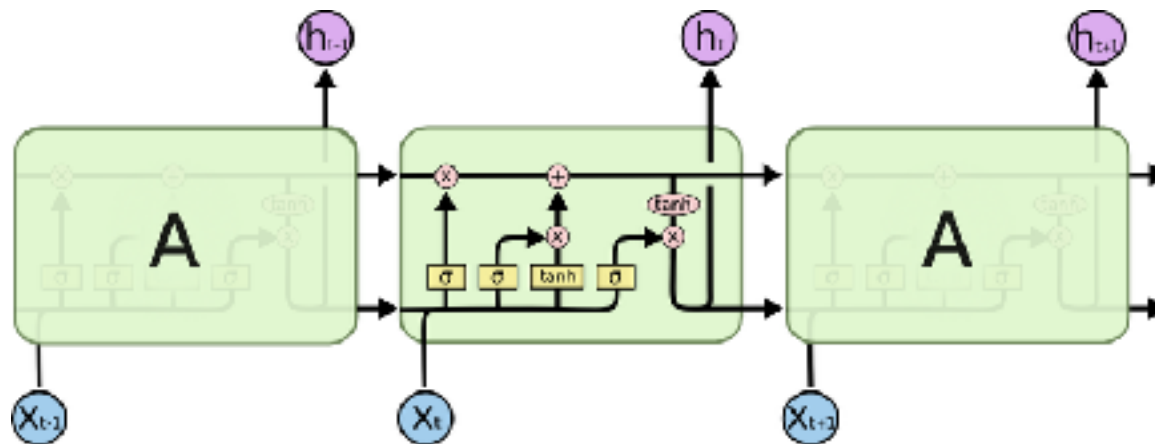
# אתגר לימוד פרוטוקול תקשורת פשוט

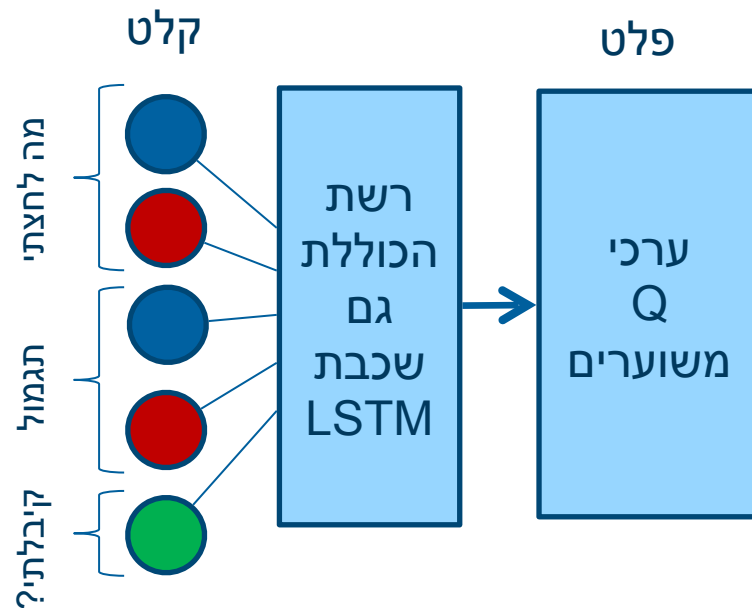
• תזכורת



## LSTM

- שכבת LSTM היא שכבה שבה כל נוירון כולל מצב פנימי שאותו הוא מעדכן על סמך הקלט.
- אנו משתמשים בה בכדי לפתור את בעיית ה partial observability וחוסר המרקוביות.

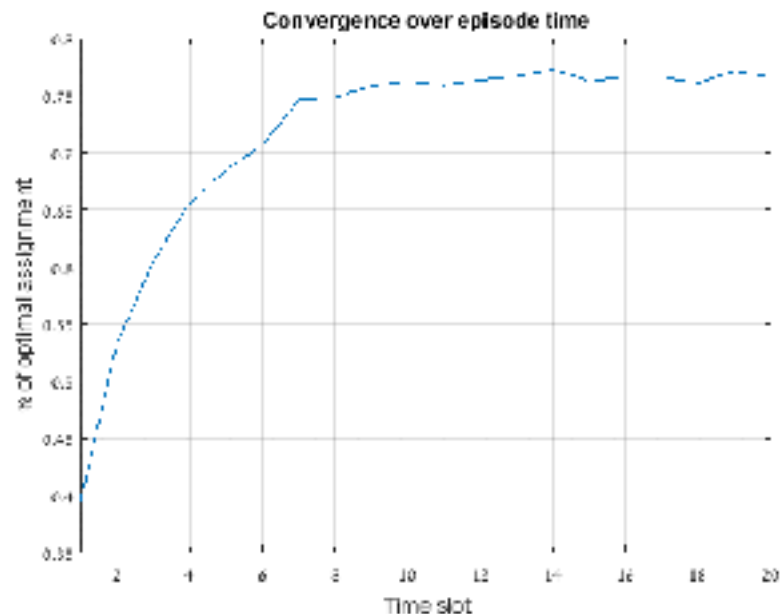
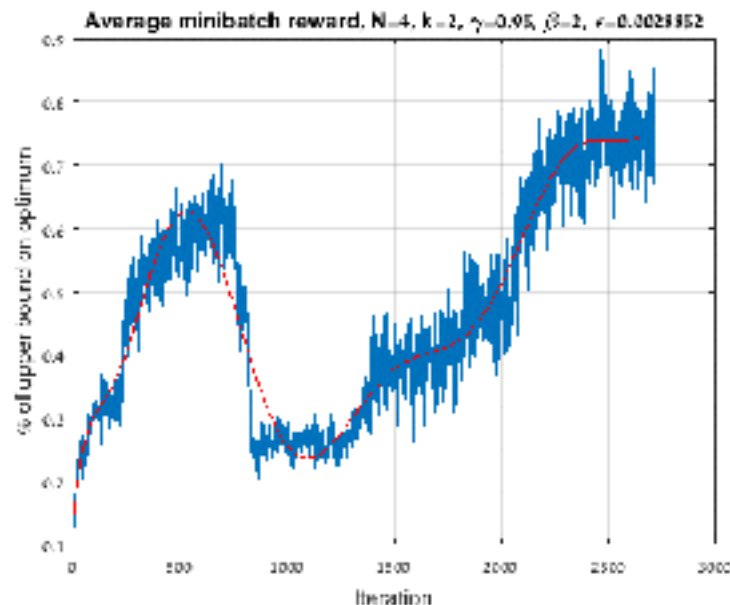




- הרשת מורכבת משכבת LSTM בכדי שתוכל ללמוד תלויות בזמן.
- לאחר ה LSTM ישנה שכבה שחוצה את ערכי ה Q – תוחלת הרווחים העתידיים.
- כל השחקנים ישתמשו תמיד באותה רשת.
- העדכון נעשה כל 30 משחקים של 20 שלבים.



• משיג כ 75% מהביצועים של חסם עליון על האופטימום.



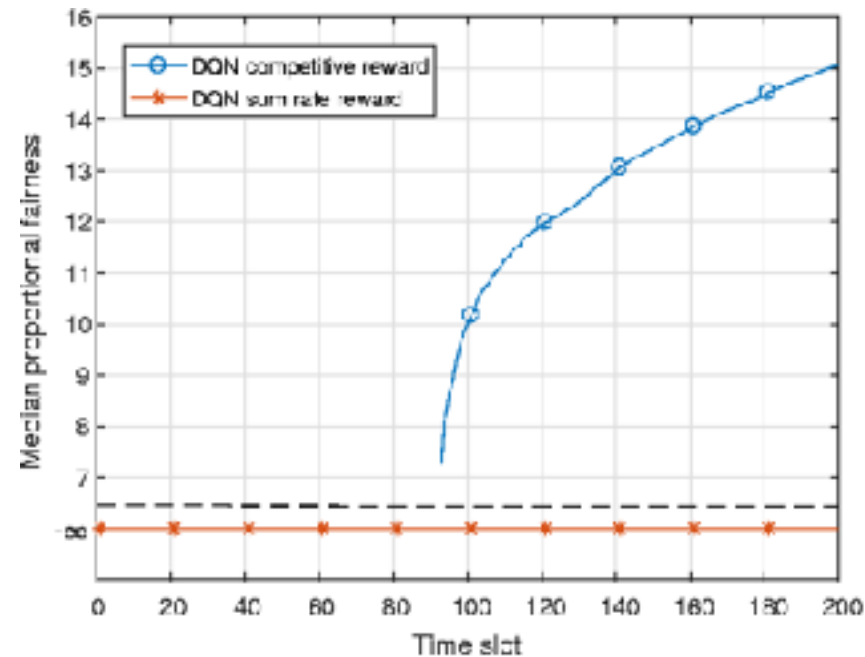
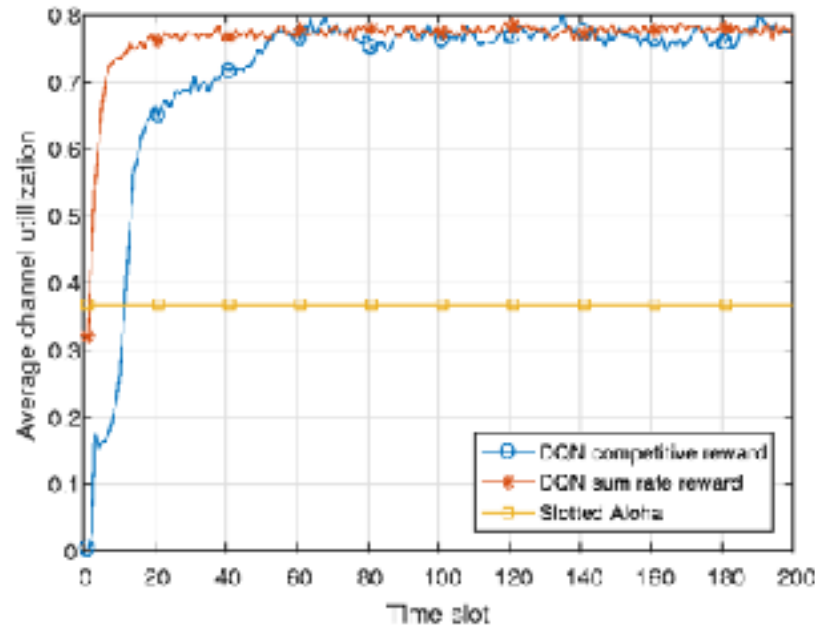
## בעיית תקשורת aloha

- זוגות משתמשים ללא יכולת העברת הודעות בין זוגות שונים מעוניינים להשתמש בערוץ N תקשורת ללא תיאום מראש.
- בכל תא זמן כל משתמש מחליט האם לשדר או לא ובאיזה ערוץ.
- אם משתמש שידר הודעה בהצלחה הוא מקבל ACK מבן הזוג שלו.
- אם יותר ממשתמש אחד משדר בחריוץ זמן אז אנו מניחים כי כל ההודעות לא הגיעו.
- אנו מעוניינים לפתח שיטות מבוזרות למיקסום פונקציית תועלת ברשת כדוגמת utilization.

בדקנו את האלגוריתם עבור שיטות תגמול שונות:

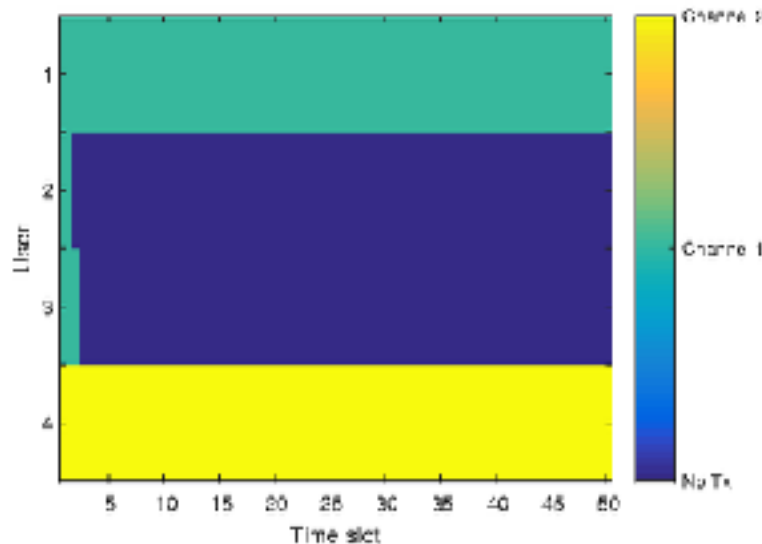
- תגמול כללי - כל המשתמשים מקבלים תגמול על כל הודעה ששודרה בהצלחה.
- תגמול פרטי - רק המשתמש ששידר בהצלחה מקבל תגמול.
- מספר  $\log$  כל משתמשים מקבלים תגמול שהוא סכום – proportional fairness – ההודעות שכל משתמש הצליח לשדר

# השוואה בין תגמול פרטי לכללי



# השוואה בין פונקציות התגמול השונות

## פעולות - תגמול כללי

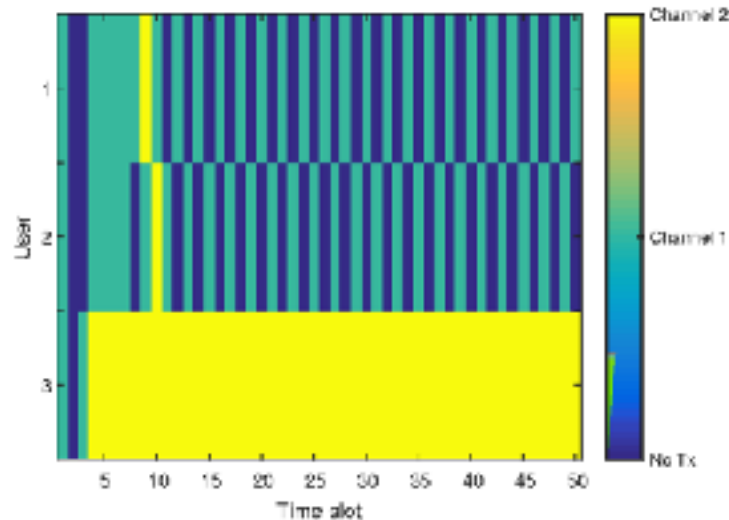


- מכיוון שכל השחקנית מתוגמלים בצורה זהה עבור כל שידור מוצלח, האסטרטגיה האופטימלית היא להגיע למצב שתמיד שניים משדרים והשאר שותקים.

## השוואה בין פונקציות התגמול השונות

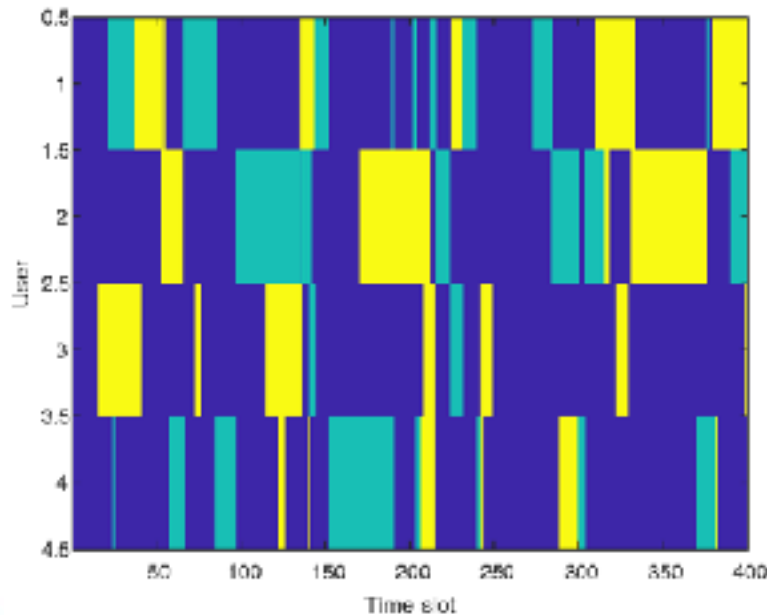
- השחקנים מתוגמלים רק על ההודעות שהם שידרו, תגמול זה מכריח אותם למצוא דרך להתחלק בזמן.

פעולות - תגמול פרטי



## השוואה בין פונקציות התגמול השונות

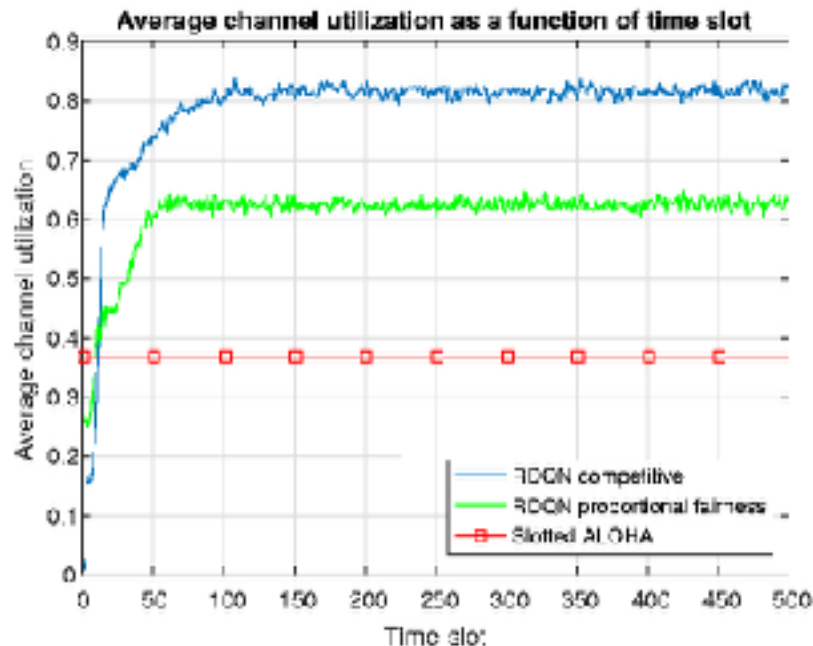
פעולות - תגמול כללי עם הוגנות



- השחקנים מתוגמלים על ידי תגמול כללי הנותן יתרון לחלוקת משאבים שווה. האסטרטגיה שנלמדה היא לשדר בערוץ עד שמישהו אחר רוצה לשדר ואז לוותר לו. כן, לא נשארים זמן ארוך מידי בשידור. למקרה שמישהו אחר ירצה לשדר.



## השוואה בין תגמול פרטי ל fairness



- שני הסכמות התגמול התכנסו לפתרונות הוגנים.
- ניצולת הערוץ הייתה גבוהה יותר עבור התגמול הפרטי.

# תודה!

לעוד פרטים ניתן לגשת למאמר ב arxiv:

**[Deep Multi-User Reinforcement Learning for Dynamic Spectrum Access in Multichannel Wireless Networks](#)**