



Πανεπιστήμιο Κύπρου - Τεχνητή Νοημοσύνη

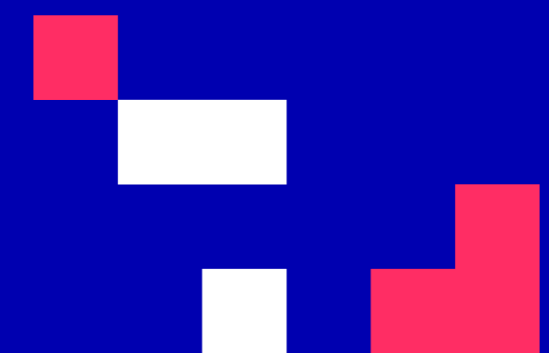
MAI612 - ΜΗΧΑΝΙΚΗ ΜΑΘΗΣΗ

Διάλεξη 17: Διαδικασίες Αποφάσεων Markov και
Δυναμικός Προγραμματισμός

Βασίλης Βασιλειάδης, PhD
Χειμερινό Εξάμηνο 2022/23



CYENS
CENTRE OF EXCELLENCE





Διάλεξη 17: Διαδικασίες Αποφάσεων Markov και Δυναμικός Προγραμματισμός

Μαθησιακά αποτελέσματα

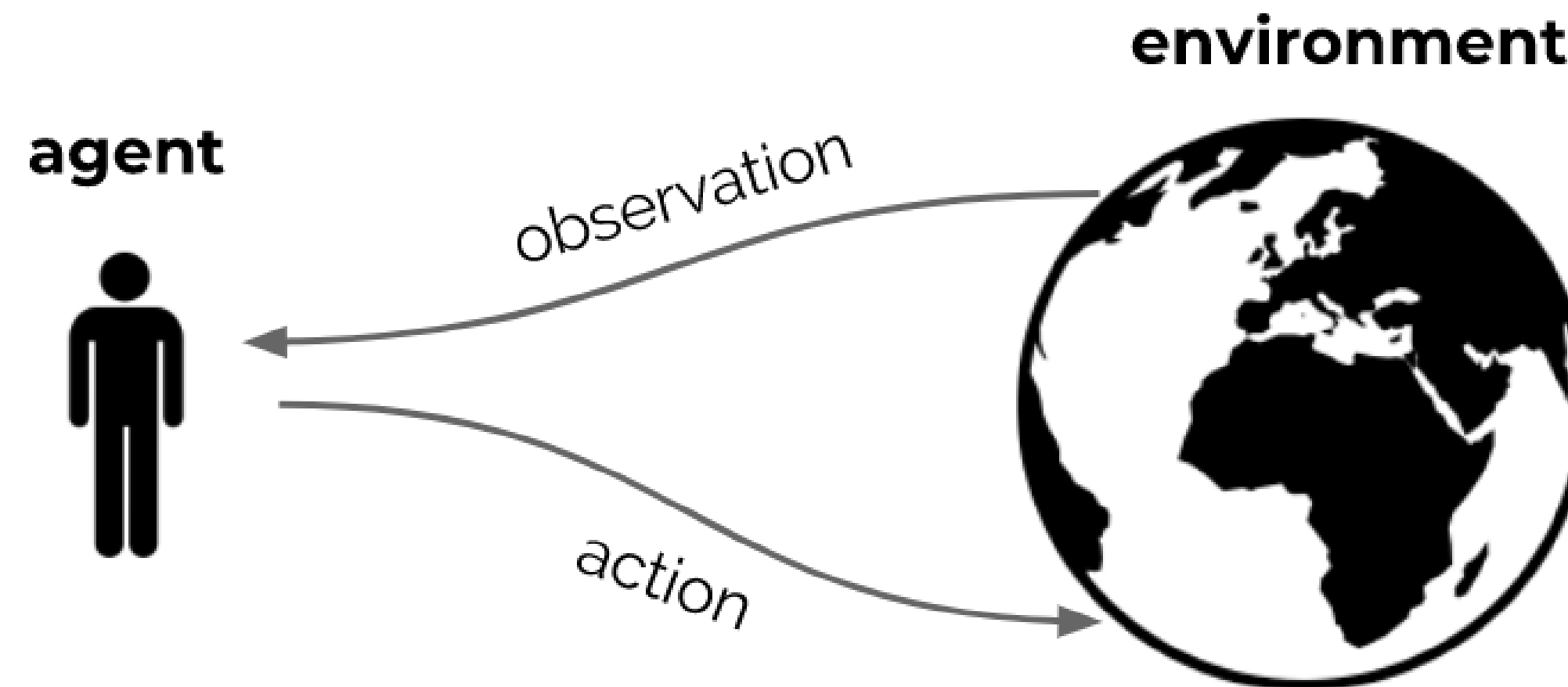
Θα καταλάβετε:

1. Πώς να επισημοποιηθεί το πρόβλημα RL χρησιμοποιώντας το πλαίσιο των διαδικασιών λήψης αποφάσεων Markov (ΔAM).
2. Τι είναι οι συναρτήσεις αξίας και πώς σχετίζονται με τους στόχους σε ένα ΔAM.
3. Δύο κατηγορίες προβλημάτων RL: πρόβλεψη και έλεγχος.
4. Πώς να χρησιμοποιήσετε τις εξισώσεις Bellman για την επίλυση προβλημάτων RL.
5. Δυναμικός προγραμματισμός και αλγόριθμοι επανάληψης αξίας και επανάληψης πολιτικής για την εξεύρεση λύσεων σε ΔAM.





Επανάληψη



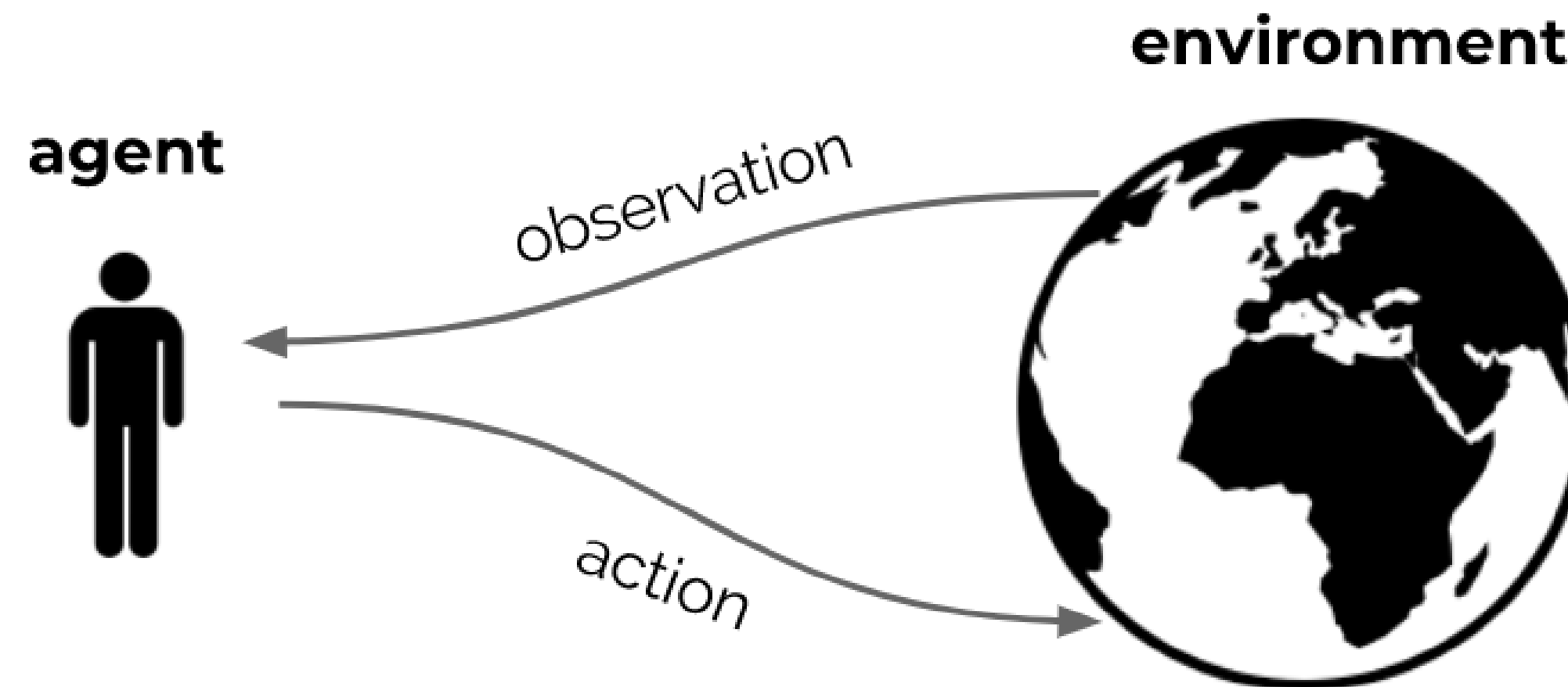
- Ενισχυτική μάθηση είναι η επιστήμη της μάθησης για τη λήψη αποφάσεων
- Οι πράκτορες μπορούν να μάθουν μια **πολιτική**, μια **συνάρτηση αξίας** ή/και ένα **μοντέλο**
- Το γενικό πρόβλημα περιλαμβάνει τη συνεκτίμηση του **χρόνου** και των **συνεπειών**
- Οι αποφάσεις επηρεάζουν την **ανταμοιβή**, το **κράτος πράκτορα** και την κατάσταση του **περιβάλλοντος**
- Η μάθηση είναι **ενεργή**: οι αποφάσεις έχουν αντίκτυπο στα δεδομένα

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Επισημοποίηση της διεπαφής EM

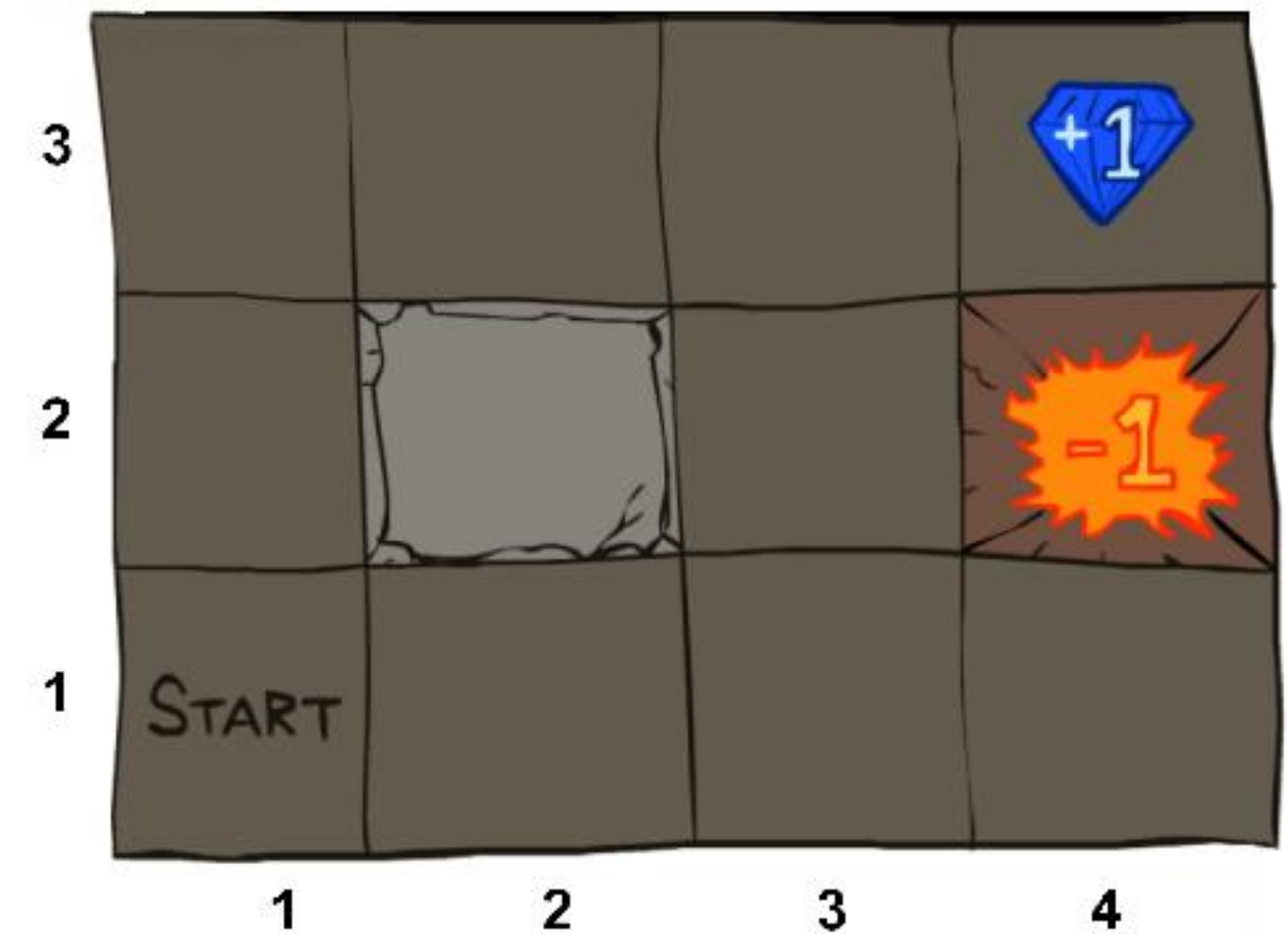


- Θα συζητήσουμε μια μαθηματική διατύπωση της αλληλεπίδρασης παράγοντα-περιβάλλοντος
- Αυτό ονομάζεται **Διαδικασία Απόφασης Markov (ΔΑΜ)**
- Μας δίνει τη δυνατότητα να μιλήσουμε ξεκάθαρα για τον **στόχο** και **τον τρόπο επίτευξής του**



Παράδειγμα: Ο Κόσμος του πλέγματος

- Ένα πρόβλημα που μοιάζει με λαβύρινθο
 - Ο πράκτορας ζει σε ένα πλέγμα
 - Τοίχοι μπλοκάρουν την πορεία του πράκτορα
- Θορυβώδης κίνηση: οι ενέργειες δεν πηγαίνουν πάντα όπως έχει προγραμματιστεί
 - 80 % του χρόνου, η δράση Βόρεια παίρνει τον πράκτορα Βόρεια (αν δεν υπάρχει τοίχος εκεί)
 - 10 % του χρόνου, αντί για βορεια ο πράκτορα πάει δυτικά και 10 % Ανατολικά
 - Αν υπάρχει τοίχος προς την κατεύθυνση που θα είχε πάρει ο πράκτορας, ο πράκτορας παραμένει στη θέση του.
- Ο πράκτορας λαμβάνει ανταμοιβές κάθε φορά που
 - Μικρή «ζωντανή» ανταμοιβή κάθε βήμα (μπορεί να είναι αρνητικό)
 - Μεγάλες ανταμοιβές έρχονται στο τέλος (καλό ή κακό)
- Ο στόχος μας: μεγιστοποιήστε το άθροισμα των ανταμοιβών

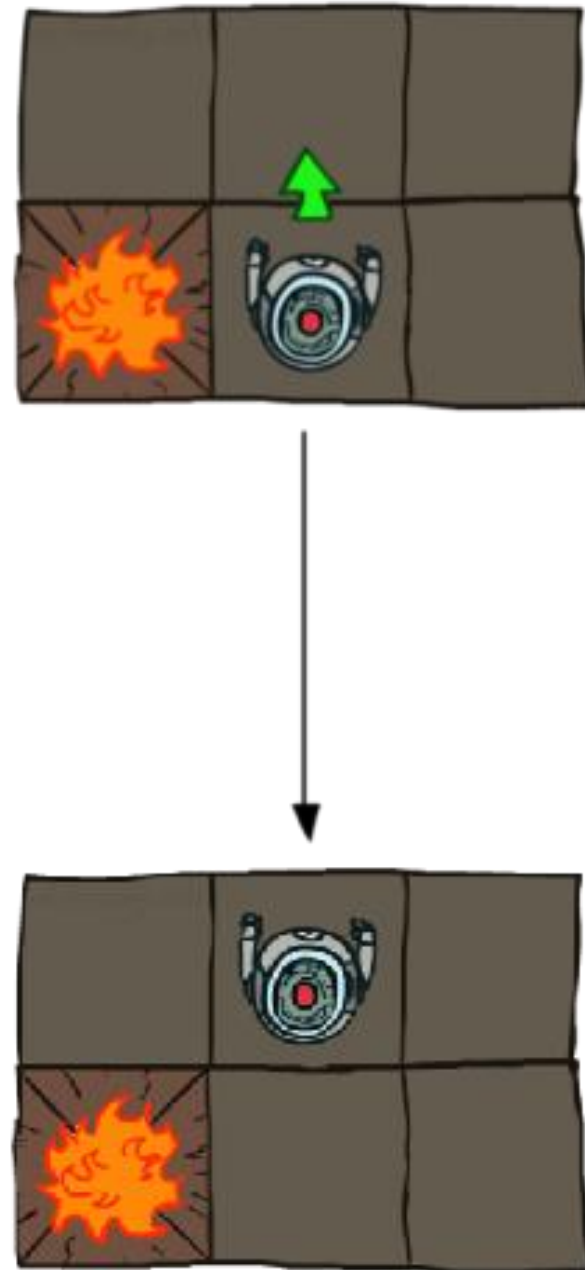


Διαφάνειες βασισμένες στο: [UC Berkeley CS188 – Intro to AI course](#)

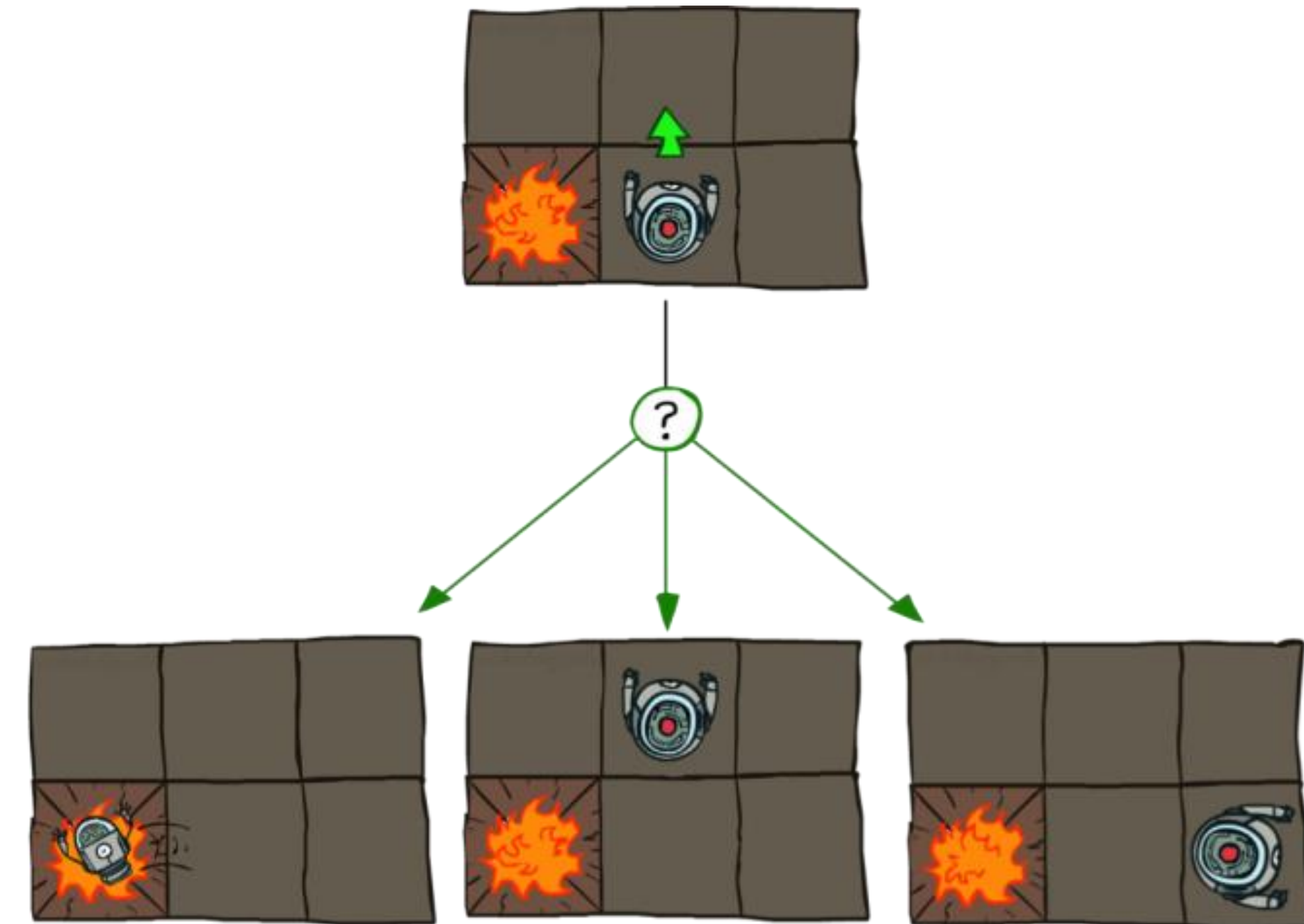


Παράδειγμα: Δράσεις στο κόσμος του πλέγματος

Ντετερμινιστικός κόσμος πλέγματος



Στοχαστικός κόσμος πλέγματος



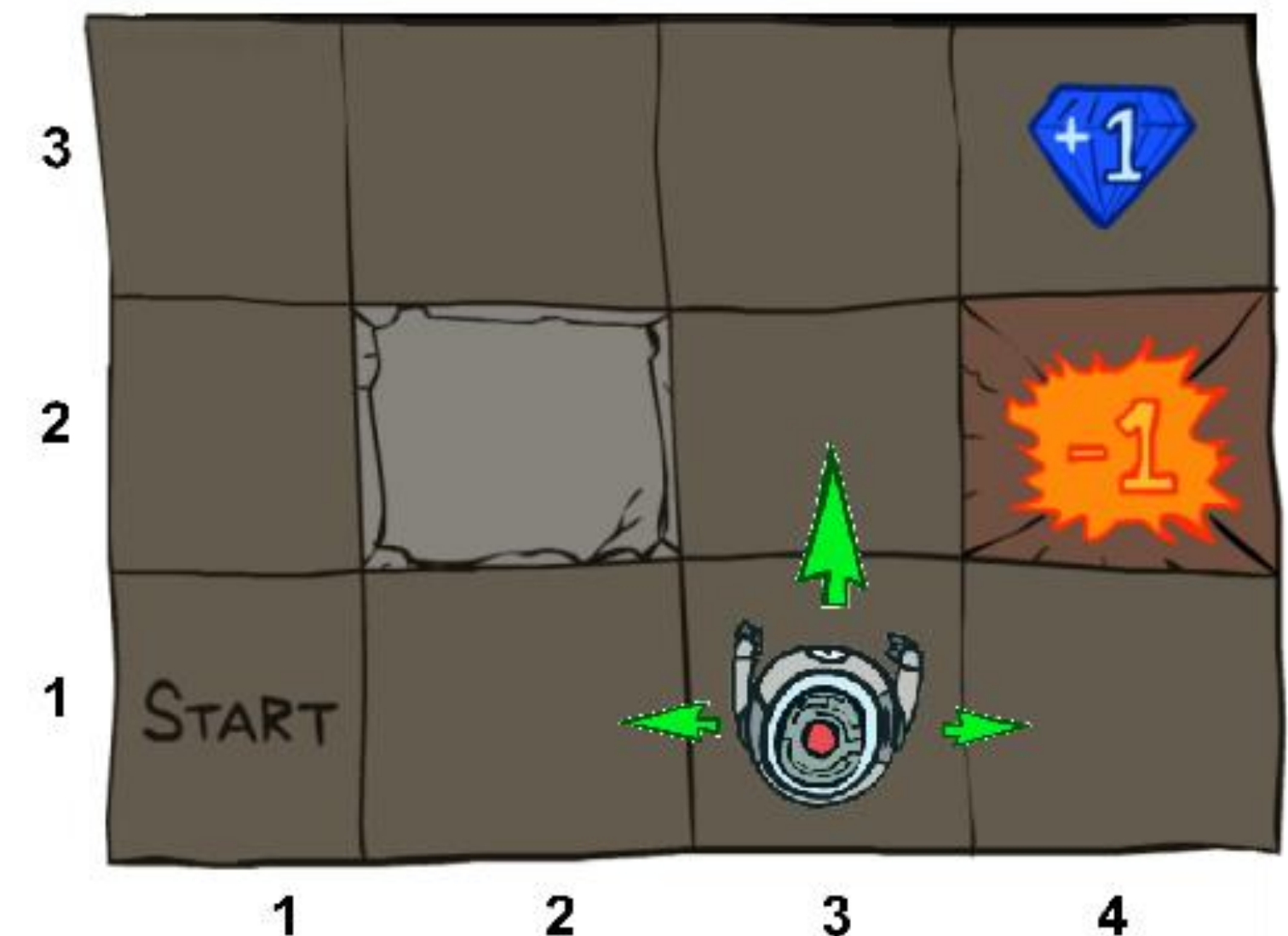
Διαφάνειες βασισμένες στο: [UC Berkeley CS188 – Intro to AI course](#)



Διαδικασία Απόφασης Markov Ορισμός της διαδικασίας λήψης αποφάσεων

Η **διαδικασία λήψης αποφάσεων Markov** είναι μια πλειάδα (S, A, T, r) όπου:

- S είναι το **σύνολο όλων των πιθανών καταστάσεων**
- A είναι το **σύνολο όλων των δυνατών ενεργειών** (π.χ. μηχανοκίνητα χειριστήρια)
- $T(s,a,s')$ είναι μια **μεταβατική συνάρτηση** που καθορίζει την πιθανότητα μετάβασης στην κατάσταση s' , δεδομένης μιας κατάστασης s και της δράσης a , δηλαδή, $p(s'|s,a)$
 - καθορίζει τη **δυναμική** του προβλήματος
- $R(s,a)$ είναι η **συνάρτηση ανταμοιβής** (μερικές φορές $r(s)$, $r(s')$ ή $r(s,a,s')$)





Ιδιότητα Markov

- «Markov» σημαίνει ότι με δεδομένο το παρόν κράτος, το μέλλον είναι ανεξάρτητο από το παρελθόν.
- Για τις ΔΑΜ, «Markov» σημαίνει ότι τα αποτελέσματα της δράσης εξαρτώνται μόνο από την τρέχουσα κατάσταση.

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1}, \dots, S_0 = s_0)$$

=

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

Σε μία ΔΑΜ **όλες οι καταστάσεις** υποτίθεται ότι έχουν την ιδιοκτησία του Markov

- Η κατάσταση συλλαμβάνει όλες τις σχετικές πληροφορίες από το ιστορικό.
- Μόλις γίνει γνωστή η κατάσταση, το ιστορικό μπορεί να πεταχτεί.
- Η κατάσταση είναι ένα επαρκές στατιστικό στοιχείο του παρελθόντος.



Andrey Markov
(1856-1922)

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)



Για παράδειγμα: ρομπότ καθαρισμού

- Ένα ρομπότ που καθαρίζει τα κουτιά σόδας
- Δύο καταστάσεις: **υψηλή** φόρτιση μπαταρίας ή **χαμηλή** φόρτιση μπαταρίας
- Δραστηριότητες: {περιμένετε, ψάξτε} ψηλά, {περιμένετε, αναζήτηση, επαναφόρτιση} στο χαμηλό
- Η δυναμική μπορεί να είναι στοχαστική
 - $P(S_{t+1} = \text{υψηλό} \mid S_t = \text{υψηλό}, A_t = \text{αναζήτηση}) = \alpha$
 - $P(S_{t+1} = \text{χαμηλό} \mid S_t = \text{υψηλό}, A_t = \text{αναζήτηση}) = 1 - \alpha$
- Η ανταμοιβή θα μπορούσε να είναι αναμενόμενος αριθμός συλλεγόμενων κουτιών (deterministic) ή πραγματικός αριθμός συλλεγόμενων κουτιών (stochastic)





Παράδειγμα: ρομπότ καθαρισμού ΔΑΜ

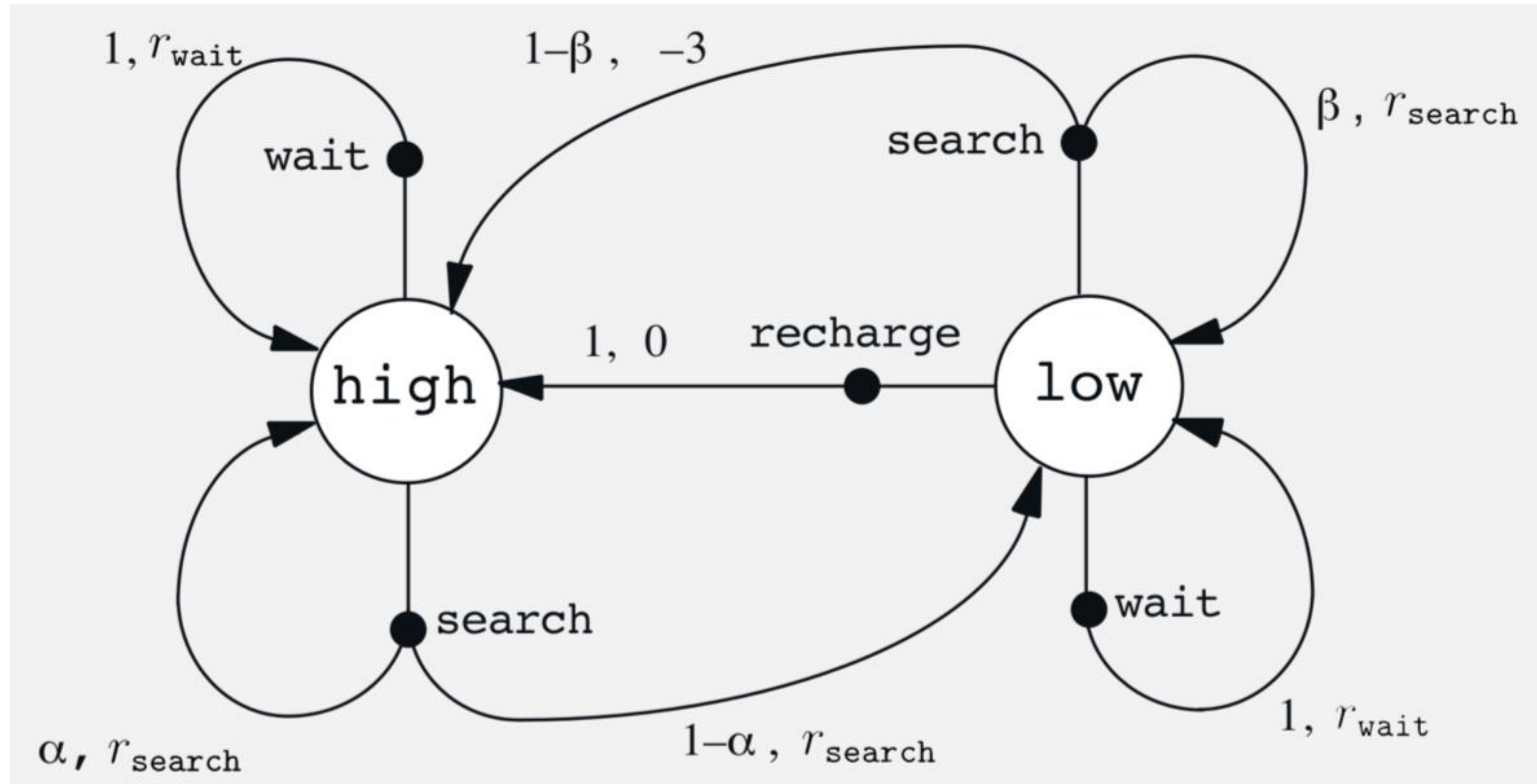
s	a	s'	$p(s' s, a)$	$r(s, a, s')$
high	search	high	α	r_{search}
high	search	low	$1 - \alpha$	r_{search}
low	search	high	$1 - \beta$	-3
low	search	low	β	r_{search}
high	wait	high	1	r_{wait}
high	wait	low	0	r_{wait}
low	wait	high	0	r_{wait}
low	wait	low	1	r_{wait}
low	recharge	high	1	0
low	recharge	low	0	0

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Παράδειγμα: ρομπότ καθαρισμού ΔΑΜ



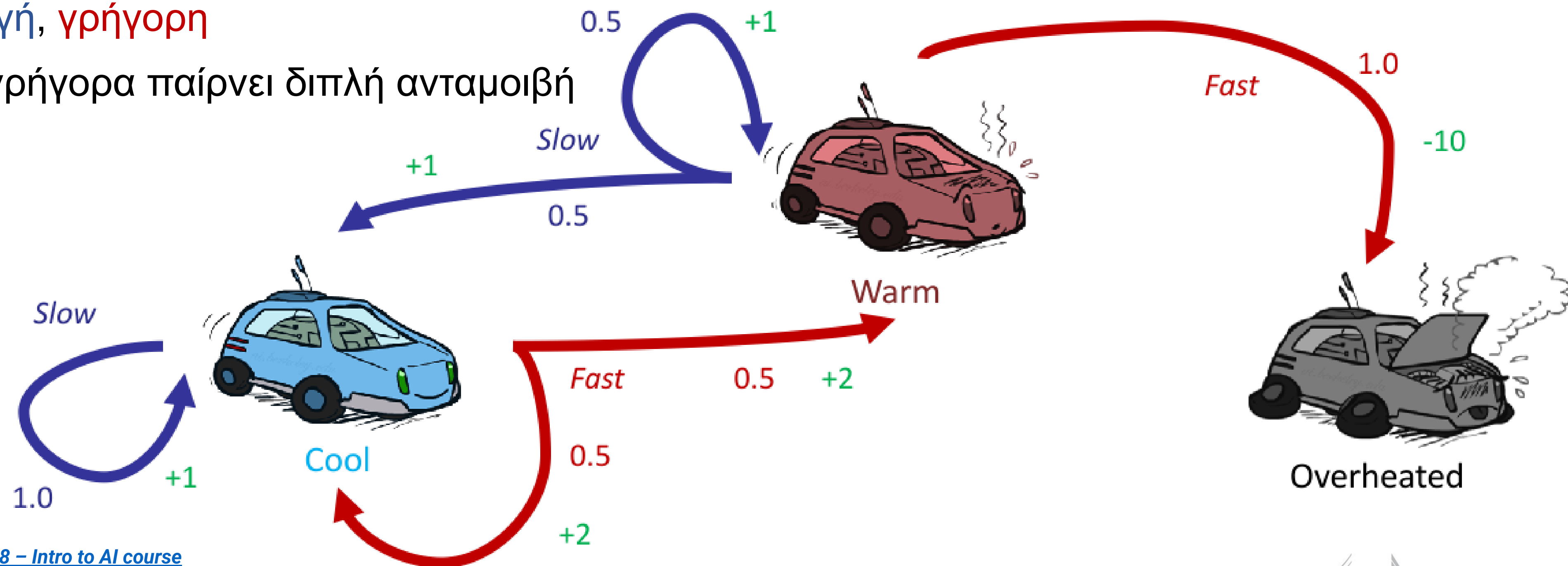
Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Παράδειγμα: Αγώνας

- Ένα ρομποτικό αυτοκίνητο θέλει να ταξιδέψει μακριά, γρήγορα
- Τρεις καταστάσεις: Δροσερό, ζεστό, υπερθερμασμένο
- Δύο καταστάσεις: Αργή, γρήγορη
- Το να πηγαίνεις πιο γρήγορα παίρνει διπλή ανταμοιβή

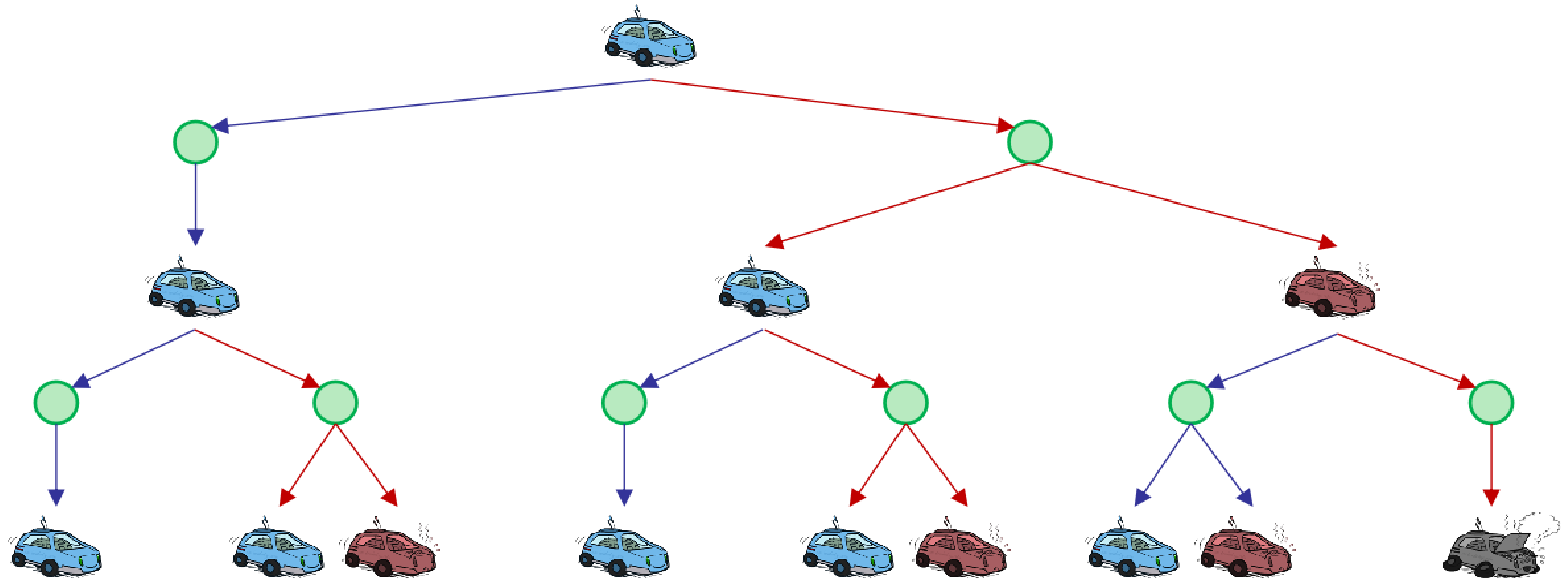


Διαφάνειες βασισμένες στο: [UC Berkeley CS188 – Intro to AI course](#)





Παράδειγμα: δέντρο αναζήτησης για τον αγώνα



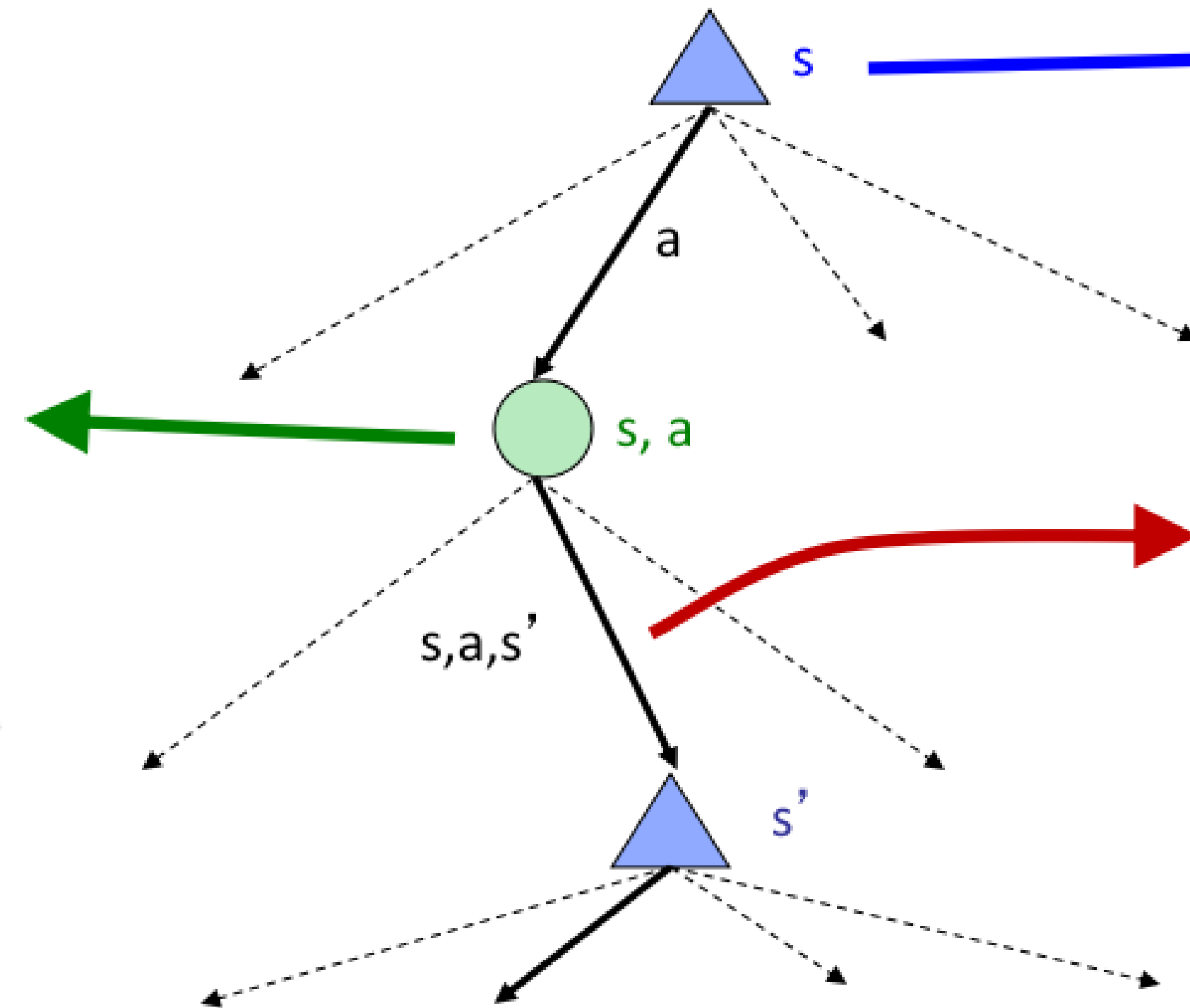
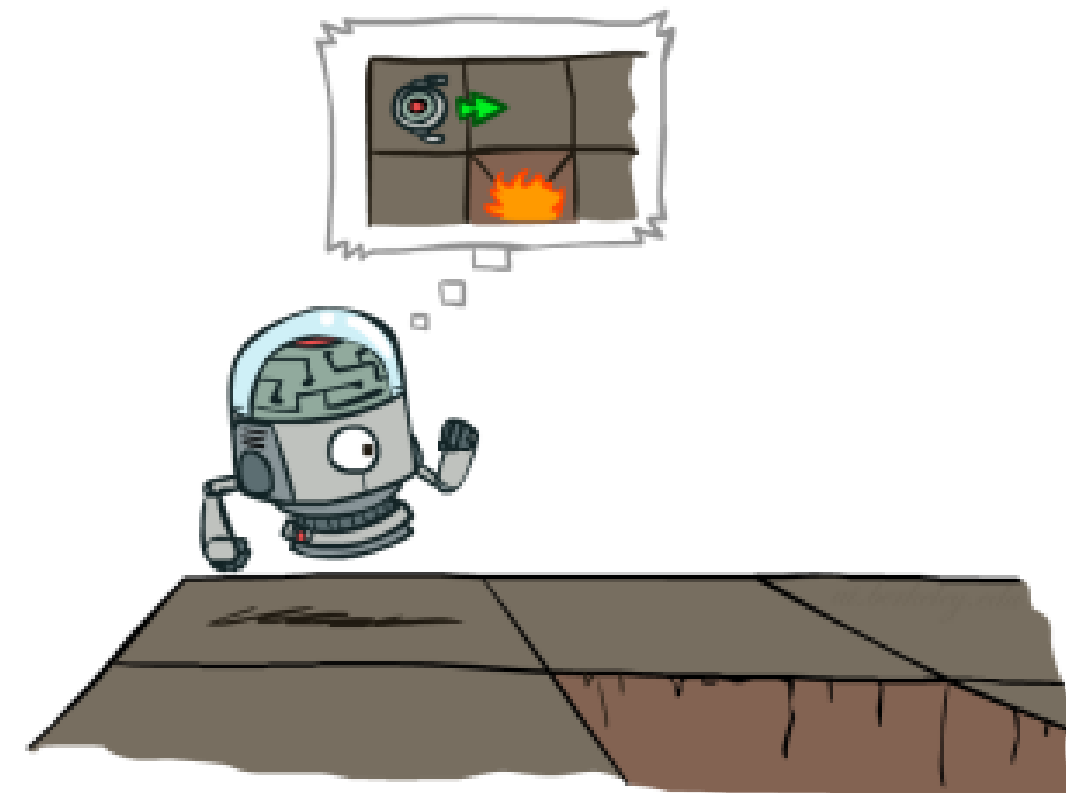
Διαφάνειες βασισμένες στο: [UC Berkeley CS188 – Intro to AI course](#)



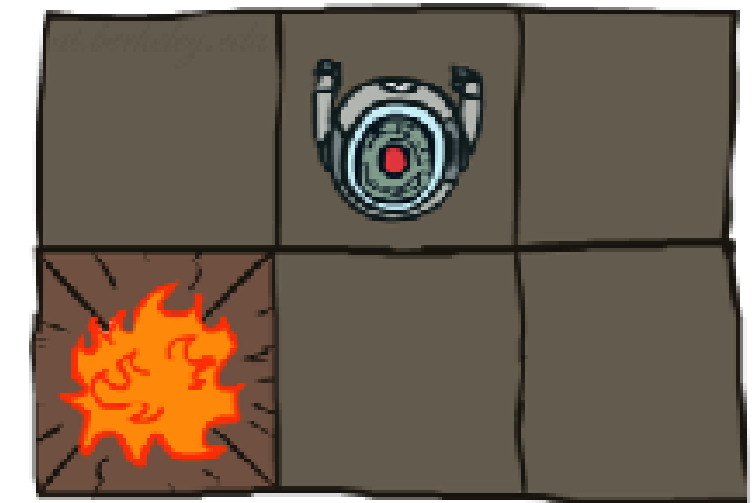


ΔΑΜ Δέντρα Αναζήτησης

- Κάθε κατάσταση ΔΑΜ προβλέπει ένα δέντρο αναζήτησης



s is a state



(s, a, s') called a transition

$$T(s, a, s') = P(s' | s, a)$$

$$R(s, a, s')$$





Επισημοποίηση του στόχου





Απόδοση

- Η δράση σε ένα ΔΑΜ έχει ως **αποτέλεσμα άμεσες ανταμοιβές** R_t , γεγονός που οδηγεί σε **απόδοση** G_t :
 - Απόδοση χωρίς έκπτωση (πρόβλημα επισοδικών/πεπερασμένων οριζόντων)

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T = \sum_{k=0}^{T-t-1} R_{t+k+1}$$

- Απόδοση με έκπτωση(πρόβλημα πεπερασμένου ή άπειρου ορίζοντα)

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t} R_T = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$$

- Μέση απόδοση (συνεχές, άπειρο πρόβλημα ορίζοντα)

$$G_t = \frac{1}{T-t-1} (R_{t+1} + R_{t+2} + \dots + R_T) = \frac{1}{T-t-1} \sum_{k=0}^{T-t-1} R_{t+k+1}$$

Σημείωση: Πρόκειται για τυχαίες μεταβλητές που εξαρτώνται από το **ΔΑΜ** και την **πολιτική**





Απόδοση με έκπτωση

- Είναι λογικό να μεγιστοποιηθεί το άθροισμα των ανταμοιβών
- Είναι επίσης λογικό να προτιμάτε τις ανταμοιβές τώρα από τις ανταμοιβές αργότερα
- Μια λύση: οι αξίες των ανταμοιβών διασπώνται εκθετικά
- Προεξοφλημένες **αποδόσεις** G_t για άπειρο ορίζοντα $T \rightarrow \rho$: $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
- Η **έκπτωση** $\gamma - [0, 1]$ είναι η παρούσα αξία των μελλοντικών ανταμοιβών
 - Η τιμή λήψης ανταμοιβής R μετά από $k + 1$ χρονικά βήματα είναι $\gamma^k R$



1

Worth Now



γ

Worth Next Step



γ^2

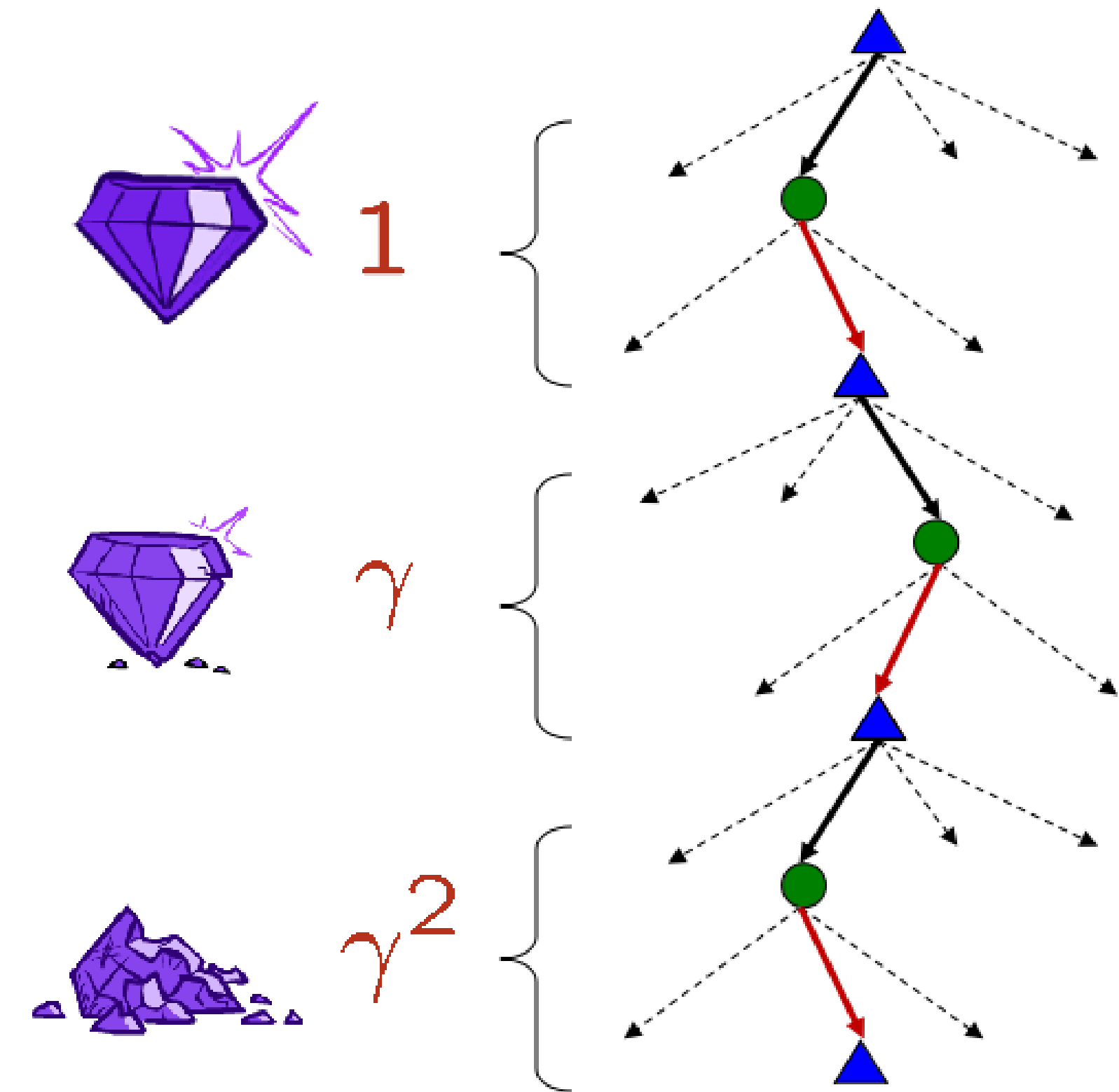
Worth In Two Steps





Έκπτωση

- Πώς να κάνετε έκπτωση;
 - Κάθε φορά που κατεβαίνουμε ένα επίπεδο, πολλαπλασιάζουμε στην έκπτωση μία φορά
- Γιατί έκπτωση;
 - Νωρίτερα οι ανταμοιβές πιθανόν να έχουν υψηλότερη αξία που αργότερα ανταμείβει
 - γ κοντά στο 0 οδηγεί σε «μυωπική» αξιολόγηση: οι άμεσες ανταμοιβές είναι πιο σημαντικές από τις καθυστερημένες ανταμοιβές
 - γ κοντά στο 1 οδηγεί σε «διορατική» αξιολόγηση
 - Μαθηματικά βολικό: αποφεύγει τις άπειρες αποδόσεις σε κυκλικές ΔΑΜ



Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)

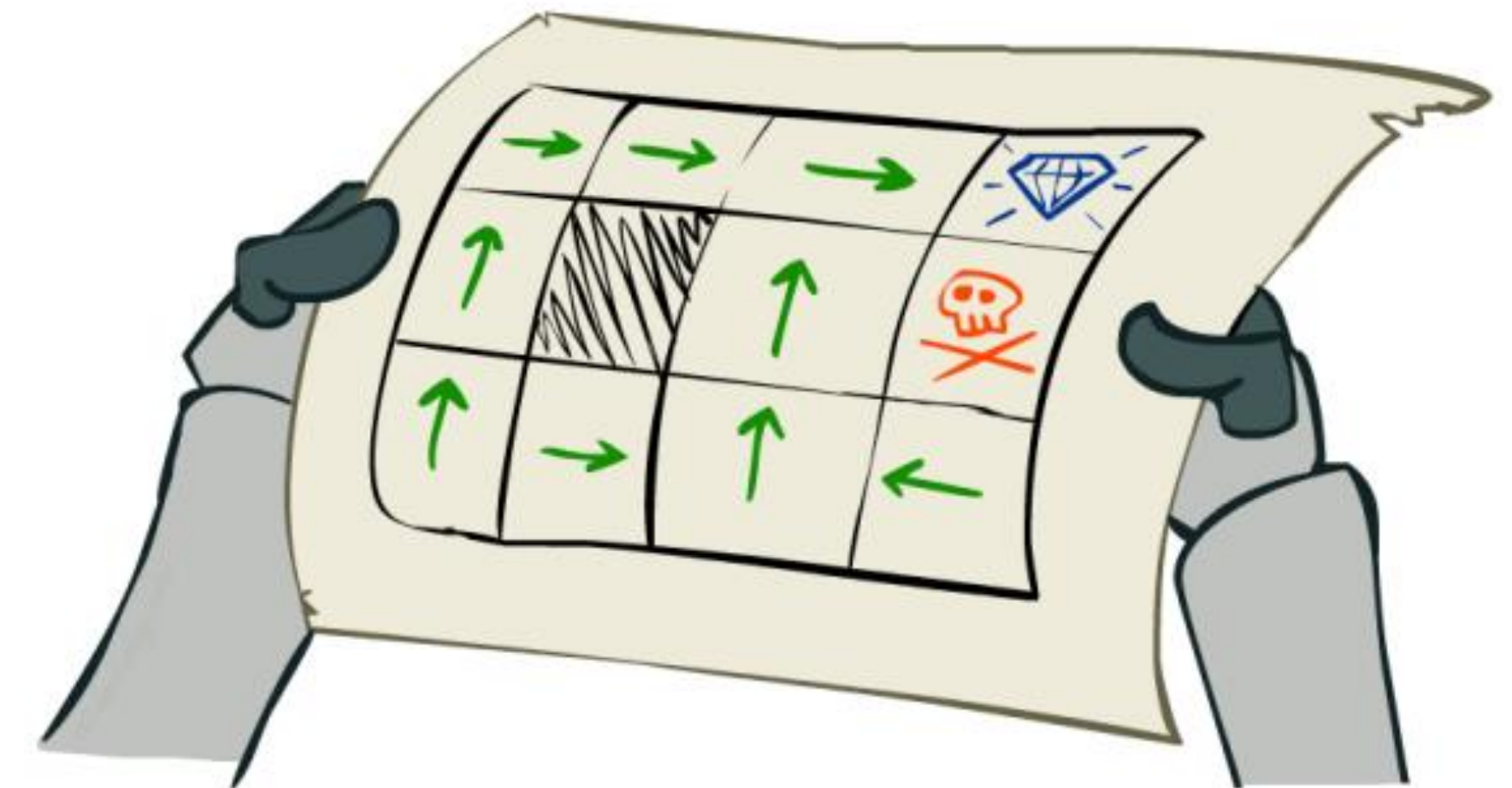


Πολιτικές

Στόχος ενός πράκτορα της RL:

βρείτε μια πολιτική συμπεριφοράς που μεγιστοποιεί την αναμενόμενη απόδοση G_t

- Μια **πολιτική** είναι μια σχέση $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ η οποία, για κάθε κράτος εκχωρεί **για κάθε ενέργεια μια A την πιθανότητα να αναλάβει αυτή τη δράση σε κατάσταση s**. Υποδεικνύεται από $\pi(a|s)$.
- Μια **βέλτιστη πολιτική** είναι αυτή που μεγιστοποιεί την αναμενόμενη απόδοση εάν ακολουθηθεί.
- Για τις ντετερμινιστικές πολιτικές, χρησιμοποιούμε μερικές φορές τον συμβολισμό $\pi_t = \pi(s_t)$ για να υποδηλώσουμε τη δράση που έχει αναλάβει η πολιτική.





Συναρτήσεις αξίας

- Η **συνάρτηση** αξίας vs *δίνει* τη μακροπρόθεσμη τιμή (αναμενόμενη απόδοση) της κατάστασης s

$$v_{\pi}(s) = \mathbb{E} [G_t \mid S_t = s, \pi]$$

- Μπορούμε να ορίσουμε τις **αξίες (κατάστασης-)δράσης**:

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

- Σύνδεση μεταξύ τους:

$$v_{\pi}(s) = \sum_a \pi(a \mid s) q_{\pi}(s, a) = \mathbb{E} [q_{\pi}(S_t, A_t) \mid S_t = s, \pi] , \forall s$$





Βέλτιστες συναρτήσεις αξίας

- Η βέλτιστη συνάρτηση κατάστασης-τιμής $v(s)$ είναι η συνάρτηση μέγιστης τιμής σε όλες τις πολιτικές

$$v^*(s) = \max_{\pi} v_{\pi}(s)$$

- Η βέλτιστη συνάρτηση τιμής δράσης $q(s, a)$ είναι η συνάρτηση μέγιστης τιμής δράσης σε όλες τις πολιτικές

$$q^*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

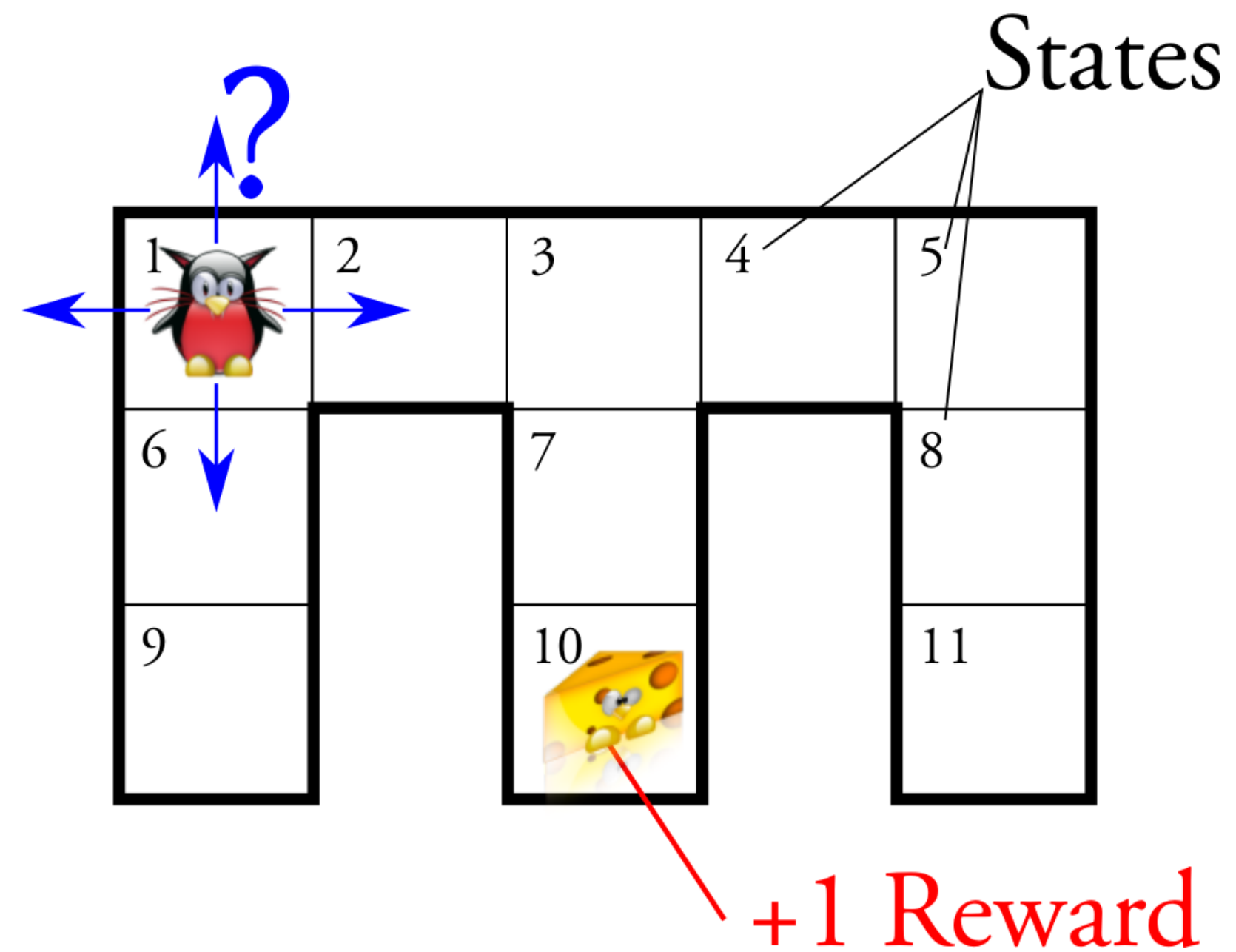
- Η συνάρτηση βέλτιστης τιμής καθορίζει την καλύτερη δυνατή απόδοση στη ΔΑΜ
- Μία ΔΑΜ «επιλύεται» όταν γνωρίζουμε τη βέλτιστη συνάρτηση τιμής





Παράδειγμα συναρτήσεων αξίας

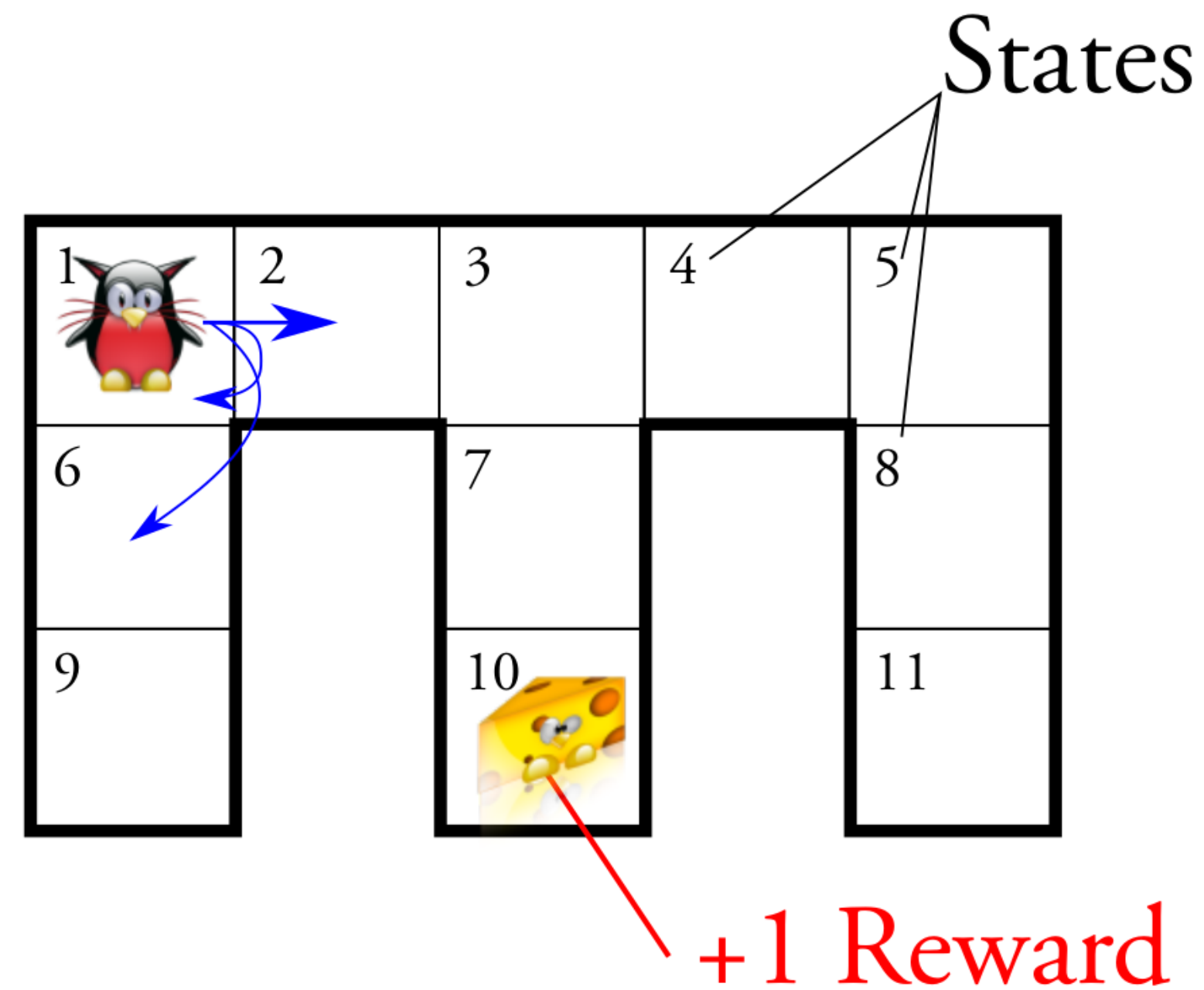
- **Στόχος:** αναζήτηση πολιτικής $\pi: S \rightarrow A$ για να μεγιστοποιήσει $\mathbb{E}_{\sim \pi} [\sum_{t=0}^{\infty} \gamma^t r_t]$





Παράδειγμα συναρτήσεων αξίας

- **Στόχος:** αναζήτηση πολιτικής $\pi: S \rightarrow A$ για να μεγιστοποιήσει $\mathbb{E}_{\sim \pi} [\sum_{t=0}^{\infty} \gamma^t r_t]$



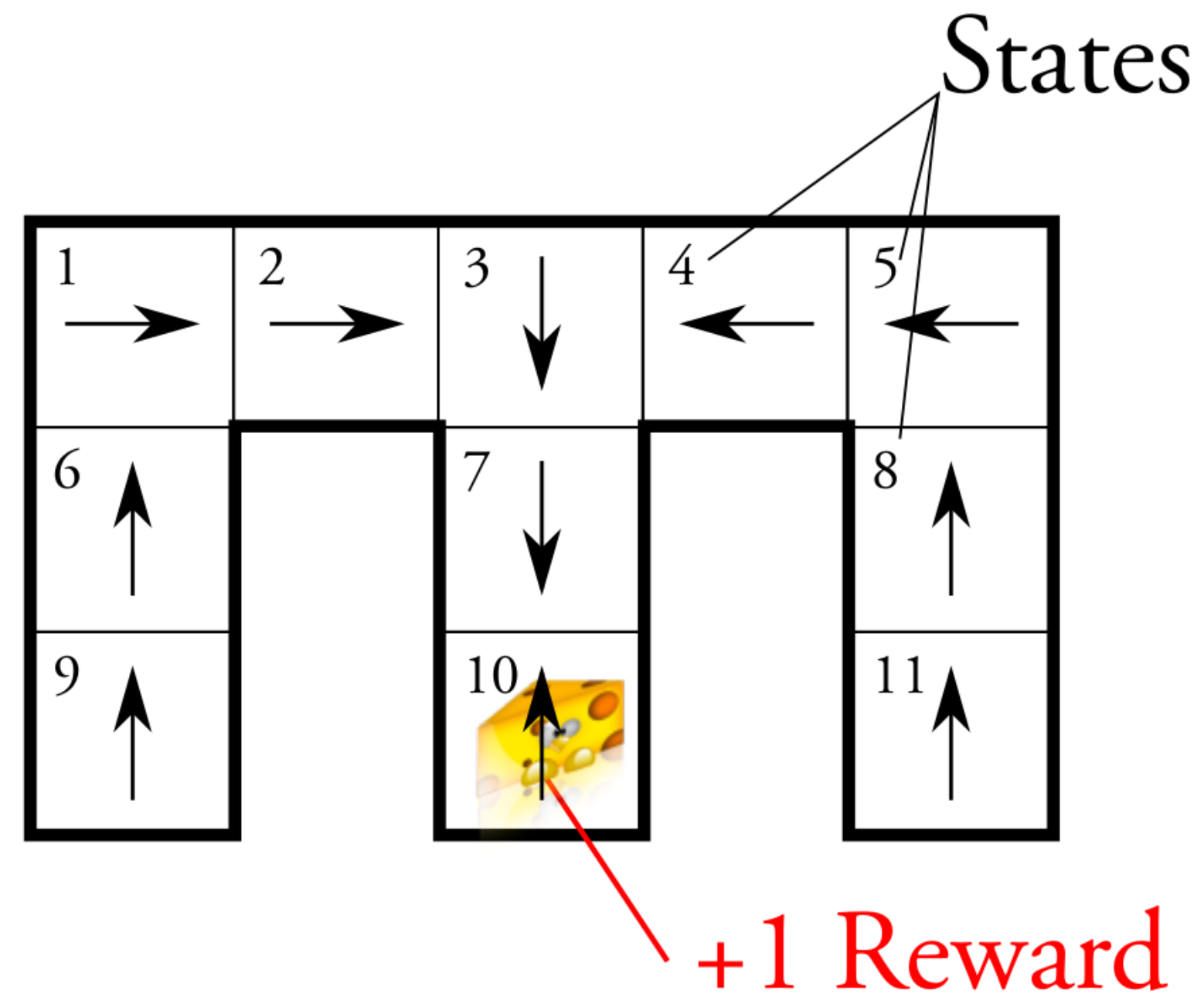
Οι ενέργειες του πράκτορα μπορεί να είναι στοχαστικές





Παράδειγμα συναρτήσεων αξίας

- **Στόχος:** αναζήτηση πολιτικής $\pi: S \rightarrow A$ για να μεγιστοποιήσει $\mathbb{E}_{\sim \pi} [\sum_{t=0}^{\infty} \gamma^t r_t]$



Αυτή είναι η βέλτιστη πολιτική

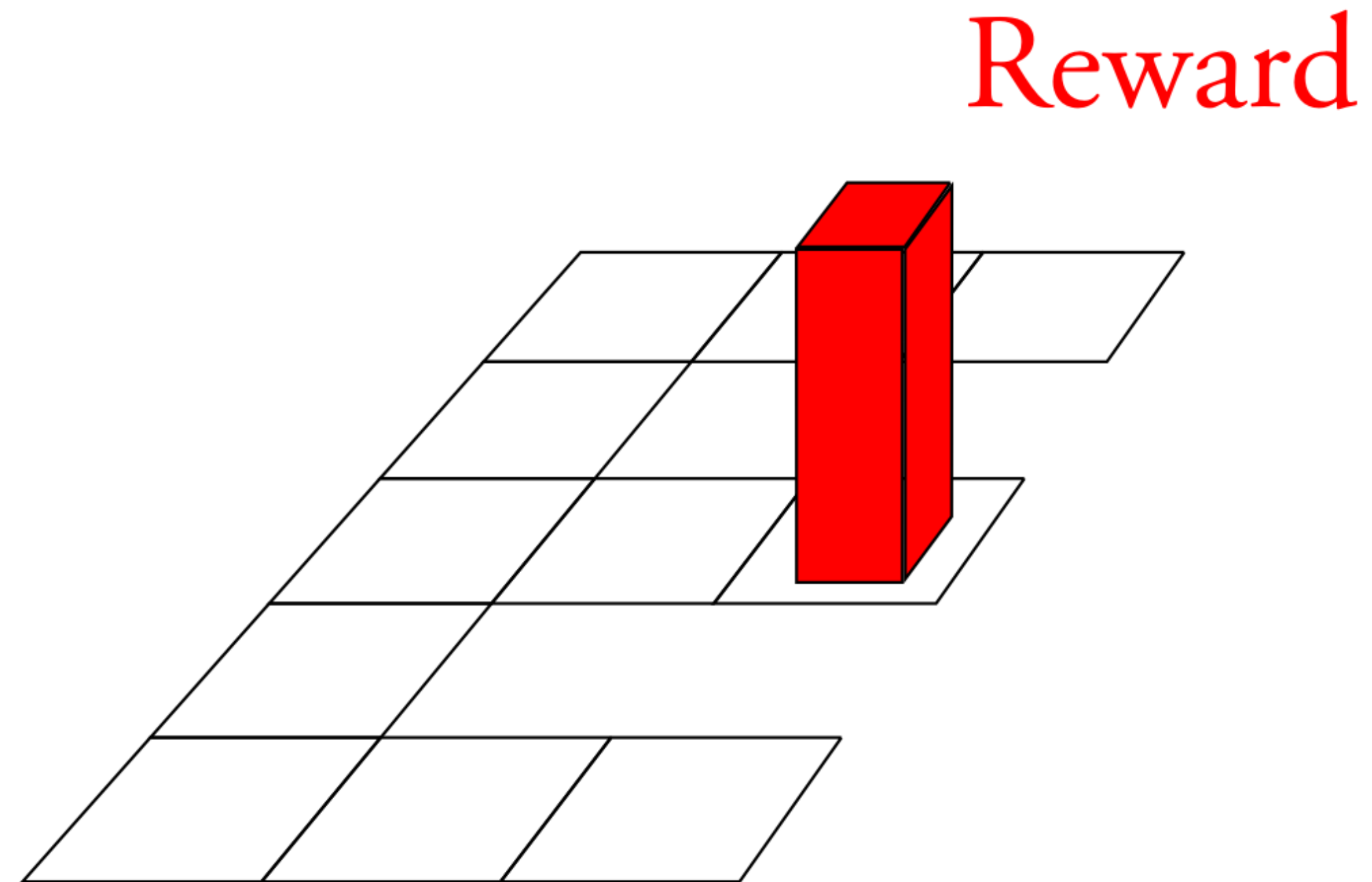
...αλλά πώς θα τη βρούμε;





Παράδειγμα συναρτήσεων αξίας

Η βέλτιστη συνάρτηση τιμής μεταμορφώνει αυτό το «αραιό τοπίο»

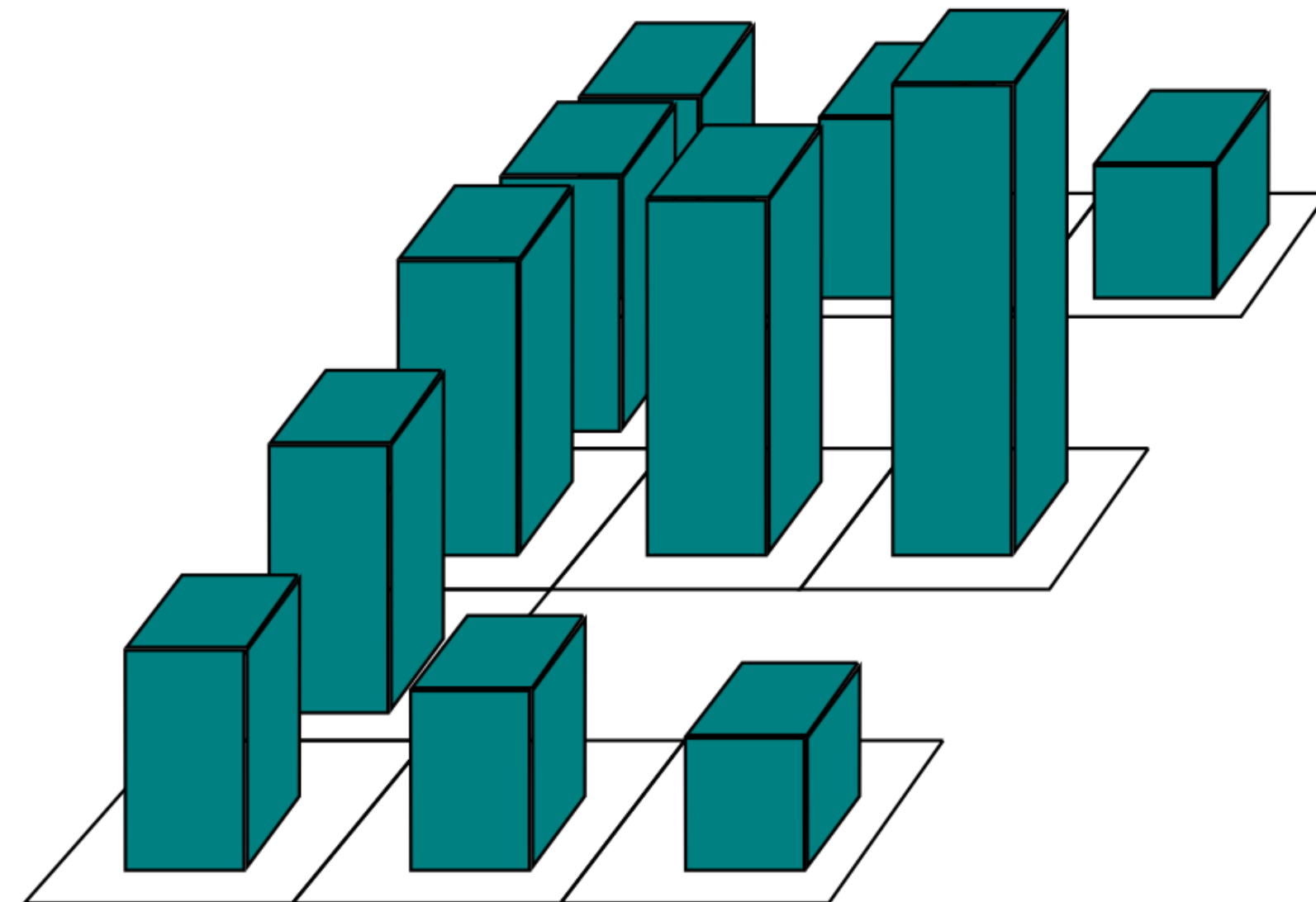




Παράδειγμα συναρτήσεων αξίας

... σε αυτό το «πυκνό» ένα, το οποίο μπορεί να χρησιμοποιηθεί για να υπολογίσει τη βέλτιστη πολιτική

Value Function





Εύρεση βέλτιστης πολιτικής

- Μια βέλτιστη πολιτική μπορεί να βρεθεί με τη μεγιστοποίηση πάνω από $q(s, a)$

$$\pi^*(s, a) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_{a \in \mathcal{A}} q^*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

- ΠΑΡΑΤΗΡΗΣΕΙΣ:

- Υπάρχει πάντα μια ντετερμινιστική βέλτιστη πολιτική για κάθε MDP
- Εάν γνωρίζουμε $q(s, a)$, έχουμε αμέσως τη βέλτιστη πολιτική
- Μπορεί να υπάρχουν πολλαπλές βέλτιστες πολιτικές
- Εάν οι πολλαπλές ενέργειες μεγιστοποιήσουν το $q(s, \cdot)$, μπορούμε επίσης να επιλέξουμε οποιαδήποτε από αυτές (συμπεριλαμβανομένης της στοχαστικής)





Εξισώσεις Bellman



Ρίσαρντ Μπέλμαν
(1920-1984)





Εξισώσεις Bellman

Συνάρτηση αξίας

- Η συνάρτηση τιμής $v(s)$ δίνει τη μακροπρόθεσμη τιμή της κατά $v_\pi(s) = \mathbb{E} [G_t \mid S_t = s, \pi]$
- Μπορεί να οριστεί αναδρομικά:

$$\begin{aligned} v_\pi(s) &= \mathbb{E} [R_{t+1} + \gamma G_{t+1} \mid S_t = s, \pi] \\ &= \mathbb{E} [R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t \sim \pi(S_t)] \\ &= \sum_a \pi(a \mid s) \sum_r \sum_{s'} p(r, s' \mid s, a) (r + \gamma v_\pi(s')) \end{aligned}$$

- Το τελευταίο βήμα γράφει ρητά την προσδοκία





Εξισώσεις Bellman

Αξίες δράσης

- Μπορούμε να ορίσουμε τιμές κρατικής δράσης: $q_\pi(s, a) = \mathbb{E} [G_t \mid S_t = s, A_t = a, \pi]$
- Αυτό συνεπάγεται

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E} [R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \mathbb{E} [R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a] \\ &= \sum_r \sum_{s'} p(r, s' \mid s, a) \left(r + \gamma \sum_{a'} \pi(a' \mid s') q_\pi(s', a') \right) \end{aligned}$$

- Σημειώστε ότι: $v_\pi(s) = \sum_a \pi(a \mid s) q_\pi(s, a) = \mathbb{E} [q_\pi(S_t, A_t) \mid S_t = s, \pi] , \forall s$





Εξισώσεις Bellman

Εξισώσεις προσδοκιών του Μπέλμαν

Δεδομένου ενός MDP, για κάθε πολιτική π , οι συναρτήσεις τιμών υπακούουν στις ακόλουθες εξισώσεις προσδοκίας:

$$v_{\pi}(s) = \sum_a \pi(s, a) \left[r(s, a) + \gamma \sum_{s'} p(s'|a, s) v_{\pi}(s') \right]$$
$$q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|a, s) \sum_{a' \in \mathcal{A}} \pi(a'|s') q_{\pi}(s', a')$$





Bellman Εξισώσεις

Εξισώσεις βελτιστοποίησης Bellman

Δεδομένου ενός ΔΑΜ, οι συναρτήσεις βέλτιστης τιμής υπακούουν στις ακόλουθες εξισώσεις προσδοκίας:

$$v^*(s) = \max_a \left[r(s, a) + \gamma \sum_{s'} p(s'|a, s) v^*(s') \right]$$
$$q^*(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|a, s) \max_{a' \in \mathcal{A}} q^*(s', a')$$





Επίλυση προβλημάτων EM χρησιμοποιώντας τις εξισώσεις Bellman





Προβλήματα στην ΕΜ

Δύο προβλήματα:

1. Η εκτίμηση v_{π} ή q_{π} ονομάζεται **αξιολόγηση πολιτικής** ή, απλά, **πρόβλεψη**
 - Λαμβάνοντας υπόψη μια πολιτική, ποια είναι η αναμενόμενη επιστροφή μου κάτω από αυτή τη συμπεριφορά;
 - Λαμβάνοντας υπόψη αυτό το πρωτόκολλο θεραπείας/τη στρατηγική συναλλαγών, ποια είναι η αναμενόμενη επιστροφή μου;
2. Η εκτίμηση v_* ή q_* ονομάζεται μερικές φορές **έλεγχος**, επειδή αυτά μπορούν να χρησιμοποιηθούν για **τη βελτιστοποίηση της πολιτικής**
 - Ποιος είναι ο βέλτιστος τρόπος συμπεριφοράς; Ποια είναι η βέλτιστη λειτουργία αξίας;
 - Ποια είναι η βέλτιστη θεραπεία; Ποια είναι η βέλτιστη πολιτική ελέγχου για την ελαχιστοποίηση του χρόνου, της κατανάλωσης καυσίμου κ.λπ.;

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Μια λύση

Εξίσωση Bellman σε μορφή πίνακα

- Η εξίσωση Bellman, για δεδομένη πολιτική π , μπορεί να εκφραστεί με τη χρήση πινάκων,

$$\mathbf{v} = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{v}$$

ΠΟΥ

$$v_i = v(s_i)$$

$$r_i^\pi = \mathbb{E} [R_{t+1} \mid S_t = s_i, A_t \sim \pi(S_t)]$$

$$P_{ij}^\pi = p(s_j \mid s_i) = \sum_a \pi(a \mid s_i) p(s_j \mid s_i, a)$$

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Μια λύση

Εξίσωση Bellman σε μορφή πίνακα

- Η εξίσωση Bellman, για δεδομένη πολιτική π , μπορεί να εκφραστεί με τη χρήση πινάκων,

$$\mathbf{v} = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{v}$$

- Αυτή είναι μια γραμμική εξίσωση που μπορεί να λυθεί ά

$$\mathbf{v} = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{v}$$

$$(\mathbf{I} - \gamma \mathbf{P}^\pi) \mathbf{v} = \mathbf{r}^\pi$$

$$\mathbf{v} = (\mathbf{I} - \gamma \mathbf{P}^\pi)^{-1} \mathbf{r}^\pi$$

- Η υπολογιστική πολυπλοκότητα είναι $O(|S|^3)$ - είναι δυνατή μόνο για μικρά προβλήματα
- Υπάρχουν **επαναληπτικές μέθοδοι** για μεγαλύτερα προβλήματα

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Μια λύση

Επίλυση της εξίσωσης της βελτιστοποίησης του Bellman

- Η εξίσωση βέλτιστης Bellman είναι μη γραμμική (λόγω του μέγιστου χειριστή)
- Δεν μπορεί να χρησιμοποιηθεί η ίδια άμεση λύση πίνακα όπως για την αξιολόγηση πολιτικής
- Πολλές επαναληπτικές μέθοδοι λύσης:
 - Χρήση μοντέλων/ **δυναμικού προγραμματισμού**
 - Επαναληπτική αξία
 - Επαναληπτική πολιτική
 - Χρήση δειγμάτων
 - Monte Carlo
 - Q-learning
 - SARSA

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Δυναμικός προγραμματισμός

Ο δυναμικός προγραμματισμός αναφέρεται σε μια συλλογή αλγορίθμων που μπορούν να χρησιμοποιηθούν για τον υπολογισμό βέλτιστων πολιτικών δεδομένου ενός τέλει μοντέλου του περιβάλλοντος ως μια διαδικασία λήψης αποφάσεων Μάρκοβ (ΔΑΜ).

Sutton & Barto 2018

- **Δυναμικός προγραμματισμός:**
 - αποσυνθέτει ένα πρόβλημα σε μικρότερα υποπρόβλήματα με αναδρομικό τρόπο
 - αποτελείται από δύο σημαντικά μέρη:

αξιολόγηση πολιτικής και βελτίωση της πολιτικής

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Αξιολόγηση πολιτικής

- Ξεκινάμε συζητώντας πώς να εκτιμήσουμε το $v_\pi(s) = \mathbb{E} [R_{t+1} + \gamma v_\pi(S_{t+1}) \mid s, \pi]$
- Ιδέα: μετατρέψτε αυτή την εξίσωση σε ενημέρωση

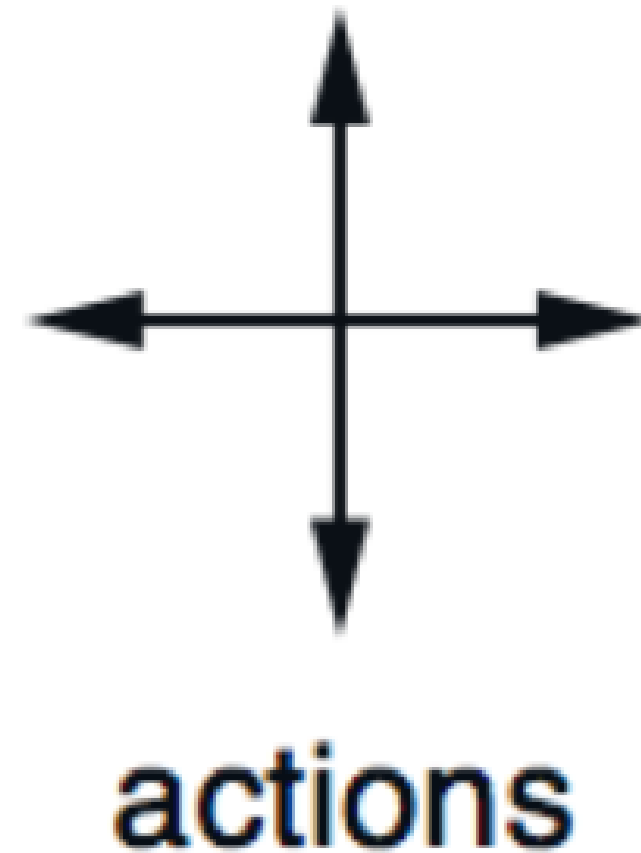
Αλγόριθμος

- Αρχικοποίηση v_0 , π.χ., έως μηδέν
- Να επαναληφθεί: $\forall s : v_{k+1}(s) \leftarrow \mathbb{E} [R_{t+1} + \gamma v_k(S_{t+1}) \mid s, \pi]$
- Σταματώντας: όποτε $v_{k+1}(s) = v_k(s)$, για όλα τα s , πρέπει να έχουμε βρει v_π





Για παράδειγμα: Αξιολόγηση πολιτικής



ΣΤΟΧΟΣ	1	2	3
4	5	6	7
8	9	10	11
12	13	14	ΣΤΟΧΟΣ

$R_t = -1$
on all transitions

$\pi = \text{τυχαία πολιτική}$
 $\gamma = 1,0$





Παράδειγμα: Αξιολόγηση πολιτικής

$$v_1(s_t) = -1 + E [v_0(s_{t+1})] = -1$$

$k = 0$

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

$$v_0 = 0$$

$$\begin{aligned} v_{k+1}(s_t) &= E [R_{t+1} + \gamma v_k(s_{t+1})] \\ &= E [-1 + v_k(s_{t+1})] \\ &= -1 + E [v_k(s_{t+1})] \end{aligned}$$

$k = 1$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

$$\begin{aligned} v_2(s_t) &= -1 + E [v_1(s_{t+1})] \\ &= -1 + (-1-1-1+0)/4 = -1.75 \end{aligned}$$

$k = 2$

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0





Παράδειγμα: Αξιολόγηση πολιτικής

$k = 3$

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0

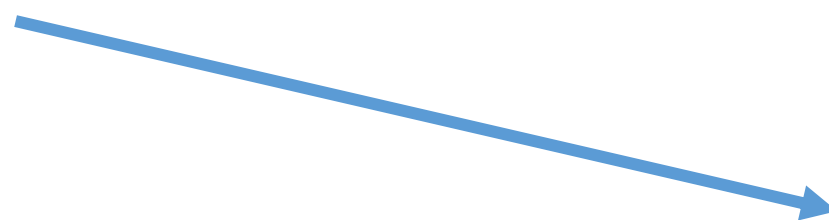
$k = 10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0

$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0

Συνάρτηση τιμής της τυχαίας πολιτικής



Η αξιολόγηση της πολιτικής μπορεί να απαιτήσει πολλά βήματα

$$V_{k+1}(s_t) = -1 + E [v_k(s_{t+1})]$$

$V_{k+1}(s_t) = v_k(s_{t+1})$ Όταν η συνάρτηση αξίας δεν αλλάζει, έχουμε ολοκληρώσει την αξιολόγηση πολιτικής

$$V_k = v_{\pi}$$





Βελτίωση της πολιτικής

- Μπορούμε να χρησιμοποιήσουμε την αξιολόγηση για να βελτιώσουμε την πολιτική μας
- Το να είμαστε άπληστοι σε σχέση με τις αξίες της τυχαίας πολιτικής μπορεί να είναι αρκετό (δεν ισχύει γενικά)

Αλγόριθμος

- Επαναλάβετε χρησιμοποιώντας:

$$\begin{aligned}\forall s : \pi_{\text{new}}(s) &= \underset{a}{\operatorname{argmax}} q_{\pi}(s, a) \\ &= \underset{a}{\operatorname{argmax}} \mathbb{E} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]\end{aligned}$$

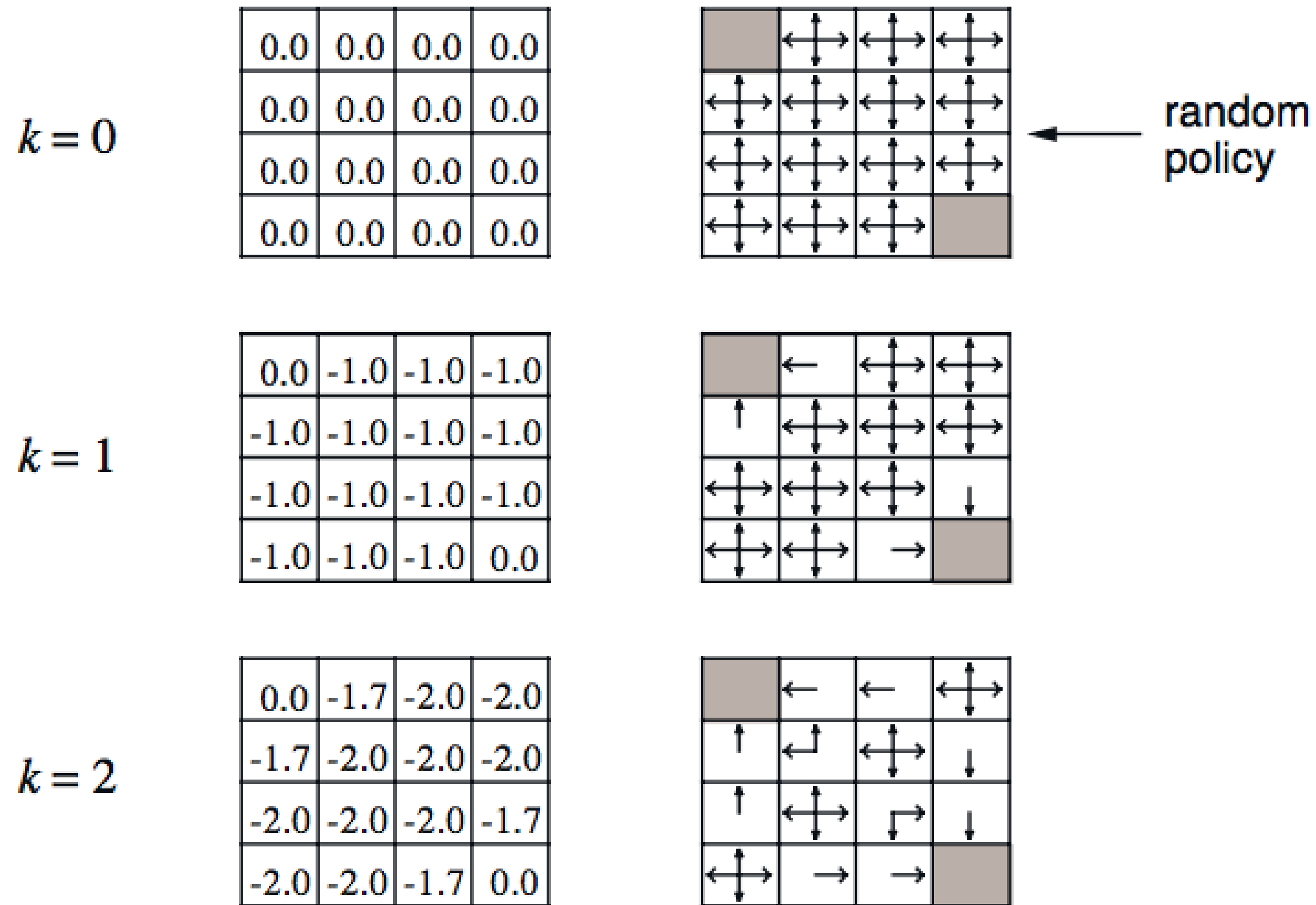
- Αξιολογήστε το νέο και επαναλάβετε

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Παράδειγμα: Αξιολόγηση Πολιτικής + Βελτίωση



Μπορούμε να πάρουμε κάποια λειτουργία αξίας και να εξάγουμε κάποια πολιτική χρησιμοποιώντας άπληστη βελτίωση

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Παράδειγμα: Αξιολόγηση Πολιτικής + Βελτίωση

Μπορούμε να πάρουμε κάποια λειτουργία αξίας και να εξάγουμε κάποια πολιτική χρησιμοποιώντας άπληστη βελτίωση

$k = 3$

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0

←	←	←	↙
↑	↖	↖	↓
↑	↗	↘	↓
↖	→	→	→

$k = 10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0

←	←	←	↙
↑	↖	↖	↓
↑	↗	↘	↓
↖	→	→	→

$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0

←	←	←	↙
↑	↖	↖	↓
↑	↗	↘	↓
↖	→	→	→

optimal policy

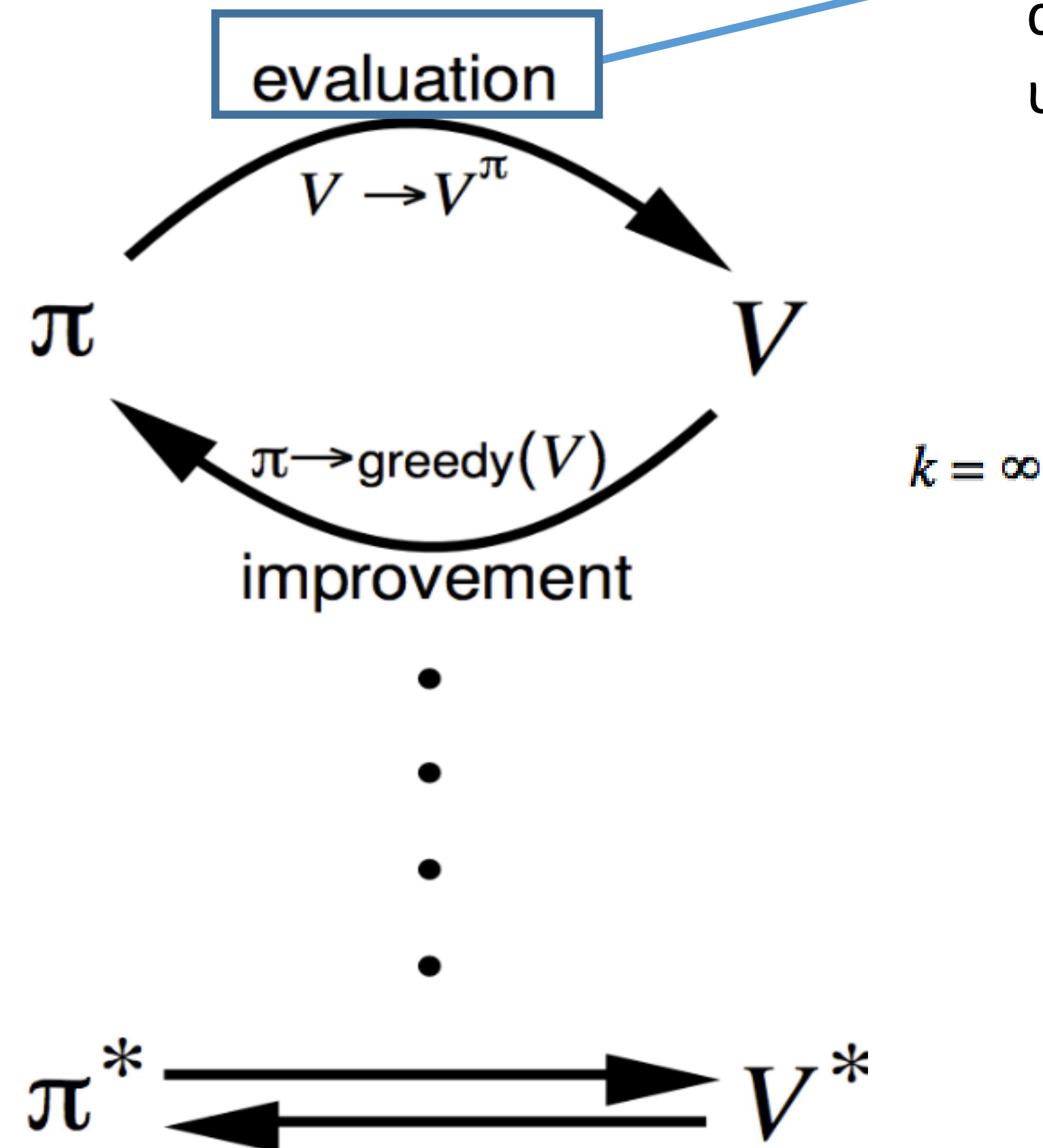
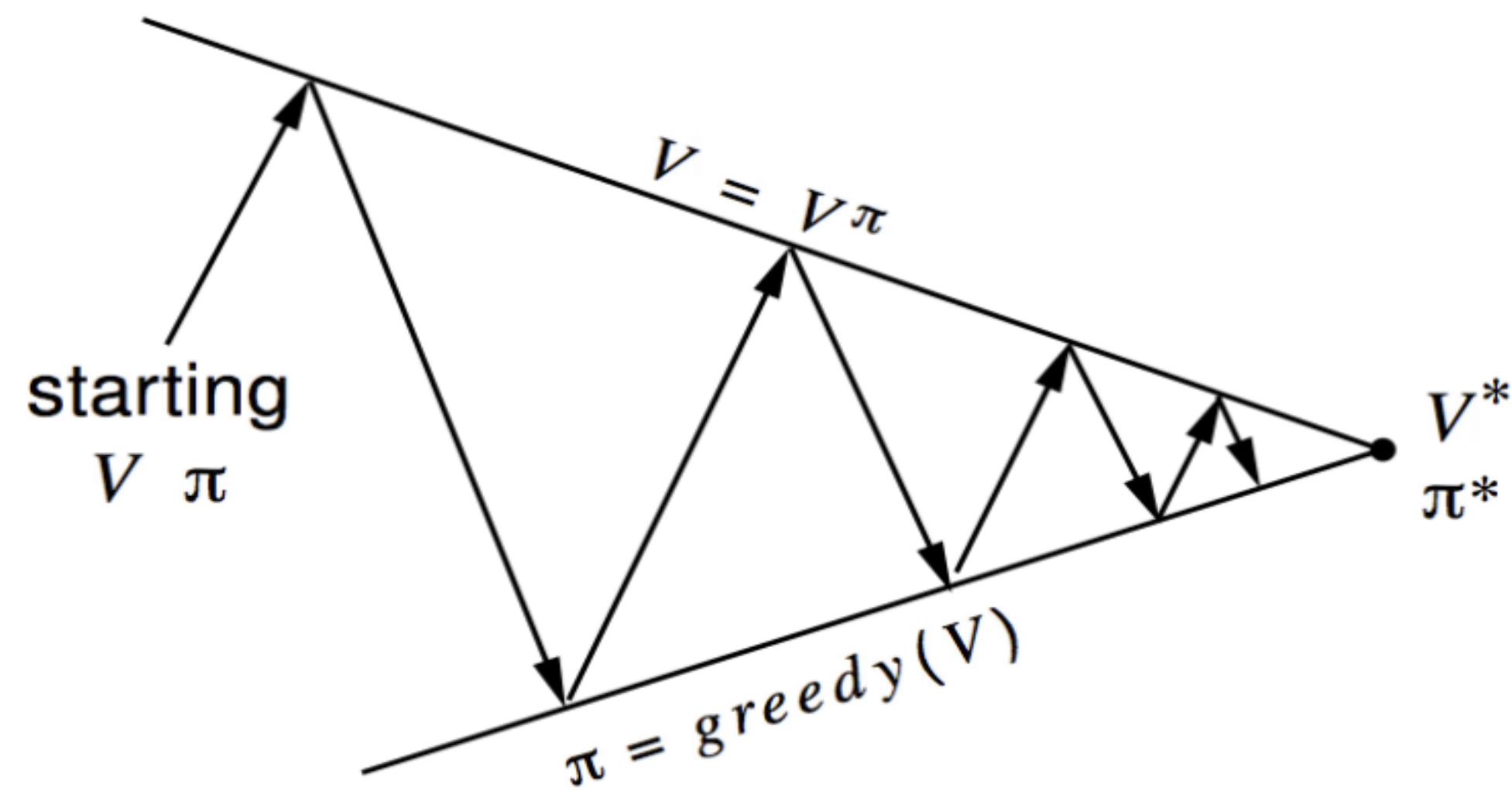
Παρατηρήστε ότι η συνάρτηση αξίας εξακολουθεί να αλλάζει, αλλά η βέλτιστη πολιτική έχει διδαχθεί από το $k=3$





Επαναληπτική μέθοδος Πολιτικής

Επαναληπτικό, δηλαδή, μπορεί να απαιτήσει πολλούς υπολογισμούς



Αξιολόγηση πολιτικής

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0

Βελτίωση της πολιτικής

	←	←	↖
↑	↖	↖	↓
↑	↗	↗	↓
↖	→	→	

Policy evaluation Estimate v^π
 Policy improvement Generate $\pi' \geq \pi$

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Επαναληπτική μέθοδος Πολιτικής

- Χρειάζεται η αξιολόγηση πολιτικής να συγκλίνει προς το π ;
- Ή μήπως πρέπει να σταματήσουμε όταν είμαστε «κοντά»; (Π.χ., με κατώτατο όριο για τη μεταβολή των τιμών)
 - Ή απλά να σταματήσουμε μετά τις επαναλήψεις της επαναληπτικής αξιολόγησης πολιτικής;
 - Στο μικρό grid world $k = 3$ ήταν αρκετό για την επίτευξη βέλτιστης πολιτικής
- **Ακραίο:** Γιατί να μην ενημερώσετε την πολιτική κάθε επανάληψης - δηλαδή σταματήστε μετά το $k = 1$;
 - Αυτό είναι ισοδύναμο με την **επανάληψη αξίας**





Επαναληπτική μέθοδος αξίας

- Θα μπορούσαμε να πάρουμε την εξίσωση βελτιστοποίησης Μπέλμαν, και να την μετατρέψουμε σε ενημέρωση.

$$\forall s : v_{k+1}(s) \leftarrow \max_a \mathbb{E} [R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = s]$$

- Αυτό ισοδυναμεί με επανάληψη πολιτικής, με $k = 1$ βήμα αξιολόγησης πολιτικής μεταξύ κάθε δύο (άπληστων) σταδίων βελτίωσης της πολιτικής

Αλγόριθμος: Εξατομίκευση αξίας

- Αρχικοποίησε v_0
- Επανάλαβε: $v_{k+1}(s) \leftarrow \max_a \mathbb{E} [R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = s]$
- Σταμάτησε: όποτε $v_{k+1}(s) = v_k(s)$, για όλα τα s , πρέπει να έχουμε βρει v^*

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Επαναληπτική μέθοδος αξίας στο παράδειγμα του αγώνα

Αλγόριθμος

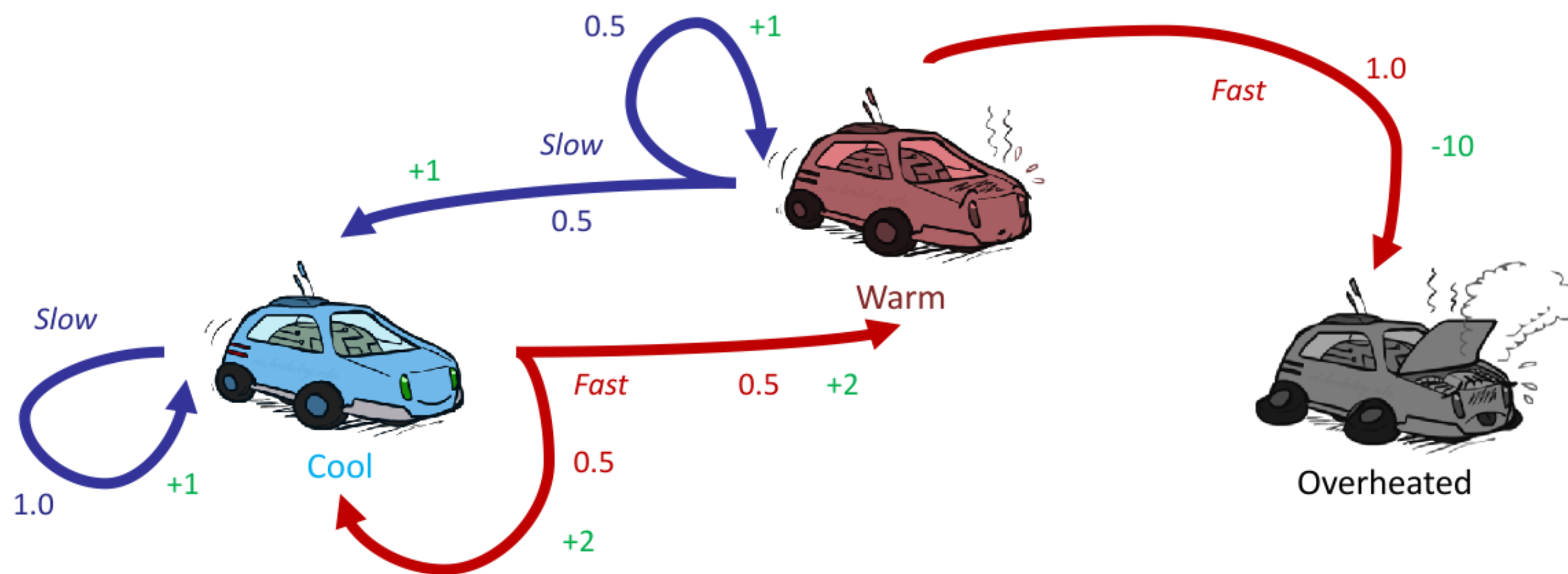
1. $\forall s \in S$, initialize $V_0(s) = 0$

2. Επαναλάβετε μέχρι τη σύγκλιση:

$$\forall s \in S, V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

Η Σύγκλιση:

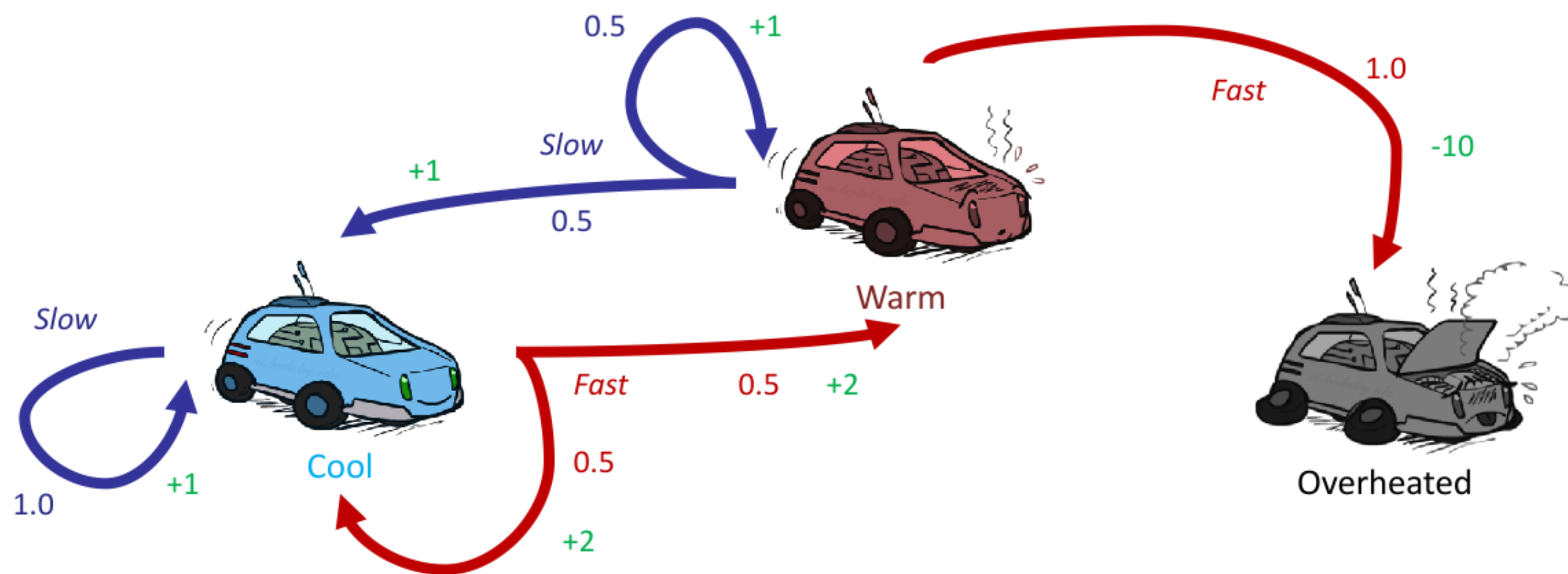
$$\forall s, V_{k+1}(s) = \bar{V}_k(s)$$





Επαναληπτική μέθοδος αξίας στο παράδειγμα του αγώνα

$$\forall s \in S, V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$



	cool	warm	overheated
V_0	0	0	0

$$\begin{aligned} V_1(\text{cool}) &= \max\{1 \cdot [1 + 0.5 \cdot 0], 0.5 \cdot [2 + 0.5 \cdot 0] + 0.5 \cdot [2 + 0.5 \cdot 0]\} \\ &= \max\{1, 2\} \\ &= \boxed{2} \end{aligned}$$

$$\begin{aligned} V_1(\text{warm}) &= \max\{0.5 \cdot [1 + 0.5 \cdot 0] + 0.5 \cdot [1 + 0.5 \cdot 0], 1 \cdot [-10 + 0.5 \cdot 0]\} \\ &= \max\{1, -10\} \\ &= \boxed{1} \end{aligned}$$

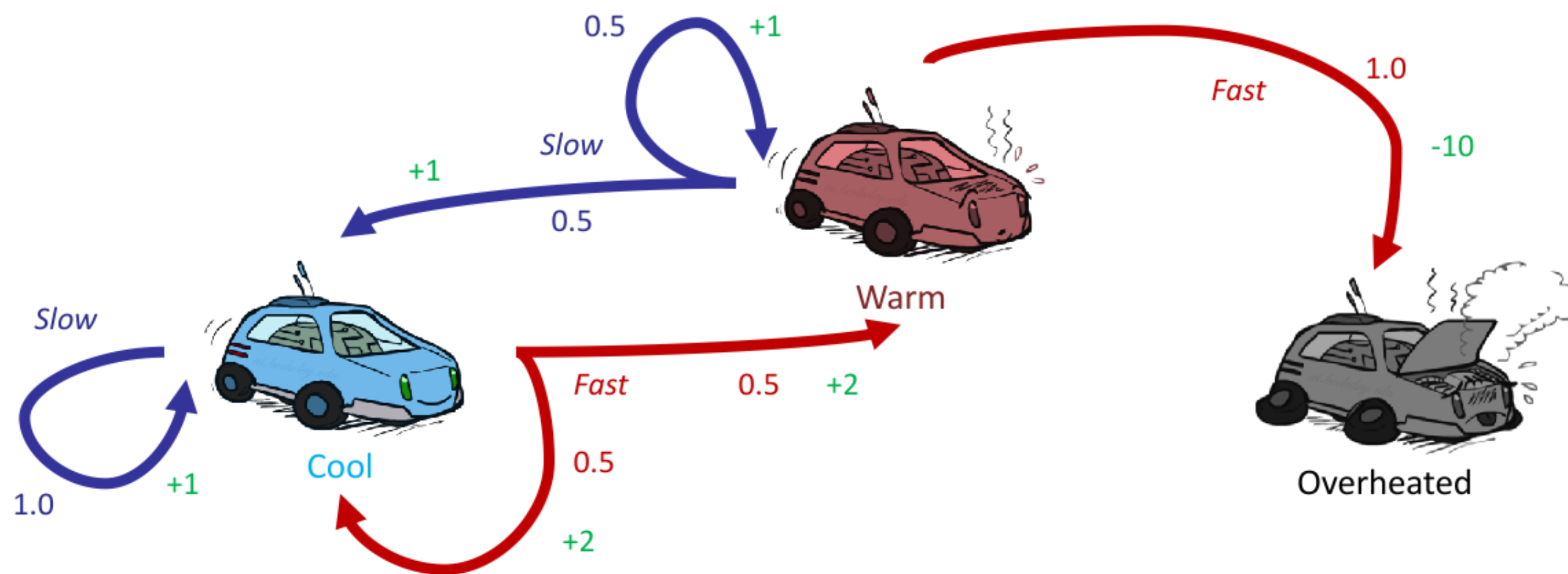
$$\begin{aligned} V_1(\text{overheated}) &= \max\{\} \\ &= \boxed{0} \end{aligned}$$





Επαναληπτική μέθοδος αξίας στο παράδειγμα του αγώνα

$$\forall s \in S, V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$



	cool	warm	overheated
V_0	0	0	0
V_1	2	1	0

$$\begin{aligned} V_2(\text{cool}) &= \max\{1 \cdot [1 + 0.5 \cdot 2], 0.5 \cdot [2 + 0.5 \cdot 2] + 0.5 \cdot [2 + 0.5 \cdot 1]\} \\ &= \max\{2, 2.75\} \\ &= \boxed{2.75} \end{aligned}$$

$$\begin{aligned} V_2(\text{warm}) &= \max\{0.5 \cdot [1 + 0.5 \cdot 2] + 0.5 \cdot [1 + 0.5 \cdot 1], 1 \cdot [-10 + 0.5 \cdot 0]\} \\ &= \max\{1.75, -10\} \\ &= \boxed{1.75} \end{aligned}$$

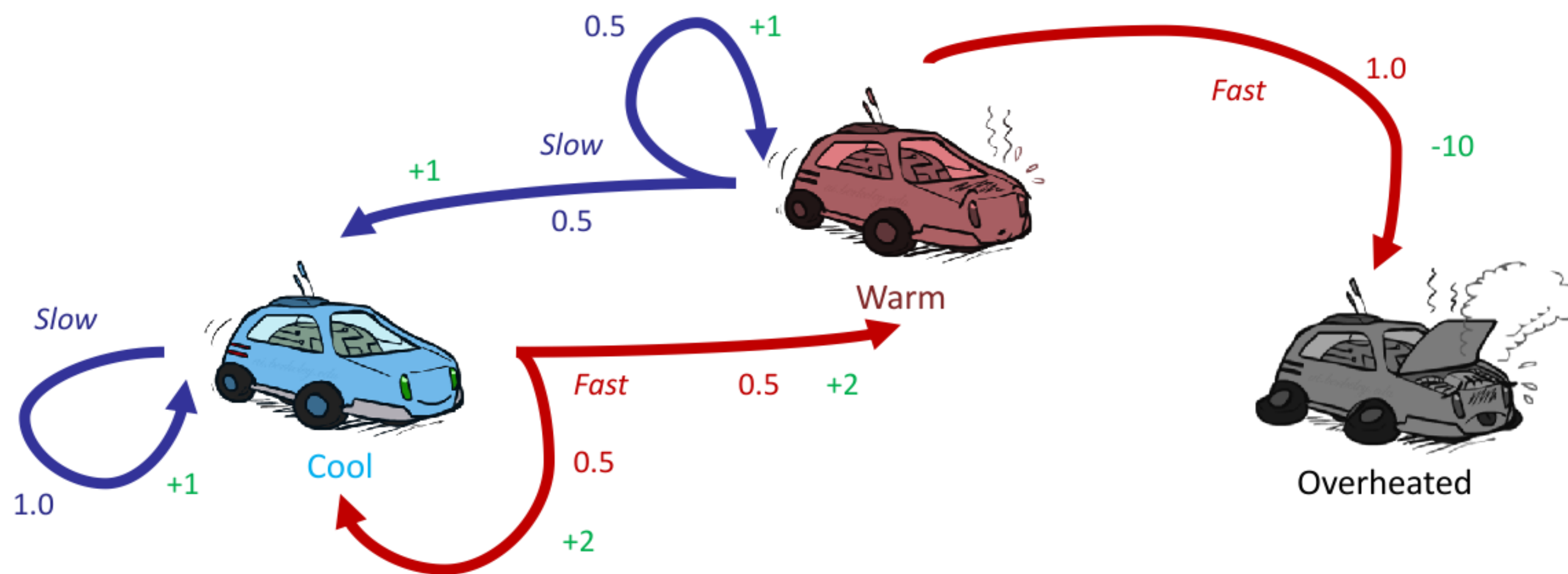
$$\begin{aligned} V_2(\text{overheated}) &= \max\{\} \\ &= \boxed{0} \end{aligned}$$





Επαναληπτική μέθοδος αξίας στο παράδειγμα του αγώνα

$$\forall s \in S, V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$



	cool	warm	overheated
V_0	0	0	0
V_1	2	1	0
V_2	2.75	1.75	0

Χρειάζονται περισσότερες επαναλήψεις για να συγκλίνουν στη βέλτιστη συνάρτηση τιμής

Διαφάνειες βασισμένες στο: [UC Berkeley CS188 – Intro to AI course](#)





Επαναληπτική μέθοδος πολιτικής Vs αξίας

- Επαναληπτική μέθοδος πολιτικής :
 - Ξεκινά με τυχαία πολιτική
 - Διενεργεί αξιολόγηση πολιτικής μέχρι τη σύγκλιση και, στη συνέχεια, ένα βήμα βελτίωσης της πολιτικής
 - Απαιτούνται λίγες επαναλήψεις για να συγκλίνουν, αλλά σε βάρος του περισσότερο χρόνου υπολογισμού λόγω της αξιολόγησης της πολιτικής
- Επαναληπτική μέθοδος αξίας:
 - Ξεκινά με μια τυχαία (ή μηδενική) συνάρτηση τιμής
 - Ακολουθεί ένα βήμα αξιολόγησης πολιτικής, ακολουθούμενο από βελτίωση της πολιτικής
 - Απαιτεί περισσότερες επαναλήψεις για να συγκλίνει στη βέλτιστη συνάρτηση τιμής





Τι έχουμε καλύψει

- Διαδικασίες λήψης αποφάσεων Markov
- Στόχοι σε ένα MDP: διαφορετική έννοια της επιστροφής
- Συναρτήσεις τιμής — αναμενόμενες αποδόσεις, συνθήκη για την κατάσταση (και δράση)
- Αρχές βελτιστοποίησης σε MDPs: βέλτιστες λειτουργίες αξίας και βέλτιστες πολιτικές
- Bellman Εξισώσεις
- Δύο κατηγορίες προβλημάτων στην RL: αξιολόγηση και έλεγχος
- Πώς να υπολογίσετε v_{π} (άλλως λύστε ένα πρόβλημα αξιολόγησης/προβλεψιμότητας)
- Πώς να υπολογίσετε τη βέλτιστη συνάρτηση τιμής μέσω δυναμικού προγραμματισμού
 - Επαναληπτική μέθοδος πολιτικής
 - Επαναληπτική μέθοδος αξίας

Διαφάνειες βασισμένες στις [Διαλέξεις Ενισχυτικής μάθησης της DeepMind](#)





Επόμενες Διαλέξεις

- Ενίσχυμένη Μάθηση χωρίς μοντέλο
- Επανάληψη



MAI4CAREU

Master programmes in Artificial
Intelligence 4 Careers in Europe



Σας ευχαριστούμε

