
Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation

Tejas D. Kulkarni*
DeepMind, London
tejasdkulkarni@gmail.com

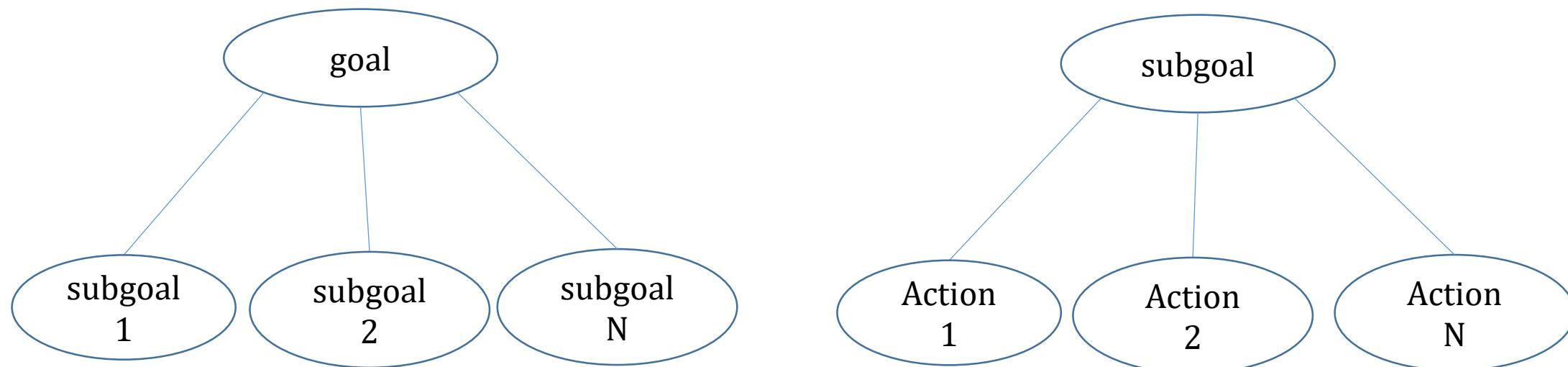
Karthik R. Narasimhan*
CSAIL, MIT
karthikn@mit.edu

Ardavan Saeedi
CSAIL, MIT
ardavans@mit.edu

Joshua B. Tenenbaum
BCS, MIT
jbt@mit.edu

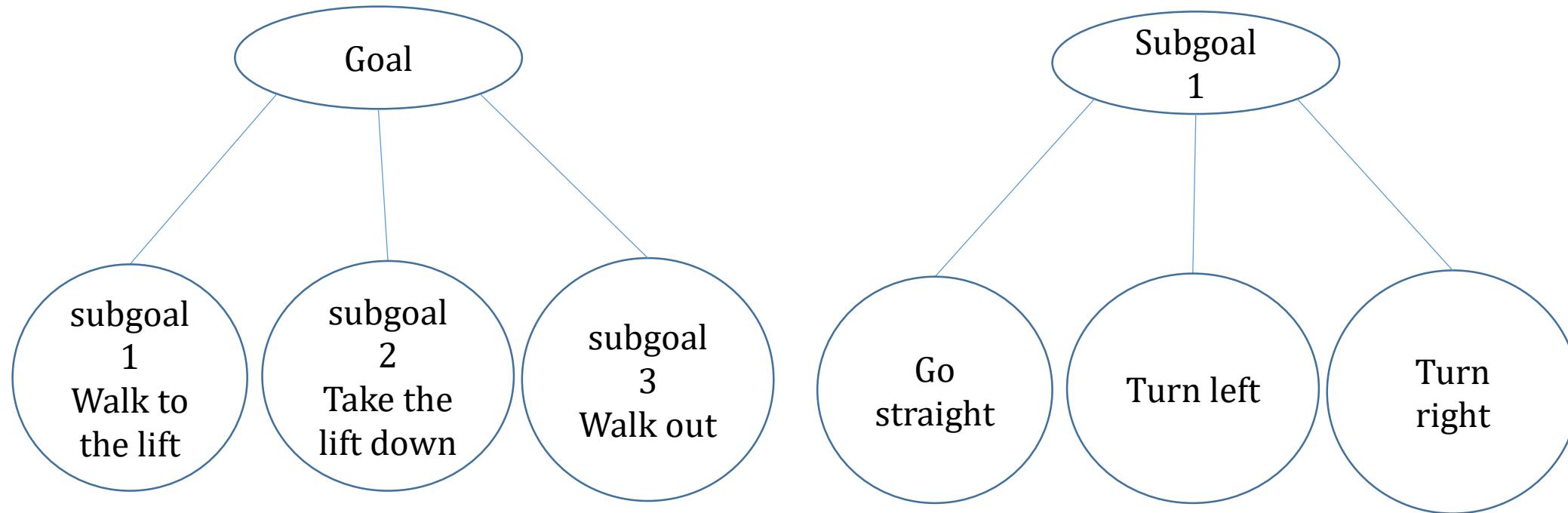
NIPS-2016

Hierarchical RL

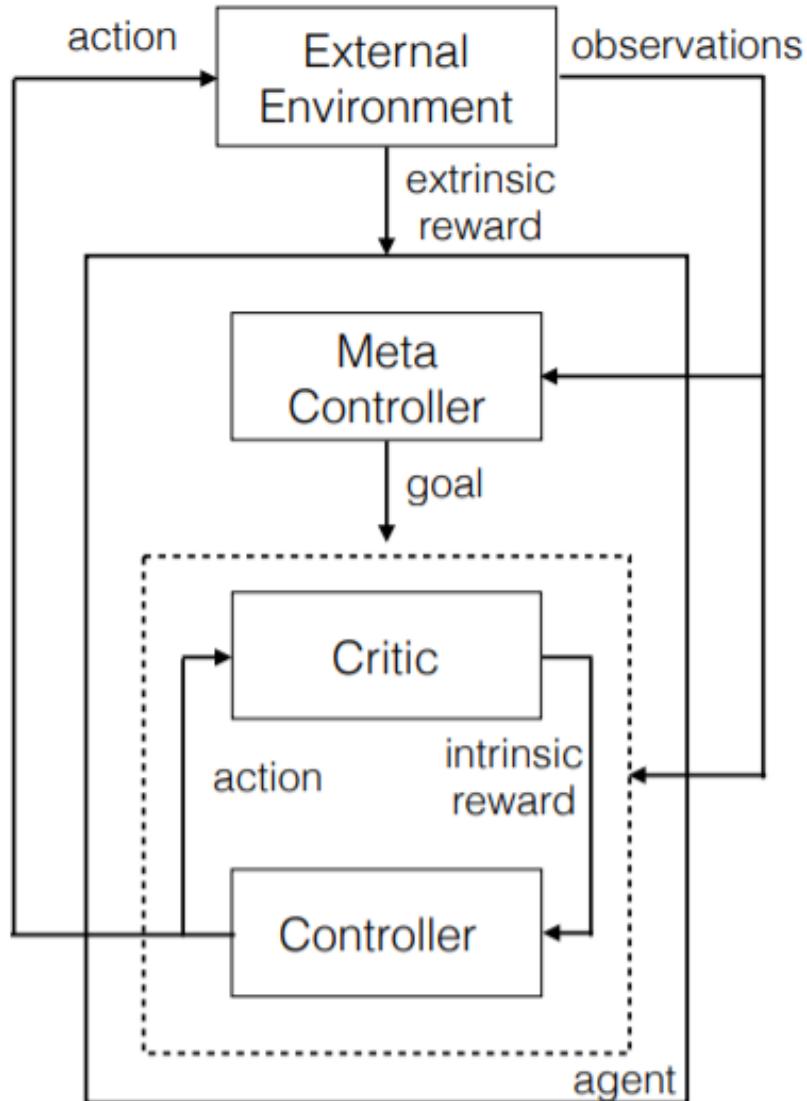


Example

- Goal: walk out of a building



Framework



- Meta controller: choose subgoals
- Controller: achieve subgoals

Maximize two total reward:

$$R_t(g) = \sum_{t'}^{\infty} \gamma^{t'-t} r_{t'}(g),$$
$$r_{t'}(g) = 1 \text{ if goal is reached and 0 otherwise}$$

$$F_t = \sum_{t'}^{\infty} \gamma^{t'-t} f_{t'}$$

Compute $f_{t'}$ using external reward

Algorithm 1 Learning algorithm for h-DQN

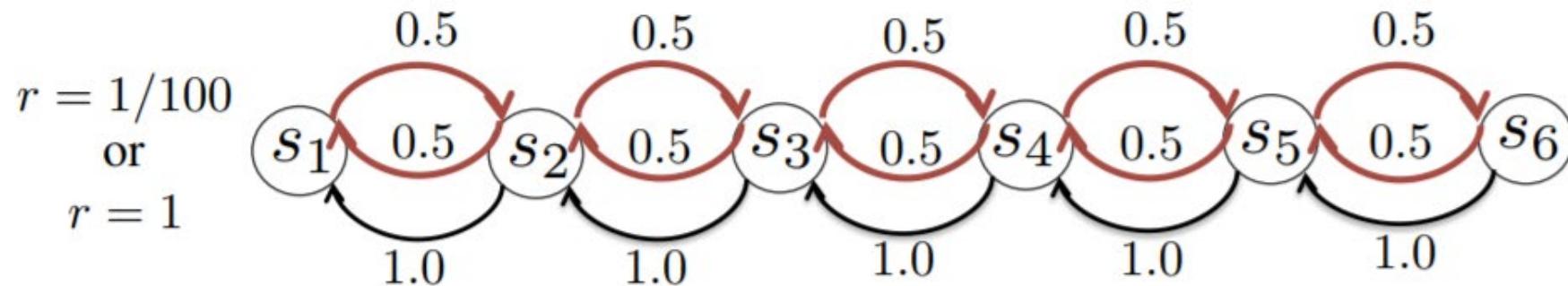
```
1: Initialize experience replay memories  $\{\mathcal{D}_1, \mathcal{D}_2\}$  and parameters  $\{\theta_1, \theta_2\}$  for the controller and meta-controller respectively.  
2: Initialize exploration probability  $\epsilon_{1,g} = 1$  for the controller for all goals  $g$  and  $\epsilon_2 = 1$  for the meta-controller.  
3: for  $i = 1, num\_episodes$  do  
4:   Initialize game and get start state description  $s$   
5:    $g \leftarrow \text{EPSGREEDY}(s, \mathcal{G}, \epsilon_2, Q_2)$  —————→ Choose subgoal  
6:   while  $s$  is not terminal do  
7:      $F \leftarrow 0$   
8:      $s_0 \leftarrow s$   
9:     while not ( $s$  is terminal or goal  $g$  reached) do  
10:       $a \leftarrow \text{EPSGREEDY}(\{s, g\}, \mathcal{A}, \epsilon_{1,g}, Q_1)$   
11:      Execute  $a$  and obtain next state  $s'$  and extrinsic reward  $f$  from environment  
12:      Obtain intrinsic reward  $r(s, a, s')$  from internal critic  
13:      Store transition  $(\{s, g\}, a, r, \{s', g\})$  in  $\mathcal{D}_1$   
14:      UPDATEPARAMS( $\mathcal{L}_1(\theta_{1,i}), \mathcal{D}_1$ )  
15:      UPDATEPARAMS( $\mathcal{L}_2(\theta_{2,i}), \mathcal{D}_2$ )  
16:       $F \leftarrow F + f$   
17:       $s \leftarrow s'$   
18:    end while  
19:    Store transition  $(s_0, g, F, s')$  in  $\mathcal{D}_2$   
20:    if  $s$  is not terminal then  
21:       $g \leftarrow \text{EPSGREEDY}(s, \mathcal{G}, \epsilon_2, Q_2)$   
22:    end if  
23:  end while  
24:  Anneal  $\epsilon_2$  and  $\epsilon_1$ .  
25: end for
```

Choose subgoal

Finish subgoal

update network

Experiment



Initial state: s_2 , terminal state: s_1

Action: left, right

Reward: 1 if s_6 is visited or $1/100$ if s_6 not visited.

Subgoal: s_4, s_5, s_6

Experiment

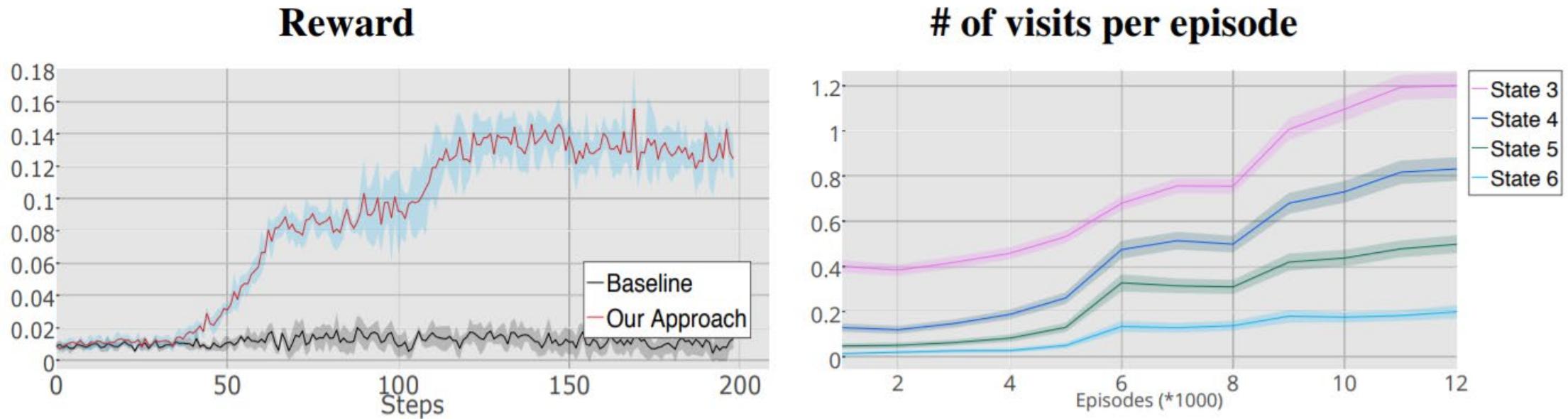
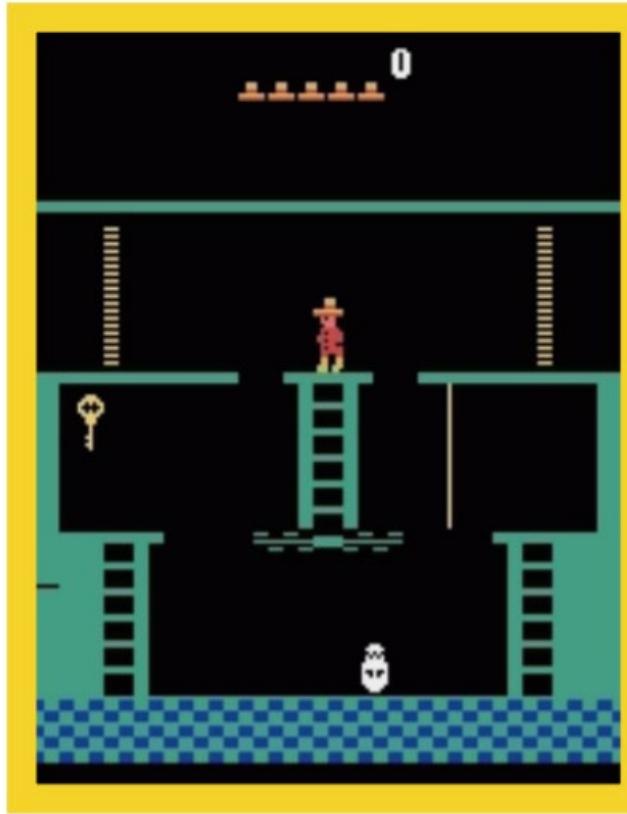


Figure 3: **(left)** Average reward (over 10 runs) of our approach compared to Q-learning. **(right)** #visits of our approach to states s_3-s_6 (over 1000 episodes). Initial state: s_2 , Terminal state: s_1 .

Experiment

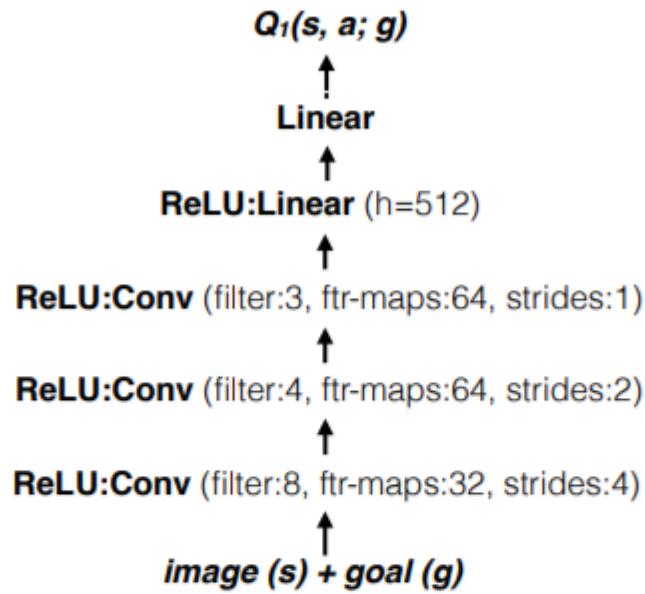


Goal: find the key and use it to open the door

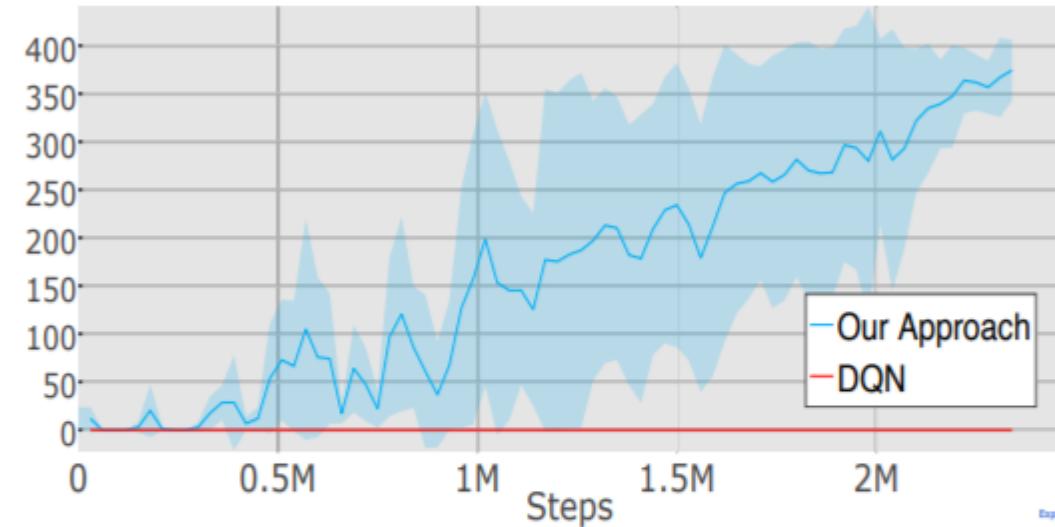
Reward: 100 if key is found and another 300 if door is opened

Subgoals: top-left door; top-right door; middle-ladder; bottom-left ladder; bottom-right ladder; key

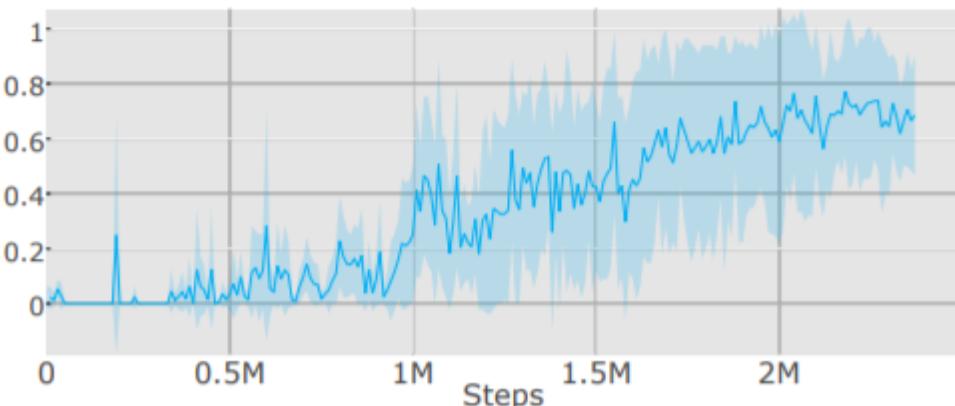
Architecture



Total extrinsic reward



Success ratio for reaching the goal 'key'



Success % of different goals over time

