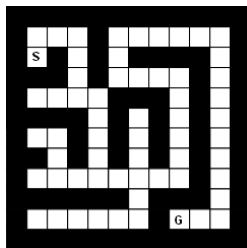


# Potential Based Reward Shaping Tutorial

ALA 2014  
Sam Devlin

# Knowledge-Based Reinforcement Learning

- ▶ Commonly, RL algorithms assume no prior knowledge
- ▶ Including domain knowledge can simplify learning



# Reward Shaping

## Q-Learning

- ▶ A popular RL algorithm
- ▶  $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q_i(s', a') - Q(s, a)]$

## Reward Shaping

- ▶ Provide heuristic knowledge by an additional reward
- ▶  $Q(s, a) \leftarrow Q(s, a) + \alpha[r + F(s, s') + \gamma \max_{a'} Q_i(s', a') - Q(s, a)]$

## Potential-Based Reward Shaping

$$F(s, s') = \gamma \Phi(s') - \Phi(s)$$

- ▶  $F(s, s')$  : Additional reward when moving from state  $s$  to  $s'$
- ▶  $\gamma$  : Discount factor
- ▶  $\Phi(s)$  : Potential of state  $s$

# Potential-Based Reward Shaping

## Formal Definition

$$\blacktriangleright F(s, s') = \gamma \Phi(s') - \Phi(s)$$

## Guarantees

- ▶ Policy invariance (optimal policy unchanged) in single agent<sup>1</sup>

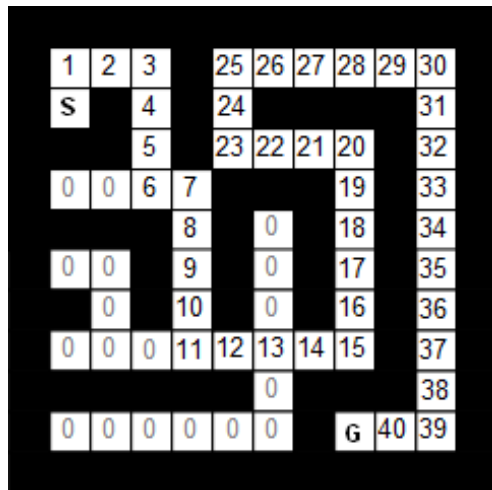
## Can

- ▶ Increase/Decrease time taken to learn optimal policy

---

<sup>1</sup> Ng, Russell and Harada. "Policy Invariance Under Reward Transformations: Theory And Application To Reward Shaping." ICML, 1999.

## An Example Potential Function



## Potential-Based Reward Shaping

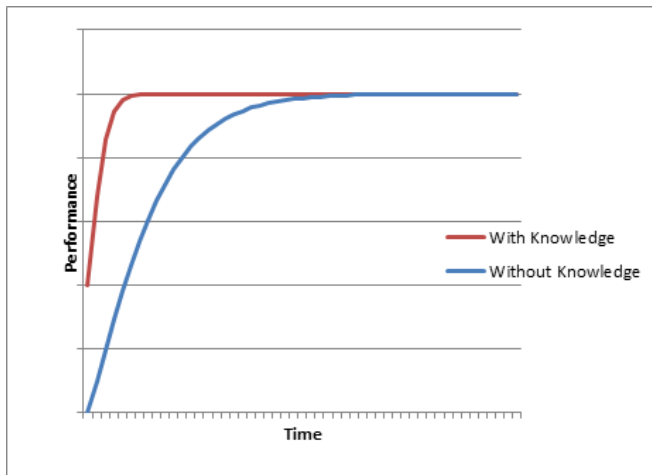


Figure : A Typical Single Agent Result

## Proof of Policy Invariance <sup>2</sup>

$$\begin{aligned}U_{\Phi}(\bar{s}) &= \sum_{j=0}^{\infty} \gamma^j (r_j + \gamma \Phi(s_{j+1}) - \Phi(s_j)) \\&= \sum_{j=0}^{\infty} \gamma^j r_j + \sum_{j=0}^{\infty} \gamma^{j+1} \Phi(s_{j+1}) - \sum_{j=0}^{\infty} \gamma^j \Phi(s_j) \\&= U(\bar{s}) + \sum_{j=1}^{\infty} \gamma^j \Phi(s_j) - \Phi(s_0) - \sum_{j=1}^{\infty} \gamma^j \Phi(s_j) \\&= U(\bar{s}) - \Phi(s_0)\end{aligned}$$

---

<sup>2</sup> Asmuth, Littman and Zinkov. "Potential-based shaping in model-based reinforcement learning." AAAI, 2008.



## Proof of Policy Invariance

$$U_{\Phi}(\bar{s}) = \sum_{j=0}^{\infty} \gamma^j (r_j + \gamma \Phi(s_{j+1}) - \Phi(s_j))$$

## Proof of Policy Invariance

$$\begin{aligned}U_{\Phi}(\bar{s}) &= \sum_{j=0}^{\infty} \gamma^j (r_j + \gamma \Phi(s_{j+1}) - \Phi(s_j)) \\&= \sum_{j=0}^{\infty} \gamma^j r_j + \sum_{j=0}^{\infty} \gamma^{j+1} \Phi(s_{j+1}) - \sum_{j=0}^{\infty} \gamma^j \Phi(s_j)\end{aligned}$$

## Proof of Policy Invariance

$$\begin{aligned}U_{\Phi}(\bar{s}) &= \sum_{j=0}^{\infty} \gamma^j (r_j + \gamma \Phi(s_{j+1}) - \Phi(s_j)) \\&= \sum_{j=0}^{\infty} \gamma^j r_j + \sum_{j=0}^{\infty} \gamma^{j+1} \Phi(s_{j+1}) - \sum_{j=0}^{\infty} \gamma^j \Phi(s_j) \\&= U(\bar{s}) + \sum_{j=1}^{\infty} \gamma^j \Phi(s_j) - \Phi(s_0) - \sum_{j=1}^{\infty} \gamma^j \Phi(s_j)\end{aligned}$$

## Proof of Policy Invariance

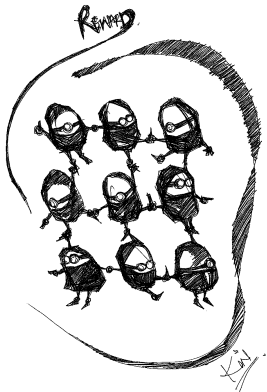
$$\begin{aligned}U_{\Phi}(\bar{s}) &= \sum_{j=0}^{\infty} \gamma^j (r_j + \gamma \Phi(s_{j+1}) - \Phi(s_j)) \\&= \sum_{j=0}^{\infty} \gamma^j r_j + \sum_{j=0}^{\infty} \gamma^{j+1} \Phi(s_{j+1}) - \sum_{j=0}^{\infty} \gamma^j \Phi(s_j) \\&= U(\bar{s}) + \sum_{j=1}^{\infty} \gamma^j \Phi(s_j) - \Phi(s_0) - \sum_{j=1}^{\infty} \gamma^j \Phi(s_j) \\&= U(\bar{s}) - \Phi(s_0)\end{aligned}$$

## Q-Table Initialization

- ▶ Wiewiora: “Potential-based shaping and Q-value initialization are equivalent.” (JAIR, 2003)
- ▶ ...If the potential function is **static**.

## Multi-Agent Reinforcement Learning

- ▶ Multiple agents learning concurrently in the same environment
- ▶ Typically learn a Nash equilibrium
- ▶ No clear notion of an optimal policy



# Multi-Agent Potential-Based Reward Shaping

## Guarantees

- ▶ Nash Equilibria not altered <sup>3</sup>

## Can

- ▶ Increase/Decrease time taken to reach a stable joint policy
- ▶ Change final joint policy

---

<sup>3</sup> Devlin and Kudenko, "Theoretical Considerations Of Potential-Based Reward Shaping For Multi-Agent Systems", AAMAS, 2011.

## Potential-Based Reward Shaping

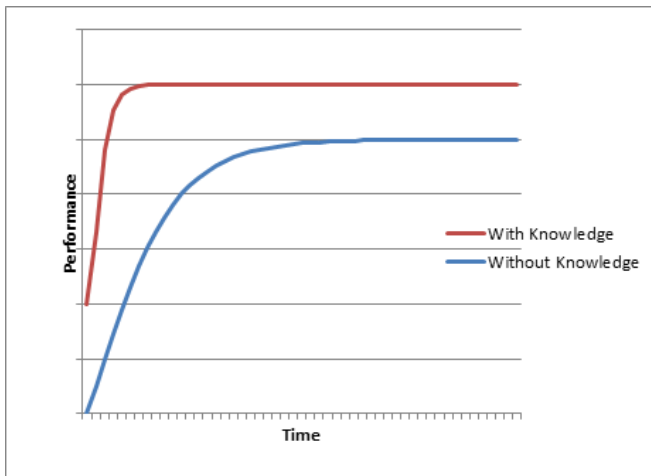
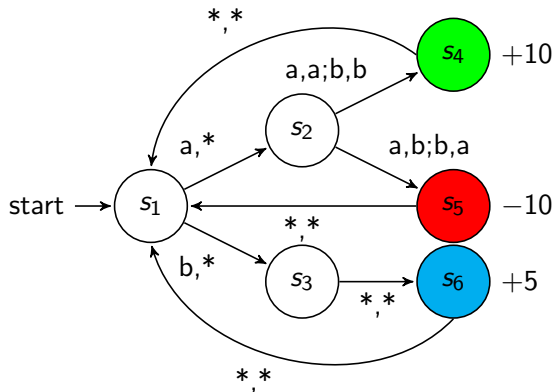


Figure : A Typical Multi-Agent Result



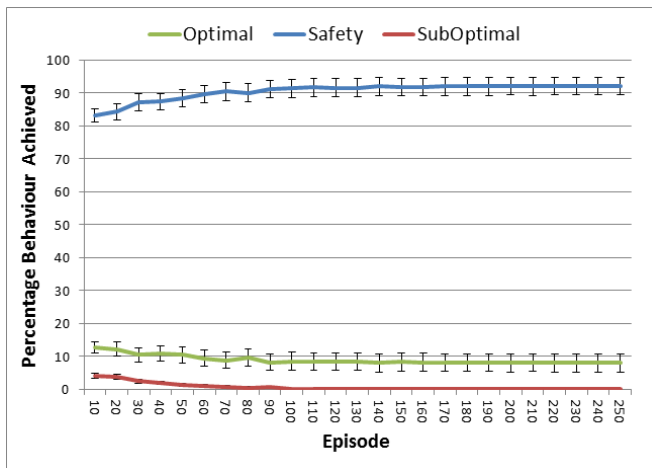
## Exploration Altered

- ▶ Reward shaping alters exploration / which actions are chosen
- ▶ In single-agent, this affects time to convergence
- ▶ In multi-agent, this may cause the agents to reach a different point of equilibrium
  - ▶ Wellman and Hu (1998) showed the joint policy converged upon in a learning MAS is highly sensitive to initial belief

Multi-Agent Example <sup>4</sup>

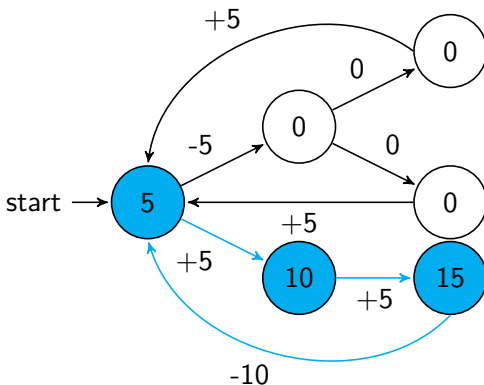
<sup>4</sup> Boutillier, "Sequential Optimality And Coordination In Multiagent Systems", IJCAI, 1999.

# Results

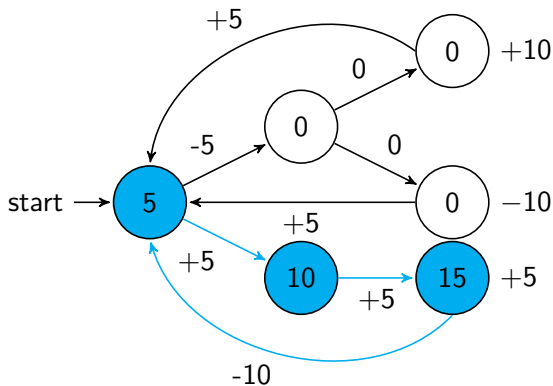


(a) Without Reward Shaping

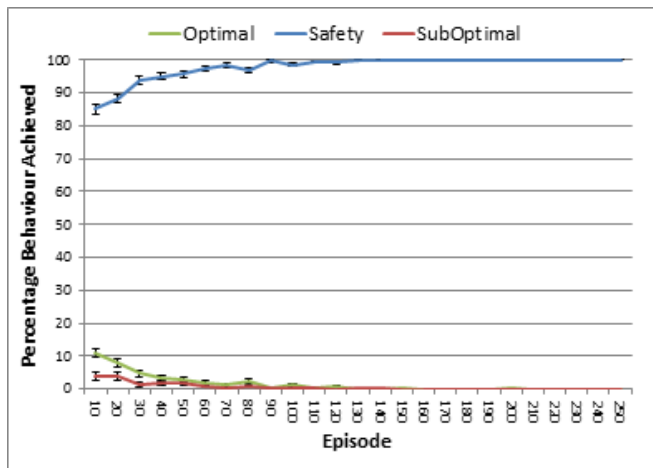
## Safe Reward Shaping



## Safe Reward Shaping

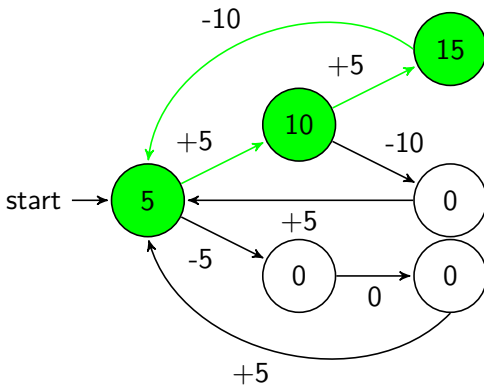


# Results

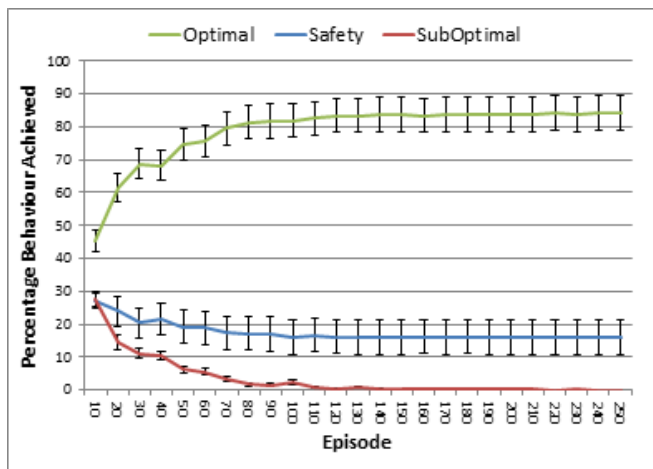


(b) With Safe Reward Shaping

## Coordinated Reward Shaping



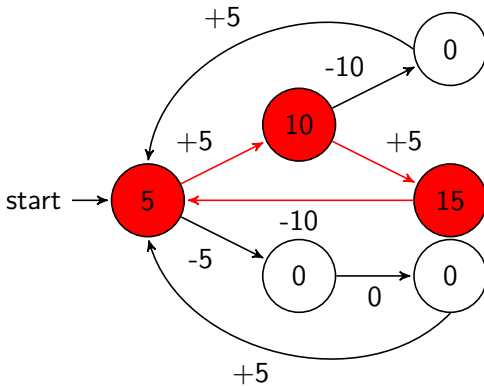
# Results



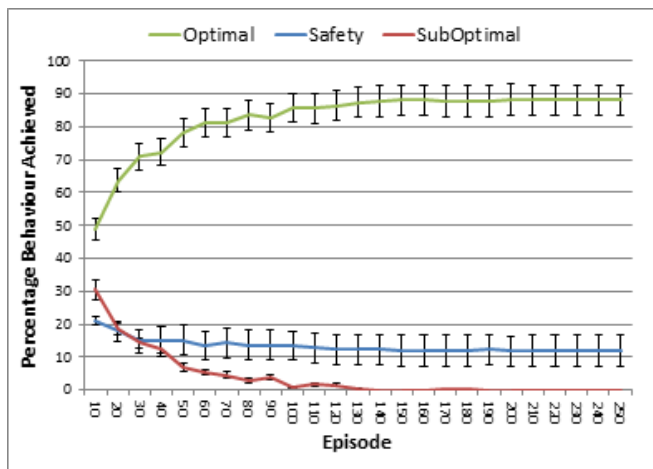
(c) With Coordinated Reward Shaping



## Miscoordinated Reward Shaping

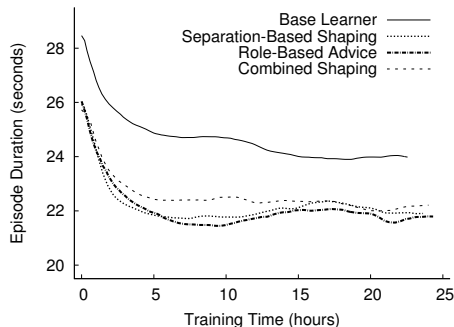
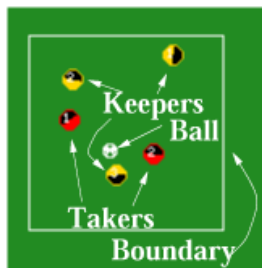


# Results



(d) With Miscoordinated Reward Shaping

## Multiagent Example 2: RoboCup KeepAway<sup>5</sup>



<sup>5</sup> Devlin, Grzes and Kudenko. "An Empirical Study Of Potential-Based Reward Shaping And Advice In Complex, Multi-Agent Systems." ACS, 2011.

## Practical Exercise: Design a Potential Function

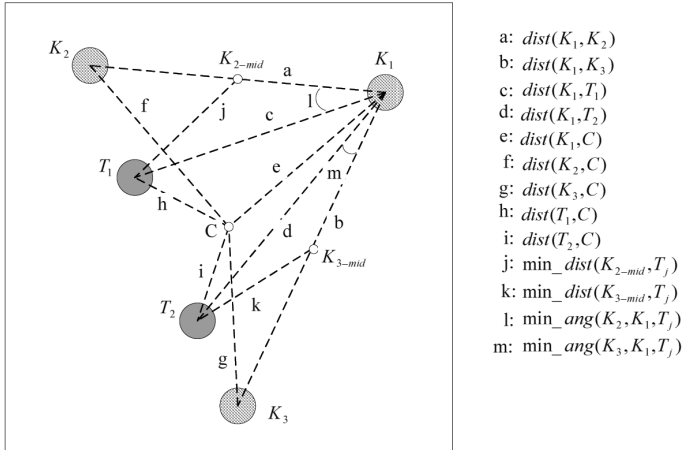


Figure : State Representation

## Past Assumptions

- ▶ Previous theoretical guarantees assume a **static** potential function
- ▶ Some claim the potential function must converge before the agent can <sup>6</sup>

---

<sup>6</sup> Laud, “Theory And Application Of Reward Shaping In Reinforcement Learning”, PhD Thesis, 2004.

## Dynamic Potential Based Reward Shaping <sup>7</sup>

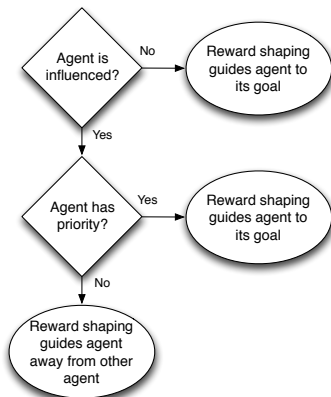
- Guarantees policy invariance or consistent Nash equilibria, provided:

$$F(s, t, s', t') = \gamma \Phi(s', t') - \Phi(s, t)$$

---

<sup>7</sup> Devlin and Kudenko. "Dynamic Potential-Based Reward Shaping."  
AAMAS, 2012.

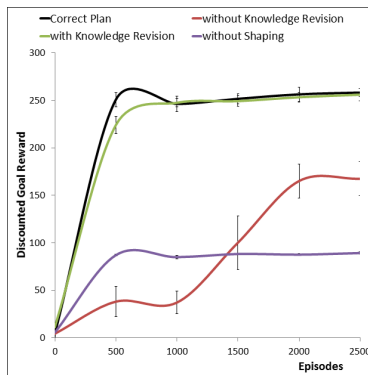
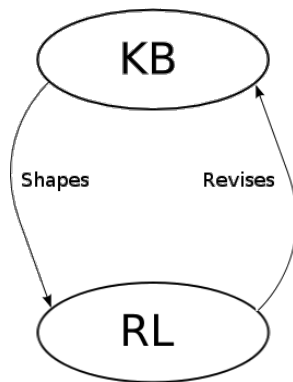
## Context Sensitive Reward Shaping<sup>8</sup>



- In different contexts we often recommend different behaviours

<sup>8</sup> De Hauwere, Devlin, Kudenko and Nowe. "Context Sensitive Reward Shaping for Sparse Interaction Multi-Agent Systems." BNAIC, 2013.

## Belief Revision <sup>9</sup>



<sup>9</sup> Efthymiadis, Devlin, and Kudenko. "Overcoming Erroneous Domain Knowledge in Plan-Based Reward Shaping." AAMAS, 2013.



## Q-Table Initialization

- ▶ Wiewiora: “Potential-based shaping and Q-value initialization are equivalent.” (JAIR, 2003)
- ▶ ...If the potential function is **static**.

# State and Action Shaping<sup>10</sup>

## Look-Ahead Advice

- ▶  $F(s, a, s', a') = \gamma \Phi(s', a') - \Phi(s, a)$
- ▶  $\pi(s) = \operatorname{argmax}_a \{Q(s, a) + \Phi(s, a)\}$
- ▶ Maintains all previous guarantees

## Look-Back Advice

- ▶  $F(s, a, s', a') = \Phi(s', a') - \gamma^{-1} \Phi(s, a)$
- ▶ No guarantees proven

---

<sup>10</sup> Wiewiora, Cottrell and Elkan. "Principled methods for advising reinforcement learning agents." ICML, 2003.

# Partial Observability

## Formal Definition

- ▶  $F(o, o') = \gamma\Phi(o') - \Phi(o)$

## Guarantees <sup>11</sup>

- ▶ Equivalence to Q-table initialisation
- ▶ Policy invariance (optimal policy unchanged) in single agent
- ▶ Consistent Nash equilibria in multi-agent systems

---

<sup>11</sup> Eck, Soh, Devlin and Kudenko. "Potential-Based Reward Shaping for Partially Observable Markov Decision Processes." AAMAS, 2013.

## Closing Remarks

## Implementation Advice

- ▶  $\gamma$  must be equal to update rule
- ▶ Use an absorbing state
- ▶ Store current potential for next iteration
- ▶ Avoid negative potentials <sup>12</sup>

---

<sup>12</sup> Grzes and Kudenko. “Theoretical and empirical analysis of reward shaping in reinforcement learning.” ICMLA, 2009.

## General Effect

- ▶ Does not modify any property of the underlying MDP or SG invariant to changes in absolute value of expected return.
- ▶ Provided a property is only reliant on the **relative difference** or **order** of expected returns, potential-based reward shaping will not affect it.

## Necessity

- ▶ For every reward shaping function that is not potential-based, there is an MDP where the optimal policy differs with and without reward shaping.<sup>13</sup>

---

<sup>13</sup> Ng, Russell and Harada. "Policy Invariance Under Reward Transformations: Theory And Application To Reward Shaping." ICML, 1999.

## References

- ▶ Ng, Russell and Harada. "Policy Invariance Under Reward Transformations: Theory And Application To Reward Shaping." ICML, 1999.
- ▶ Wiewiora. "Potential-based shaping and Q-value initialization are equivalent." JAIR, 2003
- ▶ Wiewiora, Cottrell and Elkan. "Principled methods for advising reinforcement learning agents." ICML, 2003.
- ▶ Asmuth, Littman and Zinkov. "Potential-based shaping in model-based reinforcement learning." AAAI 2008
- ▶ Devlin and Kudenko, "Theoretical Considerations Of Potential-Based Reward Shaping For Multi-Agent Systems", AAMAS, 2011.
- ▶ Devlin and Kudenko. "Dynamic Potential-Based Reward Shaping." AAMAS, 2012.



# AAMAS 2014

## Potential-Based Difference Rewards for Multiagent Reinforcement Learning.

Sam Devlin, Logan Yliniemi, Daniel Kudenko and Kagan Tumer

Learning I

Miles Davis A & B

Wednesday 09:00