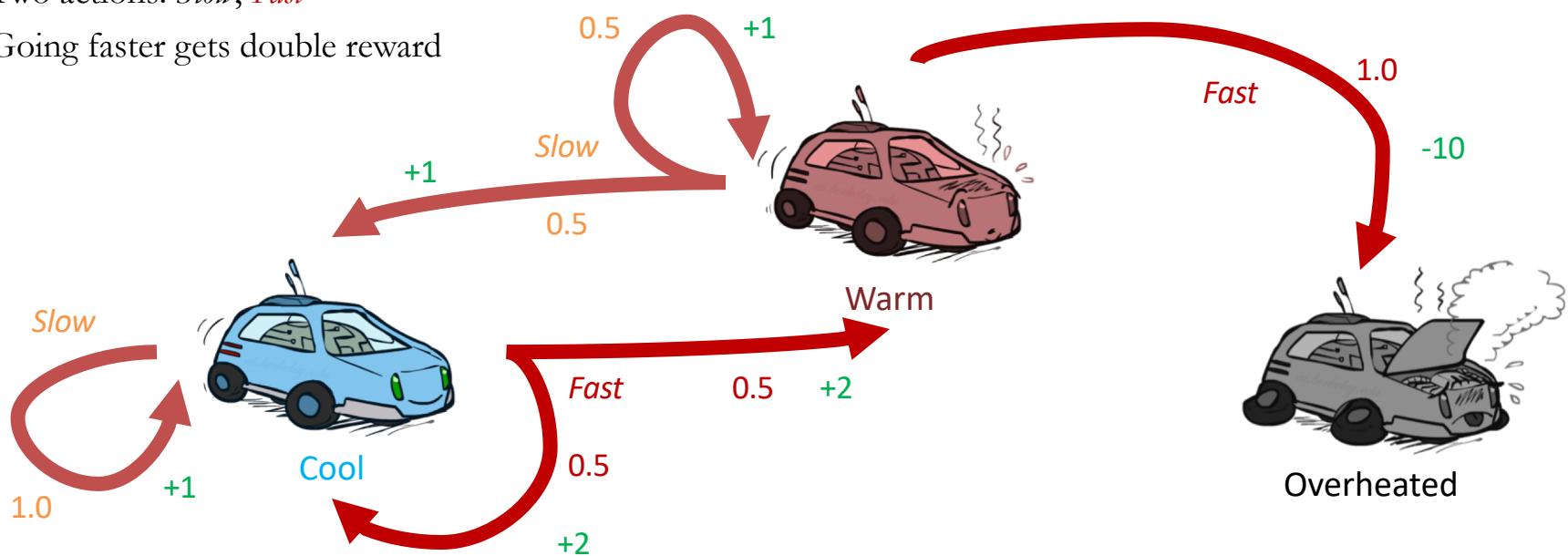


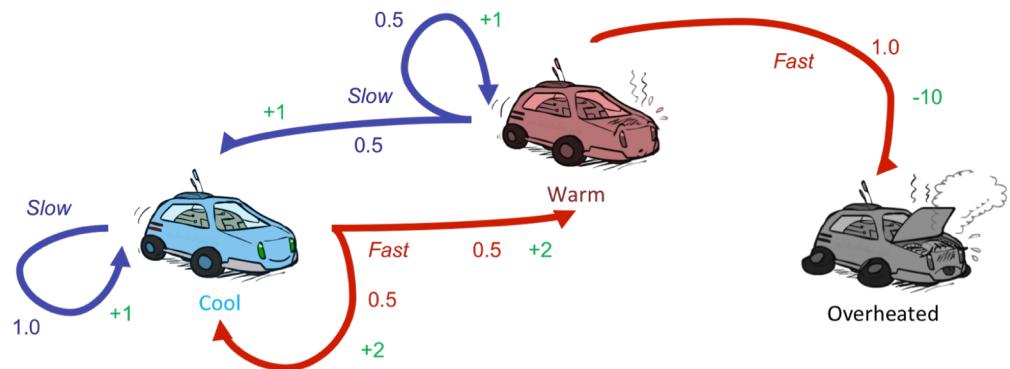
Racing

- A robot car wants to travel far, quickly
- Three states: Cool, Warm, Overheated
- Two actions: *Slow*, *Fast*
- Going faster gets double reward



Value iteration

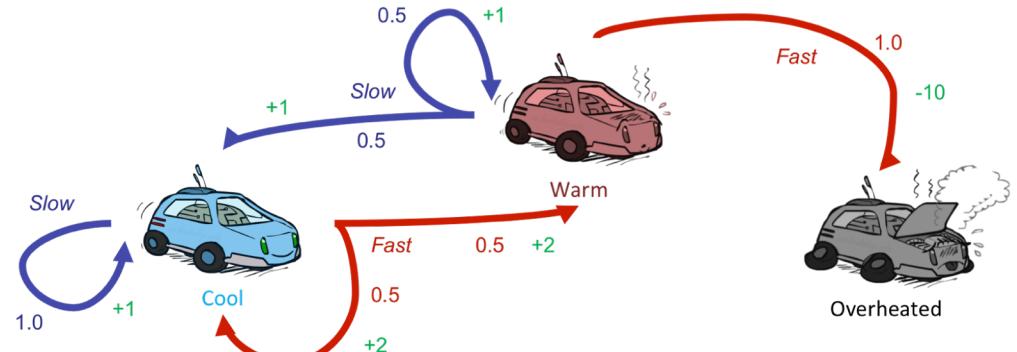
V_2	3.5	2.5	0
V_1	2	1	0
V_0	0	0	0



Assume no discount!

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

Value iteration



V_2

V_1

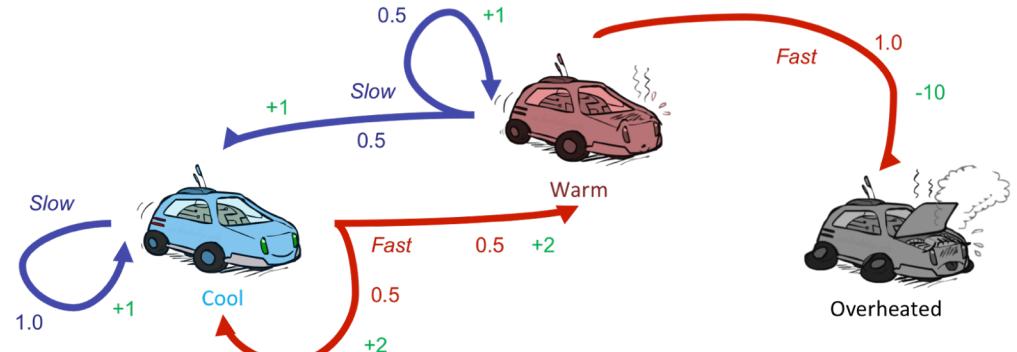
V_0

0	0	0
---	---	---

Assuming a discount factor of 0.5

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

Value iteration



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = \max_a \{ R(s, a) + \sum_{s'} [\gamma * T(s, a, s') V_k(s')] \}$$

V_2

V_1

V_0

0	0	0
---	---	---

Value iteration

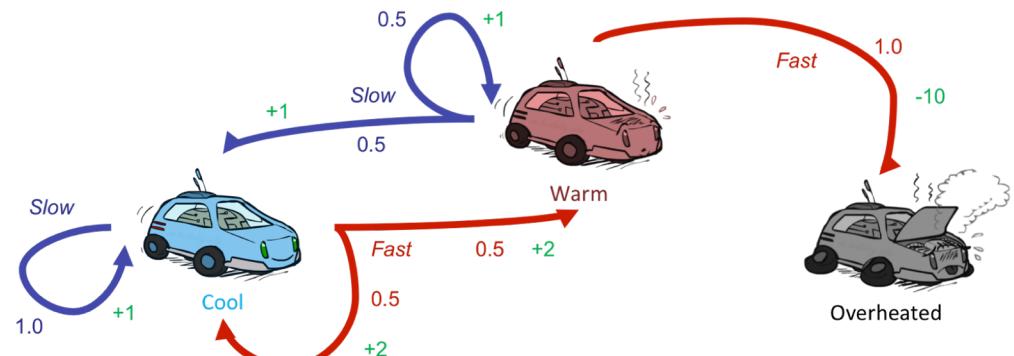
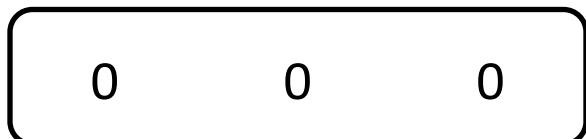
V_2



V_1



V_0



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = \max_a \{ R(s, a) + \sum_{s'} [\gamma * T(s, a, s') V_k(s')] \}$$

$$V1(\text{cool}) = \max\{$$

$$\begin{aligned} \text{slow: } & 1 + 1.0 * 0.5 * 0 = 1 \\ \text{fast: } & 2 + 0.5 * 0.5 * 0 + 0.5 * 0.5 * 0 = 2 \\ \} = & 2 \end{aligned}$$

$$V1(\text{warm}) = \max\{$$

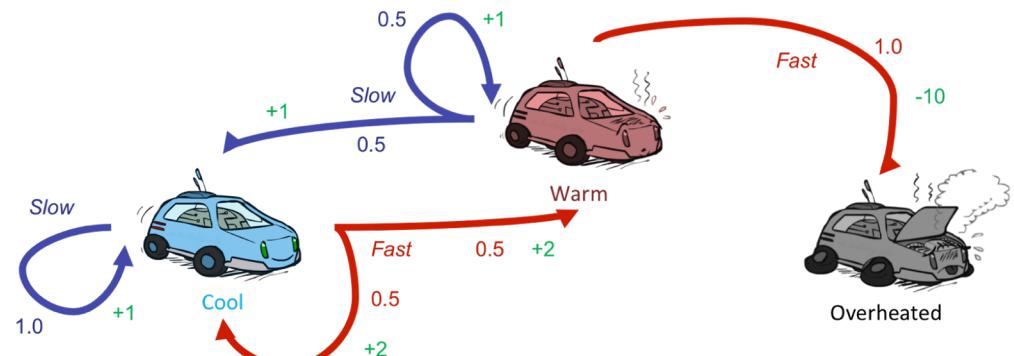
$$\begin{aligned} \text{slow: } & 1 + 0.5 * 0.5 * 0 + 0.5 * 0.5 * 0 = 1 \\ \text{fast: } & -10 + 1.0 * 0.5 * 0 = -10 \\ \} = & 1 \end{aligned}$$

$$V1(\text{overheated}) = V0(\text{overheated}) = 0$$

Value iteration

V_2		
V_1	2	1

V_0	0	0
-------	---	---



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = \max_a \{ R(s, a) + \sum_{s'} [\gamma^* T(s, a, s') V_k(s')] \}$$

Value iteration

V_2

2.75

1.75

0

V_1

2

1

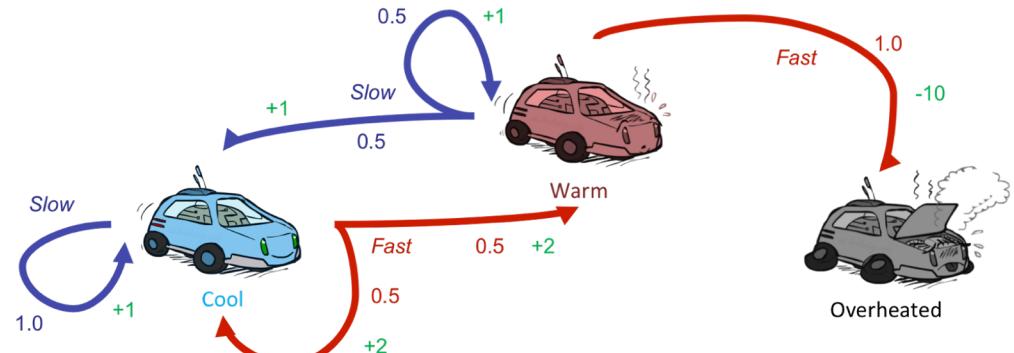
0

V_0

0

0

0



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = \max_a \{ R(s, a) + \sum_{s'} [\gamma * T(s, a, s') V_k(s')] \}$$

$$V_2(\text{cool}) = \max\{$$

$$\text{slow: } 1 + 1.0 * 0.5 * 2 = 2$$

$$\text{fast: } 2 + 0.5 * 0.5 * 2 + 0.5 * 0.5 * 1 = 2.75$$

$$\} = 2.75$$

$$V_2(\text{warm}) = \max\{$$

$$\text{slow: } 1 + 0.5 * 0.5 * 2 + 0.5 * 0.5 * 1 = 1.75$$

$$\text{fast: } -10 + 1.0 * 0.5 * 0 = -10$$

$$\} = 1.75$$

$$V_2(\text{overheated}) = V_1(\text{overheated}) = 0$$

Policy evaluation



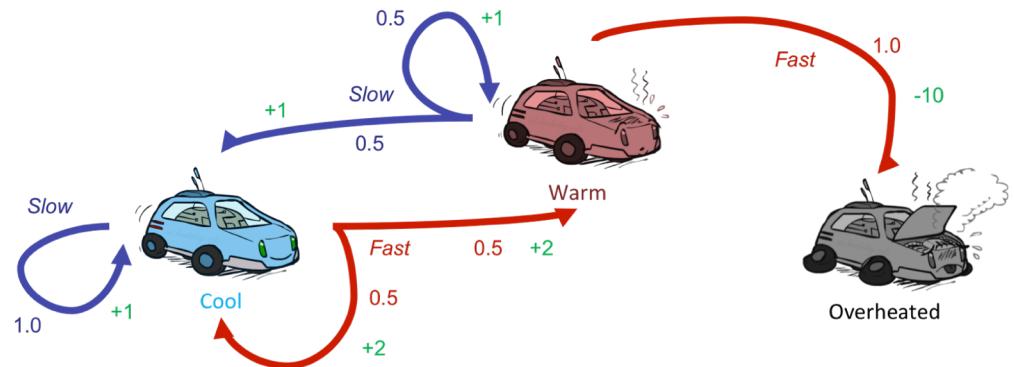
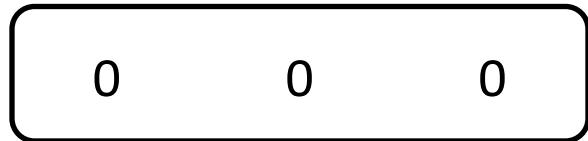
V_2



V_1



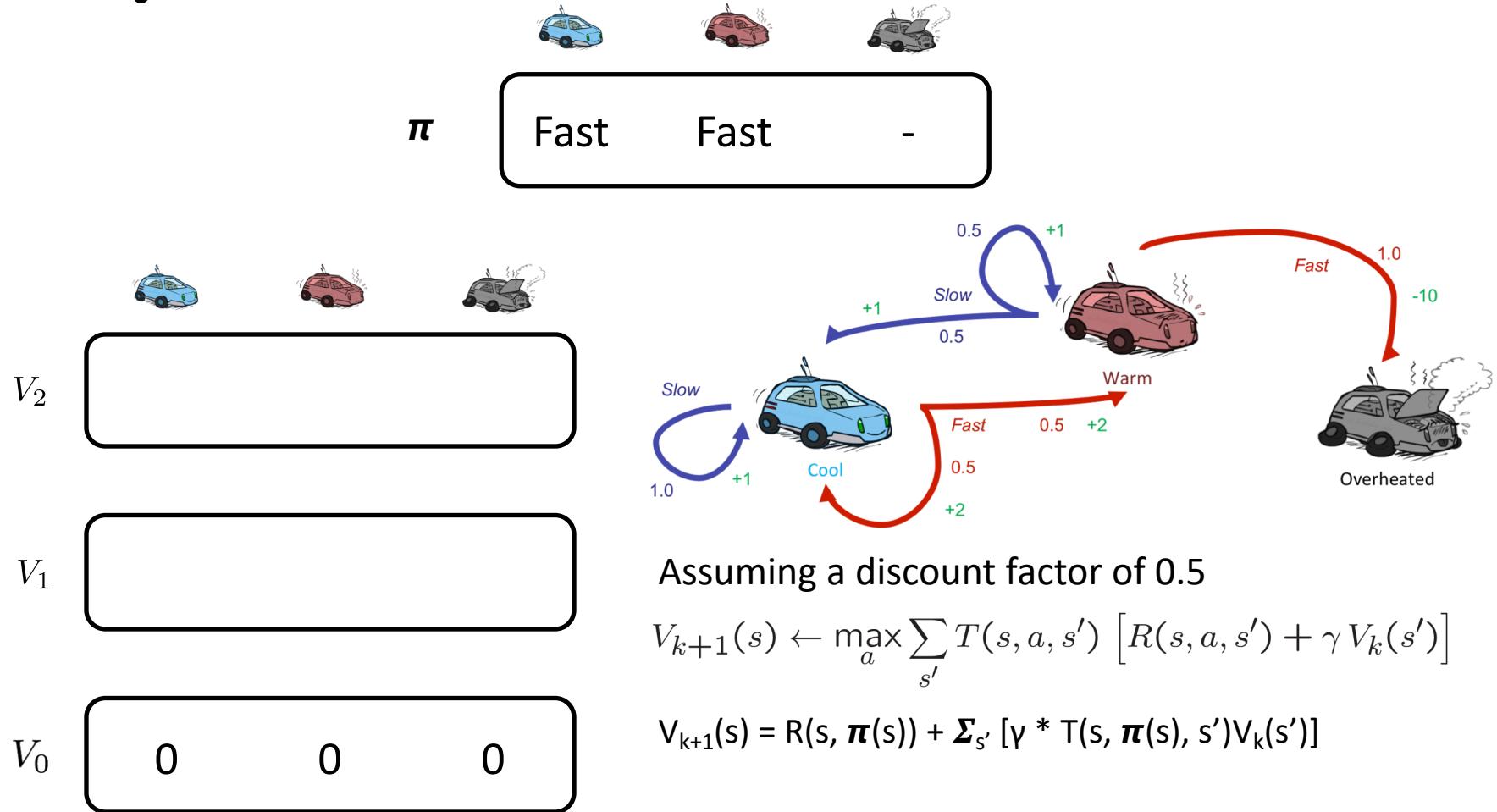
V_0



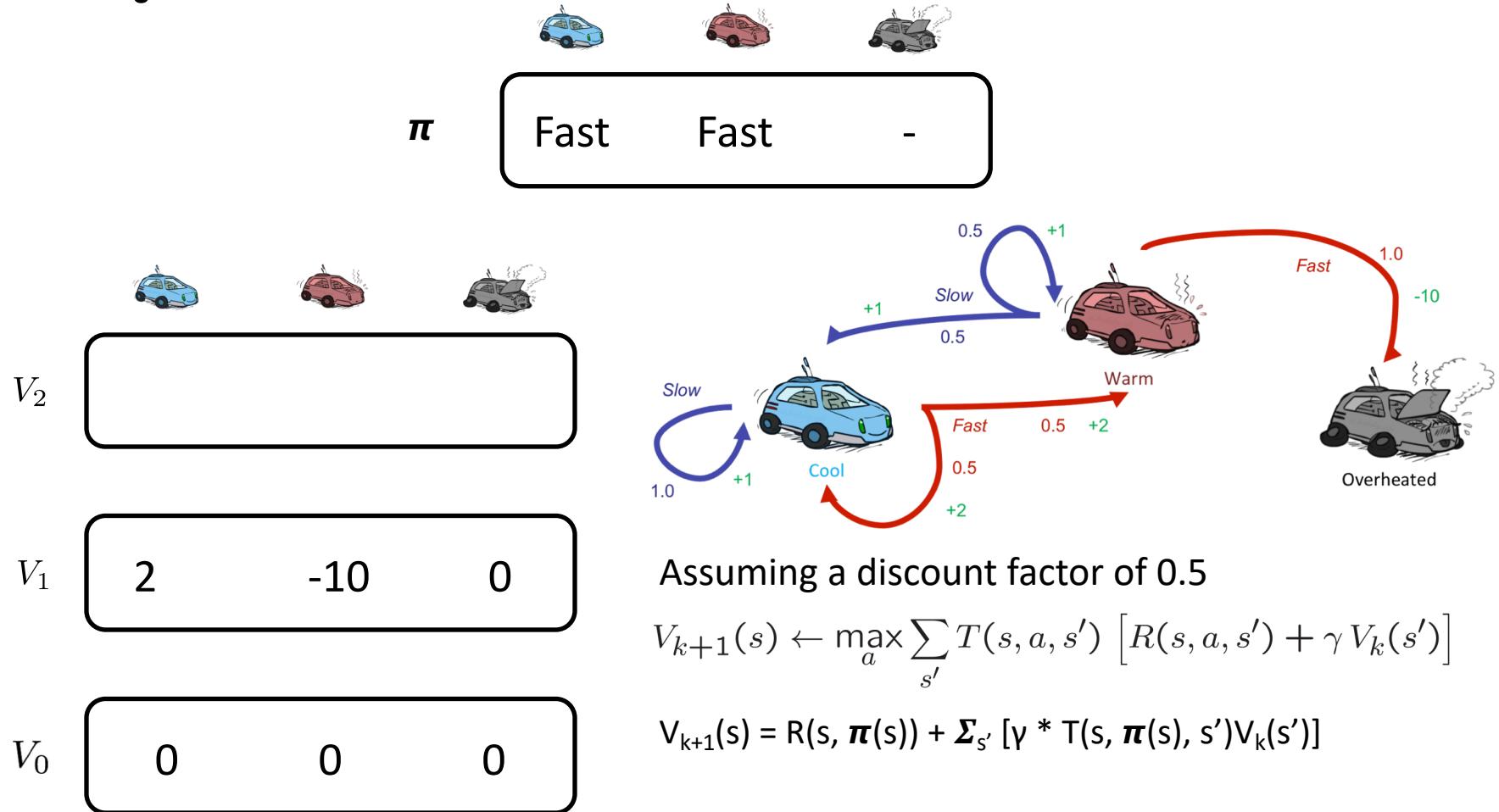
Assuming a discount factor of 0.5

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

Policy evaluation

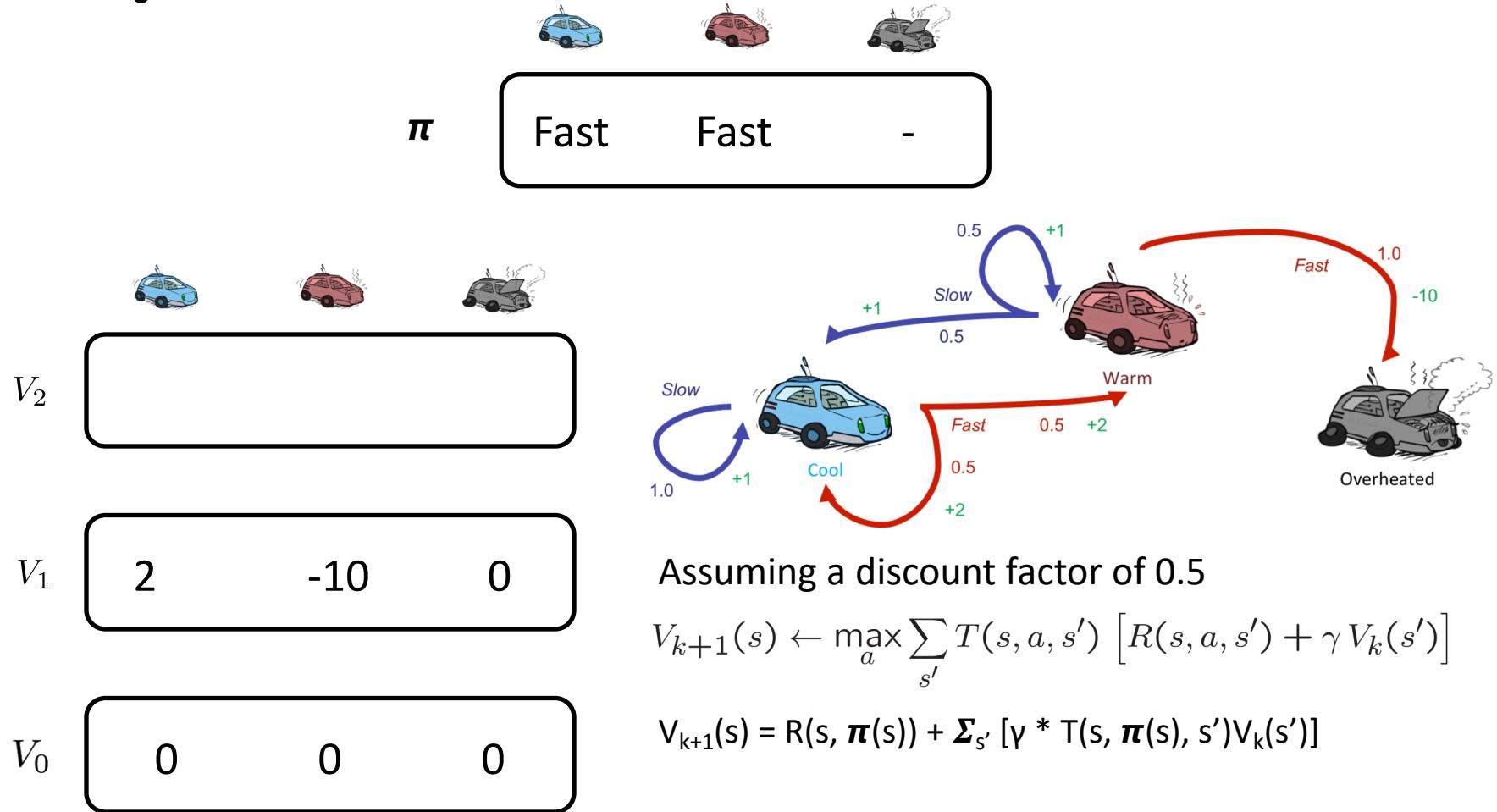


Policy evaluation

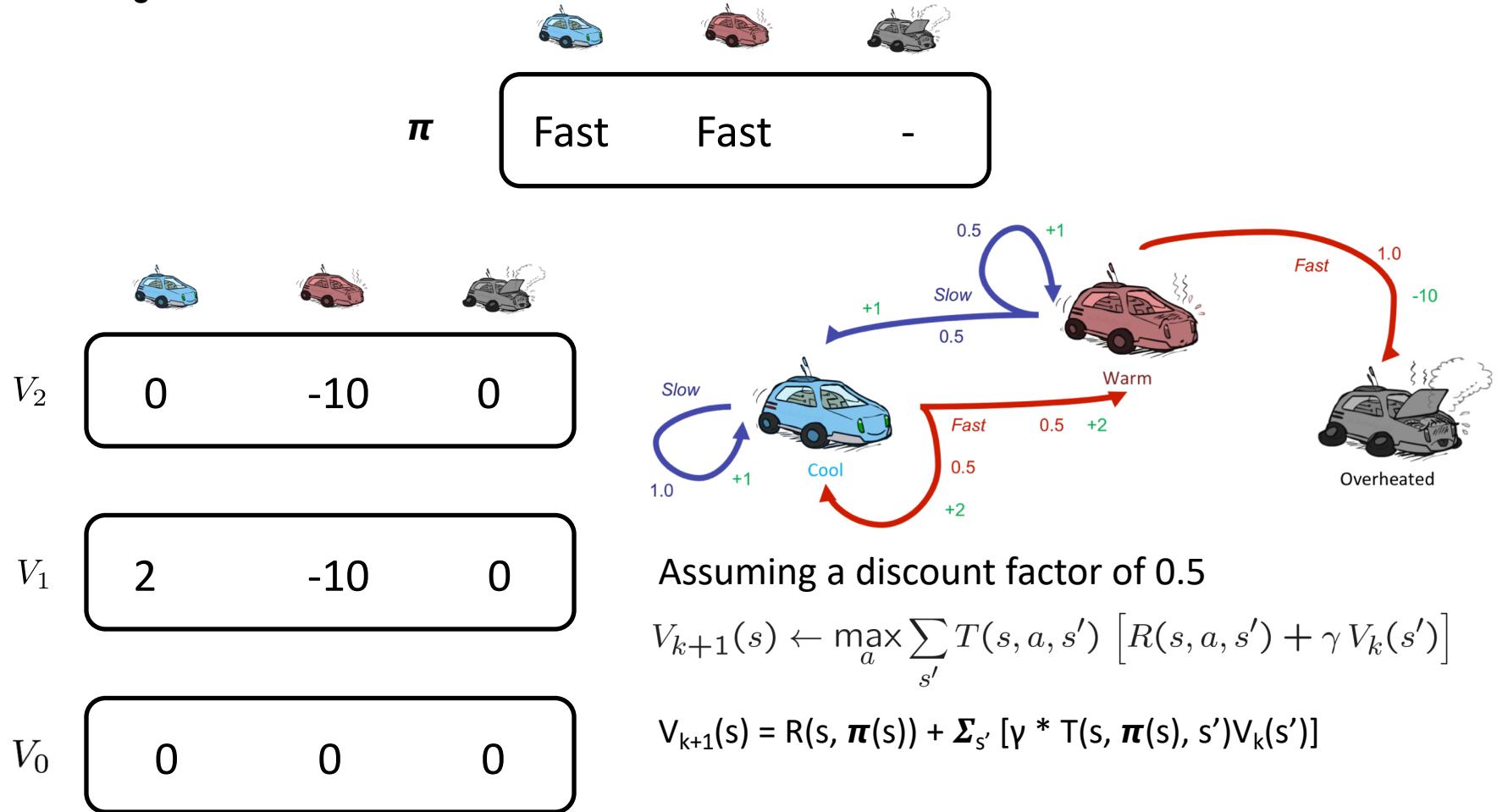


$$\begin{aligned} V1(\text{cool}) &= \text{fast: } 2 + 0.5 * 0.5 * 0 + 0.5 * 0.5 * 0 = 2 \\ V1(\text{warm}) &= \text{fast: } -10 + 1.0 * 0.5 * 0 = -10 \\ V1(\text{overheated}) &= V0(\text{overheated}) = 0 \end{aligned}$$

Policy evaluation



Policy evaluation

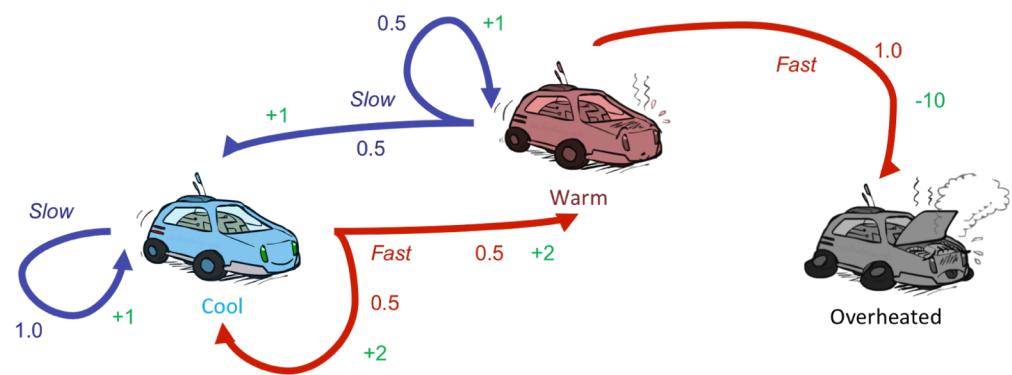
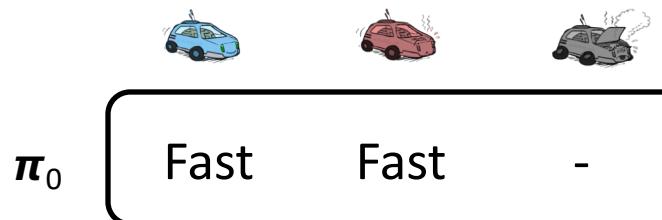


$$V2(\text{cool}) = \text{fast: } 2 + 0.5 * 0.5 * 2 + 0.5 * 0.5 * -10 = 0$$

$$V2(\text{warm}) = \text{fast: } -10 + 1.0 * 0.5 * 0 = -10$$

$$V2(\text{overheated}) = V1(\text{overheated}) = 0$$

Policy iteration



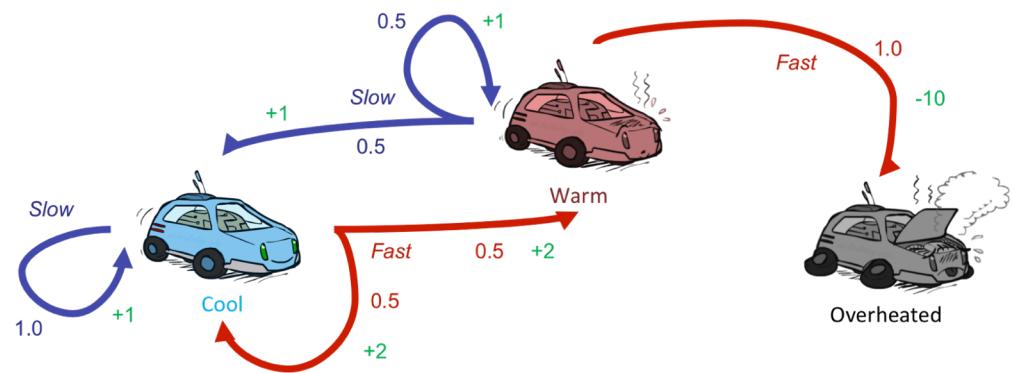
Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = R(s, \boldsymbol{\pi}(s)) + \Sigma_{s'} [\gamma * T(s, \boldsymbol{\pi}(s), s') V_k(s')]$$

Policy iteration

π_0	Fast	Fast	-
π_1			
π_2			
π_3			



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = R(s, \pi(s)) + \sum_{s'} [\gamma * T(s, \pi(s), s') V_k(s')]$$

Policy evaluation via linear programming:

$$V(\text{cool}) = 2 + 0.5 * 0.5 * V(\text{cool}) + 0.5 * 0.5 * V(\text{warm})$$

$$V(\text{warm}) = -10 + 1.0 * 0.5 * V(\text{overheat})$$

$$V(\text{overheat}) = 0$$

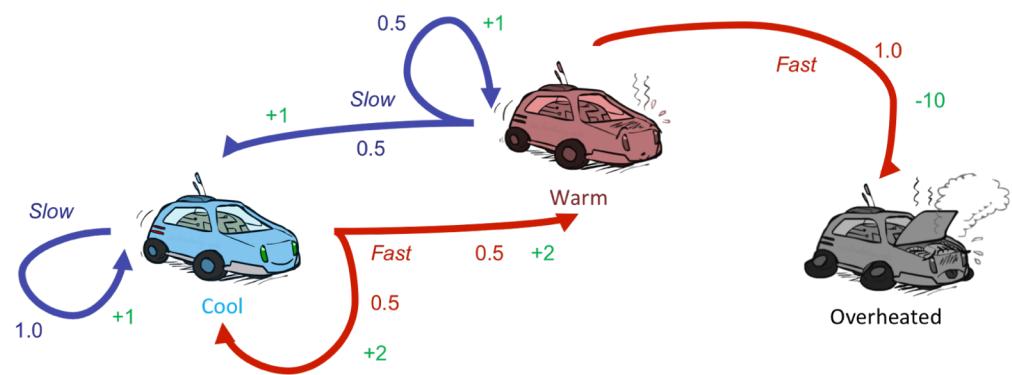
$$V(\text{cool}) = -2/3$$

$$V(\text{warm}) = -10$$

$$V(\text{overheat}) = 0$$

Policy iteration

π_0	Fast	Fast	-
π_1			
π_2			
π_3			



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = R(s, \boldsymbol{\pi}(s)) + \sum_{s'} [\gamma * T(s, \boldsymbol{\pi}(s), s') V_k(s')]$$

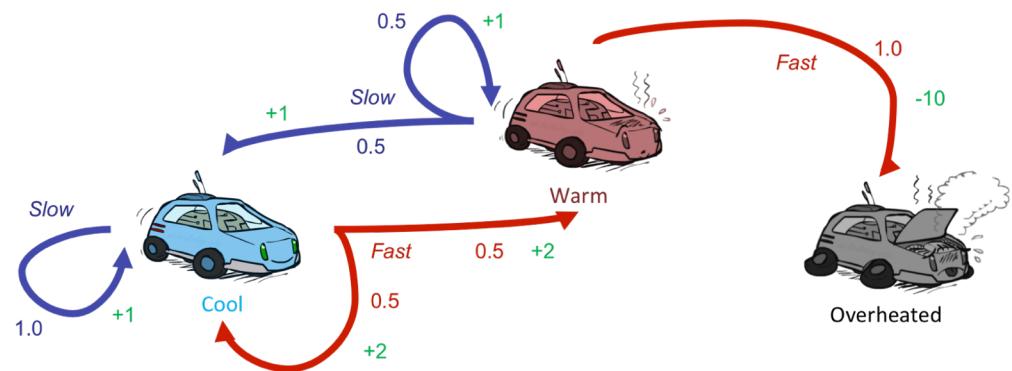
Policy improvement via policy extraction

$$V(\text{cool}) = -2/3$$

$$V(\text{warm}) = -10$$

$$V(\text{overheat}) = 0$$

Policy iteration



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = R(s, \pi(s)) + \sum_{s'} [\gamma * T(s, \pi(s), s') V_k(s')]$$

Policy improvement via policy extraction:

$$V(\text{cool}) = -2/3$$

$$V(\text{warm}) = -10$$

$$V(\text{overheat}) = 0$$

$$\pi(\text{cool}) = \text{argmax}\{$$

slow: $1 + 1.0 * 0.5 * -2/3 = 2/3$

fast: $2 + 0.5 * 0.5 * -2/3 + 0.5 * 0.5 * -10 = -2/3$
} = slow

$$\pi(\text{warm}) = \text{argmax}\{$$

$$\text{slow: } 1 + 0.5 * 0.5 * -2/3 + 0.5 * 0.5 * -10 = -5/3$$

fast: $-10 + 1.0 * 0.5 * 0 = -10$

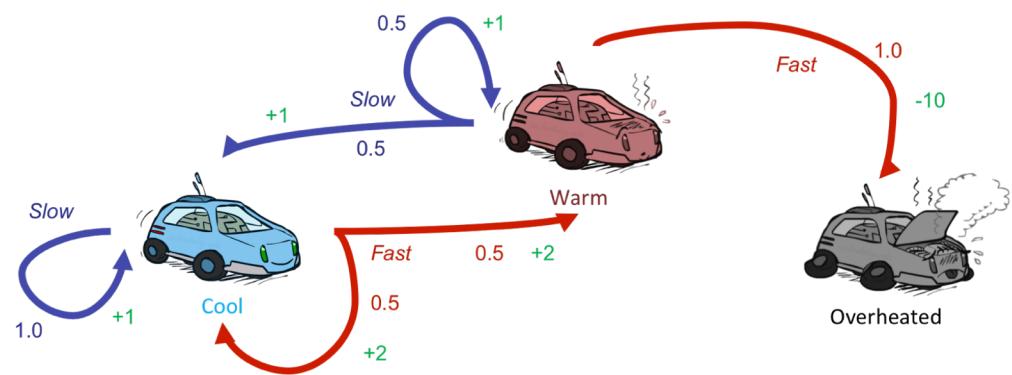
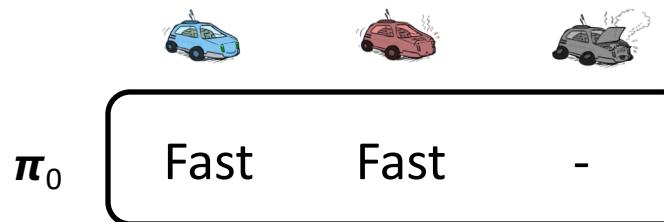
} = slow



Ira A. Fulton
Schools of Engineering

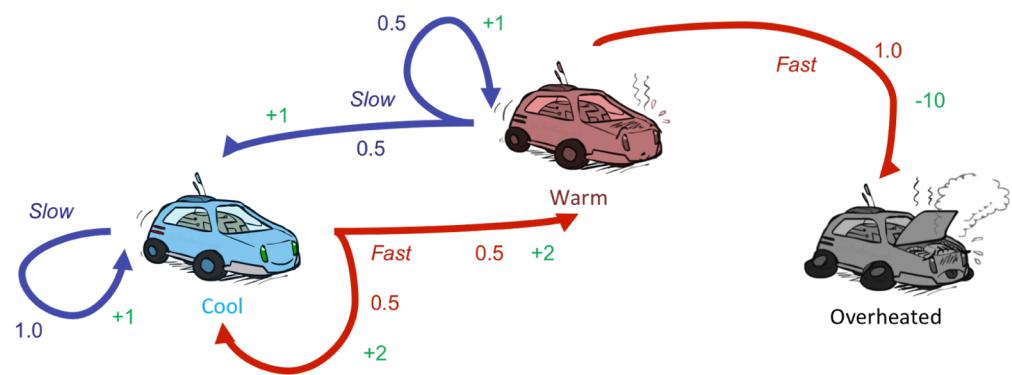
ARIZONA STATE UNIVERSITY

Policy iteration



Policy iteration

π_0	Fast	Fast	-
π_1	Slow	Slow	-
π_2			
π_3			



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = R(s, \boldsymbol{\pi}(s)) + \sum_{s'} [\gamma * T(s, \boldsymbol{\pi}(s), s') V_k(s')]$$

Policy evaluation via linear programming:

$$V(\text{cool}) = 1 + 1.0 * 0.5 * V(\text{cool})$$

$$V(\text{warm}) = 1 + 0.5 * 0.5 * V(\text{cool}) + 0.5 * 0.5 * V(\text{warm})$$

$$V(\text{overheat}) = 0$$

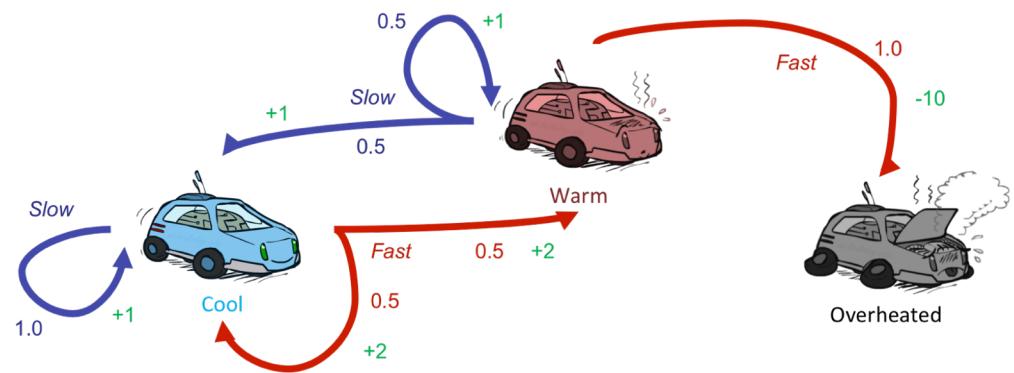
$$V(\text{cool}) = 2$$

$$V(\text{warm}) = 2$$

$$V(\text{overheat}) = 0$$

Policy iteration

π_0	Fast	Fast	-
π_1	Slow	Slow	-
π_2			
π_3			



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = R(s, \boldsymbol{\pi}(s)) + \sum_{s'} [\gamma * T(s, \boldsymbol{\pi}(s), s') V_k(s')]$$

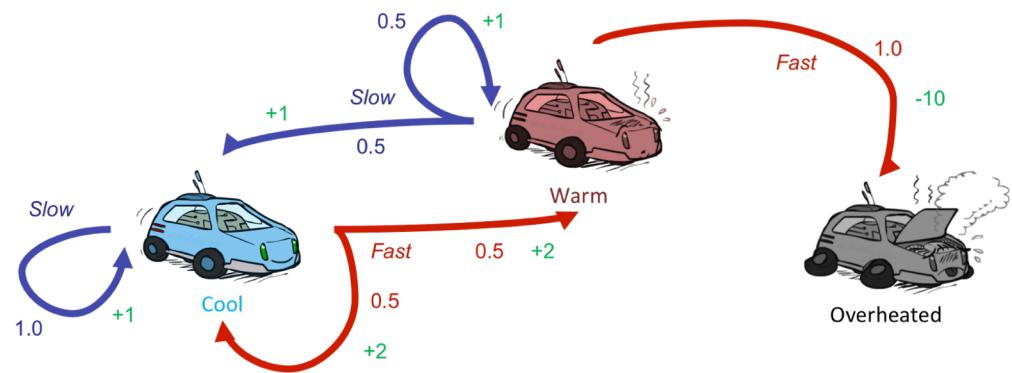
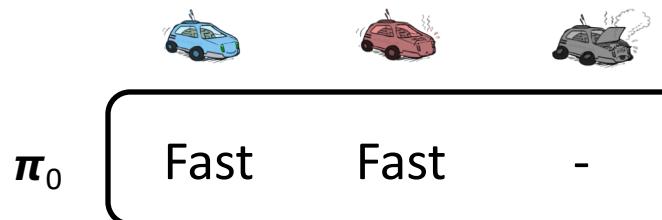
Policy improvement via policy extraction:

$$V(\text{cool}) = 2$$

$$V(\text{warm}) = 2$$

$$V(\text{overheat}) = 0$$

Policy iteration



Policy improvement via policy extraction:

$$V(\text{cool}) = 2$$

$$V(\text{warm}) = 2$$

$$V(\text{overheat}) = 0$$

$$\boldsymbol{\pi}(\text{cool}) = \text{argmax}\{$$

$$\text{slow: } 1 + 1.0 * 0.5 * 2 = 2$$

$$\text{fast: } 2 + 0.5 * 0.5 * 2 + 0.5 * 0.5 * 2 = 3$$

$$\} = \text{fast}$$

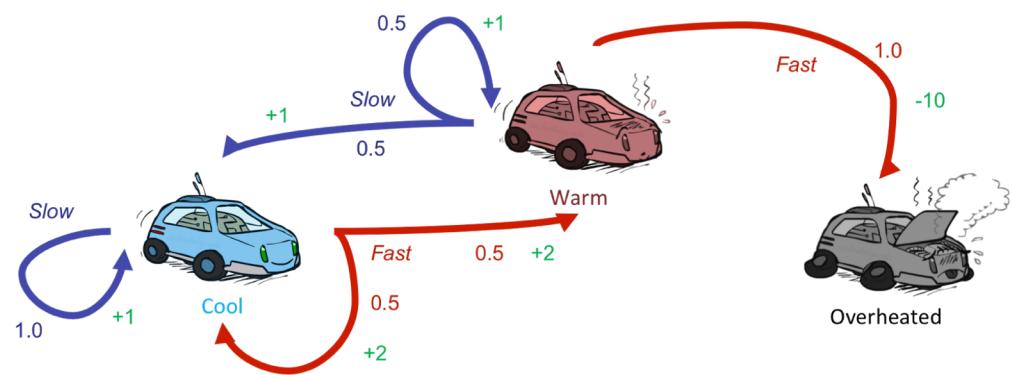
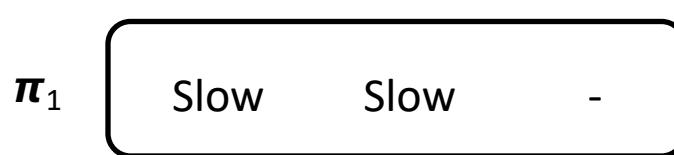
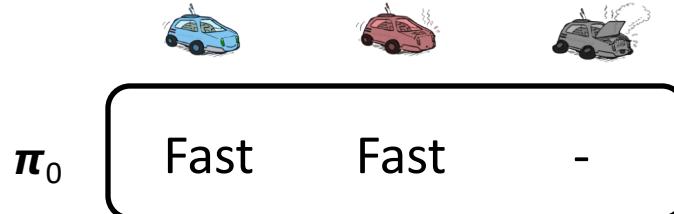
$$\boldsymbol{\pi}(\text{warm}) = \text{argmax}\{$$

$$\text{slow: } 1 + 0.5 * 0.5 * 2 + 0.5 * 0.5 * 2 = 2$$

$$\text{fast: } -10 + 1.0 * 0.5 * 0 = -10$$

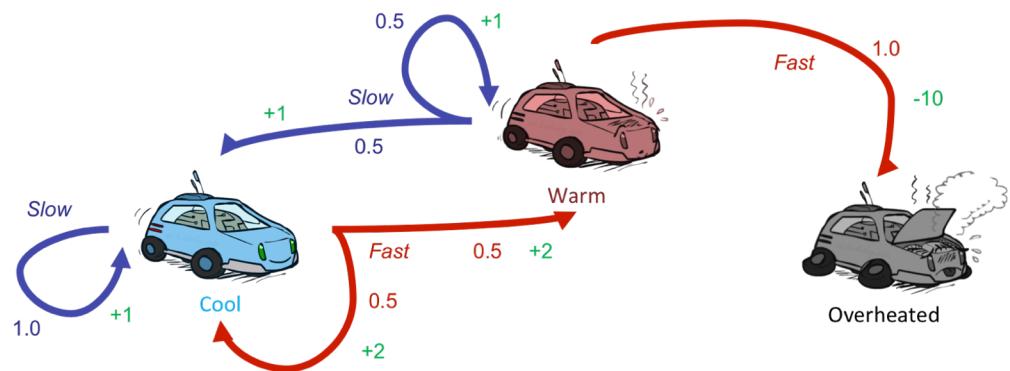
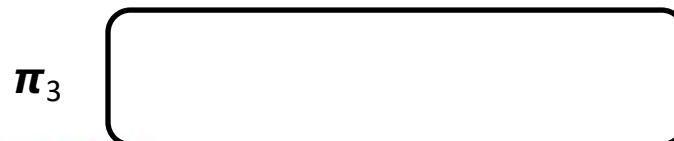
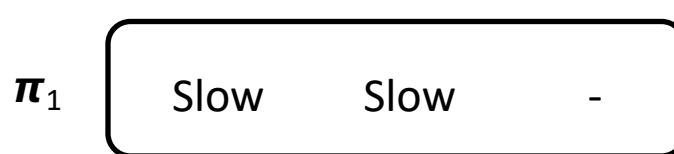
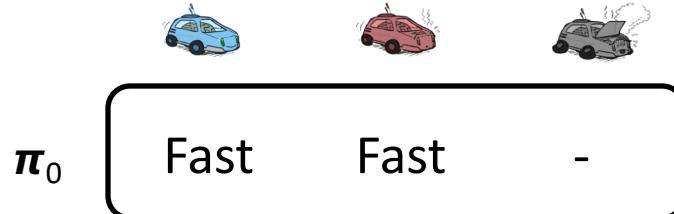
$$\} = \text{slow}$$

Policy iteration



Policy evaluation via linear programming:

Policy iteration



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = R(s, \pi(s)) + \sum_{s'} [\gamma * T(s, \pi(s), s') V_k(s')]$$

Policy evaluation via linear programming:

$$V(\text{cool}) = 2 + 0.5 * 0.5 * V(\text{cool}) + 0.5 * 0.5 * V(\text{warm})$$

$$V(\text{warm}) = 1 + 0.5 * 0.5 * V(\text{cool}) + 0.5 * 0.5 * V(\text{warm})$$

$$V(\text{overheat}) = 0$$

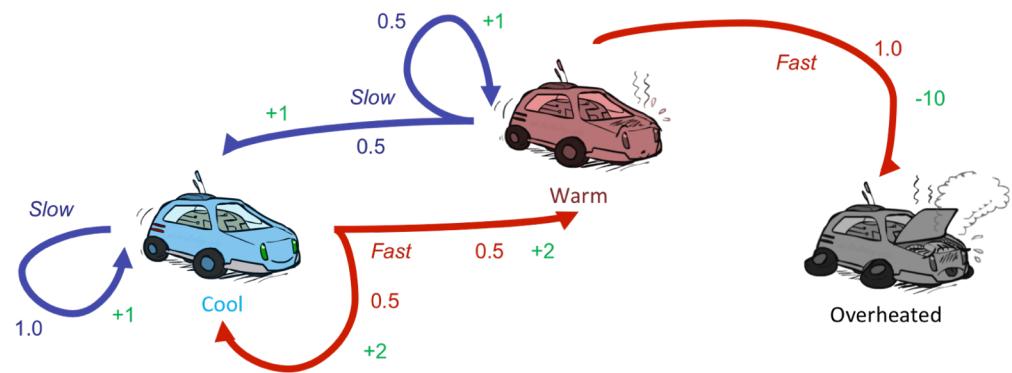
$$V(\text{cool}) = 3.5$$

$$V(\text{warm}) = 2.5$$

$$V(\text{overheat}) = 0$$

Policy iteration

π_0	Fast	Fast	-
π_1	Slow	Slow	-
π_2	Fast	Slow	-
π_3			



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = R(s, \boldsymbol{\pi}(s)) + \Sigma_{s'} [\gamma * T(s, \boldsymbol{\pi}(s), s') V_k(s')]$$

Policy evaluation via linear programming:

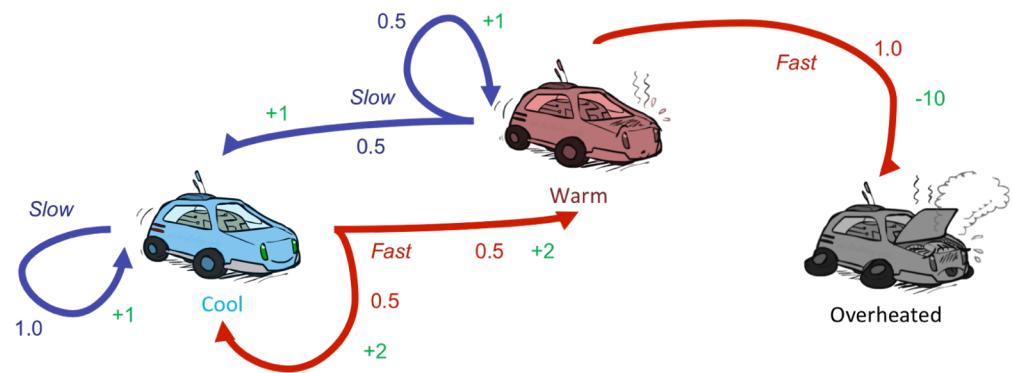
$$V(\text{cool}) = 3.5$$

$$V(\text{warm}) = 2.5$$

$$V(\text{overheat}) = 0$$

Policy iteration

π_0	Fast	Fast	-
π_1	Slow	Slow	-
π_2	Fast	Slow	-
π_3	Fast	Slow	-



Assuming a discount factor of 0.5

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

$$V_{k+1}(s) = R(s, \pi(s)) + \sum_{s'} [\gamma * T(s, \pi(s), s') V_k(s')]$$

Policy improvement via policy extraction:

$$V(\text{cool}) = 3.5$$

$$V(\text{warm}) = 2.5$$

$$V(\text{overheat}) = 0$$

$$\pi(\text{cool}) = \text{argmax}\{$$

$$\text{slow: } 1 + 1.0 * 0.5 * 3.5 = 2.75$$

$$\text{fast: } 2 + 0.5 * 0.5 * 3.5 + 0.5 * 0.5 * 2.5 = 3.5$$

$$\} = \text{fast}$$

$$\pi(\text{warm}) = \text{argmax}\{$$

$$\text{slow: } 1 + 0.5 * 0.5 * 3.5 + 0.5 * 0.5 * 2.5 = 2.5$$

$$\text{fast: } -10 + 1.0 * 0.5 * 0 = -10$$

$$\} = \text{slow}$$