# SwiftVI: Time-Efficient Planning and Learning with MDPs

**Kasper Overgaard Mortensen** · Konstantinos Skitsas · Emil Morre Christensen · Mohammad Sadegh Talebi · Andreas Pavlogiannis · Davide Mottin · Panagiotis Karras



A project based upon Emil Morre Christensen's master's thesis

## Planning







### Planning







### Planning







### Planning with MDPs



### Planning with MDPs



### Planning with MDPs using Value Iteration



- S = State space
- A = Action space
- *P* = *Transition probability funtion*
- R = Reward function
- $\gamma = Future reward discount factor$



with number of actions

transition probabilities

### A simple solution

### Beneficial to avoid updating action values! But how?

### An upper bound exists...

...such that value's updates are monotonic decreasing.







### Update only the best action

**Overhead:** Build the heaps.

Per iteration of Value Iteration:

Worst case: Update all actions. Best case: One action update per state. Approaches the best case as Value function converge.



### Learning with MDPs

#### Unknown reward function and transition probabilities.



### Learning with MDPs

Learned through a balance of exploration and exploitation.



### Learning with MDPs using Value Iteration

Substitute VI if it is used as a subroutine.





Contact me: km@cs.au.dk

Repository: github.com/constantinosskitsas/SwiftVI

Repository for Emil's master's thesis: github.com/Dugtoud/Time-Efficient-VI-for-MDPs