

# Metareasoning for Deliberation Time Distribution in the PROST Planner

Ferdinand Badenber

University of Basel

Bachelor Thesis Presentation, 2017

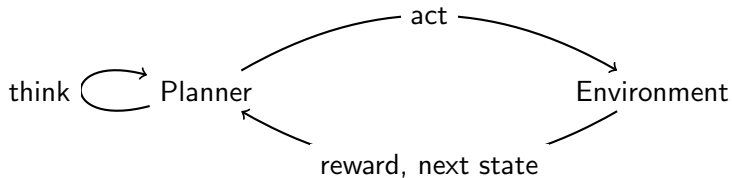
# Outline

- 1 Motivation
  - Why Metareasoning?
  - Metareasoning Problem
  
- 2 Approaches for Metareasoners
  - Hand Made Functions
  - Metareasoner of Lin. et al.
  - Improvements for the Metareasoner
  
- 3 Results
  - Results for the Hand Made Functions
  - Results for the Formal Procedure

# Table of Contents

- 1 Motivation
  - Why Metareasoning?
  - Metareasoning Problem
- 2 Approaches for Metareasoners
  - Hand Made Functions
  - Metareasoner of Lin. et al.
  - Improvements for the Metareasoner
- 3 Results
  - Results for the Hand Made Functions
  - Results for the Formal Procedure

# Cycle



# Motivation

## Why Metareasoning?

- Optimise policy in given time
- Allocate time to think where it is needed
- Act if decision is easy, clear best action
- Think if decision is difficult, multiple actions very close

# Metareasoning Problem

## Metareasoning Problem

- Steps from finite horizon MDP
- Rounds
- Limited time
- Anytime search algorithm

## Metareasoner

Decision to think or act

- Based on specific values for these factors
- After one thinking cycle of the algorithm
- Goal: only think when necessary

# Table of Contents

- 1 Motivation
  - Why Metareasoning?
  - Metareasoning Problem
- 2 Approaches for Metareasoners
  - Hand Made Functions
  - Metareasoner of Lin. et al.
  - Improvements for the Metareasoner
- 3 Results
  - Results for the Hand Made Functions
  - Results for the Formal Procedure

# Hand Made Functions

## Idea

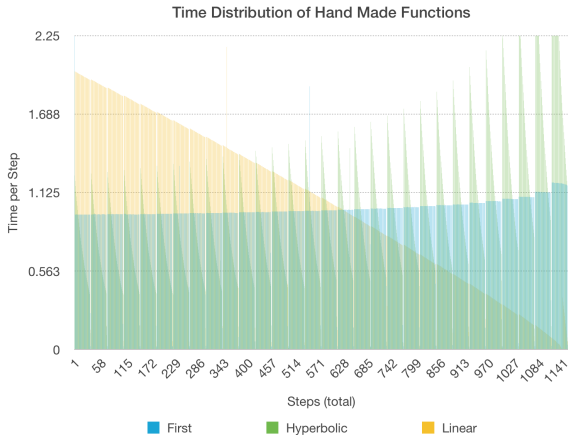
- Allocate time for each step
- Think for as long as time is left
- State of the search algorithm not considered

## Functions Tested:

- 1 Uniform (Standard)
- 2 First
- 3 Linear
- 4 Hyperbolic



# Example Time Distribution



# Formal Metareasoner of Lin et al.

## Metareasoner

- Idea: think if change of policy is likely, act if it will stay the same
- Only considers expected reward estimations (Q-values) of search algorithm
- Act if  $Q^{act} \geq Q^{think}$
- How are they calculated?

# Formal Metareasoner: $Q^{think}$ and $Q^{act}$

## $Q^{think}$

- Expected reward of the policy after another thinking cycle
- Simplification: only best action is relevant
- Estimate probability of action  $a$  being the best after the next thinking cycle
- Estimate expected reward given that action  $a$  is chosen
- Needed: next Q-values for each action

## $Q^{act}$

- Intuitive idea: Q-value of current best action
- But: average of current Q-value and next Q-value

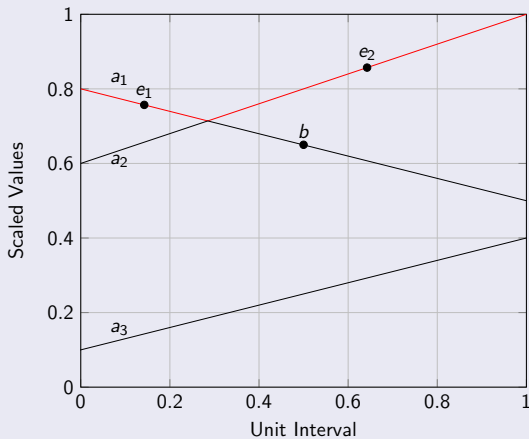
# Formal Metareasoner: Estimation of Next Q-values

## Estimating Next Q-values

- Idea: base next change in Q-values on previous change in Q-values
- Assumption: next  $\Delta Q$ -value no larger than the previous one
- Draw random  $\rho$  between 0 and 1
- $\Delta Q(a) = \hat{\Delta}Q(a) * \rho$  for all actions  $a$

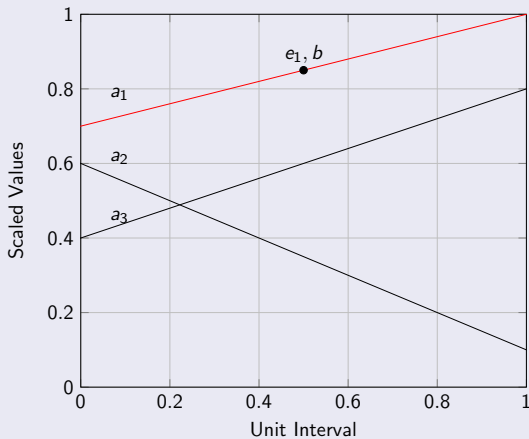
# Line Segment Example: UCT

$Q^{think} > Q^{act}$



# Line Segment Example: UCT

$$Q^{think} = Q^{act}$$

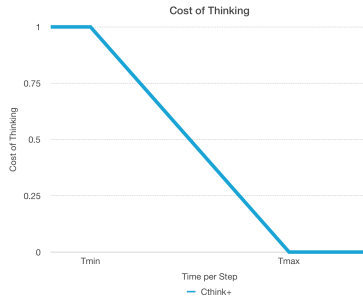


# Improvements

## Minimum Thinking Time

Problem: assumption is often not true early on

Improvement: think for at least  $T_{min}$  seconds

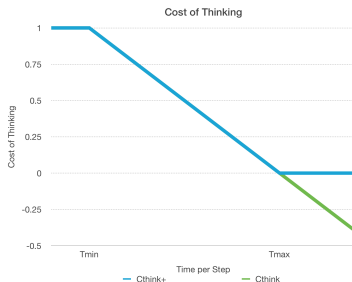
Cthink<sup>+</sup>Cthink<sup>+</sup>

Problem: time left is not considered

Improvement: subtract  $C^{think}$  from  $Q^{think}$



# Cthink



## Cthink

Problem: stopping with time left is useless

Improvement: allow a negative  $C^{think}$

# Table of Contents

- 1 Motivation
  - Why Metareasoning?
  - Metareasoning Problem
- 2 Approaches for Metareasoners
  - Hand Made Functions
  - Metareasoner of Lin. et al.
  - Improvements for the Metareasoner
- 3 Results
  - Results for the Hand Made Functions
  - Results for the Formal Procedure

# Results Hand Made Functions

## Results

Problem	Uniform	Hyperbolic	First	Linear
Wildfire	74	71	80	<b>81</b>
Triangle	72	65	72	<b>75</b>
Academic	37	37	34	<b>45</b>
Elevators	93	93	91	<b>94</b>
Tamarisk	93	<b>94</b>	92	91
Sysadmin	<b>94</b>	<b>94</b>	90	91
Recon	97	<b>99</b>	97	96
Game	<b>97</b>	93	94	93
Traffic	<b>97</b>	<b>97</b>	96	96
Crossing	87	89	91	<b>99</b>
Skill	91	91	88	<b>93</b>
Navigation	65	58	<b>83</b>	82
Total	83	82	84	<b>86</b>

# Results Formal Procedure

## Results

Problem	Uniform	Lin et al.	Minimum	Cthink <sup>+</sup>	Cthink
Wildfire	60	90	86	<b>95</b>	68
Triangle	<b>78</b>	67	62	59	68
Academic	<b>39</b>	32	36	35	38
Elevators	<b>98</b>	71	83	83	97
Tamarisk	<b>96</b>	68	86	90	92
Sysadmin	<b>100</b>	36	67	74	82
Recon	<b>98</b>	56	75	75	97
Game	<b>97</b>	64	82	86	96
Traffic	<b>99</b>	85	90	87	98
Crossing	88	58	78	83	<b>89</b>
Skill	<b>100</b>	25	71	69	86
Navigation	82	26	25	28	<b>83</b>
Total	<b>86</b>	56	70	72	83

# Summary

## Result Summary

- Hand made functions performed very well
- Default metareasoner severely underestimates thinking
- The improvements proved to be very useful

# Summary

## Outlook

- More general hand made functions
- Improve formal procedure:
  - Consider all previous  $\Delta Q$ -values
  - Replace random  $\rho$
- More sophisticated cost of thinking:  
combination of two approaches

Questions?

# BRTDP vs UCT

## BRTDP

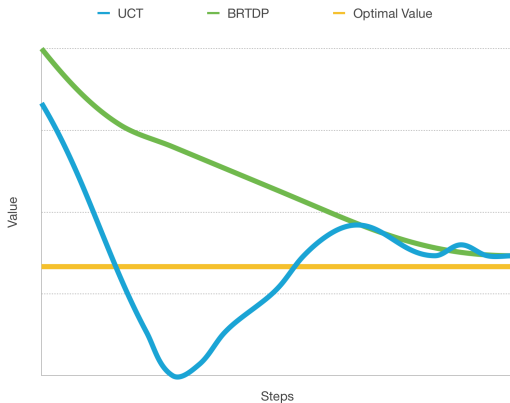
- Used in original paper
- Cost setting
- Uses upper bound of the actual Q-value
- Monotonously decreasing

## UCT

- Used by PROST planner
- Reward setting
- No guarantees

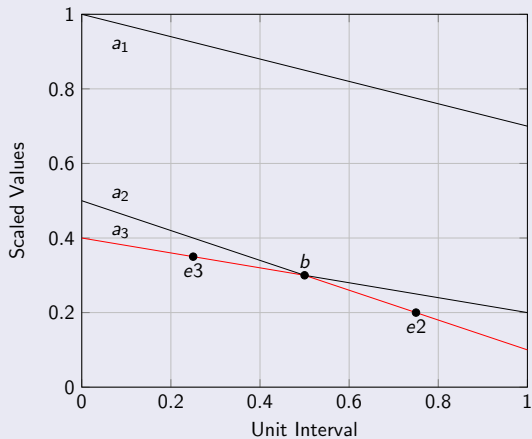


# BRTDP vs UCT: Visualisation



## Line Segment Example: BRTDP

$$Q^{think} < Q^{act}$$



# Wildfire Time per Step

