



# Bayesian Mixture Modeling and Inference based Thompson Sampling in Monte-Carlo Tree Search

Aijun Bai<sup>†</sup>, Feng Wu<sup>‡</sup>, and Xiaoping Chen<sup>†</sup>

University of Science and Technology of China<sup>†</sup>, and University of Southampton<sup>‡</sup>



## BACKGROUND

Monte-Carlo tree search (MCTS) finds near-optimal policies in domains of online planning for *Markov decision processes* (MDPs) by combining tree search methods with sampling techniques. The key idea is to iteratively evaluate each state in a best-first search tree by the mean outcome of simulation samples.

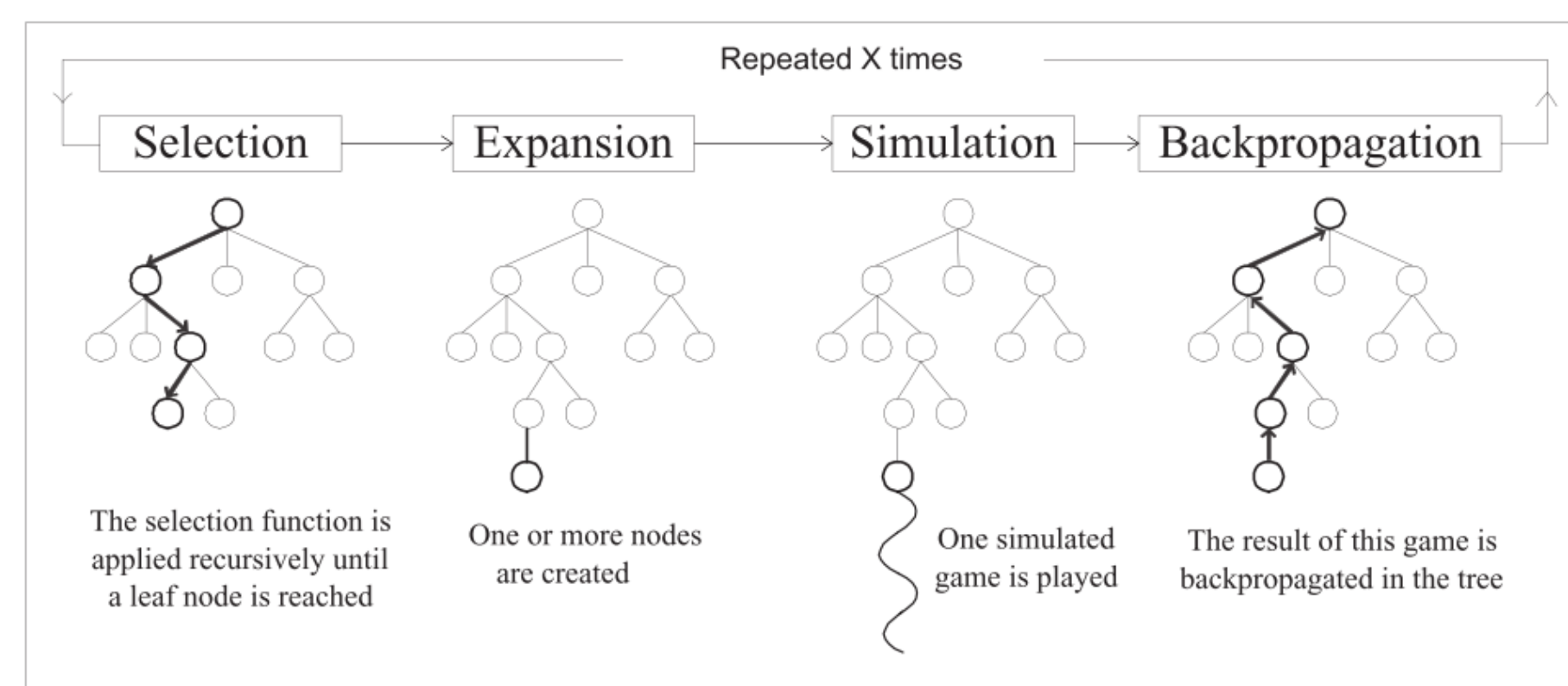


Figure 1: Monte-Carlo Tree Search [Chaslot2008]

## MOTIVATION

When applying MCTS, one of the fundamental challenges is the so-called *exploration versus exploitation* dilemma. Thompson sampling selects actions stochastically, based on the probabilities of being optimal. In this paper, we borrow the idea of Thompson sampling and propose the Dirichlet-NormalGamma MCTS (DNG-MCTS) algorithm.

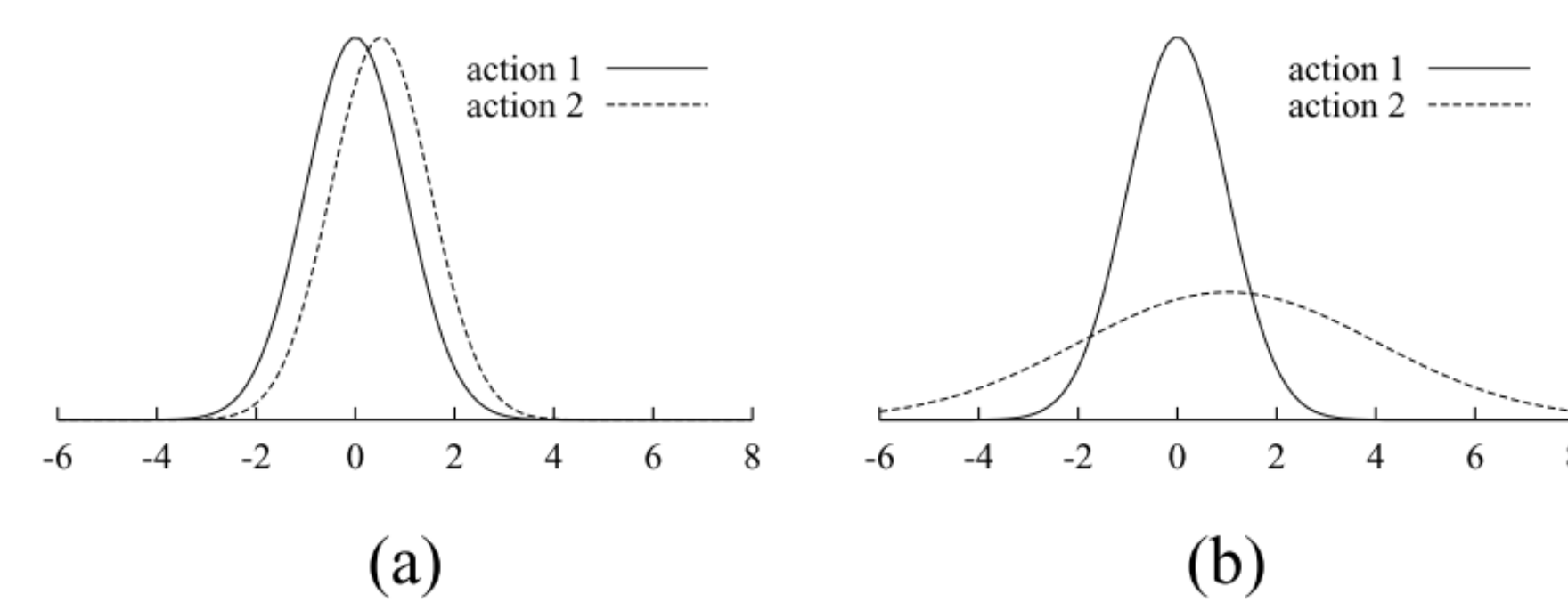


Figure 2: Thompson sampling based action selection [Dearden1998]

## ASSUMPTIONS

For a given MDP policy  $\pi$ , let random variables:

- $X_{s,\pi}$  denotes the accumulated reward of following policy  $\pi$  starting from state  $s$ ,
- $X_{s,a,\pi}$  denotes the accumulated reward of first performing action  $a$  in state  $s$  and then following policy  $\pi$  thereafter.

According to the central limit theorem on Markov chains, our assumptions are:

- $X_{s,\pi}$  is sampled from a Normal distribution,
- $X_{s,a,\pi}$  can be modeled as a mixture of Normal distributions.

## ALGORITHM

DNG-MCTS algorithm extends MCTS:

- Model the unknown distribution of  $X_{s,a,\pi}$  as a mixture of Normal distributions,
- Choose the conjugate prior in the form of a combination of Dirichlet and Normal-Gamma distributions,
- Compute the posterior distribution after each accumulated reward is observed by simulation in the search tree,
- Use Thompson sampling to guide the selection of actions at each decision node.

## CONCLUSION AND FUTURE WORK

DNG-MCTS algorithm:

- Monte-Carlo tree search framework
- Bayesian mixture modeling and inference
- Thompson sampling based action selection

- Competitive results comparing to UCT

In the future, we plan to extend our basic assumptions to using more realistic distributions and test our algorithm on real-world applications.

## EXPERIMENTS

We have tested DNG-MCTS and compared the results with UCT in three benchmark domains, namely *Canadian traveler problem*, *racetrack* and *sailing*. In each benchmark problem, we:

- Ran the algorithms for a number of iterations from the current state,
- Applied the best action based on the re-

sulted action-values,

- Repeated the loop until terminating conditions,
- Reported the total discounted cost.

DNG-MCTS produced competitive results in *CTP*, and converged faster in *racetrack* and *sailing* with respect to sample complexity.

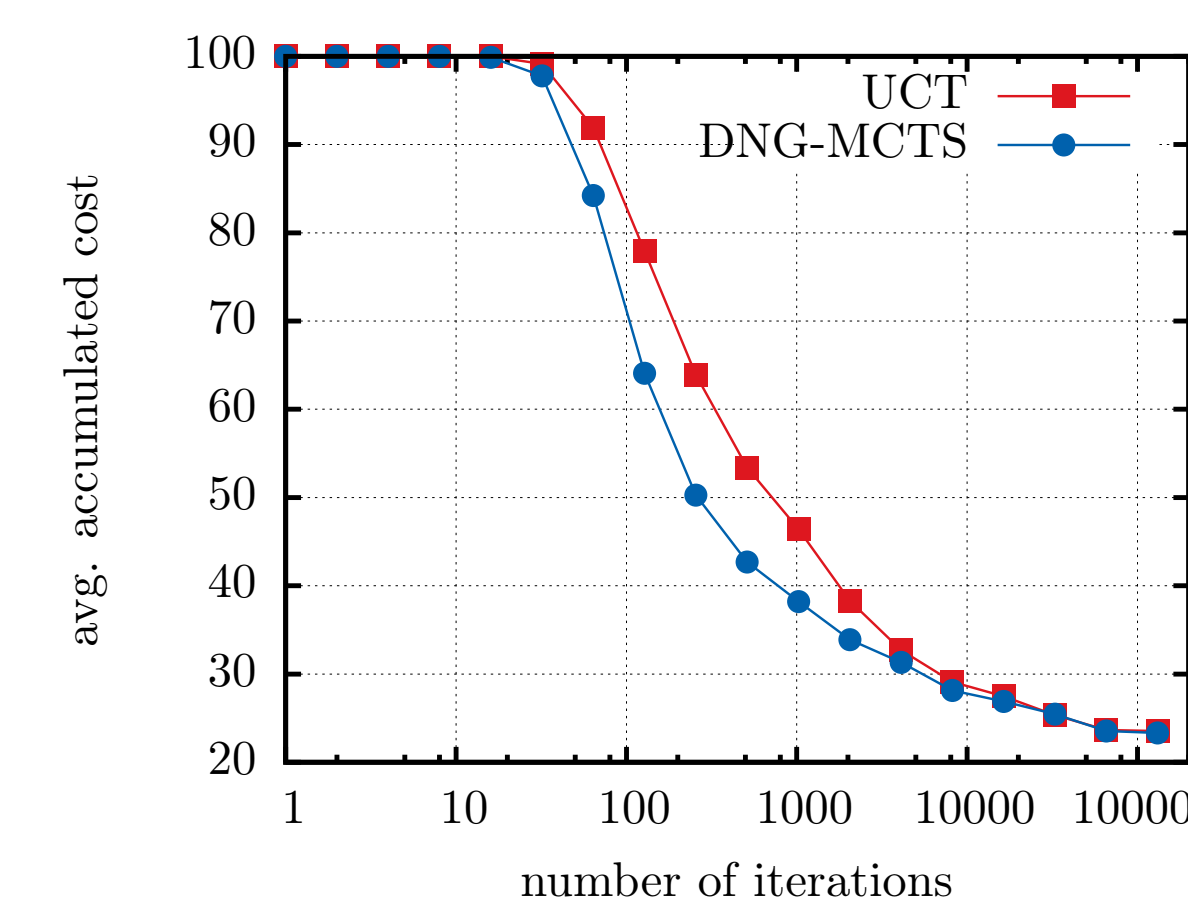


Figure 3: Racetrack-barto-big with random policy

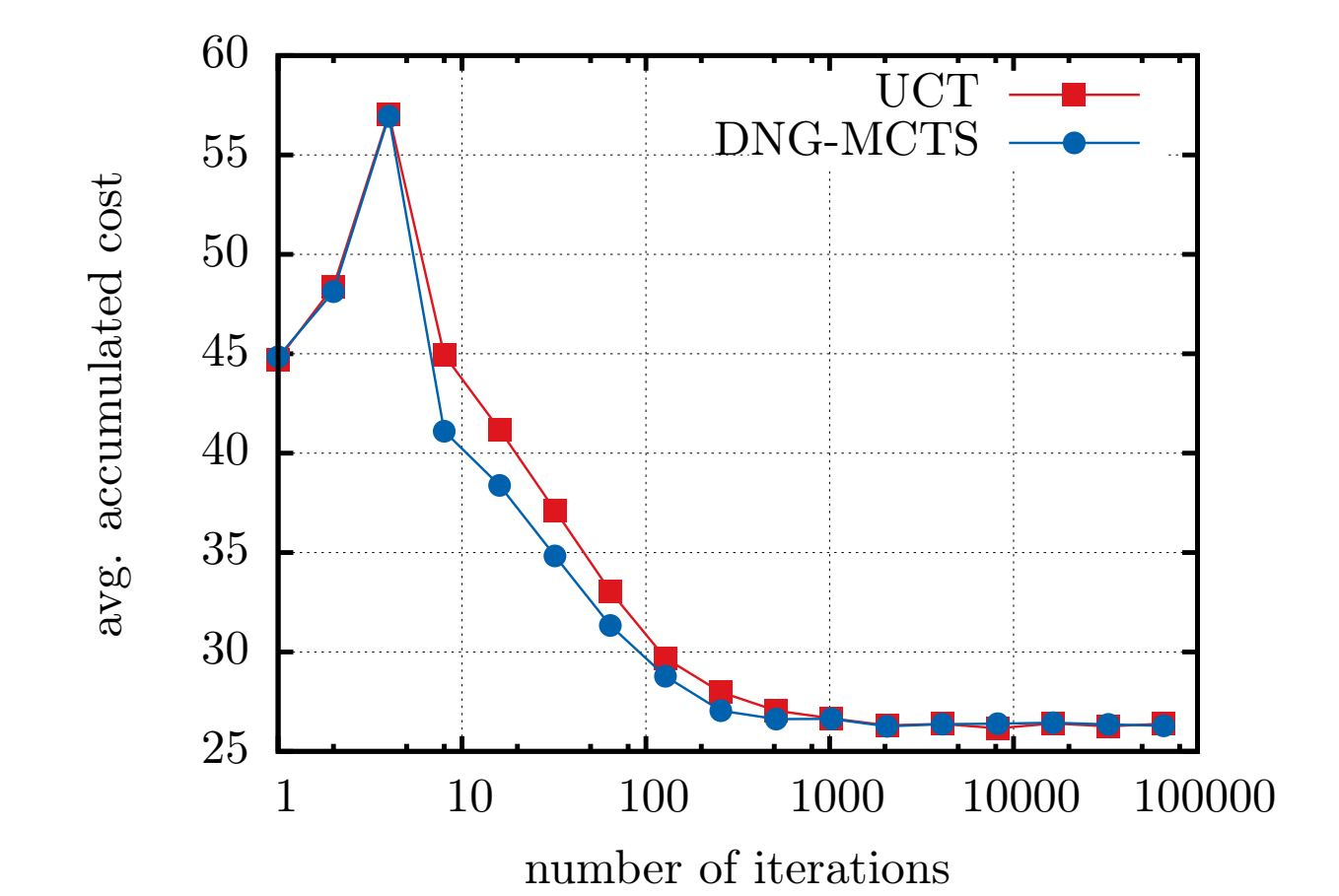


Figure 4: Sailing-100 x 100 with random policy

Table 1: CTP problems with 20 nodes. The second column indicates the belief size of the transformed MDP for each problem instance. UCTB and UCTO are the two domain-specific UCT implementations. DNG-MCTS and UCT run for 10,000 iterations. Boldface fonts are best in whole table; gray cells show best among domain-independent implementations for each group.

prob.	belief	domain-specific UCT		random rollout policy		optimistic rollout policy	
		UCTB	UCTO	UCT	DNG	UCT	DNG
20-1	$20 \times 3^{49}$	210.7±7	<b>169.0±6</b>	216.4±3	223.9±4	180.7±3	177.1±3
20-2	$20 \times 3^{49}$	176.4±4	<b>148.9±3</b>	178.5±2	178.1±2	160.8±2	155.2±2
20-3	$20 \times 3^{51}$	150.7±7	<b>132.5±6</b>	169.7±4	159.5±4	144.3±3	140.1±3
20-4	$20 \times 3^{49}$	264.8±9	<b>235.2±7</b>	264.1±4	266.8±4	238.3±3	242.7±4
20-5	$20 \times 3^{52}$	123.2±7	<b>111.3±5</b>	139.8±4	133.4±4	123.9±3	122.1±3
20-6	$20 \times 3^{49}$	165.4±6	<b>133.1±3</b>	178.0±3	169.8±3	167.8±2	141.9±2
20-7	$20 \times 3^{50}$	191.6±6	<b>148.2±4</b>	211.8±3	214.9±4	174.1±2	166.1±3
20-8	$20 \times 3^{51}$	160.1±7	<b>134.5±5</b>	218.5±4	202.3±4	152.3±3	151.4±3
20-9	$20 \times 3^{50}$	235.2±6	<b>173.9±4</b>	251.9±3	246.0±3	185.2±2	180.4±2
20-10	$20 \times 3^{49}$	180.8±7	<b>167.0±5</b>	185.7±3	188.9±4	178.5±3	170.5±3
total		1858.9	<b>1553.6</b>	2014.4	1983.68	1705.9	1647.4

## CONTACT INFORMATION

Web <http://home.ustc.edu.cn/~baj>

Email [baj@mail.ustc.edu.cn](mailto:baj@mail.ustc.edu.cn)