
Monte-Carlo Tree Search vs. Model-Predictive Controller: A Track-Following Example

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Monte-Carlo Tree Search (MCTS) has achieved remarkable success in the game
2 of Go. However, most success of MCTS is in games where actions are discrete.
3 For autonomous driving, the vehicle action such as throttle and steering angle is
4 continuous. To fill the gap, we propose an MCTS algorithm for continuous actions,
5 and used it specially for a track-following scenerio. We compared MCTS with
6 a standard Model Predictive Controller (MPC) on the Udacity simulator. Using
7 the same cost function and system model, this MCTS algorithm achieves a much
8 lower cost than MPC. MCTS drives with an adaptive speed, as well as exhibits a
9 braking behavior in sharp turns. MPC drives almost a constant speed regardless of
10 the curvy track.

11 1 Introduction

12 Autonomous driving aims to make cars safer. Nearly 1.3 million people die in road crashes each
13 year, on average 3,287 deaths a day. Road crashes cost USD \$518 billion globally, costing individual
14 countries from 1-2% of their annual GDP.¹ So far there are three major avenues for autonomous driv-
15 ing. **The classical approach** extracts perception and localization results from sensors, summarizing
16 into geometry relationship of the car with its environment. Based on the *geometry representation*
17 of the world, a controller is built. This approach is so far the most popular and widely adopted by
18 industrial leaders such as Google, Uber and Baidu. **The learning-from-demonstration approach**,
19 started from the simple full-connected neural networks [Pomerleau, 1989] in the old days to recent
20 deep convolution layers by NVIDIA [Bojarski et al., 2016], regresses the steer angle given the
21 camera view. This approach leads to a simpler architecture for autonomous driving. **The affordance**
22 **approach**[Chen et al., 2015] predicts relevant geometry features (called “affordances”) from images.
23 Based on the predicted features, a controller can be developed. This approach bears some similarity
24 to Palovian control in which animals map predictions of events into behaviors[Modayil and Sutton,
25 2014].

26 Besides these exciting progress, it is interesting to bring reinforcement learning to autonomous
27 driving. Reinforcement learning achieved remarkable success in Atari games [Mnih et al., 2015] and
28 Go [Silver et al., 2016]. Recently, Mobileye proposed an interesting architecture for autonomous
29 vehicles. Similar to the classical approach, their achitecture also has two layers. In particular, their
30 high-level path planning is implemented using a recurrent neural network over the trajectory of the
31 car [Shalev-Shwartz et al., 2016]. The low-level control is a model-based approach that learns a
32 model for the state transition in response to the car’s action. Mobileye’s efforts stand for extending
33 model-based reinforcement learning [Sutton et al., 2008, Yao and Szepesvári, 2012, Grünewälder
34 et al., 2012] to autonomous driving. Modeling the state that the car sees next turns out to be very

¹<http://asirt.org/initiatives/informing-road-users/road-safety-facts/road-crash-statistics>

35 important for cars although there has not been convincing applications published yet. However,
 36 considerable progress was made in video prediction [Zeng et al., 2017, Oh et al., 2015], which can
 37 be possibly used on cars. The reward function is also a fundamental issue to bring reinforcement
 38 learning to autonomous driving. In games, the reward signal is noise free since win or loss signals
 39 can be observed as a delayed but groundtruth reward. Although how to learn a reward function for
 40 autonomous driving is still an open problem, Hadfield-Menell et al. explored teaching a car to align
 41 with a human driver with his reward function.

42 In this paper, we study Monte-Carlo Tree Search (MCTS) for an autonomous driving setting. MCTS is
 43 especially advantageous for large and complex decision making problems, as demonstrated in the
 44 competition of AlphaGo against Mr. Lee Sedol². MCTS is well practised and relatively easy to
 45 implement. All the top Go programs have used MCTS for a decade, e.g., [Coulom, 2007, Silver,
 46 2009, Enzenberger et al., 2010]. So far the success of MCTS is largely in board games where actions
 47 are discrete. However, in autonomous driving a car’s actions like throttle, braking and steer angles are
 48 all continuous. We consider a simple motion planning setting where a car has been given a trajectory
 49 to follow, and its goal is to drive within track boundary. Note in our problem, motion planing is by no
 50 means to be realistic. Practical motion planning also considering avoiding obstacle, e.g., [Kuwata
 51 et al., 2008]. We aim to have an environment that renders a simple cost function and vehicle model
 52 under which comparing the performance of MCTS and MPC is easy.

53 2 Background

54 The classical approach is so far the most practiced and mature. A two-level architecture for au-
 55 tonomous vehicle is often used: path planning at a high level and vehicle control (with a target path
 56 and speed) at a low level [Paden et al., 2016, Berntorp, 2017]. There are a spectrum of methods for
 57 each of the problems. For example, Rapid-exploring Random Trees finds feasible trajectories for
 58 robots with high degrees of freedom [Lavelle, 1998, Kuwata et al., 2008]. MPC is classical control
 59 method [Garcia et al., 1989], and has been used for motion planning in a short time horizon [Paden
 60 et al., 2016, Kim et al., 2014, Omar et al., 1998, Yim and Oh, 2004, Raffo et al., 2009, Ng et al., 2003,
 61 Bakker et al., 1987, Kong et al., 2015, Rajamani, 2011, Besselmann and Morari, 2009, Levinson
 62 et al., 2011, Urmson et al., 2007].

63 2.1 Model Predictive Controller

64 MPC defines a cost function that considers N steps ahead. This cost function is essentially the
 65 undiscounted, N -step truncated return. With a dynamics model of the car and an initial state (and
 66 optionally some constraints), MPC produces a sequence of N actions to maximize the return by
 67 calling a large-scale nonlinear programming algorithm called *Ipopt*.

68 The state of the problem is given by $[x, y, \psi, v]$, where x, y are the x-y coordinates of the car, ψ and
 69 v are the orientation and speed of the car. Let δ denote the distance of the car center to the track axis
 70 (defined by the regressed line from the reference points), and ω denotes the angle between the car
 71 and the track axis. The reward function has the following form:

$$r(s_k, a) = -(w_{tr}\delta_{k+1}(a))^2 + w_{ang}\omega_{k+1}(a)^2 + w_v(v_{k+1}(a) - v^*)^2 + w_{sta}[steer]^2 + w_{thra}[throttle]^2 \\ + w_{steerd}(a[steer] - a_{k-1}[steer])^2 + w_{throt}(a[throttle] - a_{k-1}[throttle])^2 \quad (1)$$

72 where v^* is the target speed.

73 At each time step t , MPC receives a number of reference points which are usually provided by a
 74 high-level behavior planner. A polynomial line fitting the reference points is estimated. The model
 75 used by MPC is a system of equations:

$$x_{k+1} = x_k + v_k \cos(\psi_k) t_\delta \\ y_{k+1} = y_k + v_k \sin(\psi_k) t_\delta \\ \psi_{k+1} = \psi_k + \frac{a_k [steer] v_k}{L_f} t_\delta \\ v_{k+1} = v_k + a_k [throttle] t_\delta,$$

²<https://deepmind.com/research/alphago/>

input : An action model A and a return function that considers N steps of future rewards.

output : A policy that maximizes the return function.

Initialize the state s_0 and the action a_0 .

```

for  $t = 0, 1, \dots$  do
  Observe state  $s_t$ 
  Receive a number of reference points, and fit a polynomial line
  /* Search over  $N_p$  paths */
  for  $p = 0, \dots, N_p$  do
    Set  $\tilde{s}_0 = s_t, a' = a_t$  /* each path starts with the current state and action */
    Set  $R(p) = 0$ 
    /* planning into future  $N$  steps */
    for  $k = 0, \dots, N$  do
      Sample  $a$  from a distribution  $u(a')$ 
      Predict the next state,  $\tilde{s}_{k+1} = A(\tilde{s}_k, a)$ 
      Compute the reward  $r$  according to  $\tilde{s}_{k+1}, a, a'$ , and deviation from the reference line
      Update  $R(p) = \gamma R(p) + r$ 
      Set  $a' = a$ 
    end
  end
  Select the best path (with highest return  $R$ )
  Set  $a_t$  to the first action in the best path.
  Take action  $a_t$ 
end

```

Algorithm 1: Continuity-preserved (Monte-Carlo) Tree Search (CPTS for short). This CPTS algorithm generates a number of paths expanded from continuous actions. In expanding the child trees from a state, we enforce the continuity in the actions going down a tree. The algorithm makes sure that driving action does not change abruptly and reduces the search space significantly.

76 where t_δ is the latency between two time steps. Based on the fitted line and the model, the deviation
 77 measures δ and ω can both be computed after each of all the forward N actions. In a summary, MPC
 78 solves an optimization problem (using Ipopt) to obtain a sequence of N actions:

$$a_{0:N-1} = \arg \max_{a_{0:N-1}} R = \arg \max_{a_{0:N-1}} \sum_{k=0}^{N-1} r(\tilde{s}_k, a),$$

79 where \tilde{s}_0 was set to the current state s_t . In practice of MPC, usually only the first action a_0 was used
 80 (so MPC calls the optimization procedure every time step).

81 2.2 Monte-Carlo Tree Search

82 MCTS is a special policy search algorithm. Policy search algorithms search in the action space
 83 directly. Comparing to policy gradient methods, policy search algorithms can find global optima [??],
 84 while policy gradient methods usually suffer from local minima [??].

85 Continuity-preserved Monte-Carlo Tree Search (MCTS).

86 MCTS search exploration issue. UCT.

87 3 Experiments

88 In this experiment, we used Udacity simulator³. The simulator is developed in Unity and supports
 89 self driving car development. The reward function is computed from equation (1). In the experiment,
 90 the target speed (in the reward function) was set to 70 km/h. As shown in Figure 1 (left plot), MCTS
 91 achieved a much smaller cost than MPC. The cost function is a linear combination of seven cost
 92 components. MCTS achieved both a smaller velocity cost and a smaller trackPos cost (deviation from
 93 the track center) as shown in the second plot. The third plot in Figure 1 shows the velocity on the
 94 track. MCTS's speed is more variable due to that the track has quite a few sharp turns. MPC, on the

³<https://github.com/udacity/self-driving-car-sim>

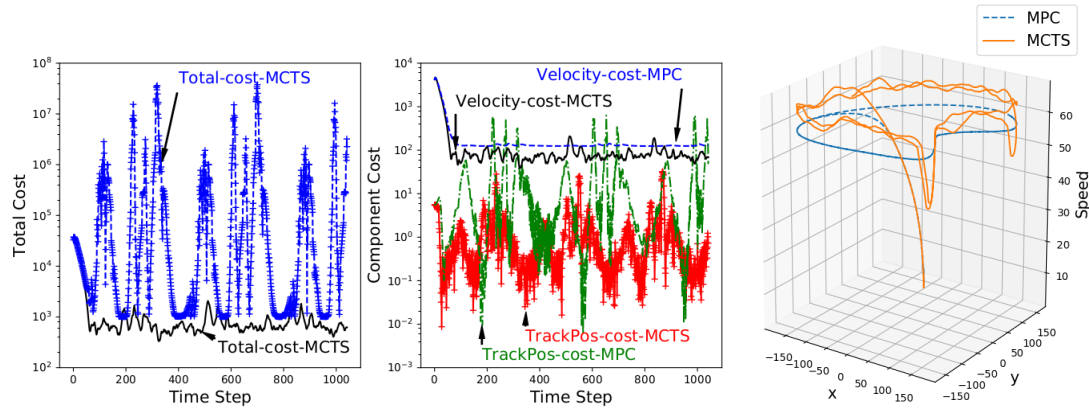


Figure 1: MCTS vs. MPC. The left plot shows the total cost over the next 8 time steps of the two algorithms. The middle plot shows the cost of the velocity component and the trackPos (distance to the center of the lane) component in the same run as the left plot. The right plot shows the velocity on the track, which shows that (a) MCTS accelerates faster in the beginning (the ascending curves from bottom), (b) MCTS drives closer to the target speed (70 km/h) than MPC most of the time. and (c) MCTS’s control is more adaptive to curvature in the track. In particular, at sharp turns we see speed dip in the orange line while MPC drives at almost a constant speed (58 km/h).



Figure 2: MCTS braking in front of a sharp turn in Udacity Simulator. MPC never shows braking behavior in the experiment.

95 other hand, drives almost a constant speed. In particular, after the beginning acceleration period, MPC
 96 drove between 57.8 km/h and 58.6 km/h with an average speed of 58.4 km/h. MCTS drove between
 97 42.8 km/h and 67.5 km/h, averaging at 62.3 km/h. Interestingly, MCTS shows braking behavior
 98 (continually negative throttles) ahead of sharp turns (Figure 2) although it was never explicitly trained
 99 to do so. In contrast, MPC never braked.

100 For MCTS, 10000 paths were generated with random samples of action (both steer angle and throttle).
 101 In generating the actions, uniform samples in the small range of last steer angle and last throttle
 102 were independently drawn. Simulation were run with lookahead depth of 8 (same as MPC). We used
 103 exactly the same model and reward function as MPC. The weights for the reward function are, $w_{tr} =$
 104 10.0 , $w_{ang} = 50.0$, $w_v = 1.0$, $w_{st} = 10.0$, $w_{thr} = 3000.0$, $w_{steerd} = 10.0$, $w_{throtld} = 3000.0$.

105 Video of MPC driving:

106 <http://youtube.com/mpc.mp4>

107 Video of MCTS driving:

108 http://youtube.com/pure_mcts_gamm0.9_steerthrotchaining.mp4.mp4

References

- 109
110 Egbert Bakker, Lars Nyborg, and Hans B Pacejka. Tyre modelling for use in vehicle dynamics studies.
111 Technical report, SAE Technical Paper, 1987.
- 112 Karl Berntorp. Path planning and integrated collision avoidance for autonomous vehicles. In *American*
113 *Control Conference (ACC)*, pages 4023–4028, 05 2017.
- 114 Thomas Besselmann and M Morari. Autonomous vehicle steering using explicit lqv-mpc. pages
115 2628–2633, 08 2009.
- 116 Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseon
117 Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao,
118 and Karol Zieba. End to end learning for self-driving cars. *CoRR*, abs/1604.07316, 2016. URL
119 <http://arxiv.org/abs/1604.07316>.
- 120 Chenyi Chen, Ari Seff, Alain L. Kornhauser, and Jianxiong Xiao. DeepDriving: Learning affordance
121 for direct perception in autonomous driving. In *ICCV*, 2015.
- 122 Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *Proceedings*
123 *of the 5th International Conference on Computers and Games, CG’06*, pages 72–83, Berlin,
124 Heidelberg, 2007. Springer-Verlag. ISBN 3-540-75537-3, 978-3-540-75537-1. URL <http://dl.acm.org/citation.cfm?id=1777826.1777833>.
- 126 Markus Enzenberger, Martin Müller, B Arneson, and Richard Segal. Fuego - an open-source
127 framework for board games and go engine based on monte carlo tree search. 2:259–270, 01 2010.
- 128 C. E. Garcia, D. M. Prett, and M. Morari. Model predictive control: Theory and practice—a
129 survey. *Automatica*, 25(3):335–348, May 1989. ISSN 0005-1098. doi: 10.1016/0005-1098(89)
130 90002-2. URL [http://dx.doi.org/10.1016/0005-1098\(89\)90002-2](http://dx.doi.org/10.1016/0005-1098(89)90002-2).
- 131 Steffen Grünewälder, Guy Lever, Luca Baldassarre, Massimiliano Pontil, and Arthur Gretton. Mod-
132 elling transition dynamics in mdps with rkhs embeddings. In *ICML*, pages 1603–1610, USA, 2012.
133 Omnipress. ISBN 978-1-4503-1285-1.
- 134 Dylan Hadfield-Menell, Anca D. Dragan, Pieter Abbeel, and Stuart J. Russell. Cooperative inverse
135 reinforcement learning. *CoRR*, abs/1606.03137, 2016. URL <http://arxiv.org/abs/1606.03137>.
- 137 E Kim, J Kim, and Myoung-ho Sunwoo. Model predictive control strategy for smooth path tracking
138 of autonomous vehicles with steering actuator dynamics. 15:1155–1164, 12 2014.
- 139 Jason Kong, Mark Pfeiffer, Georg Schildbach, and Francesco Borrelli. Kinematic and dynamic
140 vehicle models for autonomous driving control design. In *Intelligent Vehicles Symposium (IV),*
141 *2015 IEEE*, pages 1094–1099. IEEE, 2015.
- 142 Yoshiaki Kuwata, Gaston Fiore, Justin Teo, Emilio Frazzoli, and Jonathan How. Motion planning
143 for urban driving using rrt. In *2008 IEEE/RSJ International Conference on Intelligent Robots and*
144 *Systems, IROS*, pages 1681–1686, 09 2008.
- 145 Steven M. Lavalle. Rapidly-exploring random trees: A new tool for path planning. Technical report,
146 Iowa State University, 1998.
- 147 Jesse Levinson, Jake Askeland, Jan Becker, Jennifer Dolson, David Held, Soeren Kammel, J Zico
148 Kolter, Dirk Langer, Oliver Pink, Vaughan Pratt, et al. Towards fully autonomous driving: Systems
149 and algorithms. In *Intelligent Vehicles Symposium (IV), 2011 IEEE*, pages 163–168. IEEE, 2011.
- 150 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-
151 mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen,
152 Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra,
153 Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning.
154 *Nature*, 518(7540):529–533, 02 2015. URL <http://dx.doi.org/10.1038/nature14236>.
- 155 Joseph Modayil and Richard S. Sutton. Prediction driven behavior: Learning predictions that drive
156 fixed responses. In *AAAI Workshop on AI and Robotics*, 2014.

- 157 Andrew Y Ng, H Jin Kim, Michael I Jordan, Shankar Sastry, and Shiv Ballianda. Autonomous
158 helicopter flight via reinforcement learning. In *NIPS*, volume 16, 2003.
- 159 Junhyuk Oh, Xiaoxiao Guo, Honglak Lee, Richard L Lewis, and Satinder Singh. Action-Conditional
160 Video Prediction using Deep Networks in Atari Games. In C Cortes, N D Lawrence, D D Lee,
161 M Sugiyama, R Garnett, and R Garnett, editors, *NIPS 28*, pages 2845–2853. Curran Associates,
162 Inc., 2015.
- 163 T Omar, A Eskandarian, and N Bedewi. Vehicle crash modelling using recurrent neural networks.
164 *Mathematical and computer Modelling*, 28(9):31–42, 1998.
- 165 Brian Paden, Michal Cap, Sze Zheng Yong, Dmitry S. Yershov, and Emilio Frazzoli. A survey of
166 motion planning and control techniques for self-driving urban vehicles. *CoRR*, abs/1604.07446,
167 2016. URL <http://arxiv.org/abs/1604.07446>.
- 168 Dean A. Pomerleau. Alvin, an autonomous land vehicle in a neural network. Technical report,
169 Carnegie Mellon University, 1989. URL [http://repository.cmu.edu/cgi/viewcontent.
170 cgi?article=2874&context=compsci](http://repository.cmu.edu/cgi/viewcontent.cgi?article=2874&context=compsci).
- 171 G.V. Raffo, Guilherme K. Gomes, Julio Normey-Rico, Christian R. Kelber, and Leandro B. Becker.
172 A predictive controller for autonomous vehicle path tracking. 10:92 – 102, 04 2009.
- 173 Rajesh Rajamani. *Vehicle dynamics and control*. Springer Science & Business Media, 2011.
- 174 Shai Shalev-Shwartz, Nir Ben-Zrihem, Aviad Cohen, and Amnon Shashua. Long-term planning
175 by short-term prediction. *CoRR*, abs/1602.01580, 2016. URL [http://arxiv.org/abs/1602.
176 01580](http://arxiv.org/abs/1602.
176 01580).
- 177 David Silver. *Reinforcement Learning and Simulation-Based Search in Computer Go*. PhD thesis,
178 2009.
- 179 David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driess-
180 che, Julian Schrittwieser, Ioannis Antonoglou, Vedavyas Panneershelvam, Marc Lanctot, Sander
181 Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy P. Lillicrap,
182 Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the
183 game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016. doi:
184 10.1038/nature16961. URL <http://dx.doi.org/10.1038/nature16961>.
- 185 Richard Sutton, Csaba Szepesvari, Alborz Geramifard, and Michael Bowling. Dyna-style planning
186 with linear function approximation and prioritized sweeping. In *UAI*, pages 528–536, 2008.
- 187 Chris Urmson, J Andrew Bagnell, Christopher R Baker, Martial Hebert, Alonzo Kelly, Raj Rajkumar,
188 Paul E Rybski, Sebastian Scherer, Reid Simmons, Sanjiv Singh, et al. Tartan racing: A multi-modal
189 approach to the darpa urban challenge. 2007.
- 190 Hengshuai Yao and Csaba Szepesvari. Approximate policy iteration with linear action models. In
191 *AAAI*, pages 1212–1217, 2012.
- 192 Young Uk Yim and Se-Young Oh. Modeling of vehicle dynamics from real vehicle measurements
193 using a neural network with two-stage hybrid learning for accurate long-term prediction. *IEEE
194 Transactions on Vehicular Technology*, 53(4):1076–1084, 2004.
- 195 Kuo-Hao Zeng, William B. Shen, De-An Huang, Min Sun, and Juan Carlos Nieves. Visual forecasting
196 by imitating dynamics in natural sequences. *CoRR*, abs/1708.05827, 2017. URL [http://arxiv.
197 org/abs/1708.05827](http://arxiv.
197 org/abs/1708.05827).