
Monte-Carlo Tree Search vs. Model-Predictive Controller: A Lane-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 lane-following scenario. We compared MCTS with a
6 standard Model Predictive Controller (MPC) on the Udacity simulator. Using the
7 same cost function and system model, this MCTS algorithm achieves a much lower
8 cost than MPC. MCTS drives with an adaptive speed, as well as exhibits a braking
9 behavior in sharp turns. MPC drives a nearly constant speed regardless of the curvy
10 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 Pavlovian 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 architecture 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 ground truth feedback. 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. Brechtel et al. modeled the car’s environment using an
 42 Markov Decision Process (MDP) in which the state space is equidistant cells of the coordinates on the
 43 road, and the transition probabilities are approximated using a Dynamic Bayesian Networks. Their
 44 empirical studies show that the car can coordinate well when to overtake according to oncoming traffic.
 45 Exploring in a driving environment is challenging because such exploration (normally practiced in
 46 reinforcement learning without constraint) must be safety guaranteed. Recently, Mnih et al. proposed
 47 an asynchronous reinforcement learning framework that lets a number of learning agents run in
 48 parallel aiming to explore different parts of the environment. Their algorithm on a simulated driving
 49 environment achieved near to a human driver with only 12 hours of training. It is interesting to see
 50 whether this new framework can solve the specific exploration constraint in autonomous driving.

51 In this paper, we study Monte-Carlo Tree Search (MCTS) for an autonomous driving setting. MCTS
 52 is especially advantageous for large and complex decision making problems, as demonstrated in the
 53 competition of AlphaGo against Mr. Lee Sedol ². MCTS is well practiced and relatively easy to
 54 implement. All the top Go programs have used MCTS for a decade, e.g., [Coulom, 2007, Silver,
 55 2009, Enzenberger et al., 2010]. So far the success of MCTS is largely in board games where actions
 56 are discrete. However, in autonomous driving a car’s actions like throttle, braking and steer angles are
 57 all continuous. We consider a simple motion planning setting where a car has been given a trajectory
 58 to follow, and its goal is to drive within track boundary. Our treatment of motion planning is by no
 59 means to be realistic. Practical motion planning also considering avoiding obstacle, e.g., [Kuwata
 60 et al., 2008]. We aim to have an environment that renders a simple cost function and vehicle model
 61 under which evaluating performances of algorithms is relatively easy. We propose an extension of
 62 pure MCTS to continuous actions, and compare it with an off-the-shelf Model Predictive Controller
 63 (MPC).

64 2 Background

65 The classical approach is so far the most practiced and mature. A two-level architecture for au-
 66 tonomous vehicle is often used: path planning at a high level and vehicle control (with a target path
 67 and speed) at a low level [Paden et al., 2016, Berntorp, 2017]. There are a spectrum of methods for
 68 each of the problems. For example, Rapid-exploring Random Trees finds feasible trajectories for
 69 robots with high degrees with freedom [Lavalle, 1998, Kuwata et al., 2008]. MPC is classical control
 70 method [Garcia et al., 1989], and has been widely used for motion planning in a short time horizon
 71 [Paden et al., 2016, Kim et al., 2014, Omar et al., 1998, Yim and Oh, 2004, Raffo et al., 2009, Ng
 72 et al., 2003, Bakker et al., 1987, Kong et al., 2015, Rajamani, 2011, Besselmann and Morari, 2009,
 73 Levinson et al., 2011, Urmsen et al., 2007]. MPC is a major research field on its own and this section
 74 provides the application context of MPC for autonomous driving, especially lane following.

75 2.1 The Model and the Problem

76 The car’s dynamics is represented by a practical model, often referred to as the *Kinematic model*. In
 77 this model, the two wheels of the car are connected by a rigid link. The state of the car is given by
 78 $[x, y, \psi, v]$, where x, y are the x-y coordinates of the car, ψ and v are the orientation and speed of the
 79 car, respectively. The model can be expressed by,

$$\begin{aligned}
 \dot{x} &= v \cos(\psi) \\
 \dot{y} &= v \sin(\psi) \\
 \dot{\psi} &= \frac{a[\textit{steer}]v}{L_f} \\
 \dot{v} &= a[\textit{throttle}],
 \end{aligned} \tag{1}$$

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

80 where L_f is the distance between the two front wheels. This is discretized using Euler method in
 81 practice. Without loss of generality, we denote the model by an equation: $s_{k+1} = A(s_k, a_k)$, where
 82 s is a vector of the state variables and a is a vector of action variables (steer and throttle). In our
 83 problem, at each time step, an agent (MPC or MCTS) receives a number of reference coordinate
 84 points. These points are often provided by a high-level trajectory planner. The agent is also given the
 85 a distance measures δ , the distance of the car’s center to the track axis; an angle deviation measure, ω ,
 86 the difference between the car’s heading angle (ψ) and the track axis direction. In practice, both δ
 87 and ω are computed by first regressing a polynomial line from the reference points. The goal of the
 88 agents is to drive close a target speed v^* within the track. Specifically at each time step k , the agent
 89 selects an action a . Afterwards it receives a cost signal that is computed from the following equation:

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

90 2.2 Model Predictive Controller

91 MPC assumes to know the form of the cost function in equation 2. It defines a cost function that
 92 considers N steps ahead. This cost function is essentially the undiscounted, N -step truncated return.
 93 MPC produces a sequence of N actions to minimize the cost function

$$a_{0:N-1} = \arg \min_{a_{0:N-1}} R = \arg \min_{a_{0:N-1}} \sum_{t=0}^{N-1} r(\tilde{s}_t, a),$$

94 where \tilde{s}_0 was set to the current state s_k . Common practice of MPC is to use the interior point
 95 optimizer [Paden et al., 2016].

96 2.3 Monte-Carlo Tree Search

97 MCTS is a special policy search algorithm. Comparing to other reinforcement learning methods,
 98 policy search algorithms can find global optima [Valko et al., 2013, Munos, 2014]. MCTS algorithms
 99 are designed for discrete actions. In the Upper Confidence Tree (UCT) algorithm [Kocsis and
 100 Szepesvári, 2006], the actions are treated as the arms in a multi-armed bandit problem and the
 101 frequency of actions is used to measure the knowledge of the actions according to which the
 102 exploration term in selecting the action is determined. Practice and theory of MCTS for continuous
 103 actions is largely a gap. A number of recent advances aim to generalize across actions. In particular,
 104 [Couetoux et al., 2011, Yee et al., 2016] explored generalizing in actions from already exploited
 105 actions. HOOT replaces the UCB algorithm in UCT with a continuous action selection procedure
 106 [Mansley et al., 2011]. In this paper, we aimed to first extend pure MCTS for autonomous driving.

107 3 A New MCTS Algorithm

108 Besides UCT and its variants, there is a basic MCTS algorithm which works by playing a number
 109 of random games to the end; and the moves that achieves the best game scores are chosen. This
 110 algorithm is often referred to as the pure MCTS. This algorithm by definition can be extended to
 111 continuous actions in a straight forward way. However, it is very inefficient to sample the actions in a
 112 random fashion especially in real time.

113 Based on the observation that in everyday driving our steering and throttle control is continuous, we
 114 propose the following tree search algorithm. Similar to the pure MCTS, at each depth of the tree
 115 we only sample one node. This essentially searches over paths instead of trees. We will generate a
 116 number of paths expanded from continuous actions. Specifically, in expanding the path from a state,
 117 we enforce the continuity in the actions going down the path. This small trick reduces the search space
 118 significantly. The algorithm is shown in Algorithm 1. The sampling distribution $u(a')$ for the current
 119 time step is incrementally adjusted according to the last action a' . In the experiment, we used a
 120 uniform distribution over a neighborhood of a' . In general u is not limited to the uniform distribution
 121 but it can incorporate knowledge we gained through exploitation and be learned automatically.

input : A system model A and a cost function that considers N steps of future costs.
output : A policy that minimize the cost 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  $\tilde{s}_0 = s_t, a' = a_{t-1}$  /* each path starts with the current state and last 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 cost  $r$  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 lowest cost  $R$ 
  Set  $a_t$  to the first action in the best path.
  Take action  $a_t$ 
end

```

Algorithm 1: Continuity-preserved (Monte-Carlo) Tree Search for following lane.

122 4 Experiments

123 In the experiment, we used the Udacity simulator³. The simulator is developed by Unity to support
124 self-driving car development. Both algorithms used the same model in equation 1 and the same
125 cost function in equation 2. Simulation for both algorithms was run with lookahead depth of 8.
126 The weights for the cost function are, $w_{tr} = 10.0, w_{ang} = 50.0, w_v = 1.0, w_{st} = 10.0, w_{thr} =$
127 $3000.0, w_{steerd} = 10.0, w_{throtld} = 3000.0$.

128 For both MPC and MCTS, only the first action a_0 was used although N actions were produced at a
129 single time step. In the experiment, the target speed was set to 70 km/h.

130 As shown in Figure 1 (left plot), MCTS achieved a much smaller cost than MPC. The cost function
131 is a linear combination of seven cost components. MCTS achieved both a smaller speed cost and a
132 smaller “trackPos” cost (deviation from the track center) than MPC most of the time as shown in the
133 second plot.

134 The third plot in Figure 1 is the speeds of the agents on the track, which shows that (a) MCTS
135 accelerates faster in the beginning (the ascending curves from bottom); (b) MCTS drives closer to the
136 target speed (70 km/h) than MPC most of the time; (c) MCTS’s control is more adaptive to curvature
137 in the track. In particular, at sharp turns we see speed dip in the orange line while MPC drives at
138 almost a constant speed. After the beginning acceleration period, MPC drove between 57.8 km/h
139 and 58.6 km/h with an average speed of 58.4 km/h. MCTS drove between 42.8 km/h and 67.5 km/h,
140 averaging at 62.3 km/h.

141 Interestingly, in the experiment we observed that MCTS showed braking behavior (continually
142 negative throttles) ahead of sharp turns (Figure 2) although it was never explicitly trained to do so.
143 In contrast, MPC never braked. For MCTS, 10,000 paths were generated. The $u(a')$ is a uniform
144 distribution over a close neighborhood of last steer angle and last throttle. In particular, the steer angle
145 and throttle were independently drawn uniformly from the intervals, $(a'[steer] - 0.02, a'[steer] +$
146 $0.02)$ and $(a'[throttle] - 0.2, a'[throttle] + 0.2)$, respectively.

147 We produced videos of MCTS driving

148 <https://youtu.be/YP7qPJSJAVU>

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

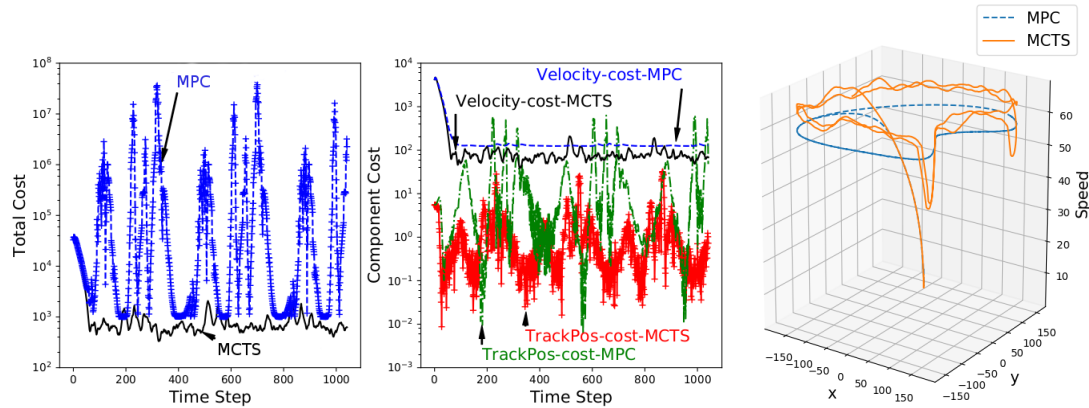


Figure 1: MCTS vs. MPC. The left plot compares the total cost over the next 8 time steps of the two algorithms. The middle plot compares the cost of the speed component and the “trackPos” (distance to the center of the lane) component. The right plot compares the speeds of the two agents driving on the track.



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

149 and MPC driving:

150 <https://youtu.be/SL150wMenyY>

151 5 Conclusion

152 In this paper, we proposed a primitive Monte-Carlo Tree Search algorithm for following lane in
 153 autonomous driving. The algorithm is inspired by that in driving action change is smooth: the next
 154 steering or throttle is a small modification to the existing values. In a simulated driving environment,
 155 the algorithm achieves much smaller cost than a standard off-the-shelf Model Predictive Controller
 156 under the same setting. We hope this incremental step could help the efforts of bringing reinforcement
 157 learning to autonomous driving. We plan to do a more comprehensive study of Monte-Carlo Tree
 158 Search algorithms for continuous actions and use it for autonomous driving.

159 References

- 160 Egbert Bakker, Lars Nyborg, and Hans B Pacejka. Tyre modelling for use in vehicle dynamics studies.
 161 Technical report, SAE Technical Paper, 1987.
- 162 Karl Berntorp. Path planning and integrated collision avoidance for autonomous vehicles. In *American
 163 Control Conference (ACC)*, pages 4023–4028, 05 2017.
- 164 Thomas Besselmann and M Morari. Autonomous vehicle steering using explicit lqv-mpc. pages
 165 2628–2633, 08 2009.

- 166 Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon
167 Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao,
168 and Karol Zieba. End to end learning for self-driving cars. *CoRR*, abs/1604.07316, 2016. URL
169 <http://arxiv.org/abs/1604.07316>.
- 170 S. Brechtel, T. Gindele, and R. Dillmann. Probabilistic mdp-behavior planning for cars. In *14th*
171 *International Conference on Intelligent Transportation Systems, IEEE*, page 1537–1542, 2011.
- 172 Chenyi Chen, Ari Seff, Alain L. Kornhauser, and Jianxiong Xiao. DeepDriving: Learning affordance
173 for direct perception in autonomous driving. In *ICCV*, 2015.
- 174 Adrien Couetoux, Jean-Baptiste Hoock, Nataliya Sokolovska, Olivier Teytaud, and Nicolas Bonnard.
175 Continuous Upper Confidence Trees. In *LION'11: Proceedings of the 5th International Conference*
176 *on Learning and Intelligent Optimization*, page TBA, Italy, January 2011. URL [https://hal.](https://hal.archives-ouvertes.fr/hal-00542673)
177 [archives-ouvertes.fr/hal-00542673](https://hal.archives-ouvertes.fr/hal-00542673).
- 178 Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *Proceedings*
179 *of the 5th International Conference on Computers and Games, CG'06*, pages 72–83, Berlin,
180 Heidelberg, 2007. Springer-Verlag. ISBN 3-540-75537-3, 978-3-540-75537-1. URL [http:](http://dl.acm.org/citation.cfm?id=1777826.1777833)
181 [/dl.acm.org/citation.cfm?id=1777826.1777833](http://dl.acm.org/citation.cfm?id=1777826.1777833).
- 182 Markus Enzenberger, Martin Müller, B Arneson, and Richard Segal. Fuego - an open-source
183 framework for board games and go engine based on monte carlo tree search. 2:259–270, 01 2010.
- 184 C. E. Garcia, D. M. Prett, and M. Morari. Model predictive control: Theory and practice—a
185 survey. *Automatica*, 25(3):335–348, May 1989. ISSN 0005-1098. doi: 10.1016/0005-1098(89)
186 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).
- 187 Steffen Grünewälder, Guy Lever, Luca Baldassarre, Massimiliano Pontil, and Arthur Gretton. Mod-
188 elling transition dynamics in mdps with rkhs embeddings. In *ICML*, pages 1603–1610, USA, 2012.
189 Omnipress. ISBN 978-1-4503-1285-1.
- 190 Dylan Hadfield-Menell, Anca D. Dragan, Pieter Abbeel, and Stuart J. Russell. Cooperative inverse
191 reinforcement learning. *CoRR*, abs/1606.03137, 2016. URL [http://arxiv.org/abs/1606.](http://arxiv.org/abs/1606.03137)
192 [03137](http://arxiv.org/abs/1606.03137).
- 193 E Kim, J Kim, and Myoung-ho Sunwoo. Model predictive control strategy for smooth path tracking
194 of autonomous vehicles with steering actuator dynamics. 15:1155–1164, 12 2014.
- 195 Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *Proceedings of the*
196 *17th European Conference on Machine Learning, ECML'06*, pages 282–293, Berlin, Heidelberg,
197 2006. Springer-Verlag. ISBN 3-540-45375-X, 978-3-540-45375-8. doi: 10.1007/11871842_29.
198 URL http://dx.doi.org/10.1007/11871842_29.
- 199 Jason Kong, Mark Pfeiffer, Georg Schildbach, and Francesco Borrelli. Kinematic and dynamic
200 vehicle models for autonomous driving control design. In *Intelligent Vehicles Symposium (IV),*
201 *2015 IEEE*, pages 1094–1099. IEEE, 2015.
- 202 Yoshiaki Kuwata, Gaston Fiore, Justin Teo, Emilio Frazzoli, and Jonathan How. Motion planning
203 for urban driving using rrt. In *2008 IEEE/RSJ International Conference on Intelligent Robots and*
204 *Systems, IROS*, pages 1681–1686, 09 2008.
- 205 Steven M. Lavalle. Rapidly-exploring random trees: A new tool for path planning. Technical report,
206 Iowa State University, 1998.
- 207 Jesse Levinson, Jake Askeland, Jan Becker, Jennifer Dolson, David Held, Soeren Kammel, J Zico
208 Kolter, Dirk Langer, Oliver Pink, Vaughan Pratt, et al. Towards fully autonomous driving: Systems
209 and algorithms. In *Intelligent Vehicles Symposium (IV), 2011 IEEE*, pages 163–168. IEEE, 2011.
- 210 Christopher R. Mansley, Ari Weinstein, and Michael L. Littman. Sample-based planning for continu-
211 ous action markov decision processes. In *21st International Conference on Auto- mated Planning*
212 *and Scheduling, ICAPS*, 2011.

- 213 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-
214 mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen,
215 Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra,
216 Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning.
217 *Nature*, 518(7540):529–533, 02 2015. URL <http://dx.doi.org/10.1038/nature14236>.
- 218 Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim
219 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement
220 learning. *CoRR*, abs/1602.01783, 2016. URL <http://arxiv.org/abs/1602.01783>.
- 221 Joseph Modayil and Richard S. Sutton. Prediction driven behavior: Learning predictions that drive
222 fixed responses. In *AAAI Workshop on AI and Robotics*, 2014.
- 223 Rémi Munos. From bandits to monte-carlo tree search: The optimistic principle applied to opti-
224 mization and planning. *Foundations and Trends in Machine Learning*, 7(1):1–129, 2014. ISSN
225 1935-8237. doi: 10.1561/22000000038. URL <http://dx.doi.org/10.1561/22000000038>.
- 226 Andrew Y Ng, H Jin Kim, Michael I Jordan, Shankar Sastry, and Shiv Ballianda. Autonomous
227 helicopter flight via reinforcement learning. In *NIPS*, volume 16, 2003.
- 228 Junhyuk Oh, Xiaoxiao Guo, Honglak Lee, Richard L Lewis, and Satinder Singh. Action-Conditional
229 Video Prediction using Deep Networks in Atari Games. In C Cortes, N D Lawrence, D D Lee,
230 M Sugiyama, R Garnett, and R Garnett, editors, *NIPS 28*, pages 2845–2853. Curran Associates,
231 Inc., 2015.
- 232 T Omar, A Eskandarian, and N Bedewi. Vehicle crash modelling using recurrent neural networks.
233 *Mathematical and computer Modelling*, 28(9):31–42, 1998.
- 234 Brian Paden, Michal Cáp, Sze Zheng Yong, Dmitry S. Yershov, and Emilio Frazzoli. A survey of
235 motion planning and control techniques for self-driving urban vehicles. *CoRR*, abs/1604.07446,
236 2016. URL <http://arxiv.org/abs/1604.07446>.
- 237 Dean A. Pomerleau. Alvin, an autonomous land vehicle in a neural network. Technical report,
238 Carnegie Mellon University, 1989. URL <http://repository.cmu.edu/cgi/viewcontent.cgi?article=2874&context=compsci>.
- 240 G.V. Raffo, Guilherme K. Gomes, Julio Normey-Rico, Christian R. Kelber, and Leandro B. Becker.
241 A predictive controller for autonomous vehicle path tracking. 10:92 – 102, 04 2009.
- 242 Rajesh Rajamani. *Vehicle dynamics and control*. Springer Science & Business Media, 2011.
- 243 Shai Shalev-Shwartz, Nir Ben-Zrihem, Aviad Cohen, and Amnon Shashua. Long-term planning
244 by short-term prediction. *CoRR*, abs/1602.01580, 2016. URL <http://arxiv.org/abs/1602.01580>.
- 246 David Silver. *Reinforcement Learning and Simulation-Based Search in Computer Go*. PhD thesis,
247 2009.
- 248 David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driess-
249 che, Julian Schrittwieser, Ioannis Antonoglou, Vedavyas Panneershelvam, Marc Lanctot, Sander
250 Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy P. Lillicrap,
251 Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the
252 game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016. doi:
253 10.1038/nature16961. URL <http://dx.doi.org/10.1038/nature16961>.
- 254 Richard Sutton, Csaba Szepesvari, Alborz Geramifard, and Michael Bowling. Dyna-style planning
255 with linear function approximation and prioritized sweeping. In *UAI*, pages 528–536, 2008.
- 256 Chris Urmson, J Andrew Bagnell, Christopher R Baker, Martial Hebert, Alonzo Kelly, Raj Rajkumar,
257 Paul E Rybski, Sebastian Scherer, Reid Simmons, Sanjiv Singh, et al. Tartan racing: A multi-modal
258 approach to the darpa urban challenge. 2007.
- 259 Michal Valko, Alexandra Carpentier, and Rémi Munos. Stochastic Simultaneous Optimistic Opti-
260 mization. In *International Conference on Machine Learning*, Atlanta, United States, June 2013.
261 URL <https://hal.inria.fr/hal-00789606>.

- 262 Hengshuai Yao and Csaba Szepesvári. Approximate policy iteration with linear action models. In
263 *AAAI*, pages 1212–1217, 2012.
- 264 Timothy Yee, Viliam Lisy, and Michael Bowling. Monte carlo tree search in continuous action spaces
265 with execution uncertainty. In *Proceedings of the Twenty-Fifth International Joint Conference on*
266 *Artificial Intelligence, IJCAI'16*, pages 690–696. AAAI Press, 2016. ISBN 978-1-57735-770-4.
267 URL <http://dl.acm.org/citation.cfm?id=3060621.3060718>.
- 268 Young Uk Yim and Se-Young Oh. Modeling of vehicle dynamics from real vehicle measurements
269 using a neural network with two-stage hybrid learning for accurate long-term prediction. *IEEE*
270 *Transactions on Vehicular Technology*, 53(4):1076–1084, 2004.
- 271 Kuo-Hao Zeng, William B. Shen, De-An Huang, Min Sun, and Juan Carlos Niebles. Visual forecasting
272 by imitating dynamics in natural sequences. *CoRR*, abs/1708.05827, 2017. URL <http://arxiv.org/abs/1708.05827>.
273