

Agile Catching with Whole-Body MPC and Blackbox Policy Learning

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** We address a benchmark task in agile robotics: catching objects thrown
2 at high-speed. This is a challenging task that involves tracking, intercepting, and
3 cradling a thrown object with access only to visual observations of the object and
4 the proprioceptive state of the robot, all within a fraction of a second. We present
5 the relative merits of two fundamentally different solution strategies: (i) Model
6 Predictive Control using accelerated constrained trajectory optimization, and (ii)
7 Reinforcement Learning using zeroth-order optimization. We provide insights
8 into various performance trade-offs including sample efficiency, sim-to-real transfer,
9 robustness to distribution shifts, and whole-body multimodality via extensive
10 on-hardware experiments. We conclude with proposals on fusing “classical” and
11 “learning-based” techniques for agile robot control.



Figure 1: Mobile Manipulator with Lacrosse Head catching a ball within a second. (right) Automatic ball thrower with controllable yaw angles and speed of around 5m/s.

12 1 Introduction

13 Chasing a ball in flight and completing a dramatic diving catch is a memorable moment of athleti-
14 cism - a benchmark of human agility - in several popular sports. In this paper, we consider the task
15 of tracking, intercepting and catching balls moving at high speeds on a mobile manipulator platform
16 (see Figure 1), whose end-effector is equipped with a Lacrosse head. Within a fraction of a second,
17 the robot must start continuously translating visual observations of the ball into feasible whole body
18 motions, controlling both the base and the arm in a coordinated fashion. In the final milliseconds, the
19 control system must be robust to perceptual occlusions while also executing a cradling maneuver to
20 stabilize the catch and prevent bounce-out. The physics of this task can be surprisingly complex: de-
21 spite its geometric simplicity, a ball in flight can swing and curve in unpredictable ways due to drag
22 and Magnus effects [1]; furthermore, the contact interaction between the ball and the deformable
23 end-effector involves complex soft-body physics which is challenging to model accurately.

24 In this paper, we study the relative merits of synthesizing high speed visual feedback controllers for
25 this task from two ends of a design spectrum: Model Predictive Control (MPC) [2, 3] represent-
26 ing a “pure control” strategy, and Blackbox policy optimization [4] representing a “pure learning”
27 approach. MPC optimizes robot trajectories in real time in response to state uncertainty - it is
28 nearly “zero-shot” in terms of data requirements and gracefully handles kinematics, dynamics and
29 task-specific constraints, but can be computationally expensive and sensitive to errors in dynamics
30 modeling. On the other hand, policy learning via blackbox or RL (Reinforcement Learning) meth-
31 ods can be extremely data inefficient, but can adapt, in principle, to complex and unknown real
32 world dynamics. Our primary contribution is to provide insights into subtle trade-offs in reaction
33 time, sample efficiency, robustness to distribution shift, and versatility in terms of whole-body mul-
34 timodal behaviors in a unified experimental evaluation of robot agility. We conclude the paper with
35 proposals to combine the “best of both worlds” in future work.

36 **Related Work:** Both classes of techniques have been previously applied to the robotic catching
37 task. Examples of optimization-based control for ball catching include [5, 6, 7, 8, 9]; [10] and [11]
38 present an unified approach subsuming catch point selection, catch configuration computation and
39 path generation in a single, nonlinear optimization problem (also see, [12], [13]). Several papers
40 utilize human demonstration and machine learning for parts of the control stack. [14] probabilis-
41 tically predict various feasible catching configurations and develop controllers to guide hand–arm
42 motion, which is learned from human demonstration. [15] also learn motion primitives from hu-
43 man demonstration and generate new movements. [16] use bi-level motion planning plus a learning
44 based tracking controller. Some papers aim for soft catching explicitly. [17] extend [14] further,
45 offering a soft catching procedure that is more resilient to imprecisions in controlling the arm and
46 desired time of catch. [18] extend [10] further for enabling soft landing. [8, 19] add heuristics for
47 soft catching, moving the hand along the predicted path of the ball, while decreasing its velocity to
48 allow the dissipation of the impact energy.

49 2 Problem formulation and proposed solution

50 We describe the trajectory of the object to be caught by a function F_o , which maps a query time
51 $t \in \mathbb{R}_{\geq 0}$ to the object’s position and velocity at time t , i.e., $(p_o(t), v_o(t)) \in \mathbb{R}^{3 \times 3}$. Depending on the
52 aerodynamic and inertial properties of the object, F_o may be highly non-trivial. Our knowledge of
53 F_o is encoded via a known \hat{F}_o which maps a query time $t \in \mathbb{R}_{\geq 0}$ and a set of parameters $\theta_o \in \mathbb{R}^d$ to a
54 prediction for the object position and velocity at time t , i.e., $(\hat{p}_o(t; \theta_o), \hat{v}_o(t; \theta_o))$. For this work, we
55 limit our scope to spherical, rigid balls and implement \hat{F}_o via classical Newtonian physics; catching
56 objects with non-trivial aerodynamics and non-uniform shapes is left to future work. However, we
57 only observe the ball position and velocity indirectly via two fixed cameras, and use θ_o to encode
58 our vision system’s current position and velocity estimate.

59 For the robot, we let $q \in \mathbb{R}^7$ denote the joint configuration vector, where $q_1 \in \mathbb{R}$ corresponds to the
60 translational base joint, and $q_{2:7} \in \mathbb{R}^6$ represent the arm joint angles.

61 2.1 Catching via optimal control

62 We assume that there exists a lower-level position and/or velocity controller that compensates for
63 the arm’s nonlinear manipulator dynamics. Abstracting away the closed-loop behavior of this lower-
64 level control system, we plan for the motion of the arm by assuming second-order integrator dynam-
65 ics¹ for q , i.e., $\ddot{q}(t) = u_a(t) \in \mathbb{R}^7$.

66 With this assumption, the optimal catching problem (OCP) can be formalized as a *free-end-time*
67 *constrained optimal control problem* over the function $u_a(\cdot)$ and catching time t_f :

$$\underset{u_a(\cdot), t_f}{\text{minimize}} J(u_a, t_f) := \int_0^{t_f} (\lambda + \|u_a(\tau)\|^2) d\tau + \Psi(q(t_f), \dot{q}(t_f), t_f), \quad (2.1)$$

68 where $\lambda \in \mathbb{R}_{>0}$ is a weighting constant and $\Psi : \mathbb{R}^7 \times \mathbb{R}^7 \times \mathbb{R}_{>0} \rightarrow 0$ is a terminal cost; subject to
69 the second-order integrator dynamics $\ddot{q}(t) = u_a(t)$, and the following constraints:

$$\forall \tau \in [0, t_f], \quad u_a(\tau) \in [\underline{u}_a, \bar{u}_a], \quad q(\tau) \in [\underline{q}, \bar{q}], \quad \dot{q}(\tau) \in [\underline{\dot{q}}, \bar{\dot{q}}], \quad c(q(t_f), t_f) \geq 0. \quad (2.2)$$

70 The first three constraints capture limits on the control effort and the joint configurations and ve-
71 locities. The terminal cost Ψ and endpoint constraint function c capture two desirable properties:
72 (i) $SE(3)$ pose alignment of the lacrosse head with the ball’s position and velocity direction at the
73 catching instant, and (ii) minimizing any residual velocity of the lacrosse head perpendicular to the
74 ball’s velocity vector. We provide details on these functions within the appendix.

75 **Conversion to Multi-Stage Trajectory Optimization:** The OCP is a non-trivial problem which
76 could be solved by leveraging the necessary conditions of optimality for free end-time problems
77 and using boundary-value-problems solvers. However, this would entail optimizing over control,
78 state, and co-state trajectories using dense discretization of the dynamics and inequality constraints
79 (e.g., via collocation). Instead, we simplify the computational burden by optimizing over a restricted
80 class of solutions – a sequence of acceleration and coasting phases, and in the process, convert the
81 problem into a *multi-stage* discrete-time trajectory optimization problem that is subsequently solved

¹Note that the lower-level control system may have some non-trivial closed-loop response characteristics, including delays. However, these can be pre-compensated for by adjusting the *commanded* (q, \dot{q}) setpoints from the planned (q, \dot{q}) trajectory.

82 using a state of the art shooting-based Sequential Quadratic Programming (SQP) solver [20]. We
83 describe this conversion in the appendix.

84 **Asynchronous Implementation:** Running concurrently to the catching controller is an estimator
85 that generates updates of the predictor parameters θ_o , necessitating online re-planning. We achieve
86 this via an asynchronous implementation where the optimization problem is continually re-solved in
87 a separate thread, using the latest estimate for θ_o and the current robot state (q, \dot{q}) . The commanded
88 (q, \dot{q}) for the robot’s lower-level PD controllers are computed by decoding the most recent stage-
89 wise solution to a continuous-time trajectory, thereby guaranteeing a consistent control rate.

90 **Cradling:** Following the intercept of the ball, we use a simple cradling motion primitive, modeled
91 as 2nd-order ODE in q , to slow the lacrosse head and simultaneously rotate the net to point upwards.

92 2.2 Blackbox Gradient Sensing Optimization

93 The catching problem can also be formulated as a Partially-Observable Markov Decision Process
94 (POMDP), and solved via Blackbox policy optimization [21, 22, 23]. In this setting, consider a
95 POMDP (S, O, A, R, P) where S is the state space partially observed by O , the observation space,
96 A is the action space, $R : S \times A \mapsto \mathbb{R}$ is the reward function and $P : S \times A \mapsto S$ is the dynamics
97 function. The optimization objective is to learn a parameterized policy $\pi_\theta : O \mapsto A$ that maximizes
98 the expected total episode return, $J(\theta) = \mathbb{E}_{\tau=(s_0, a_0, \dots, s_T)} \sum_{t=0}^T r(s_t, \pi_\theta(o_t))$.

99 **Reward function:** The reward function is different for training in sim vs. real due to differences
100 in quality of data from each. In both cases we reward the net getting close to the ball during the
101 episode. In sim, we additionally reward orientation alignment before the catch + a stability reward
102 for keeping the ball in the net; in real, we use a flat reward for successful catches (detected by a
103 sensor). Finally, we discourage excessive motion via penalizing position/velocity/acceleration/jerk
104 in sim, and hardware limit violation in real. See Appendix B for more details.

105 **Policy Network:** We use a two-tower CNN neural network. The first tower process the histor-
106 ical joint positions represented as an image of size $(n_{\text{hist}}, 7)$, where n_{hist} is the number of past
107 timesteps. The second CNN tower process the predicted ball trajectory represented as an image of
108 size $(n_{\text{pred}}, 6)$, where n_{pred} is the number of predicted timesteps. The output of the two towers is
109 concatenated into a single tensor, which is fed into two fully-connected layers. The final output is
110 then taken as the commanded joint velocities. In total, our policy network has 3255 parameters.

111 **Blackbox Gradient Sensing and Sim-to-Real Finetuning:** We apply Blackbox Gradient Sens-
112 ing (BGS) [24] for optimizing the policy neural network parameters θ . The algorithm opti-
113 mizes a smoothed version $J_\sigma(\theta)$ of the original total-reward objective $J(\theta)$, given as: $J_\sigma(\theta) =$
114 $\mathbb{E}_{\delta \sim \mathcal{N}(0, \mathbf{I}_d)} [J(\theta + \sigma\delta)]$, where $\sigma > 0$ controls the precision of the smoothing, and δ is an isotropic
115 random Gaussian vector. We first train in a simulation environment implemented in PyBullet [25].
116 Once the policy performs well in simulation, we transfer the policy to the real robot and run further
117 BGS finetuning steps using the mechanical thrower.

118 3 Experiments

119 We evaluate both our SQP and blackbox (BB) agents in simulation, on the real robot, and also ex-
120 plore performance under various distribution shifts of the thrower. Our SQP agent uses a state of
121 the art SQP solver [20] built on top of trajax [26], a JAX library for differentiable optimal con-
122 trol. For our BB agent, we use a distributed BGS library [4] with policy networks implemented in
123 Tensorflow Keras. The robot used is a combination of an ABB IRB 120T 6-DOF arm mounted on
124 a one-dimensional Festo linear actuator, creating a 7-DOF system. The ball location is determined
125 using a stereo pair of Ximea MQ013CG-ON cameras running with a trained recurrent tracker model.

126 **Error bars:** We show catch success for the real robot with error bars which give at least 95%
127 coverage, by using the Clopper–Pearson method to compute binomial confidence intervals.

128 **Inference speed and Reaction time:** The BB agent computes a single policy action in time 7.253
129 ms (std. 0.160 ms), whereas SQP takes 43.046 ms to solve (std. 21.255 ms). Recall that the SQP
130 runs asynchronously, so this solve time does not block the agent; the synchronous part runs in 2.139
131 ms (std. 0.212 ms). Vision/hardware joint data processing takes about 5 ms. Overall agents are set to
132 synchronously run at 75Hz. The mechanical thrower is 3.9 meters away from the robot and imparts
133 4.5 m/s horizontal velocity alone; including z-component the speed is ~ 5.5 m/s at catch time.

134 **Simulation to Reality Transfer:** Figure 2 highlights the real robot catch performance of both SQP
135 and BB agents. First, we see that while BB performance in sim is mostly monotonically increasing

136 (Figure 2, left), this does not necessary translate to monotonic improvement on the real hardware
 137 (Figure 2, middle). Secondly, we see that SQP suffers less performance degradation compared to
 138 BB when transferring to real. Finally, we see that it takes 40 iterations of fine-tuning on real (30
 139 ball throws per iteration) in order for the fine-tuned BB agent to match SQP’s real performance (and
 140 eventually exceed it). Both methods achieve about 80 to 85% success on mechanical ball throws.

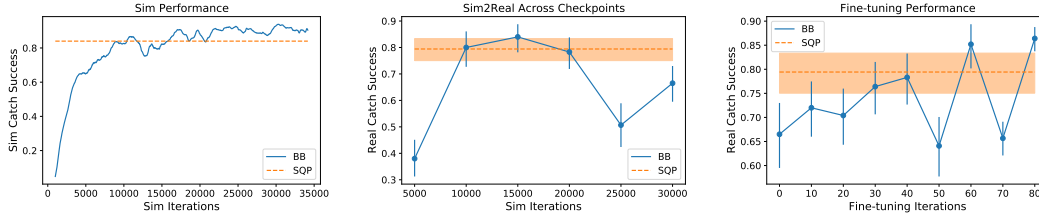


Figure 2: **(Left)** Performance of agents in sim. **(Middle)** Performance of agents on real without fine-tuning. **(Right)** Performance of sim2real transfer after fine-tuning the BB agent starting from the 30k iteration sim checkpoint. Note that each iteration corresponds to 30 mechanical ball throws.

141 **Robustness to Distribution Shifts:** Next, we look at the robustness of both agents to out-of-
 142 distribution throws. We consider three different distribution shifts: (i) varying the speed of the
 143 thrower, (ii) varying the yaw angle of the thrower, and (iii) throwing balls by hand instead of using
 144 the mechanical thrower. The first two distribution shifts are plotted in Figure 3. In Figure 3 (left), we
 145 see that while BB is reasonably robust to faster throws, its performance significantly degrades for
 146 slower throws. This is in contrast to the SQP agent, which moderately degrades in performance for
 147 faster throws (most likely due to computational bottlenecks), but is quite robust to slower throws. In
 148 Figure 3 (middle), we see that both agents have similar performance across the in-distribution yaw
 149 angles, but for out-of-distribution angles SQP maintains its performance better relative to BB.

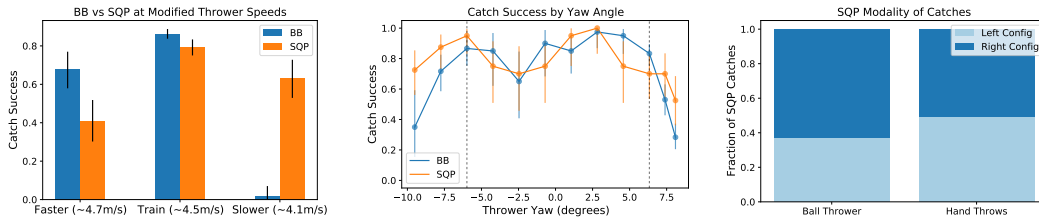


Figure 3: **(Left)** Catch performance as thrower speed varies between faster (~ 4.7 m/s), training (~ 4.5 m/s), and slower (~ 4.1 m/s) throws. **(Middle)** Catch performance as the thrower yaw angle varies from -9.5° to 8° . Note that the training distribution varies between -6° and 6.3° (marked by the dashed vertical black line). **(Right)** Distribution of left and right catches by the SQP agent on both mechanical ball throws and hand throws. Note that the BB agent catches to the right 100% of the time, likely due to the learning bias from the ball throw distribution.

150 Our last distribution shift involves hand throws (lobs) to the thrower instead of using the mechanical
 151 thrower. Using hand throws, the SQP agent has a 68.9% catch success (over 196 throws), whereas
 152 the BB agent catch performance degrades to 2.0% (over 150 throws).

153 **Multimodality:** In Figure 3 (right), we demonstrate that the SQP agent is able to catch balls in both
 154 a left and right pose configuration at fairly even rates matching the bias of the thrower. On the other
 155 hand, the BB agent is only able to catch to the right, since the ball thrower distribution is biased
 156 (60/40%) towards throwing to the right.

157 4 Conclusion and future work

158 While the fine-tuned blackbox agent has the highest catching success performance, the SQP agent
 159 is much more robust to distribution shifts in the thrower. To obtain the “best of both”, we plan to
 160 investigate two different strategies: (i) use BGS to learn the various cost parameters of SQP which
 161 we currently tune by hand, and (ii) design a mixed policy which uses SQP to move the end-effector
 162 to the ball, and then hands over control to the blackbox agent for final cradling motion. Future
 163 extensions include handling multiple non-spherical objects with adaptive dynamics prediction.

References

- [1] Aerodynamics of sports balls. *Annual Review of Fluid Mechanics*, 17(1):151–189, 1985.
- [2] F. Borrelli, A. Bemporad, and M. Morari. *Predictive control for linear and hybrid systems*. Cambridge University Press, 2017.
- [3] J. B. Rawlings. Tutorial overview of model predictive control. *IEEE control systems magazine*, 20(3):38–52, 2000.
- [4] K. Choromanski, M. Rowland, V. Sindhvani, R. Turner, and A. Weller. Structured evolution with compact architectures for scalable policy optimization. In *International Conference on Machine Learning*, pages 970–978. PMLR, 2018.
- [5] H. Yu, D. Guo, H. Yin, A. Chen, K. Xu, Y. Wang, and R. Xiong. Neural motion prediction for in-flight uneven object catching. *ArXiv*, abs/2103.08368, 2021.
- [6] U. Frese, B. Bäuml, S. Haidacher, G. Schreiber, I. Schäfer, M. Hähnle, and G. Hirzinger. Off-the-shelf vision for a robotic ball catcher. *Proceedings 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems. Expanding the Societal Role of Robotics in the the Next Millennium (Cat. No.01CH37180)*, 3:1623–1629 vol.3, 2001.
- [7] J. Kober, M. Glisson, and M. Mistry. Playing catch and juggling with a humanoid robot. *2012 12th IEEE-RAS International Conference on Humanoid Robots (Humanoids 2012)*, pages 875–881, 2012.
- [8] W. Hong and J. Slotine. Experiments in hand-eye coordination using active vision. In *ISER*, 1995.
- [9] B. Hove and J. Slotine. Experiments in robotic catching. *1991 American Control Conference*, pages 380–386, 1991.
- [10] B. Bäuml, T. Wimböck, and G. Hirzinger. Kinematically optimal catching a flying ball with a hand-arm-system. *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2592–2599, 2010.
- [11] R. Lampariello, D. Nguyen-Tuong, C. Castellini, G. Hirzinger, and J. Peters. Trajectory planning for optimal robot catching in real-time. *2011 IEEE International Conference on Robotics and Automation*, pages 3719–3726, 2011.
- [12] O. Koç, G. J. Maeda, and J. Peters. Online optimal trajectory generation for robot table tennis. *Robotics Auton. Syst.*, 105:121–137, 2018.
- [13] Y. Jia, M. Gardner, and X. Mu. Batting an in-flight object to the target. *The International Journal of Robotics Research*, 38:451 – 485, 2019.
- [14] S. Kim, A. Shukla, and A. Billard. Catching objects in flight. *IEEE Transactions on Robotics*, 30:1049–1065, 2014.
- [15] M. Riley and C. Atkeson. Robot catching: Towards engaging human-humanoid interaction. *Autonomous Robots*, 12:119–128, 2002.
- [16] K. Dong, K. Pereida, F. Shkurti, and A. P. Schoellig. Catch the ball: Accurate high-speed motions for mobile manipulators via inverse dynamics learning. *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6718–6725, 2020.
- [17] S. S. M. Salehian, M. Khoramshahi, and A. Billard. A dynamical system approach for softly catching a flying object: Theory and experiment. *IEEE Transactions on Robotics*, 32:462–471, 2016.
- [18] B. Bäuml, O. Birbach, T. Wimböck, U. Frese, A. Dietrich, and G. Hirzinger. Catching flying balls with a mobile humanoid: System overview and design considerations. *2011 11th IEEE-RAS International Conference on Humanoid Robots*, pages 513–520, 2011.
- [19] V. Lippiello and F. Ruggiero. 3d monocular robotic ball catching with an iterative trajectory estimation refinement. *2012 IEEE International Conference on Robotics and Automation*, pages 3950–3955, 2012.

- 212 [20] S. Singh, J.-J. Slotine, and V. Sindhvani. Optimizing trajectories with closed-loop dynamic
213 SQP. In *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022.
- 214 [21] K. Choromanski, M. Rowland, V. Sindhvani, R. E. Turner, and A. Weller. Structured Evolution
215 with Compact Architectures for Scalable Policy Optimization. In *Proceedings of the 35th*
216 *International Conference on Machine Learning*, pages 969–977. PMLR, 2018.
- 217 [22] K. Choromanski, A. Pacchiano, J. Parker-Holder, Y. Tang, D. Jain, Y. Yang, A. Iscen, J. Hsu,
218 and V. Sindhvani. Provably robust blackbox optimization for reinforcement learning. In L. P.
219 Kaelbling, D. Kragic, and K. Sugiura, editors, *3rd Annual Conference on Robot Learning,*
220 *CoRL 2019, Osaka, Japan, October 30 - November 1, 2019, Proceedings*, volume 100 of
221 *Proceedings of Machine Learning Research*, pages 683–696. PMLR, 2019. URL [http://](http://proceedings.mlr.press/v100/choromanski20a.html)
222 proceedings.mlr.press/v100/choromanski20a.html.
- 223 [23] H. Mania, A. Guy, and B. Recht. Simple random search provides a competitive approach to
224 reinforcement learning. *NeurIPS*, 2018.
- 225 [24] S. Abeyruwan, L. Graesser, D. B. D’Ambrosio, A. Singh, A. Shankar, A. Bewley, D. Jain,
226 K. Choromanski, and P. R. Sanketi. i-sim2real: Reinforcement learning of robotic policies in
227 tight human-robot interaction loops. *arXiv preprint arXiv:2207.06572*, 2022.
- 228 [25] E. Coumans and Y. Bai. Pybullet, a python module for physics simulation for games, robotics
229 and machine learning. <http://pybullet.org>, 2016–2021.
- 230 [26] R. Frostig, V. Sindhvani, S. Singh, and S. Tu. trajax: differentiable optimal control on accel-
231 erators, 2021. URL <http://github.com/google/trajax>.
- 232 [27] R. Sarkar, D. U. Patil, and I. N. Kar. Characterization of minimum time-fuel optimal control
233 for lti systems. *arXiv preprint arXiv:2102.10831*, 2021.
- 234 [28] B. Bäuml, T. Wimböck, and G. Hirzinger. Kinematically optimal catching a flying ball with
235 a hand-arm-system. In *2010 IEEE/RSJ International Conference on Intelligent Robots and*
236 *Systems*. IEEE, 2010.

237 A Catching via Optimal Control

238 We provide some more details on the formulation of the Optimal Catching Problem. For what fol-
 239 lows, let $\text{FK} : q \in \mathbb{R}^7 \mapsto \text{FK}(q) = (p_h(q), R_h(q)) \in \mathbb{R}^3 \times \text{SO}(3)$ denote the forward-kinematics
 240 transform that maps the joint configuration vector q to the lacrosse head's $SE(3)$ pose. This trans-
 241 form may be computed in a differentiable manner, e.g., by using a product-of-exponentials method.

242 **Desired Catching Properties.** The endpoint catching constraints capture the requirement that the
 243 lacrosse head must be positioned and oriented correctly to accept the incoming projectile. In partic-
 244 ular, let $(p_o(t_f), v_o(t_f))$ be the true 3D position and velocity of the object at the catching time t_f .
 245 Then, we require:

$$\|p_h(q(t_f)) - p_o(t_f)\| \leq \epsilon_p, \quad \text{and} \quad (R_h(q(t_f))e_2)^T \frac{v_o(t_f)}{\|v_o(t_f)\|} \geq \cos \epsilon_r, \quad (\text{A.1})$$

246 where $\epsilon_p, \epsilon_r \in \mathbb{R}_{>0}$ are prescribed tolerances on the position and angular errors, respectively, and
 247 $e_2 = (0, 1, 0)^T$. The second constraint above encourages the local \hat{y} -axis on the lacrosse head,
 248 which is orthogonal to the net's catching plane, to be aligned with the ball's velocity vector at t_f .

249 The constraint above is written assuming access to the ball's true 3D position and velocity. However,
 250 since we only have access to a prediction of these quantities via the parametric predictor $\hat{F}_o(\cdot; \theta_o)$,
 251 we enforce the above constraints w.r.t. the predicted quantities $\hat{p}_o(t_f; \theta_o), \hat{v}_o(t_f; \theta_o)$, making the
 252 endpoint catching constraint function $c(q(t_f), t_f; \theta_o)$ parametric in θ_o .

253 In conjunction with the hard constraints above, the terminal cost Ψ takes the following form:

$$\Psi(q(t_f), \dot{q}(t_f), t_f; \theta_o) := w_p \psi_p(q(t_f), t_f; \theta_o) + w_v \psi_v(q(t_f), \dot{q}(t_f)) \quad (\text{A.2})$$

$$\psi_p(q(t_f), t_f; \theta_o) := \|p_h(q(t_f)) - \hat{p}_o(t_f; \theta_o)\|^2 + \left(1 - (R_h(q(t_f))e_2)^T \frac{\hat{v}_o(t_f; \theta_o)}{\|\hat{v}_o(t_f; \theta_o)\|}\right) \quad (\text{A.3})$$

$$\psi_v(q(t_f), \dot{q}(t_f)) := \left\| v_h(q(t_f), \dot{q}(t_f)) - \begin{bmatrix} 0 \\ v_c \\ 0 \end{bmatrix} \right\|^2, \quad (\text{A.4})$$

254 where $w_p, w_v \in \mathbb{R}_{\geq 0}$ are constant weights, and $v_h(q(t_f), \dot{q}(t_f)) \in \mathbb{R}^3$ is the lacrosse head local
 255 body-frame translational velocity, computed via the Jacobian-vector product $\partial_q p_h(q) \dot{q}$. The constant
 256 $v_c \in \mathbb{R}$ is a desired catching *speed*. Thus, the terminal cost Ψ penalizes the catching-time pose
 257 errors, as defined within (A.1), as well as the motion of the lacrosse head perpendicular to the ball's
 258 velocity vector at the catching instant.

259 The overall OCP is thus parametric in θ_o , the parameters of the ball's 3D predictor function \hat{F}_o , and
 260 problem parameters $\{\epsilon_p, \epsilon_r, v_c, w_p, w_v, \lambda\}$.

261 A.1 Conversion to Multi-Stage Trajectory Optimization

262 To begin, we assume that the acceleration limits are given by symmetric intervals $[-\ddot{q}_a, \ddot{q}_a]$, where
 263 $\ddot{q}_a \in \mathbb{R}_{>0}^7$ is a fixed vector. Then, we can define an N -stage discrete-time trajectory optimization
 264 problem, where each "stage" is composed of a constant acceleration phase followed by constant
 265 cruise phase. Formally, stage- k for $k \in \{0, \dots, N-1\}$ lasts for $\delta t[k]$ seconds, where $\delta t[k] \in \mathbb{R}_{\geq 0}$.
 266 Then, within the acceleration phase of stage- k , joint $i \in \{1, \dots, 7\}$ accelerates at $\pm \ddot{q}_a$ starting at
 267 $(q_i, \dot{q}_i)[k]$ to achieve a net velocity change of $\delta \dot{q}_i[k]$. In the cruise phase, the joint moves at a constant
 268 rate of $\dot{q}_i[k] + \delta \dot{q}_i[k]$ for $\delta t[k] - (|\delta \dot{q}_i[k]| / \ddot{q}_{a_i})$ seconds.

269 We can summarize the stage transition above by defining a composite state $x[k] :=$
 270 $(q[k], \dot{q}[k], t[k]) \in \mathbb{R}^{15}$, and control $u[k] := (\delta \dot{q}[k], \delta t[k]) \in \mathbb{R}^8$. Then, the stage-"dynamics"
 271 are written as:

$$x[k+1] = \begin{bmatrix} q[k+1] \\ \dot{q}[k+1] \\ t[k] \end{bmatrix} = \begin{bmatrix} q[k] + (\dot{q}[k] + \delta \dot{q}[k])\delta t[k] - (1/2)\Delta_{\ddot{q}_a}^{-1}(\delta \dot{q}[k] \circ |\delta \dot{q}[k]|) \\ \dot{q}[k] + \delta \dot{q}[k] \\ t[k] + \delta t[k] \end{bmatrix} \quad (\text{A.5})$$

272 where \circ denotes the Hadamard product, and Δ_v is the diagonal matrix form of the vector v .

273 Let $\mathbf{u} := (u[0], \dots, u[N-1])$. The stage-equivalent discrete-time objective is given as:

$$J(\mathbf{u}) = \sum_{k=0}^{N-1} (\lambda \delta t[k] + \|\delta \dot{q}[k]\|^2) + \Psi(x[N]). \quad (\text{A.6})$$

274 **Remark A.1.** Note that the exact conversion of the integral objective in (2.1) to the stage-wise
 275 discrete-time objective would result in a stage-cost of the form $\lambda\delta t[k] + \ddot{q}_a^T |\delta\dot{q}[k]|$. However, this
 276 was found to be numerically less robust than the C^2 smooth objective used above.

277 The terminal cost and endpoint catching inequality constraints from (A.1) carry over directly, and are
 278 applied to $x[N] = (q(t_f), \dot{q}(t_f), t_f)$, where $t_f = \sum_{k=0}^{N-1} \delta t[k]$. We now tackle the limit constraints
 279 on (q, \dot{q}, \ddot{q}) . For acceleration, we require:

$$|\delta\dot{q}[k]| \leq \ddot{q}_a \delta t[k], \quad k = 0, \dots, N-1. \quad (\text{A.7})$$

280 Since $\dot{q}(t)$ linearly interpolates between the stage-values $\dot{q}[k]$, the velocity limit constraints need
 281 only be enforced at the stage values:

$$\underline{\dot{q}} \leq \dot{q}[k] \leq \bar{\dot{q}}, \quad k = 0, \dots, N. \quad (\text{A.8})$$

282 Finally, to handle the limit constraints on $q(\tau)$ for all $\tau \in [0, t_f]$, we must account for both the
 283 parabolic (constant acceleration) and linear (cruise) profiles within each stage. There exist two
 284 cases:

- 285 • Case 1: $\dot{q}_i[k](\dot{q}_i[k] + \delta\dot{q}_i[k]) \geq 0$. In this case $q_i(\tau)$ interpolates in-between $\{q_i[k], q_i[k+1]\}$
 286 for all $\tau \in [t[k], t[k+1]]$. Thus, we need only apply the limit constraints on the endpoints
 287 $q_i[k], q_i[k+1]$.
- 288 • Case 2: $\dot{q}_i[k](\dot{q}_i[k] + \delta\dot{q}_i[k]) < 0$. In this case, there is a local max/min for $q_i(\tau)$ within
 289 $[t[k], t[k+1]]$ where $\dot{q}_i(\tau) = 0$. Denote this max/min as $\hat{q}_i[k]$. Then, in addition to
 290 enforcing the limit constraints at $q_i[k], q_i[k+1]$, we must also enforce the constraint on
 291 $\hat{q}_i[k]$. The expression for $\hat{q}_i[k]$ is given by:

$$\hat{q}_i[k] = q_i[k] + \begin{cases} \frac{\dot{q}_i[k]^2}{2\ddot{q}_{a_i}} & \text{if } \dot{q}_i[k] > 0 \\ -\frac{\dot{q}_i[k]^2}{2\ddot{q}_{a_i}} & \text{if } \dot{q}_i[k] < 0 \end{cases}.$$

292 Given the discrete-time “stage”-dynamics, optimization objective, and constraints, we can use any
 293 off-the-shelf constrained discrete-time trajectory optimization solver. In this work, we leverage
 294 Dynamic Shooting SQP, introduced in [20].

295 **Remark A.2.** Note that the combination of max/min acceleration and cruise phases within each
 296 stage reflects the nature of mixed control-effort/minimum-time optimal control solutions, colloquially
 297 characterized as the ‘bang-off-bang’ strategy. Further, recent work [27] has shown that for LTI
 298 systems with a single control input, the optimal solution to a mixed control-effort/minimum-time
 299 problem with an endpoint reachability constraint is a sequence of “bang-off” stages. This justifies
 300 our use of such a stage-wise reduction of the original continuous-time OCP, and is similar in spirit
 301 to previous works on catching using trapezoidal velocity profiles [28].

302 B Detailed Reward Functions

303 We provide detailed descriptions of the reward functions used for training the Blackbox policy and
 304 how they are computed below.

305 **Object Position Reward.** This reward is based on the closest distance the end-effector comes to
 306 the object during the episode. The closest distance is scaled on an exponential curve with a cutoff
 307 at 20cm scoring 1.0 for any episodes that get closer than this. This reward is used both in sim and
 308 real. This reward is useful for the Blackbox Policy to learn to get the net close to the ball.

309 **Object Orientation Reward.** This reward is based on the orientation of the net right when the ball
 310 gets to within 20cm of the net. The score is computed as a dot product of the velocity vector of the
 311 object and the axis of the net, scaled between 0 to 1 as a reward. This reward is only used in sim and
 312 encourages the policy to point the net in the right direction for a catch.

313 **Object Stability Reward.** This reward is based on how stable the object remains after it is close
 314 (defined as within 20cm of the net). Entering the close criteria and staying there through the end
 315 of the episode provides a flat 0.2 reward. The remaining 0.8 part of the stability reward is given
 316 by measuring the speed of the ball while it’s close for 0.25s. Each time-step during this duration
 317 contributes equally and is scored on an exponential curve based on object speed, capping out at
 318 speeds less than 0.2m/s scoring full for that timestep. A full score would be keeping the speed less
 319 than 0.2m/s for the full 0.25s. This reward is only used in SIM as the precision of ball tracking is
 320 difficult when the ball is in the net or obscured.

321 **Object Catch Reward.** This reward is only used in real and is measured using a proximity sensor
322 attached close to the net that can reliably detect whether a ball is in the net or not. The ball is
323 declared as caught if the sensor detects a ball continuously in the net for greater than $0.25s$. This
324 provides a flat 0 or 1 reward.

325 **Penalties for exceeding Robot dynamic constraints.** In Sim, there are multiple penalty rewards
326 used to ensure the policy learns to operate within the robot constraints such as joint position, velocity,
327 acceleration and jerk limits. The penalty rewards are implemented as a flat 1.0 if the agent actions
328 stay within constraints and reduces to 0 depending on how much it violates them. The reward is
329 reduced depending on how many timesteps and by how much it exceeds them. In Real the hardware
330 produces a fault error code and freezes when movements exceed constraints. So in real the penalty
331 is just scored based on whether or not the hardware encounters the fault code.