
Supplementary Materials for: Responsive Safety in Reinforcement Learning by PID Lagrangian Methods

1. Derivations

Here we derive the effects of the proportional and derivative control terms on the Differential Multiplier Method.

1.1. Traditional Lagrangian method

Start from continuous Lagrangian dynamics for objective f and constraint g :

$$\mathcal{L} = f(\mathbf{x}) + \lambda g(\mathbf{x}) \quad (1)$$

$$\dot{x}_i = -\frac{\partial f}{\partial x_i} - \lambda \frac{\partial g}{\partial x_i} \quad (2)$$

$$\dot{\lambda} = \alpha g(\mathbf{x}) \quad (3)$$

with α a scalar. Taking the time-derivative of \dot{x}_i :

$$\ddot{x}_i = -\frac{d}{dt} \left(\frac{\partial f}{\partial x_i} \right) - \dot{\lambda} \frac{\partial g}{\partial x_i} - \lambda \frac{d}{dt} \left(\frac{\partial g}{\partial x_i} \right) \quad (4)$$

$$\frac{d}{dt} \left(\frac{\partial f}{\partial x_i} \right) = \sum_j \frac{\partial}{\partial x_j} \left(\frac{\partial f}{\partial x_i} \right) \frac{dx_j}{dt} = \sum_j \frac{\partial^2 f}{\partial x_i \partial x_j} \dot{x}_j \quad (5)$$

$$\frac{d}{dt} \left(\frac{\partial g}{\partial x_i} \right) = \sum_j \frac{\partial^2 g}{\partial x_i \partial x_j} \dot{x}_j \quad (6)$$

and substituting for $\dot{\lambda}$:

$$\ddot{x}_i = -\sum_j \left(\frac{\partial^2 f}{\partial x_i \partial x_j} + \lambda \frac{\partial^2 g}{\partial x_i \partial x_j} \right) \dot{x}_j - \alpha g(\mathbf{x}) \frac{\partial g}{\partial x_i} \quad (7)$$

Defining the damping matrix, A :

$$A_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j} + \lambda \frac{\partial^2 g}{\partial x_i \partial x_j} \quad (8)$$

yields the oscillator equation for the i -th component:

$$\ddot{x}_i + \sum_j A_{ij} \dot{x}_j + \alpha g(\mathbf{x}) \frac{\partial g}{\partial x_i} = 0 \quad (9)$$

Written in vector form:

$$\ddot{\mathbf{x}} + A\dot{\mathbf{x}} + \alpha g(\mathbf{x})\nabla g = 0 \quad (10)$$

$$A = \nabla^2 f + \lambda \nabla^2 g \quad (11)$$

with $\nabla^2 f$ denoting the Hessian of f .

1.2. Proportional-Integral Method

We amend the differential equation for λ as:

$$\dot{\lambda} = \alpha g(\mathbf{x}) + \beta \dot{g}(\mathbf{x}) = \alpha g(\mathbf{x}) + \beta \sum_j \frac{\partial g}{\partial x_j} \dot{x}_j \quad (12)$$

Substitution into \ddot{x} yields:

$$\ddot{x}_i = - \sum_j \left(\frac{\partial^2 f}{\partial x_i \partial x_j} + \lambda \frac{\partial^2 g}{\partial x_i \partial x_j} \right) \dot{x}_j - \alpha g(\mathbf{x}) \frac{\partial g}{\partial x_i} - \beta \left(\sum_j \frac{\partial g}{\partial x_j} \dot{x}_j \right) \frac{\partial g}{\partial x_i} \quad (13)$$

$$= - \sum_j \left(\frac{\partial^2 f}{\partial x_i \partial x_j} + \lambda \frac{\partial^2 g}{\partial x_i \partial x_j} + \beta \frac{\partial g}{\partial x_i} \frac{\partial g}{\partial x_j} \right) \dot{x}_j - \alpha g(\mathbf{x}) \frac{\partial g}{\partial x_i} \quad (14)$$

giving the same dynamics as the traditional Lagrangian method except with damping matrix modified by addition of a positive-semi-definite term:

$$\ddot{\mathbf{x}} + (A + \beta \nabla g \nabla^\top g) \dot{\mathbf{x}} + \alpha g(\mathbf{x}) \nabla g = 0 \quad (15)$$

1.3. Integral-Derivative Method

We use the following differential equation for λ :

$$\dot{\lambda} = \alpha g(\mathbf{x}) + \gamma \ddot{g}(\mathbf{x}) \quad (16)$$

Expanding the second-derivative:

$$\ddot{g}(x) = \frac{d}{dt} \left(\frac{d}{dt} g(\mathbf{x}) \right) \quad (17)$$

$$= \frac{d}{dt} \left(\sum_j \frac{\partial g}{\partial x_j} \dot{x}_j \right) \quad (18)$$

$$= \sum_j \left(\frac{d}{dt} \left(\frac{\partial g}{\partial x_j} \right) \dot{x}_j + \frac{\partial g}{\partial x_j} \frac{d}{dt} (\dot{x}_j) \right) \quad (19)$$

$$= \sum_j \left(\dot{x}_j \sum_k \frac{\partial^2 g}{\partial x_k \partial x_j} \dot{x}_k + \frac{\partial g}{\partial x_j} \ddot{x}_j \right) \quad (20)$$

Substituting the full expression for $\dot{\lambda}$ into \ddot{x}_i yields:

$$\ddot{x}_i = - \sum_j \left(\frac{\partial^2 f}{\partial x_i \partial x_j} + \lambda \frac{\partial^2 g}{\partial x_i \partial x_j} \right) \dot{x}_j - \alpha g(\mathbf{x}) \frac{\partial g}{\partial x_i} - \gamma \frac{\partial g}{\partial x_i} \left(\sum_j \left(\sum_k \dot{x}_j \frac{\partial^2 g}{\partial x_j \partial x_k} \dot{x}_k + \frac{\partial g}{\partial x_j} \ddot{x}_j \right) \right) \quad (21)$$

where the terms with coefficient γ are due to the derivative update term in λ . Now included are terms of second-order in the velocity and additional mixing terms on the acceleration. It is instructive to consider the resulting vector equation in full:

$$0 = \ddot{\mathbf{x}} + A \dot{\mathbf{x}} + \alpha g(\mathbf{x}) \nabla g + \gamma \nabla g \dot{\mathbf{x}}^\top \nabla^2 g \dot{\mathbf{x}} + \gamma \nabla g \nabla^\top g \ddot{\mathbf{x}} \quad (22)$$

$$= (I + \gamma \nabla g \nabla^\top g) \ddot{\mathbf{x}} + A \dot{\mathbf{x}} + (\alpha g(\mathbf{x}) + \gamma \dot{\mathbf{x}}^\top \nabla^2 g \dot{\mathbf{x}}) \nabla g \quad (23)$$

The acceleration terms are coupled together by a positive definite matrix (identity plus the outer product of the vector ∇g with itself). Therefore a form without coupling of acceleration can be restored by left-multiplying by the inverse of this matrix, letting $B = (I + \gamma \nabla g \nabla^\top g)$:

$$\ddot{\mathbf{x}} + B^{-1} A \dot{\mathbf{x}} + (\alpha g(\mathbf{x}) + \gamma \dot{\mathbf{x}}^\top \nabla^2 g \dot{\mathbf{x}}) B^{-1} \nabla g = 0 \quad (24)$$

The effects of the new terms from the derivative update rule are discussed in the main text.

110 **1.4. Proportional-Integral-Derivative Method**

111 Finally, the full PID-Lagrangian method has dynamics which combines the independent effects of the proportional and
 112 derivative methods:

113
 114
$$\ddot{\mathbf{x}} + B^{-1} (A + \beta \nabla g \nabla^T g) \dot{\mathbf{x}} + (\alpha g(\mathbf{x}) + \gamma \dot{\mathbf{x}}^T \nabla^2 g \dot{\mathbf{x}}) B^{-1} \nabla g = 0$$
 (25)
 115

116
 117 **2. Training Details**

118 All of our experiments began with randomly initialized agents. The reward-value and cost-value estimators shared parameters
 119 with the policy. We used Generalized Advantage Estimation for both reward and cost advantages. The control input is
 120 updated once per iteration, which in our settings included multiple gradient updates on the policy. Training batches typically
 121 included the end of several trajectories, allowing an estimate of the average episodic sum of costs at each iteration.
 122

123
 124 *Table 1.* Experiment hyperparameters.

125
 126

HYPERPARAMETER	VALUE
LEARNING RATE	3×10^{-4}
NN HIDDEN LAYER SIZE	512
NN NONLINEARITY	<i>tanh</i>
BATCH DIMENSION, TIME	128
BATCH DIMENSION, ENV	104
PPO EPOCHS	8
PPO MINIBATCHES	1
PPO RATIO-CLIP	0.1
DISCOUNT, γ	0.99
λ_{GAE}	0.97
K_I	1
K_P, K_D	AS SPECIFIED
COST SCALING	1/10
EXPONENTIAL MOVING AVERAGE, K_P	0.95
EXPONENTIAL MOVING AVERAGE, K_D	0.9
DIFFERENCE ITERATES DELAY, K_D	15
OBSERVATION NORMALIZATION	TRUE
$\hat{\beta}$	1 (UNLESS SPECIFIED)
EXPONENTIAL MOVING AVERAGE, $\hat{\beta}$	0.9

127
 128
 129
 130
 131
 132
 133
 134
 135
 136
 137
 138
 139
 140
 141
 142
 143
 144
 145
 146
 147
 148
 149
 150
 151
 152
 153
 154
 155
 156
 157
 158
 159
 160
 161
 162
 163
 164

3. Additional Learning Curves

3.1. Cost and Reward Curves for Additional Environments

This section contains learning curves from the experiments used to make the figures showing cost figure-of-merit versus Lagrange multiplier learning rate, K_I , in the main text. The main text includes curves from DOGGOGOAL2, the other environments are shown here.

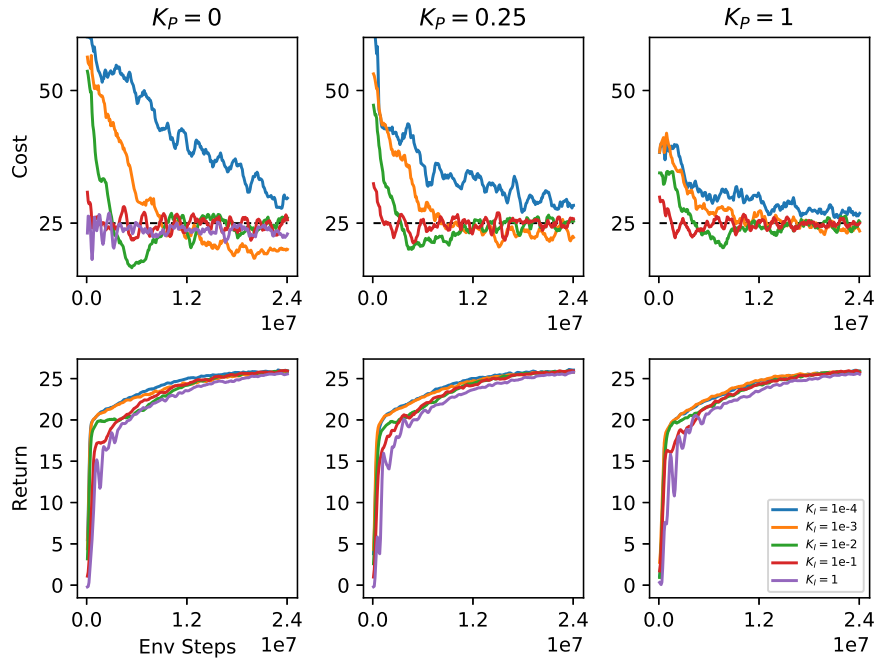


Figure 1. Costs and returns with varying Lagrange multiplier learning rate, K_I , and proportional control coefficient, K_P , in POINTGOAL1, cost-limit=25.

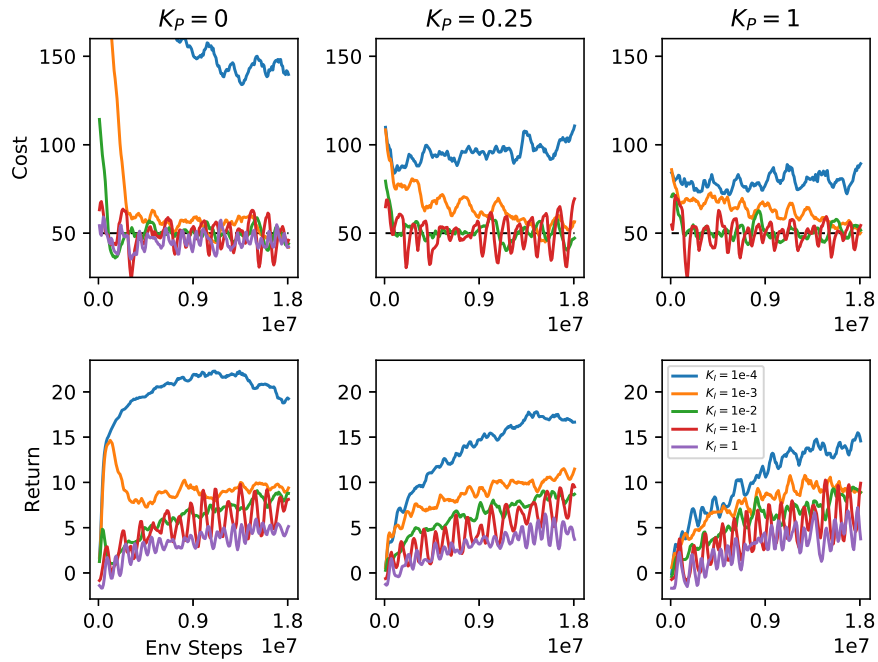


Figure 2. Costs and returns with varying Lagrange multiplier learning rate, K_I , and proportional control coefficient, K_P , in CARBUTON1, cost-limit=50.

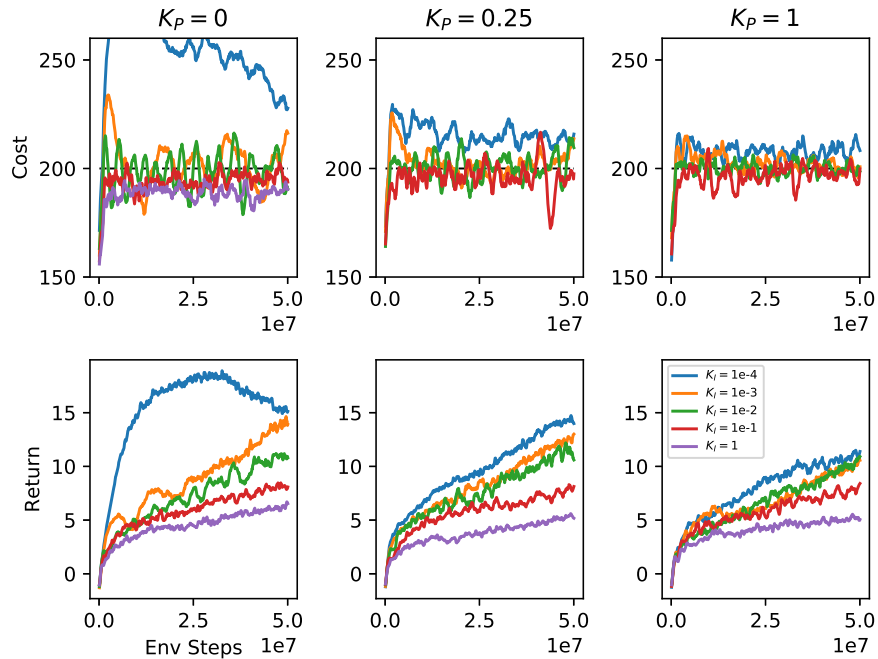


Figure 3. Costs and returns with varying Lagrange multiplier learning rate, K_I , and proportional control coefficient, K_P , in DOGGOBUTTON1, cost-limit=200.

3.2. Derivative-Control Learning Curves

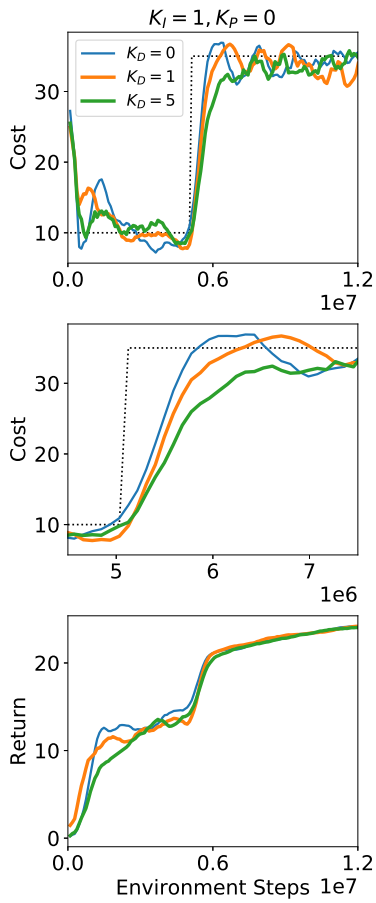


Figure 4. Derivative control slows the increase in cost, causing it to rise more gradually, and reducing overshoot. The cost limit was increased from 10 to 35 at 5M environment steps. Environment: POINTGOAL1.

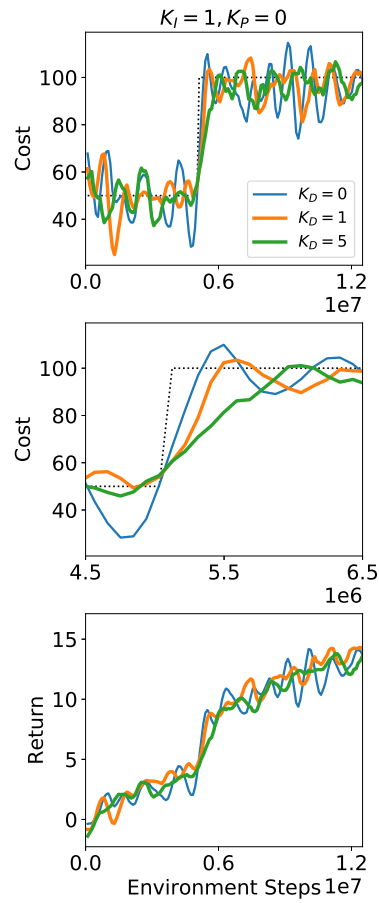


Figure 5. Derivative control slows the increase in cost, causing it to rise more gradually, and reducing overshoot. The cost limit was increased from 50 to 100 at 5M environment steps. Environment: CARBUTTON1.

3.3. Comparison to Unconstrained Algorithms

This section shows learning curves for the same four environments referenced in the main text. These figures demonstrate that the unconstrained algorithm (PPO) does not satisfy the cost constraints, and as a result it achieves higher rewards.

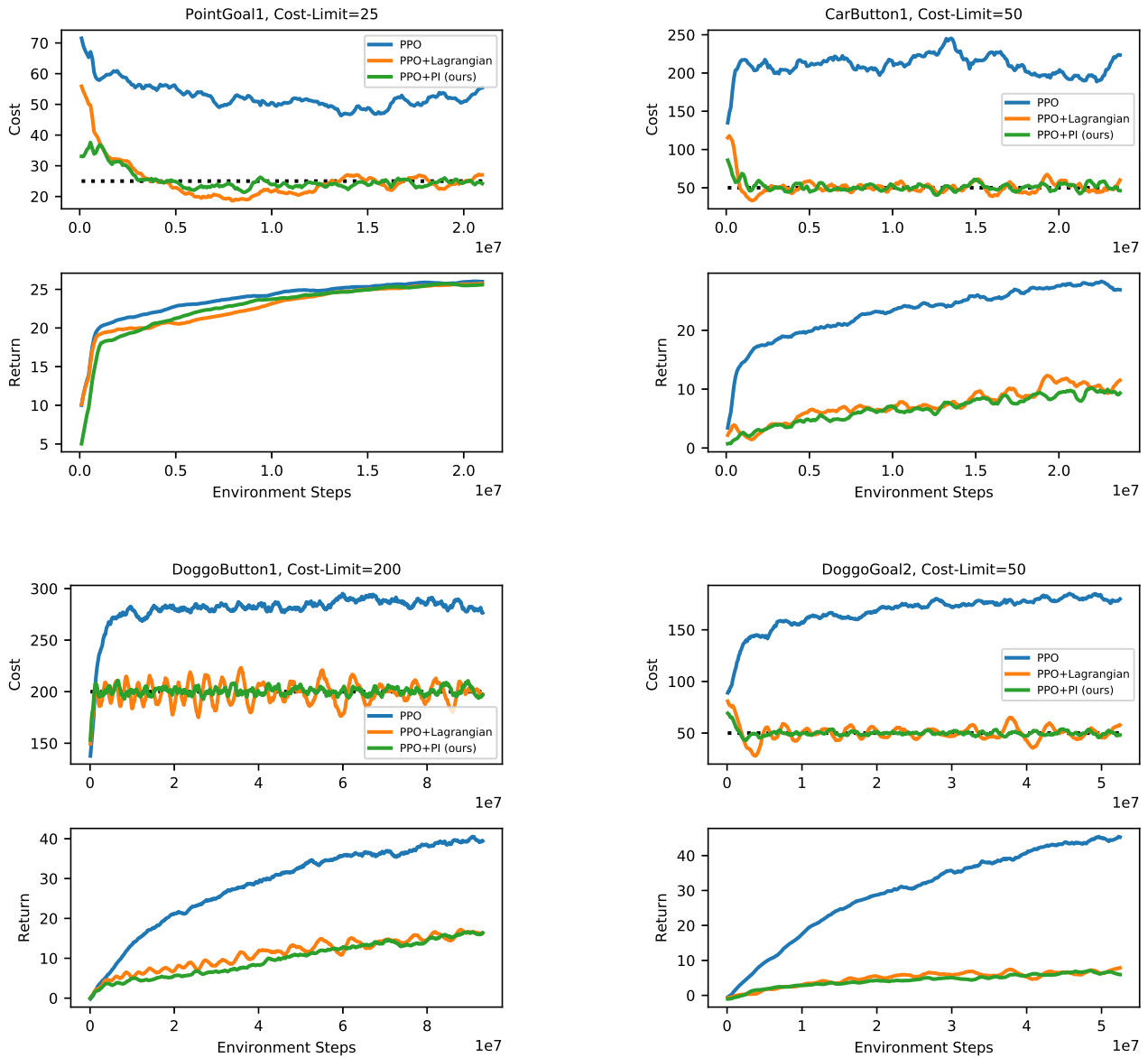


Figure 6. Cost and reward curves for three variants of PPO: unconstrained, Lagrangian, and PI-Controlled. The unconstrained algorithm wildly violates all cost-limits used in our experiments. PPO+Lagrangian and PPO+PI use the same Lagrange multiplier learning rate, K_I .

4. Adaptive Objective-Balancing

Alternative, KL-Based Estimator The magnitudes of the gradients in θ -space might not fully reflect the relative impacts of reward- and cost-learning on agent behavior. Reinforcement learning offers an alternative grounding in policy-space, for example using:

$$\beta_{KL} = \frac{D_{KL}(\pi_{\theta_k} || \pi_{\theta_{k+1}}^R)}{D_{KL}(\pi_{\theta_k} || \pi_{\theta_{k+1}}^C)} \quad (26)$$

Here, $\pi_{\theta_{k+1}}^R$ and $\pi_{\theta_{k+1}}^C$ are hypothetical new policies found using only the reward-objective or only the (un-scaled) cost-objective, respectively. We experimented with this estimator and found it to work, although not as well as the gradient-norm estimator in our cases. Some results are included in the figures below.

Figures We include figures for POINTGOAL1 using I-control and PI-control, and the more challenging DOGGOGOAL2 using I-control. Observe in the un-balanced case ($\hat{\beta} = 1$) that as the controller settings scale with the reward, the learning dynamics remained the same. For example, see $(K_I = 0.1, \rho = 10)$, $(K_I = 0.01, \rho = 1)$, and $(K_I = 0.001, \rho = 0.1)$. Using gradient-based objective balancing (β_{∇}), however, the learning dynamics were roughly the same across all reward scales, for given controller settings. Alternatively, KL-based objective balancing (β_{KL}) was also effective, but did not produce dynamics as uniform as the gradient-based method.

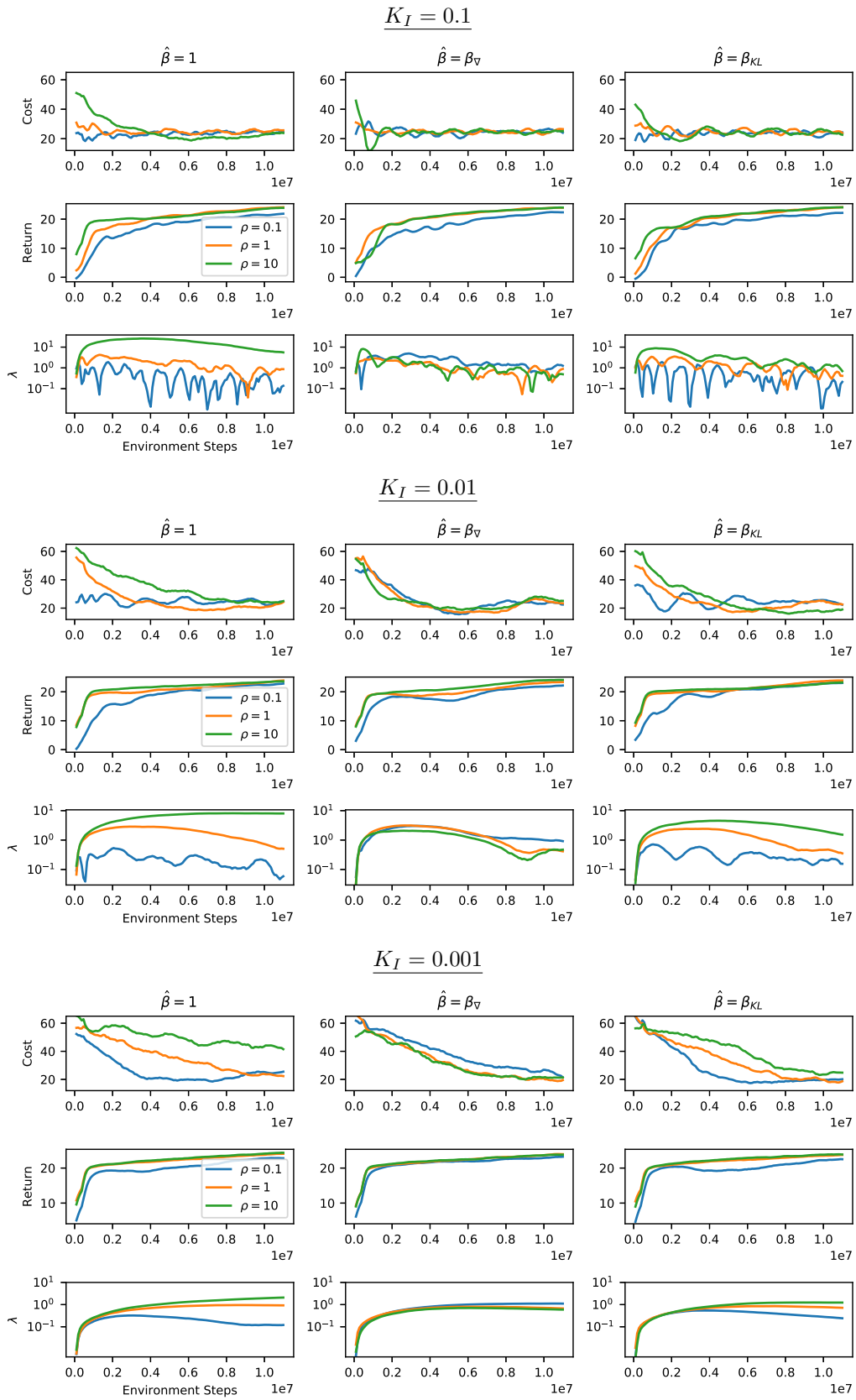


Figure 7. Reward scaling, I-control, POINTGOAL1, cost-limit=25.

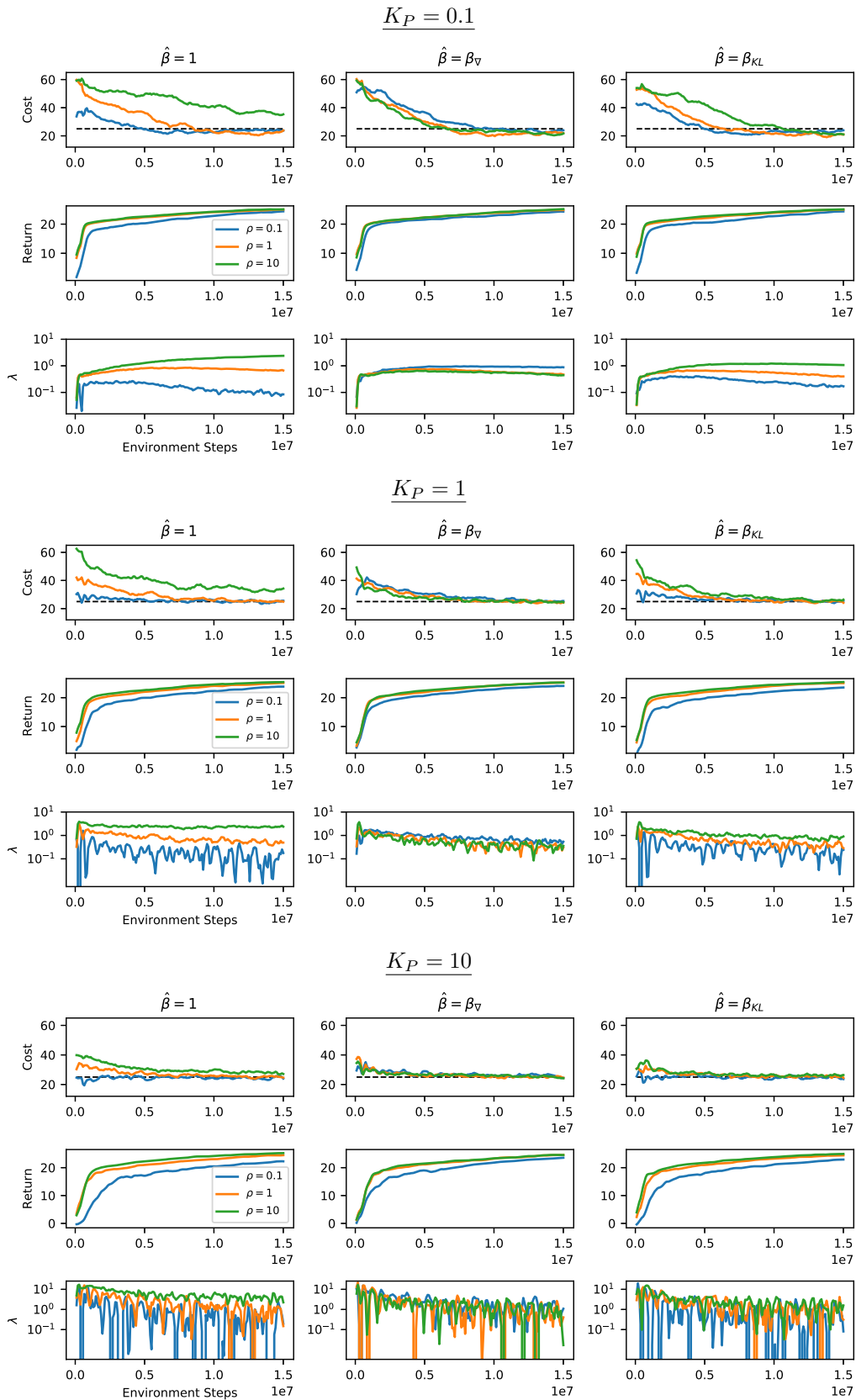


Figure 8. Reward scaling, PI-control with $K_I = 0.001$, POINTGOAL1, cost-limit=25.

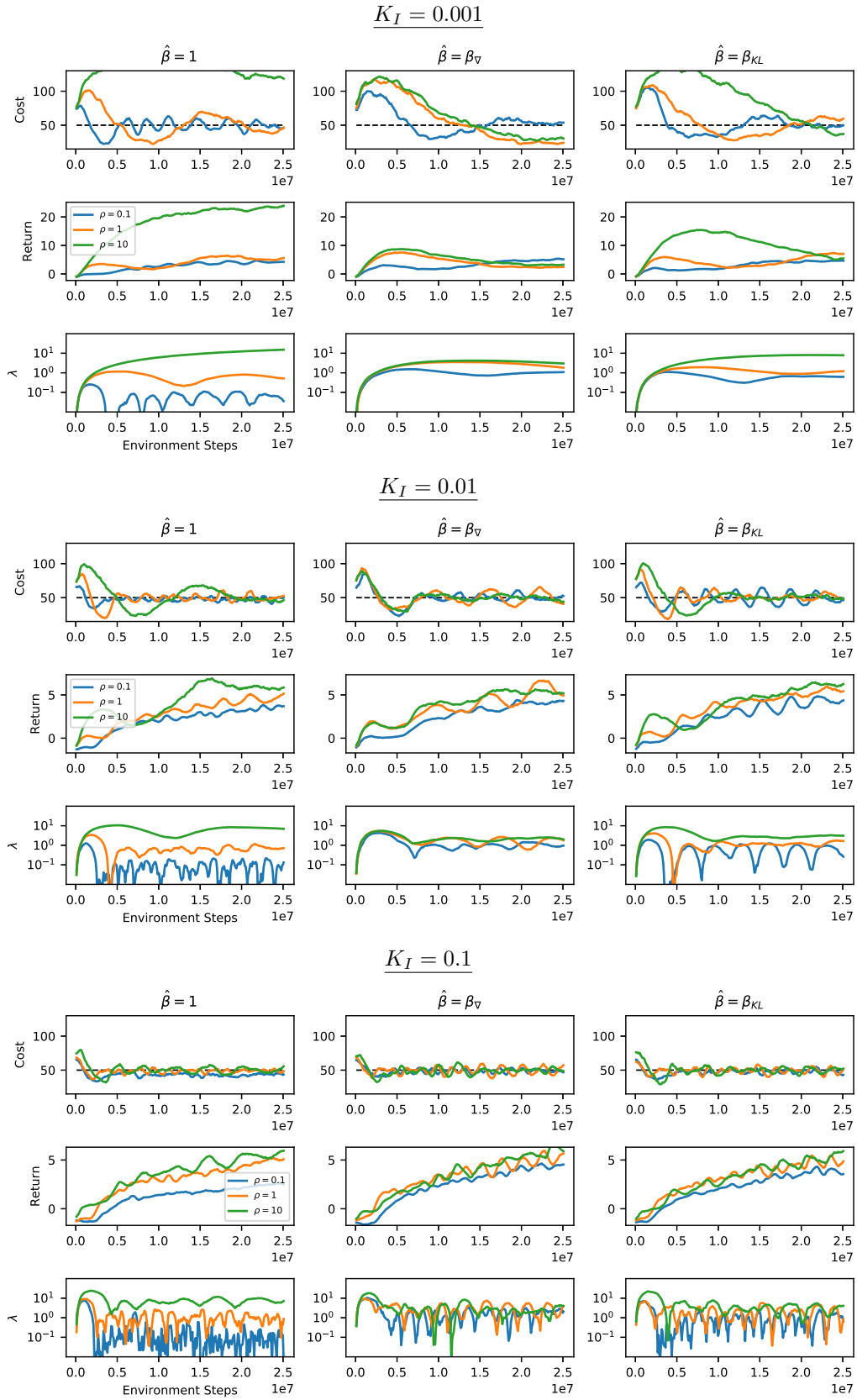


Figure 9. Reward scaling, I-control, DOGGOGOAL2, cost-limit=50.

550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604