

A Multiplicative Value Function for Safe and Efficient Reinforcement Learning

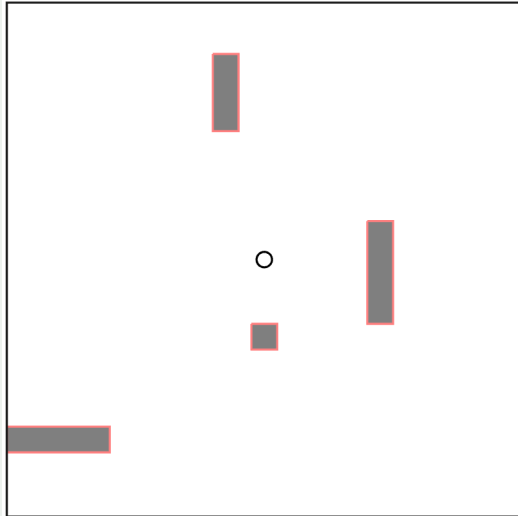


Nick Bührer ¹, Zhejun Zhang ¹, Alexander Liniger ¹,
Fisher Yu ¹, Luc Van Gool ^{1,2}.

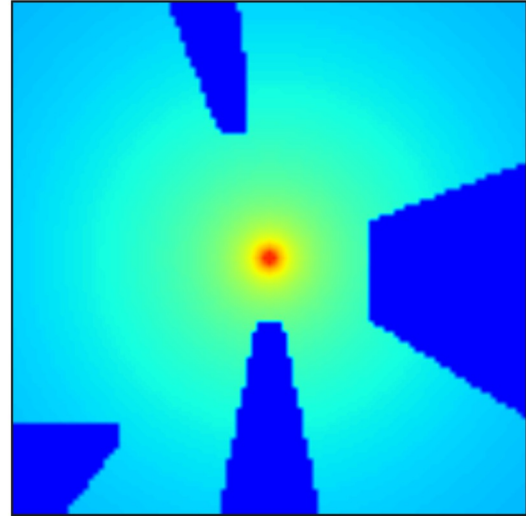
¹ Computer Vision Lab, ETH Zurich, Switzerland. ² PSI, KU Leuven, Belgium.



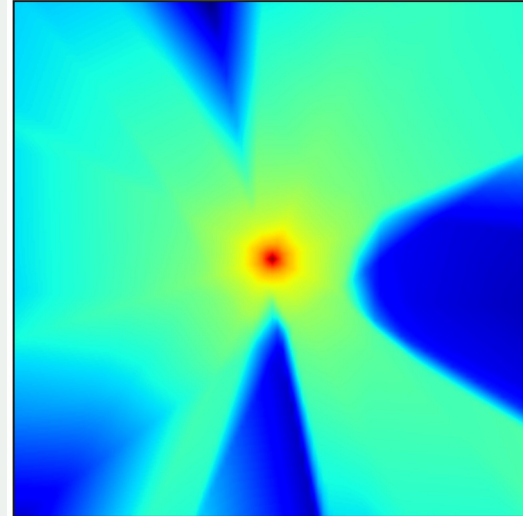
Motivation



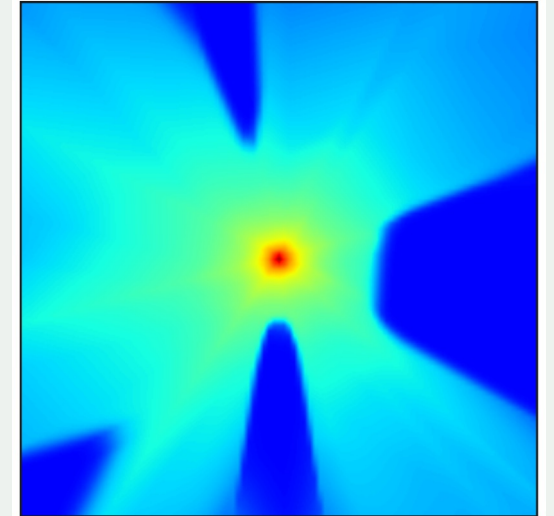
(A) Environment



(B) Ground-Truth Value



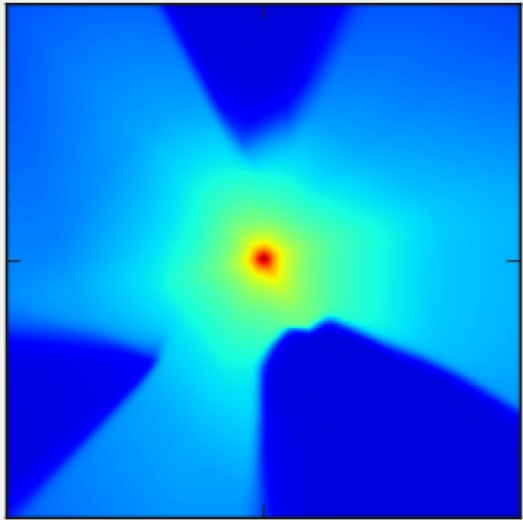
(C) Regular Value Function



(D) Multiplicative Value Function

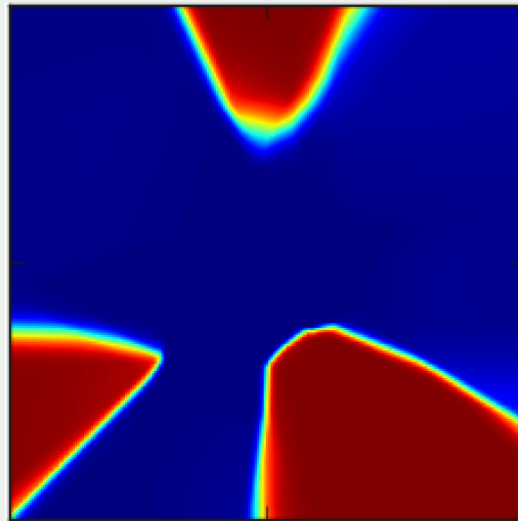
Multiplicative Value Function

Multiplicative Value
Function V_{mult}^{π}

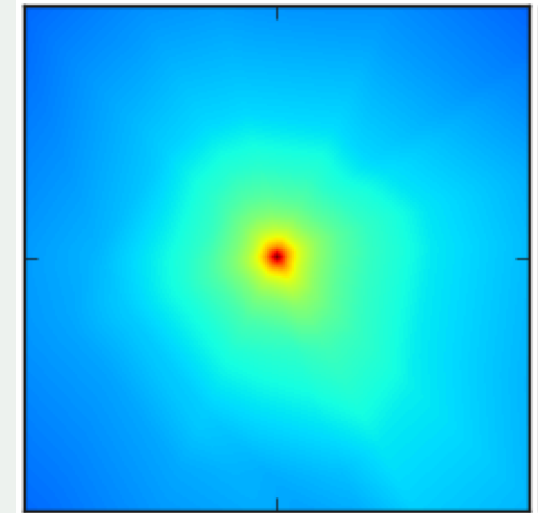


=

Probabilistic
Safety Critic Ψ^{π}



Constraint-Free
Reward Critic \bar{V}^{π}



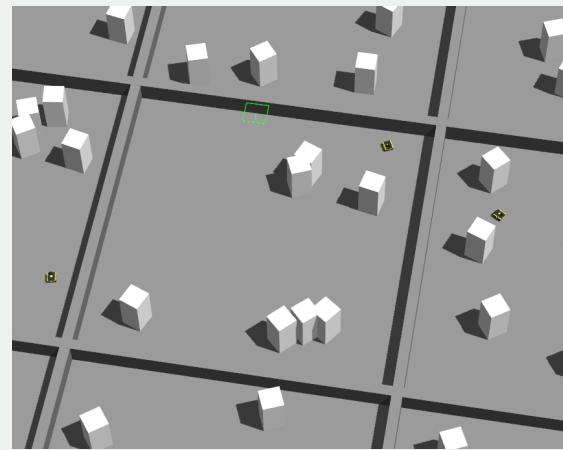
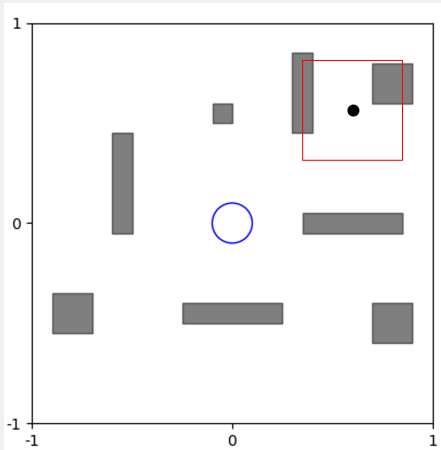
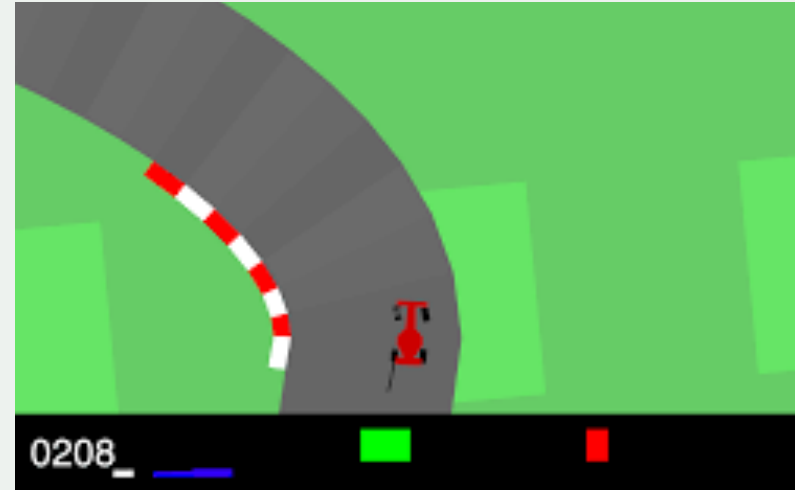
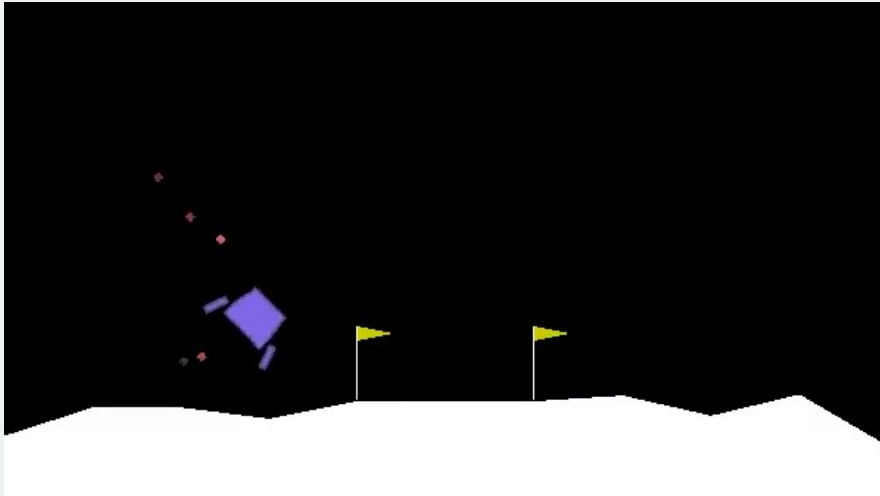
$$V_{\text{mult}}^{\pi}(s) = (\bar{V}^{\pi}(s) - \bar{v}_{\min}) \cdot (1 - \Psi^{\pi}(s)) + \bar{v}_{\min}, \quad \bar{v}_{\min} := \min_s \bar{V}^{\pi}(s)$$

Apply to SAC and PPO

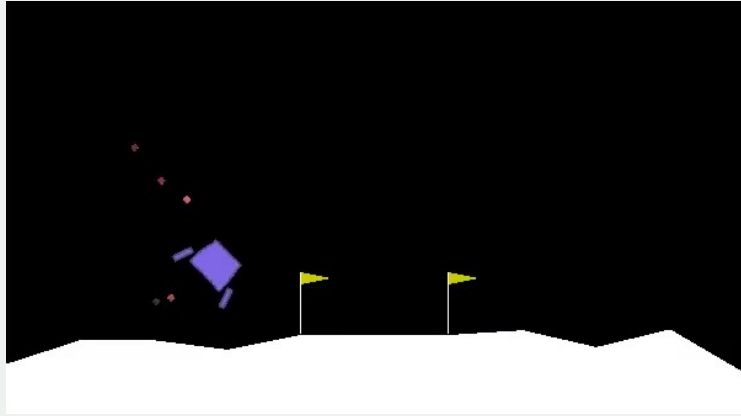
- **Q Function:** $Q_{\text{mult}}^{\pi}(s, a) = (\bar{Q}^{\pi}(s, a) - \bar{q}_{\text{min}}) \cdot (1 - \Psi^{\pi}(s, a)) + \bar{q}_{\text{min}}$
- **Advantage**
 - V1 bootstrap Q: $A_{\text{mult}}^{\pi}(s_t, a_t) = [\bar{r}_t + \gamma V_{\text{mult}}^{\pi}(s_{t+1})] - V_{\text{mult}}^{\pi}(s_t)$
 - V2 without bootstrap: $A_{\text{mult}}^{\pi}(s_t, a_t) = Q_{\text{mult}}^{\pi}(s_t, a_t) - V_{\text{mult}}^{\pi}(s_t)$
 - V3 bootstrap the safety critic inside $Q_{\text{mult}}^{\pi}(s_t, a_t)$ of V2
- **SAC:** $\max_{\theta} \mathbb{E}_{a_{\theta} \sim \pi_{\theta}} [Q_{\text{mult}}^{\pi_{\theta}}(s, a_{\theta}) - \alpha \log \pi_{\theta}(s_{\theta} | x)]$
- **PPO:** $\max_{\theta} \mathbb{E}_{a \sim \pi_{\theta}} \left[\min \left\{ \frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)} A_{\text{mult}}^{\pi_{\theta_k}}(s, a), g \left(\epsilon, A_{\text{mult}}^{\pi_{\theta_k}}(s, a) \right) \right\} \right]$



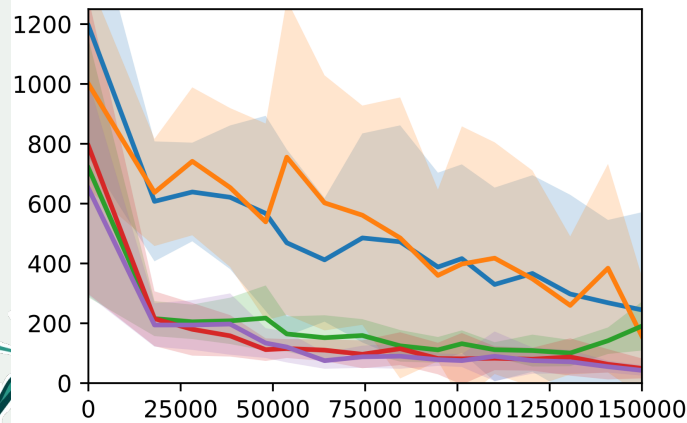
Safety-Focused Environments



Quantitative Results



Value Loss of Lunar Lander Safe

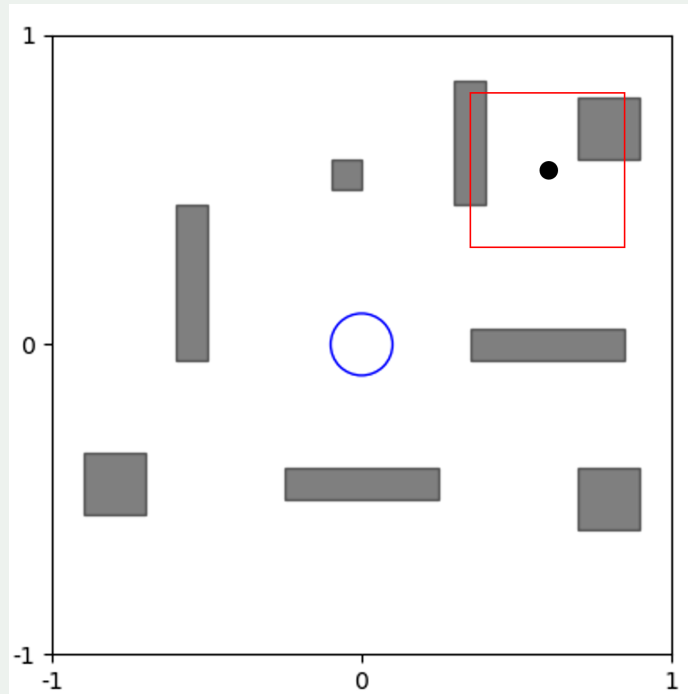


Blue: PPO base
Orange: Lagrange
Green: V1
Violet: V2
Red: V3

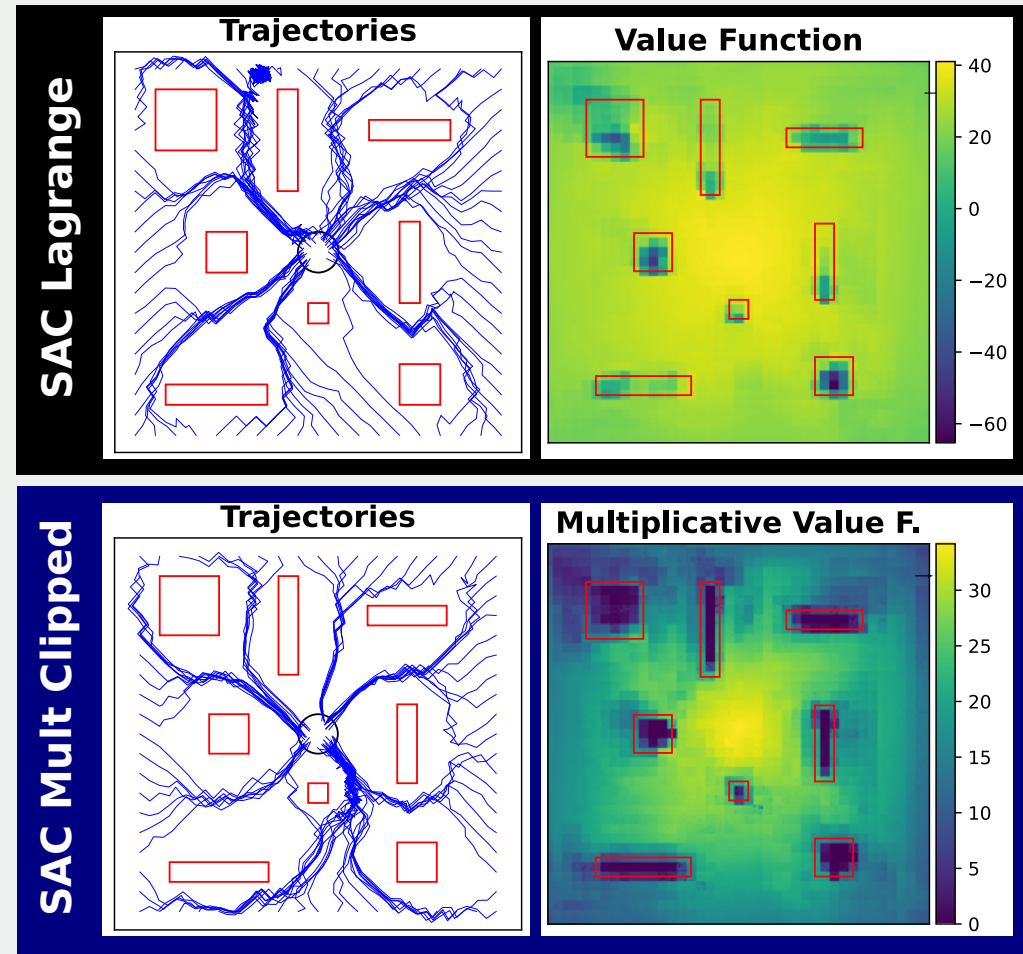
| | Reward ↑ | % Constraint violations ↓ | Reward ↑ | % Constraint violations ↓ |
|-------------------|-----------------|------------------------------|-----------------|------------------------------|
| Lunar Lander Safe | | | | |
| SAC | | | | |
| | 50k | | 150k | |
| SAC base | 90 ± 108 | 10 ± 16 | 181 ± 117 | 3 ± 6 |
| Lagrange | 111 ± 105 | 17 ± 13 | 184 ± 128 | 2 ± 3 |
| Mult | -35 ± 27 | 3 ± 5 | -34 ± 22 | 3 ± 4 |
| Mult Clipped | 134 ± 94 | 14 ± 13 | 243 ± 49 | 8 ± 15 |
| Mult Lagrange | 125 ± 59 | 29 ± 15 | 251 ± 20 | 2 ± 2 |
| PPO | | | | |
| | 50k | | 150k | |
| PPO base | -126 ± 158 | 77 ± 29 | 225 ± 100 | 10 ± 30 |
| Lagrange | -24 ± 146 | 54 ± 39 | 204 ± 116 | 12 ± 24 |
| V1 | 101 ± 84 | 41 ± 19 | 205 ± 78 | 7 ± 16 |
| V2 | 89 ± 122 | 44 ± 34 | 251 ± 28 | 5 ± 9 |
| V3 | 144 ± 4 | 26 ± 22 | 264 ± 5 | 1 ± 2 |
| FOCOPS | -129 ± 21 | 64 ± 24 | 117 ± 80 | 30 ± 19 |



Qualitative Results



Point Robot Navigation
Via local occupancy grid and
a vector pointing to the goal.



Simulation and Real-world

- Differential drive robot with 1D-Lidar
- Gazebo simulation
- Zero-shot sim-to-real



Summary

- Multiplicative value function
 - Constraint-free reward critic \otimes Probabilistic safety critic
- Integration into SAC and PPO
 - Increased sample efficiency and learning stability.
- Experiments
 - Safety-focused RL environments and real-world robot navigation.
- Future works
 - Theoretical justification.

