# EigenSafe

*A spectral framework for learning-based stochastic safety filtering*

**Inkyu Jang**[1], Jonghae Park[1], Chams E. Mballo[2], Sihyun Cho[3], Claire J. Tomlin[2], and H. Jin Kim[13]

[1] Department of Aerospace Engineering, Seoul National University
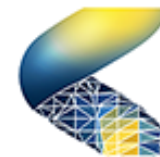[2] Electrical Engineering & Computer Sciences, University of California, Berkeley
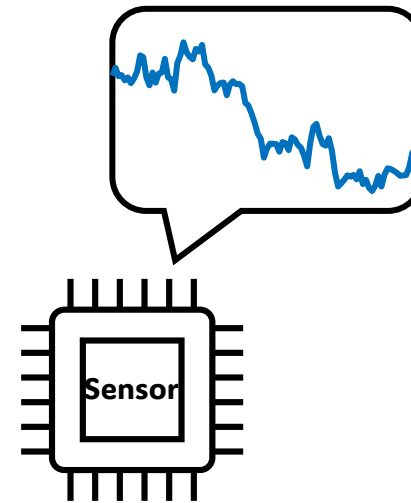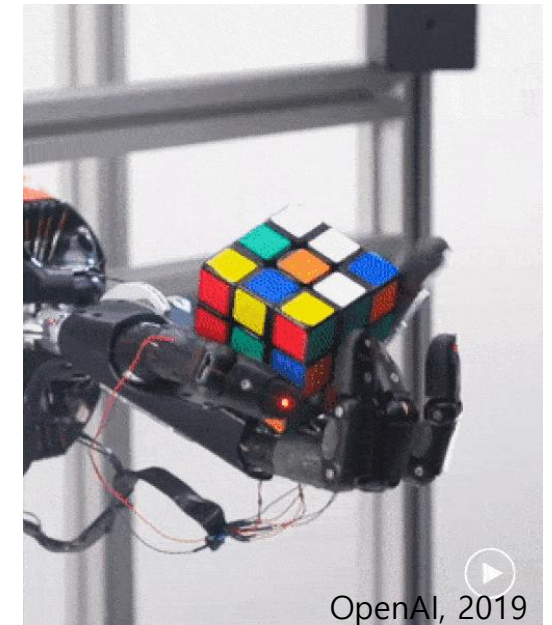[3] Graduate School of Artificial intelligence, Seoul National University

# Real-World Robots are Stochastic



Inherent stochasticity



Limited observability



OpenAI, 2019

System complexity

# Safety probability

$$Z_\pi(t, x, u) = \mathbb{P}_\pi[x_\tau \text{ safe } \forall \tau \in [0, t] \mid x_0 = x, \, u_0 = u]$$

## The law of total probability and Markov property give the dynamic programming principle

$$Z_\pi(t + 1, x, u) = \mathbb{E}_{x^+ \sim P, u^+ \sim \pi}[Z_\pi(t, x^+, u^+)]$$
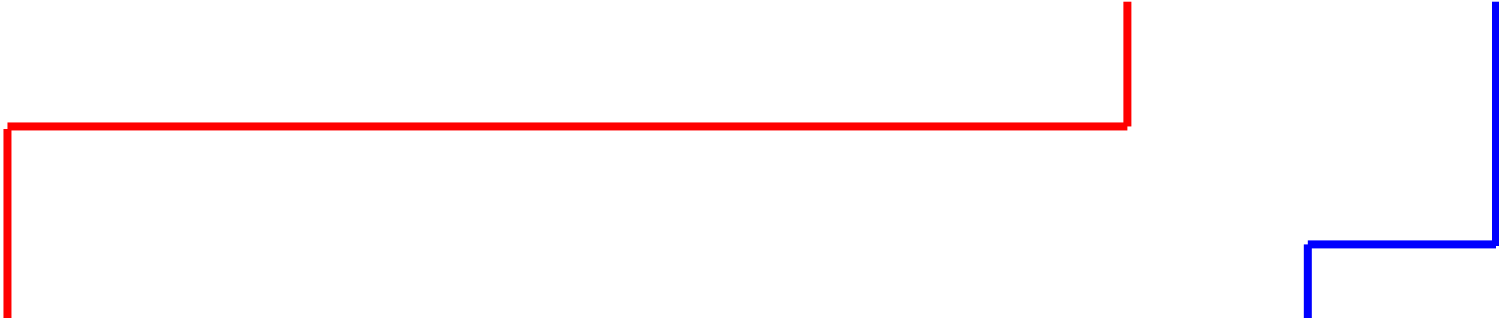
$$Z_\pi(0, x, u) = 1_{\text{safe}}(x, u)$$

# The Operator-Theoretic Perspective

Define a *linear* operator $A_\pi$

$$A_\pi \beta(x, u) := \begin{cases} \mathbb{E}_\pi[\beta(x^+, u^+)] & (x, u) \text{ safe} \\ 0 & (x, u) \text{ unsafe} \end{cases}$$

$$Z_\pi(t, x, u) = A_\pi^t 1_{\text{safe}}(x, u) = \underbrace{A_\pi \circ \cdots \circ A_\pi}_{t \text{ times}} 1_{\text{safe}}(x, u)$$

$1_{\text{safe}}$ is the safety indicator function: it returns 0 if $(x, u)$ is already unsafe, 1 otherwise.

$$Z_\pi(t, x, u) = A_\pi^t 1_{\text{safe}}(x, u) \approx c \cdot \textcolor{red}{\psi_\pi(x, u)} \cdot \textcolor{blue}{\gamma_\pi^t}$$

The dominant eigenfunction $\textcolor{red}{\psi_\pi}$

- Measures safety of each state-action pair $(x, u)$

- Always nonnegative

- **A valid stochastic CBF** that can be used in safety filtering

The dominant eigenvalue $\textcolor{blue}{\gamma_\pi}$

- Safety of the overall closed-loop system

- Always between 0 and 1

# Learning

## 1) Eigenpair learning

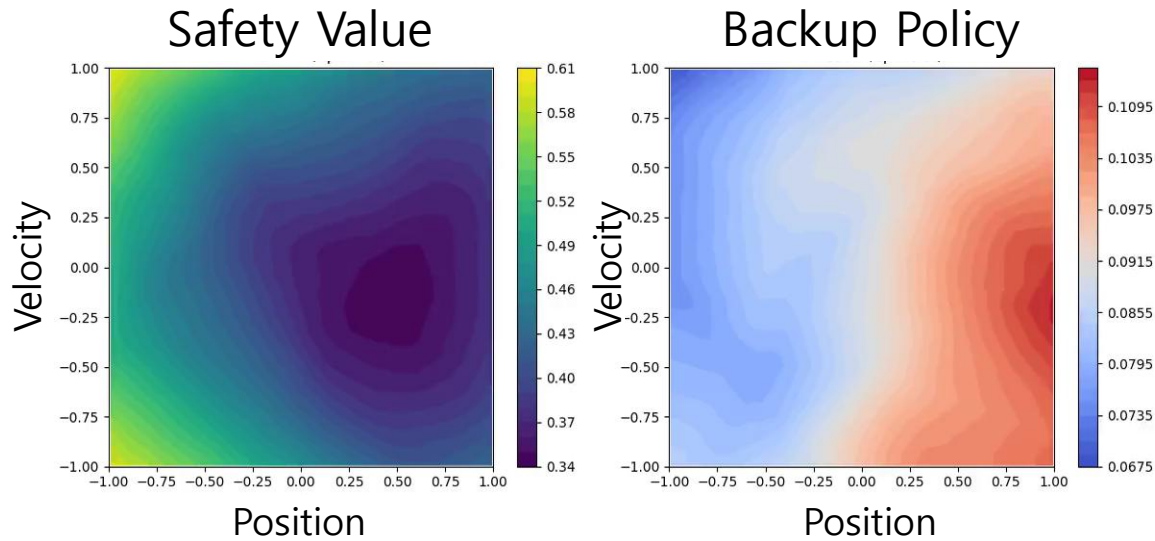$$J_{\text{eig}}[\psi, \lambda] = W_\lambda \cdot \mathbb{E}_{(x,u,x^+)\sim\mathcal{D},u^+\sim\pi}\left[\left(\psi(x^+,u^+) - \lambda\psi(x,u)\right)^2\right] + W_n \cdot \left(1 - \mathbb{E}_{(x,u,\cdot)\sim\mathcal{D}}\psi(x,u)\right)^2$$

Eigenpair loss · Normalization loss

## 2) Backup policy learning (DDPG-style)

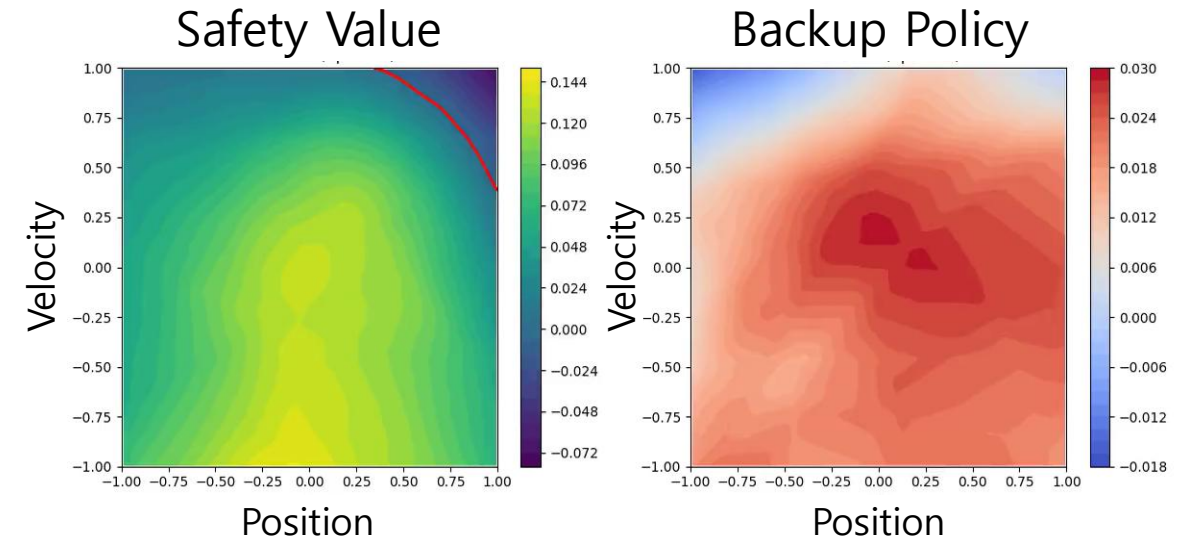$$J_{\text{policy}}[\pi] = -\mathbb{E}_{(x,\cdot,\cdot)\sim\mathcal{D},u\sim\pi}[\psi(x,u)]$$

# Stochastic Double Integrator (Gaussian noise on acceleration)

## EigenSafe
### (ours)

Safety Value

Backup Policy



The Eigenfunction correctly evaluates
relative safety across state-action pairs

## Hamilton-Jacobi Reachability + RL
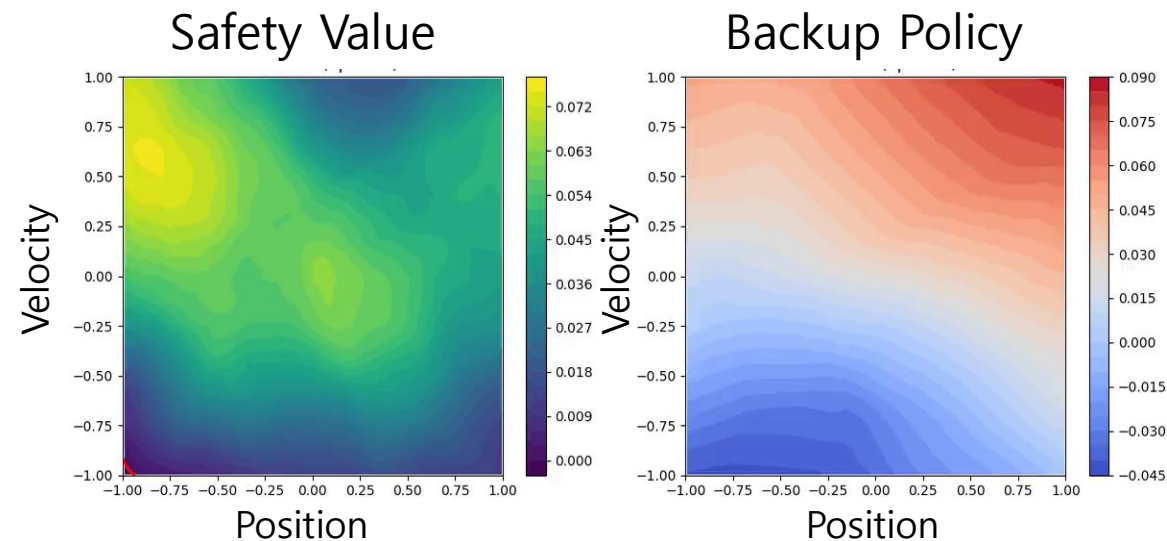### Fisac et al., ICRA 2019

Safety Value

Backup Policy



HJ reachability is not designed for stochastic
systems and fails to compute a nonempty safe set

# **EigenSafe Generalizes to Deterministic Systems** (Deterministic Double Integrator)

## **EigenSafe**
(ours)

Safety Value

Backup Policy



The Eigenfunction is the indicator for
the biggest control invariant set.

## Hamilton-Jacobi Reachability + RL
Fisac et al., ICRA 2019

Safety Value

Backup Policy

cheetah_flip
The cheetah should stay upright
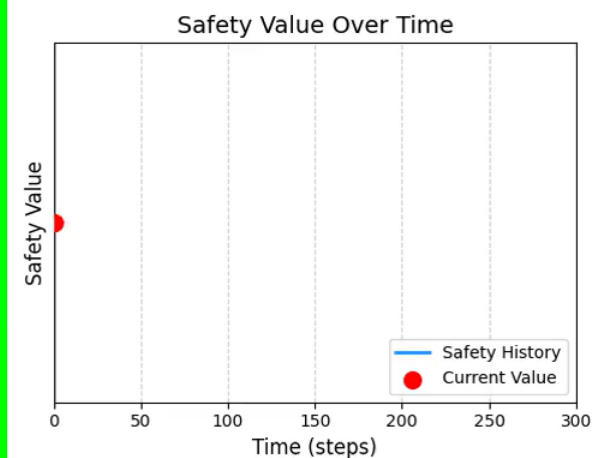
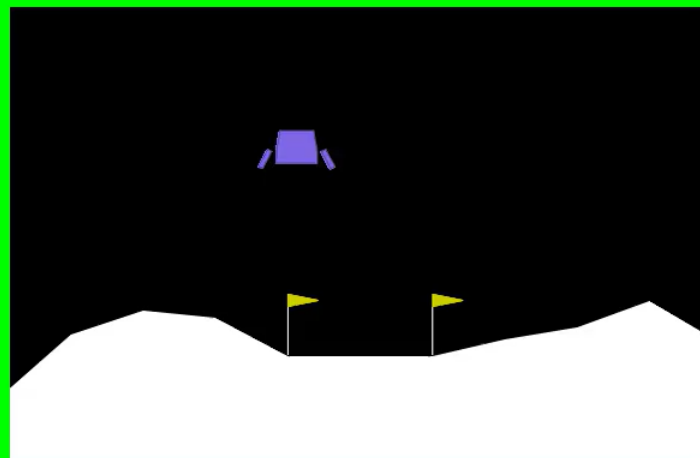cheetah_run
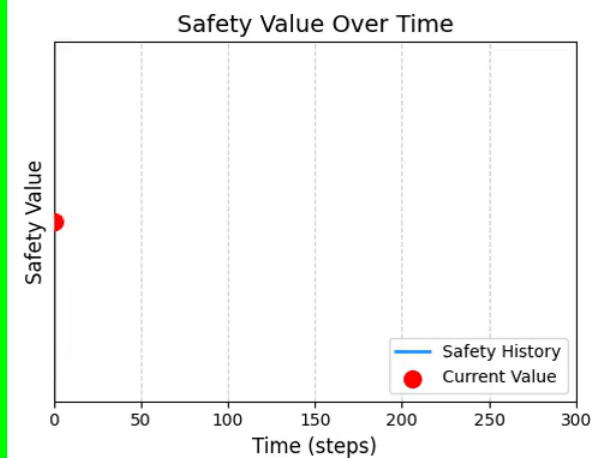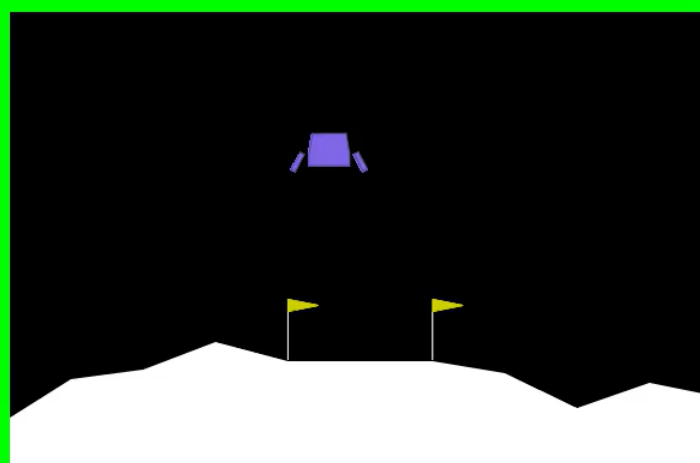The cheetah should move forwards

# lunar_lander

The lander should properly land

# Thank you!

## Contact

**Inkyu Jang**
Department of Aerospace Engineering, Automation and Systems Research Institute
Seoul National University, Korea
janginkyu.larr@gmail.com

**Jonghae Park**
PhD Student, SNU

**Chams Mballo**
Postdoc, UC Berkeley

**Sihyun Cho**
PhD Student, SNU

**Claire Tomlin**
Professor, UC Berkeley

**H. Jin Kim**
Professor, SNU