# What is holding us back?



Four Ingredients for Safety
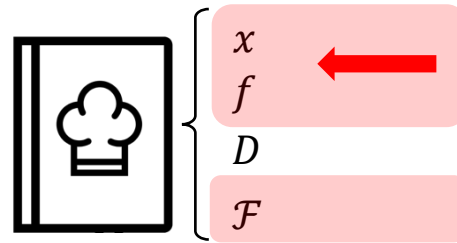
$$\begin{cases} x \\ f \\ D \\ \mathcal{F} \end{cases}$$

What is *state space, X*?

What are the *dynamics, f*?

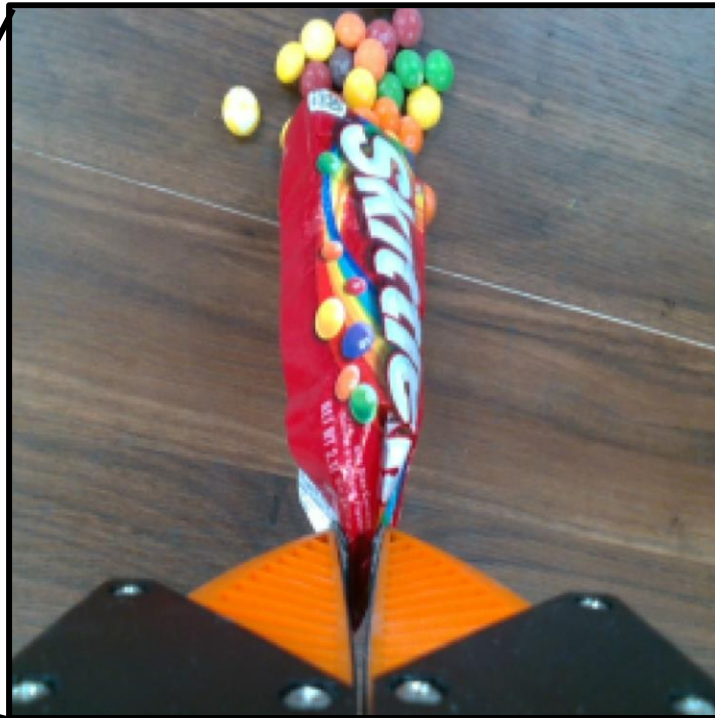What is *failure set, $\mathcal{F} \subset X$*?

How can safety go beyond collision-avoidance?

We need to reason about safety
*beyond* hand-designed state and dynamics

# Although we don't know how to write $x$ or $\mathcal{F}$ by hand...

state    failure



Failure is observable from high-dimensional observations!

But how can we *predict* if the robot's actions will result in these (failure) observations?

🔑 **Idea:** Compute a safety filter in the *latent state space* learned by generative world models

**Before**

$$x' = f_x(x, a)$$

$$\mathcal{F} \subset X$$

**Ours**

$$z' = f_z(z, a)$$

$$\mathcal{F} \subset Z$$

K. Nakamura, L. Peters, A. Bajcsy. "Generalizing Safety Beyond Collision-Avoidance via Latent-Space Reachability Analysis." RSS, 2025.

🔑 **Idea:** Compute a safety filter in the *latent state space* learned by generative world models

*Before*

$$x' = f_x(x, a)$$

$$\mathcal{F} \subset X$$

*Ours*

$$z' = f_z(z, a)$$

$$\mathcal{F} \subset Z$$

**States** (*pose, joint angles, velocities*)

**Safety Spec** (*robot falls*)

$$x = \left[ \mathbf{p}, \dot{\mathbf{p}}, \boldsymbol{\theta}, \dot{\boldsymbol{\theta}}, \boldsymbol{\theta}_{\mathrm{J}}, \dot{\boldsymbol{\theta}}_{\mathrm{J}} \right]$$

$$g(x) = \min_i \left\{ p_g^i - \bar{p}_g^{\,i} \right\},$$

**Dynamics** (*first-principles model, physics simulator*)

**States** (*embedding of image(s)*)

**Safety Spec** (*classifier on embedding*)

$z$    $\mathcal{E}$   $z$     $z$   $\mathcal{F}$   ✅ / ❌

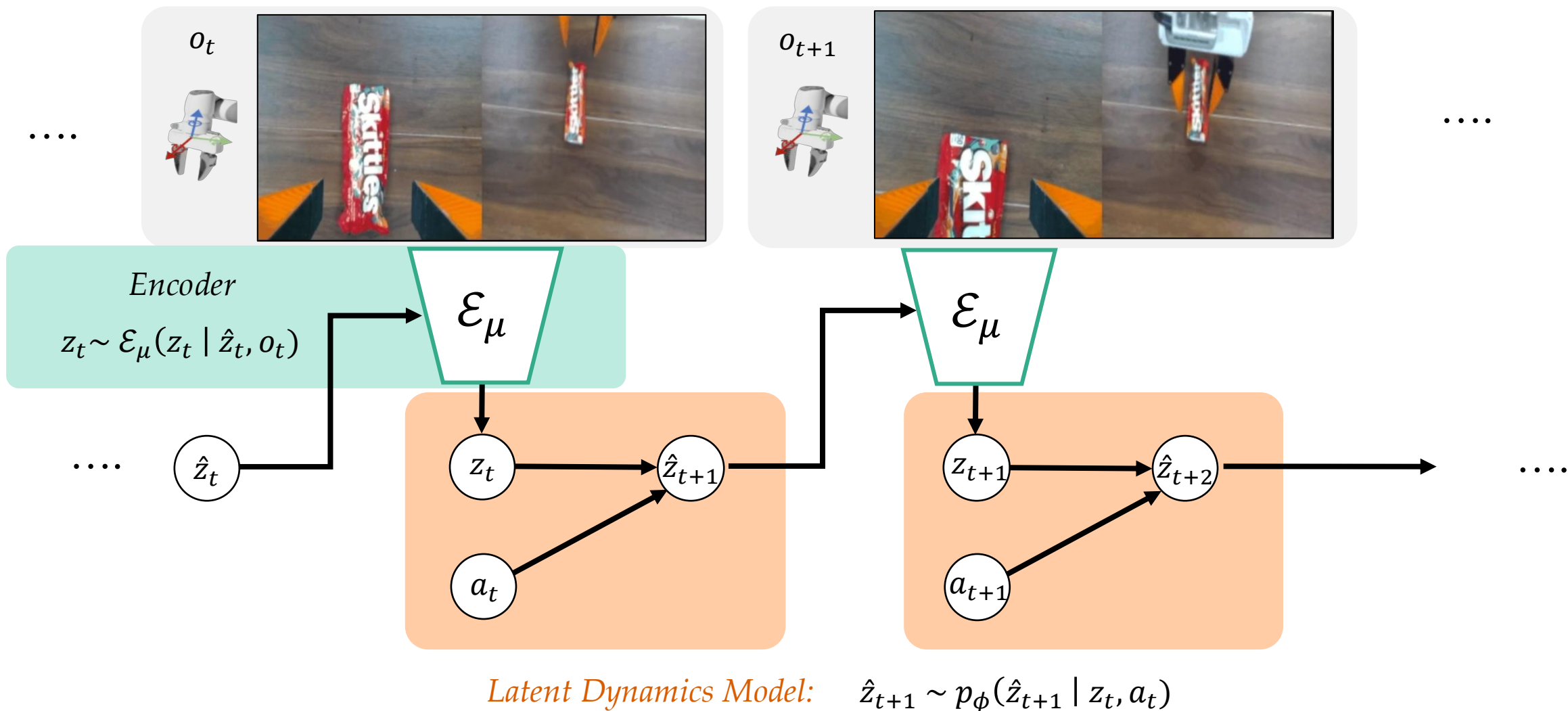**"World Model" Dynamics** (*operates on embedding*)

$$z' = f_z(z, a)$$

# **World Model Training** Time

# World Model Training Time



$o_t$

$o_{t+1}$

*Encoder*

$z_t \sim \mathcal{E}_\mu(z_t \mid \hat{z}_t, o_t)$

$\mathcal{E}_\mu$

$\mathcal{E}_\mu$

$\hat{z}_t$

$z_t$

$\hat{z}_{t+1}$

$a_t$

$z_{t+1}$

$\hat{z}_{t+2}$

$a_{t+1}$

*Latent Dynamics Model:*    $\hat{z}_{t+1} \sim p_\phi(\hat{z}_{t+1} \mid z_t, a_t)$

*Training loss*: reconstruction or teacher-forcing (min. diff. between $\hat{z}_t$ and $z_t$) + auxiliary losses

# **World Model Training** Time



$o_t$

$o_{t+1}$

....

*Encoder*

$z_t \sim \mathcal{E}_\mu(z_t \mid \hat{z}_t, o_t)$

$\mathcal{E}_\mu$

$\mathcal{E}_\mu$

....

$\hat{z}_t$

$z_t$

$\hat{z}_{t+1}$

$z_{t+1}$

$a_t$

$a_{t+1}$

*Latent Dynamics Model:* $\hat{z}_{t+1} \sim p_\phi(\hat{z}_{t+}$

*Training loss:* reconstruction or teacher-forcing (min. diff. between $\hat{z}_t$ and $z_t$) + auxiliary losses

# **Safety Analysis** Time

Latent Safety Problem

$$V(z_0) := \max_{\pi} \left( \min_{t \geq 0} \ \ell_\theta(z_t^\pi) \right)$$

Initial Observation



Imagined Failure



*Latent Dynamics Model*

$$z' = f_z(z, a)$$

$z_0$

$z$

$\mathcal{Z}$

*Latent Failure Set* ("spilled skittles")

*World Model*

# **Safety Analysis** Time

## Latent Hamilton-Jacobi Safety Bellman Equation

$$V(z) = \min\{\ell_\theta(z), \max_{a \in \mathcal{A}} \mathbb{E}_{\hat{z}' \sim p_\phi(\cdot|z,a)}[V(\hat{z}')]\}$$



Imagined Failure

Robot is **doomed to** fail

*Latent Unsafe Set*
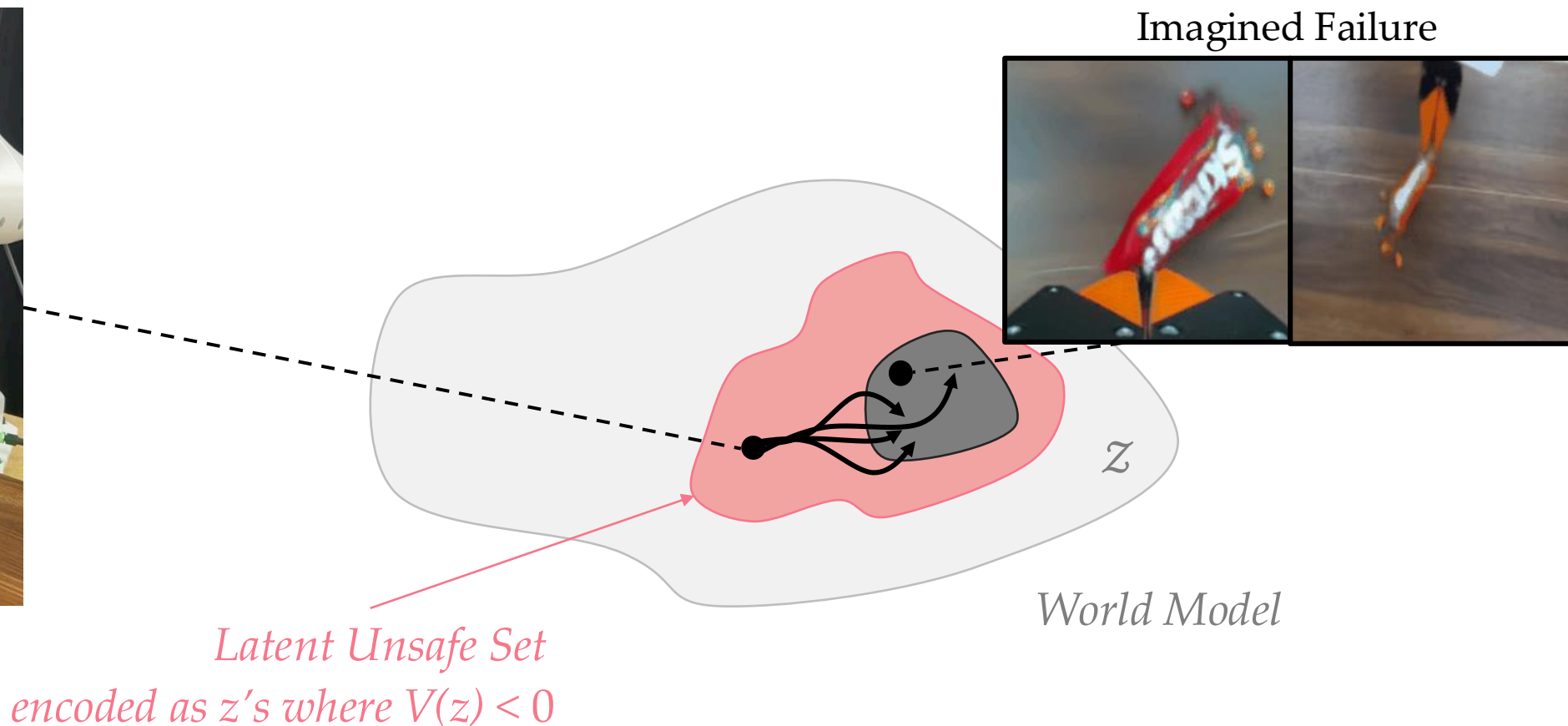*encoded as z's where V(z) < 0*

*World Model*

$z$

# **Safety Analysis** Time

Latent Hamilton-Jacobi Safety Bellman Equation

$$V(z) = \min\{\ell_\theta(z), \max_{a \in \mathcal{A}} \mathbb{E}_{\hat{z}' \sim p_\phi(\cdot | z, a)}[V(\hat{z}')]\}$$

⚠️

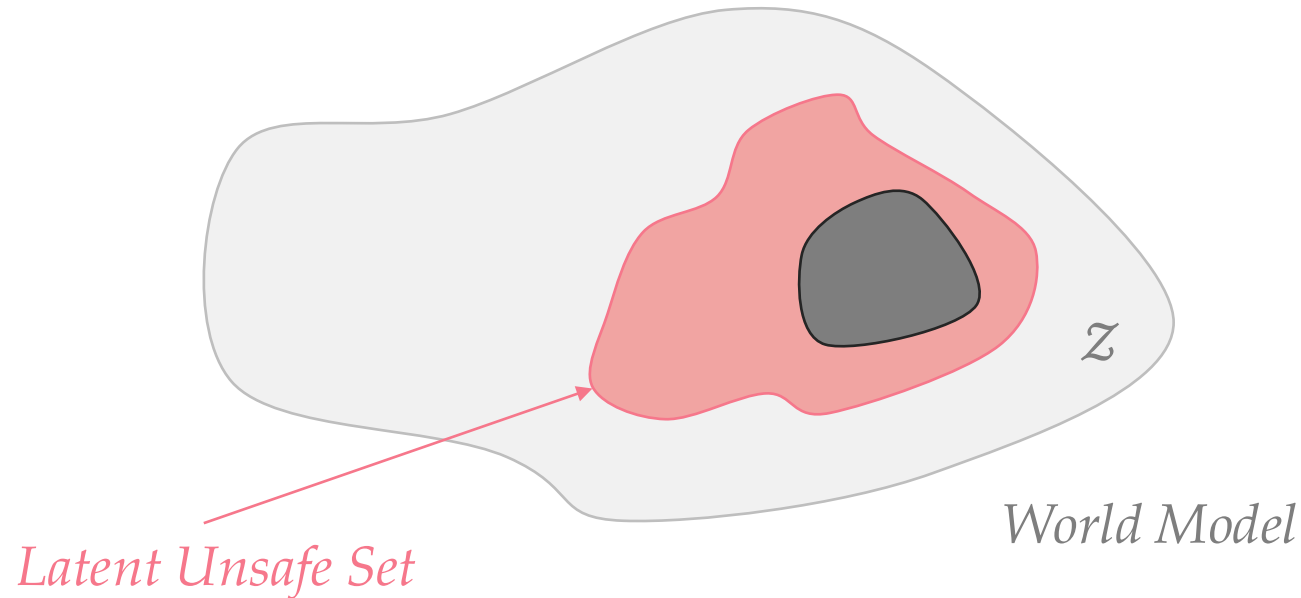*Challenge:*
latent state (z) is
high-dimensional!

e.g., **512-D!**

Note: still better
than space of all
*image observations:*
2 x 3 × 128 × 128 =
**98,304-D**



$z$

*World Model*

*Latent Unsafe Set*

**Safety Analysis** Time

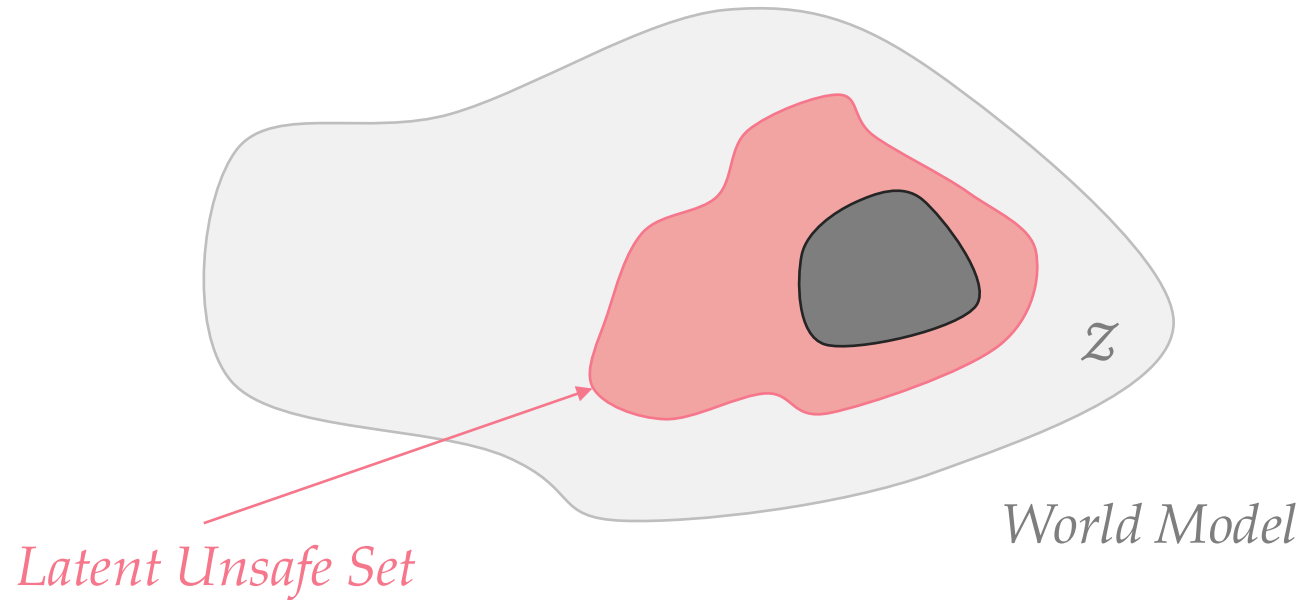Approximation via Reinforcement Learning in World Model

$$V(z) = (1 - \gamma)\ell_\theta(z) + \gamma \min\{\ell_\theta(z), \max_{a \in \mathcal{A}} \mathbb{E}_{\hat{z}' \sim p_\phi(\cdot|z,a)}[V(\hat{z}')]\}$$
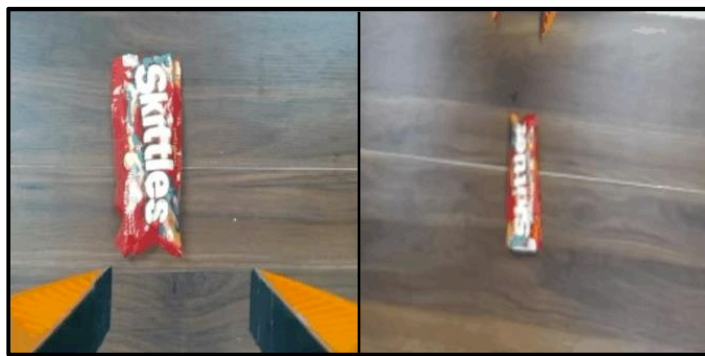


*Latent Unsafe Set*

*World Model*

$z$

**Related Work:**
[Fisac*, Lugovoy* et al. "Bridging Safety Analysis and RL", ICRA 2019]

Deployment-time Guardrail

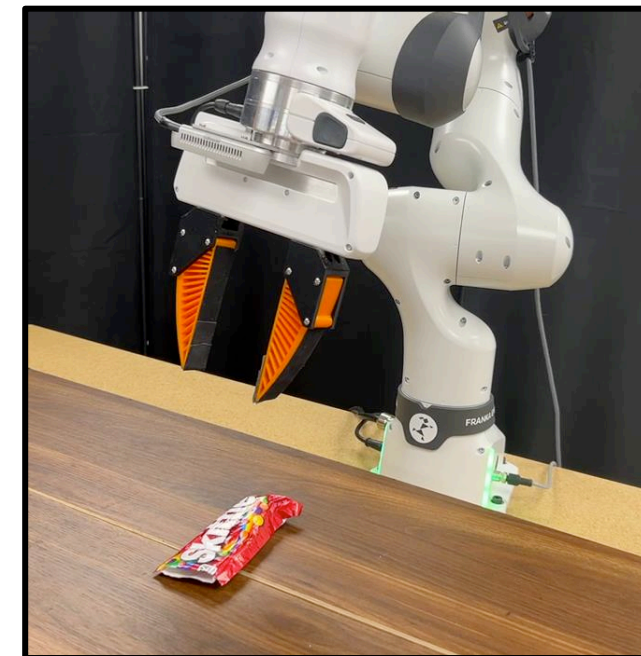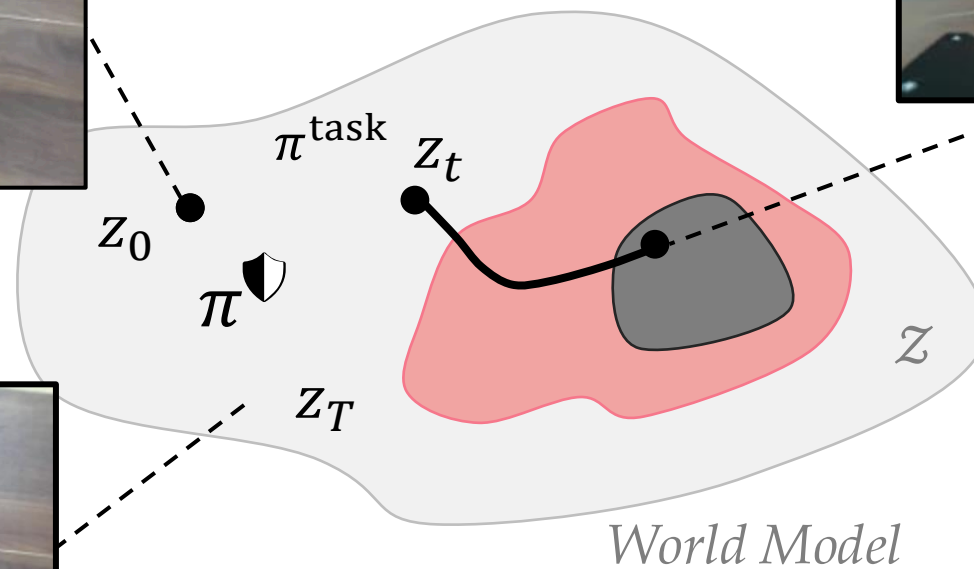Safety **Monitor** $V^{\mathbb{O}}(z)$  Safety **Policy** $\pi^{\mathbb{O}}(z)$

Imagined Failure

Observation

Safety Filtered

$\pi^{\text{task}}$ $z_t$

$z_0$

$\pi^{\mathbb{O}}$

$z_T$

$\mathcal{Z}$

*World Model*

No Safety Filter
*(direct teleop.)*

🚫

*Failure
happens (spill)!*

**Safety Computation Time** *(Offline)*

$\mathcal{U}_Z$

$\mathcal{F}_Z$

$\mathcal{Z}$

# Robot Safeguards a Human Teleoperator

*Robot allows a safe grasp…*      *Sliding motion is **filtered to slow** …*      *But unsafe pickup is **filtered to stop**!*
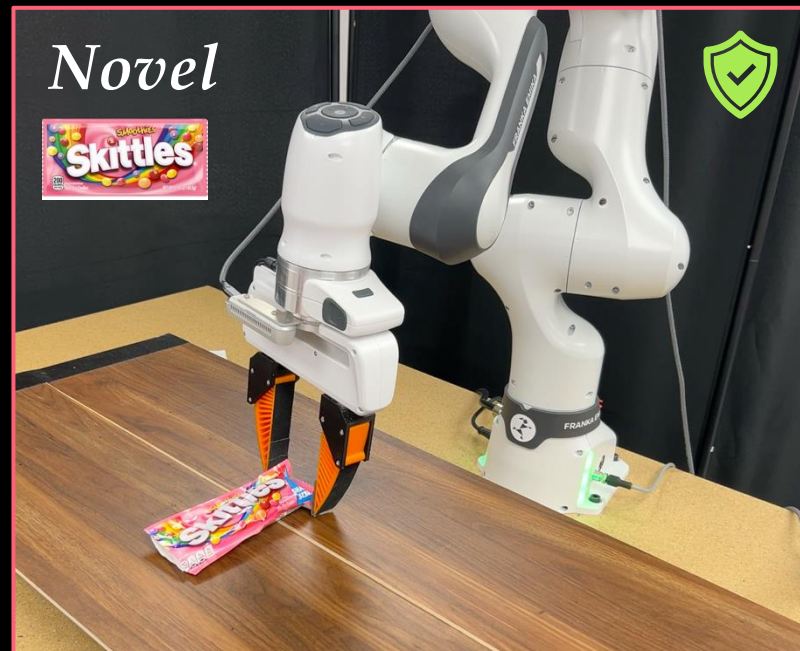
***Green border** = Robot is in control!*

*Robot POV*

$V(z')$

# Dirty Laundry

*Failure* – Observability (of bag opening)   *Failure* – OOD Dynamics (peanut M&Ms)

# More sophisticated safety filtering can "remove" unsafe modes but maintain task performance

Base Diffusion Policy

"Switching" Safety

Optimization-based Safety



*Executes unsafe interaction mode*

*Stops unsafe lifting motion*

*Guides base policy to safer grasp*

[Nakamura, Bishop, Man, Johnson, Manchester, Bajcsy. *L4DC* 2026]