

Last Time

- what makes safety hard?
- safety filters

Lecture 3

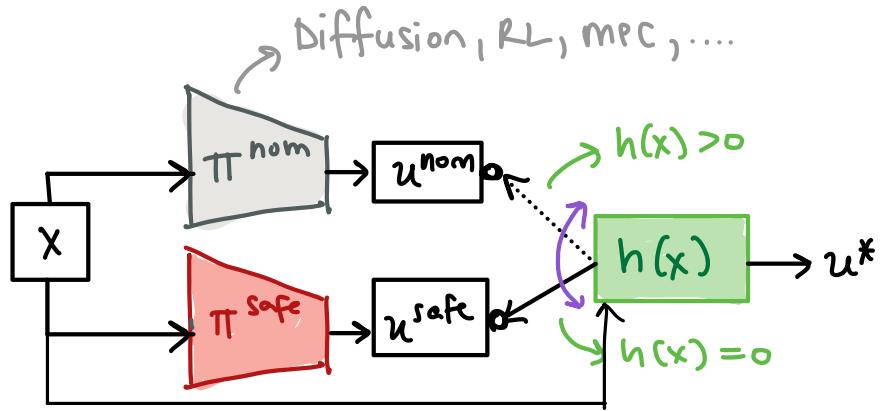
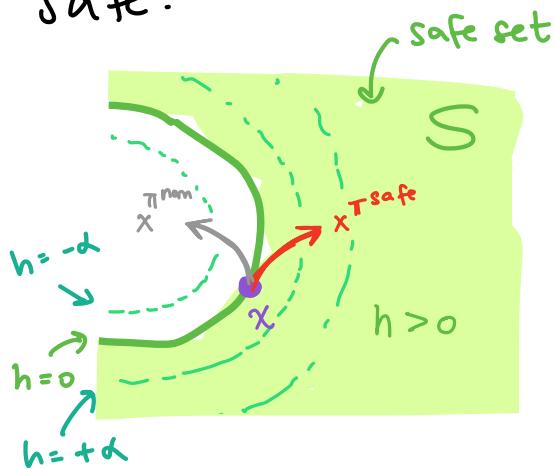
EAIS SP'25

Andrea Bojcsy

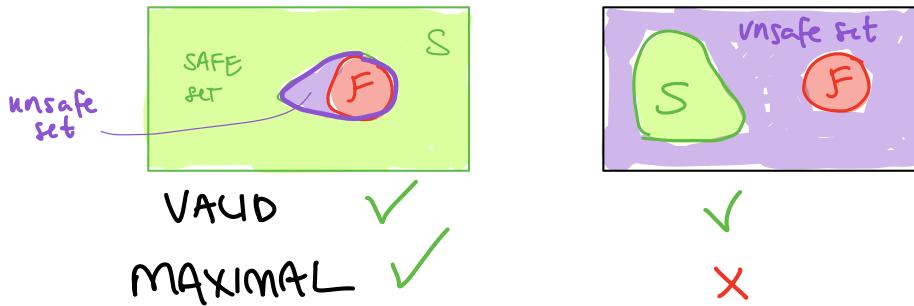
This Time:

- formalizing safe sets
- computing safety filters

Recall that last lecture we had this idea of a safety filter. It monitors and minimally modifies a base policy to keep it safe:



**!** Obtaining a VALID and not overly conservative safe set  $S \triangleq h(x)$  is really challenging for general systems



Today, we will tackle this challenge head-on.

synthesizing (i.e computing) safe sets  $\triangleq$  safety filters

↳ we will talk about a general, computational framework for getting  $S$  and  $\pi^{\text{safe}}$  and  $h(\cdot)$  fast

**TODAY** → ① is guaranteed to be VALID & MAXIMAL

**NEXT** → ② naturally handles disturbances robustly

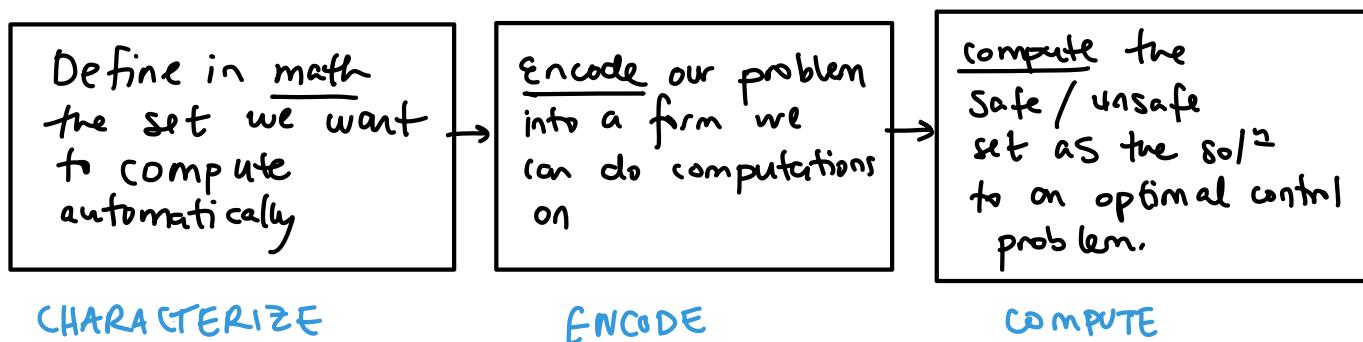
**NEXT-NEXT** → ③ is compatible w/ modern comp. tools { RL  
SSL !

## Formalize safety via reachability

We want to compute optimal controllers that ensure our robot never enters failure, AND figure out from which initial conditions is the robot doomed to fail in the future. These questions / objective fall under something called "reachability analysis".

This is the fundamental problem of identifying "if a certain state of a system is reachable from an initial state of the system."

## Safety Analysis Roadmap.



While there are many ways to compute the safe/unsafe set, we will primarily study

### Hamilton-Jacobi (HJ) Reachability

Some nice properties of this paradigm:

- 1) encode control bounds & state constraints
- 2) automatically gives both safe sets AND safe policy
- 3) general nonlinear dyn. systems
- 4) robustness to uncertainty / other agents + adversaries.

Define in math the set we want to compute automatically

Encode our problem into a form we can do computations on

compute the safe / unsafe set as the sol<sup>n</sup> to an optimal control problem.

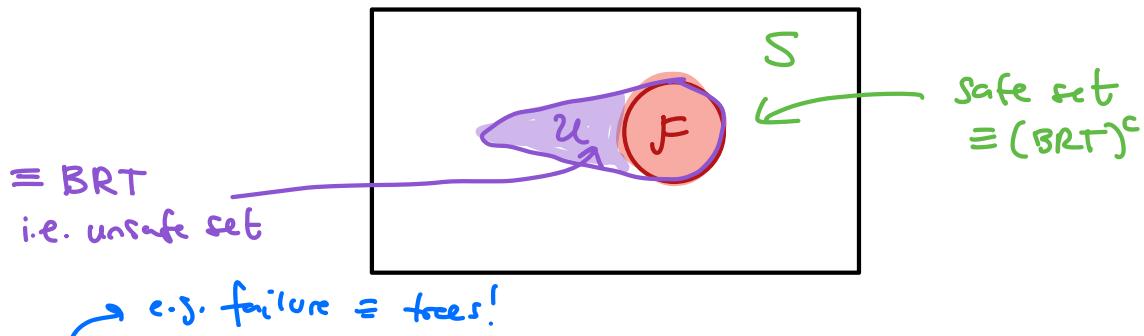
CHARACTERIZE

ENCODE

COMPUTE

How do we mathematically describe the safe set  $S$  and/or the unsafe set  $\mathcal{U}$ ?

The BACKWARDS REACHABLE TUBE (BRT) of a (failure) set  $\mathcal{F}$  and a dynamical system  $\dot{x} = f(x, u)$  is precisely the set of all states that will eventually reach  $\mathcal{F}$  despite the robot's best control efforts.



let  $\mathcal{F} \subset \mathcal{X}$  be the set of states we want to do analysis over.

let  $BRT(t) \subseteq \mathcal{X}$  at time  $t$  (typically unsafe set):

BACKWARDS REACHABLE TUBE (BRT) of set  $\mathcal{F} \subset \mathcal{X}$  and system

$\dot{x} = f(x, u)$  is:

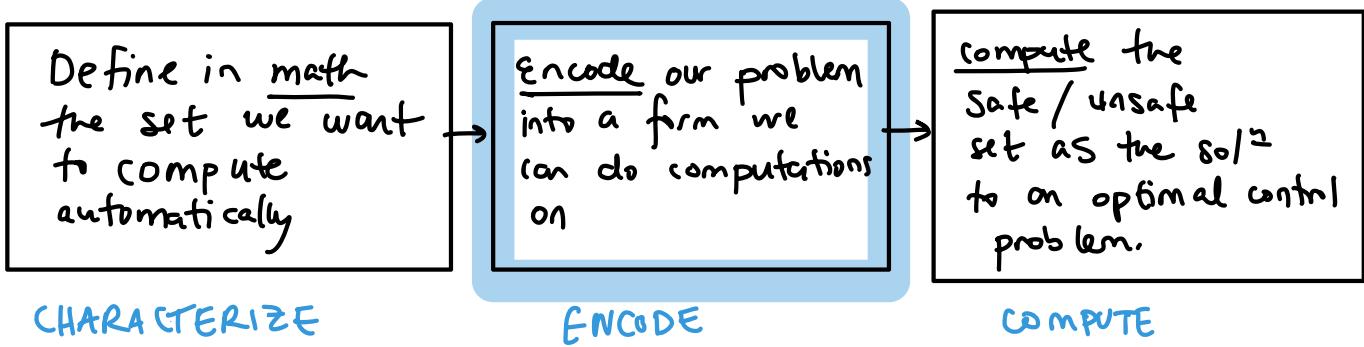
$$BRT(t) = \left\{ \underbrace{x \in \mathcal{X}}_{\substack{\text{set of initial} \\ \text{states s.t....}}} : \forall \underbrace{u(\cdot) \in U_t^T}_{\substack{\text{initial} \\ \text{states} \\ \text{for all ctrl.} \\ \text{signals}}} \underbrace{x_{x,t}^{u(\cdot)}(\tau) \in \mathcal{F}}_{\substack{\text{the state traj.} \\ \text{enters the set } \mathcal{F} \\ \text{at some pt. in time } \tau}} \text{ for some } \tau \in [t, T] \right\}$$

this is the math. def<sup>n</sup> of being "doomed"! ;)

Highlights:

Failure set ( $\mathcal{F}$ ) = safety constraint

Unsafe set ( $\mathcal{U}$ ) = BRT



## HJ Reachability:

We know connection btwn. the BRT & the failure set; let's talk about computing the BRT!

HJ Reachability uses LEVEL SET METHODS to convert the BRT / constraint satisfaction problem into an optimal ctrl. problem.

Here's the process..

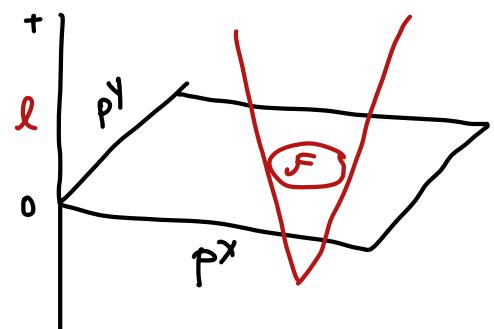
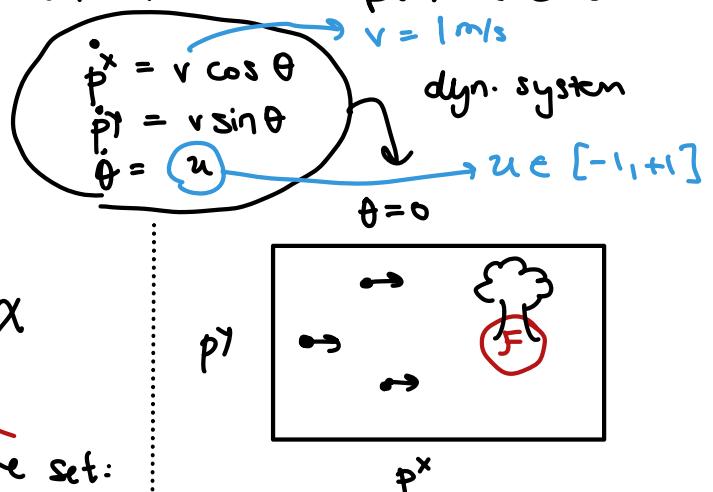
(1) we have the failure set  $\mathcal{F} \subset X$

(2) Define a function  $l(x) : X \rightarrow \mathbb{R}$  to implicitly represent this failure set:

$$l(x) < 0 \iff x \in \mathcal{F}$$

e.g. signed dist. func. is what we use in practice!

- signed-dist.  $> 0$  when  $x$  outside  $\mathcal{F}$
- signed-dist.  $< 0$  when  $x$  inside  $\mathcal{F}$
- signed-dist.  $= 0$  when  $x$  on  $\partial\mathcal{F}$



(3) Now, we want to optimize  $u(\cdot)$

with respect to  $\ell(x)$  since this is our optimal control cost function!

$$J(x, u(\cdot), t) := \min_{\tau \in [t, T]} \ell(x_{x,t}^u(\tau))$$

("closest our system got to failure when applying  $u(\cdot)$  and starting from  $x$ ")

④ By looking @ the sign of the cost  $J(\cdot, \cdot, \cdot)$  we can tell if the traj. ever entered  $\mathcal{F}$  given  $u(\cdot)$ !

If we want to stay safe, control should maximize  $J$ !

$$V(x, t) := \underset{u(\cdot) \in \mathcal{U}_t^T}{\text{maximize}} J(x, u(\cdot), t)$$

$$= \underset{u(\cdot) \in \mathcal{U}_t^T}{\text{maximize}} \left[ \min_{\tau \in [t, T]} \ell(x_{x,t}^u(\tau)) \right]$$

- If  $V(x, t) < 0$  for some state  $x_0$ , then this means that the controller  $u(\cdot)$  tried, but failed, to prevent failure despite its best efforts.  $\Rightarrow x_0 \in \text{BRT}!$

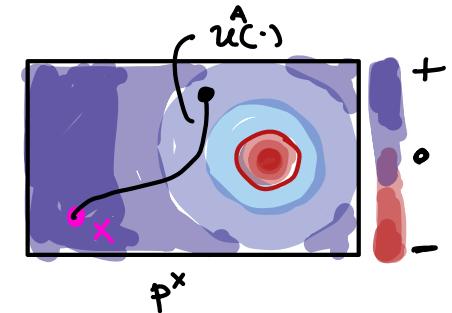
- If  $V(x, t) \geq 0$  for some  $x_0$ , then this means there exists a ctrl signal  $u(\cdot)$  that can prevent failure  $\Rightarrow x_0 \notin \text{BRT}!$

**!** Once we obtain  $V(x, t)$ , we also obtain the unsafe set (BRT)

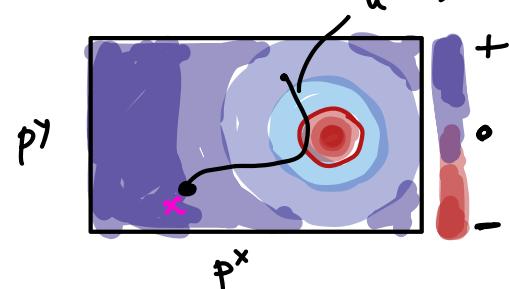
this is the  
unsafe set!

$$\text{BRT}(t) = \{x : V(x, t) < 0\}$$

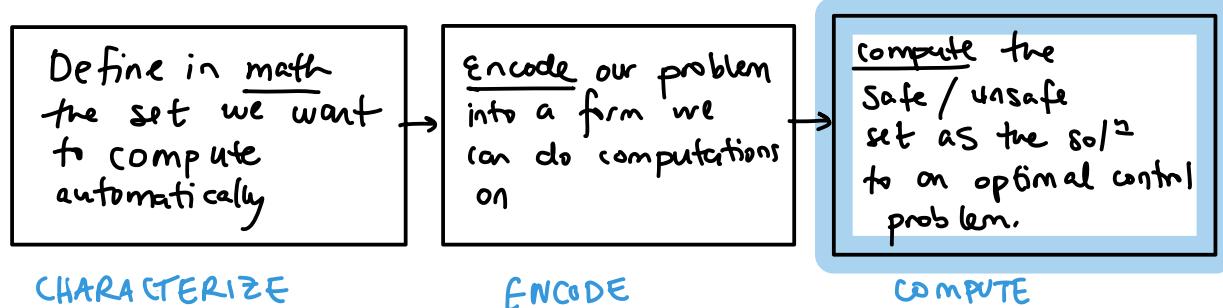
$$J(x, u^*(\cdot), t) > 0$$



$$J(x, u^*(\cdot), t) < 0$$



Now, we have an optimal ctrl. problem whose solution will automatically give us unsafe set (ERT) — how do we solve?



Hmm, but the min-over-time is not the usually running cost...  
Good news — Principle of dyn. programming still works!

$$\begin{aligned}
 V(x, t) &:= \max_{u(\cdot)} \min_{\tau \in [t, T]} l(x(\tau)) \\
 &= \max_{u(\cdot)} \min \left\{ \min_{\tau \in [t, t+\delta]} l(x(\tau)), \min_{s \in [t+\delta, T]} l(x(s)) \right\} \\
 &\quad \text{"either min happens now" } \uparrow \quad \text{... or min happens in future" } \uparrow \\
 &= \max_{u(\cdot)} \min \left\{ \min_{\tau \in [t, t+\delta]} l(x(\tau)), \underbrace{J(x(t+\delta), u(\cdot), t+\delta)}_{\text{cost in future}} \right\} \\
 &= \max_{u(\cdot)} \min \left\{ \min_{\tau \in [t, t+\delta]} l(x(\tau)), \underbrace{V(x(t+\delta), t+\delta)}_{\text{by principle of optimality}} \right\} \\
 &= \max_{u(\cdot)} \min \left\{ \min_{\tau \in [t, t+\delta]} l(x(\tau)), V(x(t+\delta), t+\delta) \right\} \quad \text{(*)}
 \end{aligned}$$

We ultimately can recover a very similar Bellman-like equation! But the key difference is we do a **min** with our  $l(x)$  function @ each backup

Safety-centric dynamic programming equation(s):

discrete time ( $t \in \mathbb{Z}$ )

Hamilton - Jacobi - Bellman Equation:

$$V_t(x) = \min \left\{ l(x), \max_{u \in U} V_{t+1}^{\stackrel{x}{\leftarrow}}(f(x, u)) \right\}$$

$$V_T(x) = l(x)$$

"remember if failing now" ... "try best not to fail in future"

continuous time ( $t \in \mathbb{R}$ )

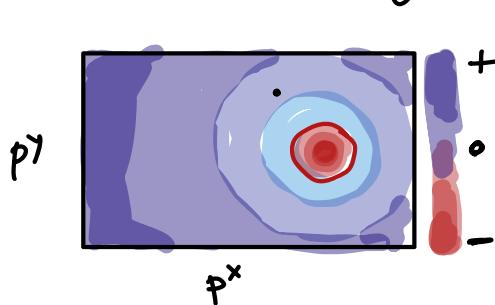
HJ variational Inequality (HJ-VI)

$$\min \left\{ l(x) - V(x, t), \frac{\partial V}{\partial t} + \max_{u \in U} \frac{\partial V}{\partial x} \cdot f(x, u) \right\} = 0$$

$$V(x, T) = l(x)$$

"keeps track of entry F" "best value change via action  $u^*$ "  
 "change in value over time"  $\dot{x}$

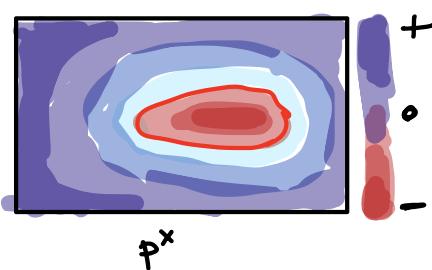
$$V(x, T) = l(x)$$



Apply Dyn. Progr.

$T-1, T-2, \dots, 0$

$$V(x, 0)$$



$$\begin{aligned} V(x, 0) \geq 0 &\Leftrightarrow x \in S \\ V(x, 0) < 0 &\Leftrightarrow x \in U \end{aligned}$$

We computed our safety filter ingredients!

**NOTE** If you are interested in derivation of HJ-VI, lets start from  $\star$

With this reformulation, we can focus on studying what happens when  $t \rightarrow 0$ ? i.e. How does value function @ current state change when we make small changes in our decisions?

As before, do TSE of  $V(x(t+\delta), t+\delta)$  around  $(x, t)$

$$V(x(t+\delta), t+\delta) \approx V(x, t) + \underbrace{\frac{\partial V}{\partial x} \cdot dx}_{= f(x, u) \cdot \delta} + \underbrace{\frac{\partial V}{\partial t} \cdot \delta}_{\text{ignore}} + h.b.t.$$

$$= \max_{u(\cdot)} \min \left\{ \underbrace{l(x(t))}_{\approx \text{for very small } \delta}, V(x, t) + \frac{\partial V}{\partial x} f(x, u) \delta + \frac{\partial V}{\partial t} \delta \right\}$$

$$\stackrel{\substack{\text{only} \rightarrow \\ \text{opt. cur. ctrl.}}}{=} \max_{u(t)} \min \left\{ l(x(t)), V(x, t) + \frac{\partial V}{\partial x} f(x, u) \delta + \frac{\partial V}{\partial t} \delta \right\} \stackrel{\substack{\text{ctrl. only influence this}}}{=}$$

$$V(x, t) = \min \left\{ l(x(t)), V(x, t) + \frac{\partial V}{\partial t} \delta + \max_{u(t)} \frac{\partial V}{\partial x} f(x, u) \delta \right\}$$

↓ subtract  $V(x, t)$  from each side

$$\Rightarrow \min \left\{ \underbrace{l(x(t)) - V(x, t)}_{\text{Case 1: this is active}}, \delta \left( \frac{\partial V}{\partial t} + \max_{u(t)} \frac{\partial V}{\partial x} \cdot f(x, u) \right) \right\} = 0$$

Case 1: this is active:

$$V(x, t) = l(x)$$

this is just HJ-Bellman PDE!

Since this statement is true for all  $\delta > 0$ , we can scale the RHS by a pos. number  $\hat{\delta}$ : it doesn't affect the minimization comparison. Thus, we can remove the  $\delta$  and achieve:

### HJ variational Inequality (HJ-VI)

$$\min \left\{ l(x) - V(x, t), \frac{\partial V}{\partial t} + \max_{u \in U} \frac{\partial V}{\partial x} \cdot f(x, u) \right\} = 0$$

$$V(x, T) = l(x)$$