Last time:

- ☐ HJI Equation
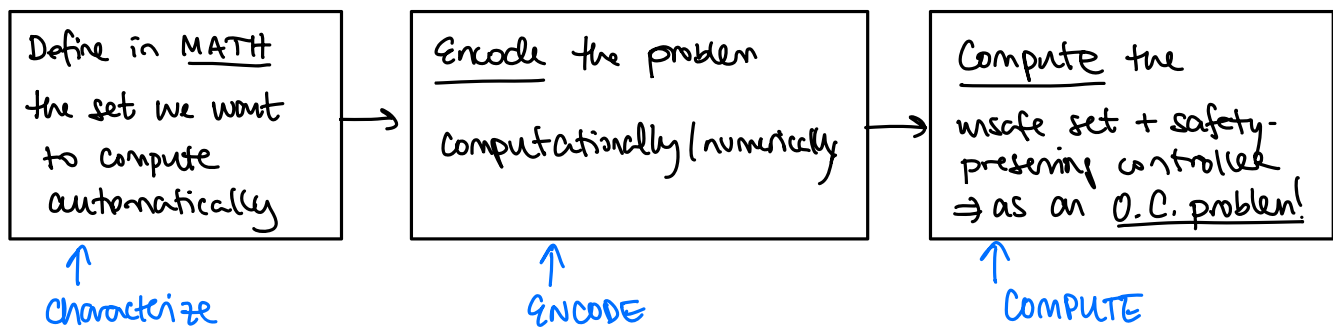- ☐ start of Safety Analysis

Today:

- ☐ HJ Safety Analysis!
- ☐ safety filters

# Formalizing safety via reachability

We now have a handle on how to solve general robust optimal control problems w/ potentially multi-cpts. But what if we wanted to ensure that our system abides by some state constraints? For example, what if we want to synthesize an optimal control that guarantees that our robot never hits an obstacle? What are the initial conditions from which robot is doomed to collide? These questions fall under reachability analysis which is a fundamental problem of identifying "if a certain state of a system is reachable from an initial state of the system":
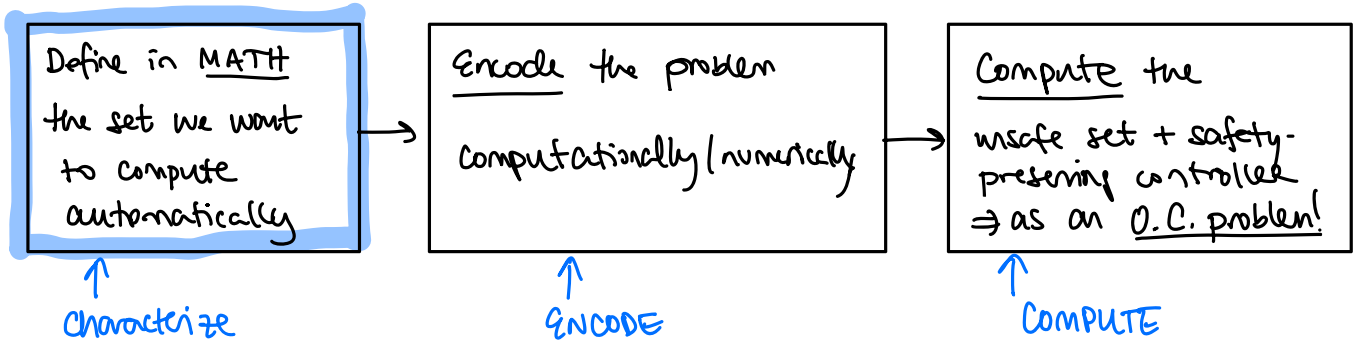
$\Rightarrow$ fundamental to program analysis, to dynamical systems, to biology!

## Safety Analysis Roadmap

| Define in <u>MATH</u> the set we want to compute automatically | Encode the problem computationally/numerically | <u>Compute</u> the unsafe set + safety-preserving controller $\Rightarrow$ as an <u>O.C. problem!</u> |
|---|---|---|
| ↑ Characterize | ↑ ENCODE | ↑ COMPUTE |

While there are many ways to compute safe/unsafe sets, we will study Hamilton-Jacobi Reachability analysis of safe sets & controllers. Why HJ?

1) automatically handles control bounds / state constraints

2) —"— synthesizes safe set <u>AND</u> safety controller

3) general nonlinear systems $\dot{x} = f(x, u)$

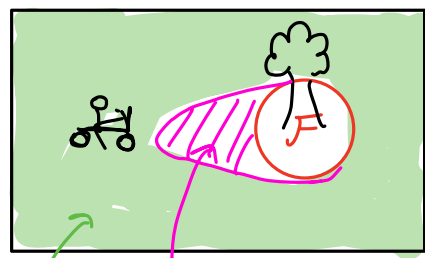4) multi-agent interactions (non-deterministic uncertainty)

Reachability is very expressive framework for defining $\underbrace{\text{safety}}_{\text{"collision"}}$ & $\underbrace{\text{liveness}}_{\text{"goal-reach"}}$

Let's focus on  BACKWARDS REACHABLE TUBES (BRT) of dyn. sys.

BRT is the set of all states of a system that will eventually reach some "target" set despite the robot's best control effort.

Let $F$ be the failure set; then the BRT represents the potential unsafe set of states for the system & thus should be avoided.



BRT (unsafe set)

BRT$^c$ (safe set)

Let $\mathcal{L} \subseteq \mathbb{R}^n$ be the  set of states we are interested in performing analysis on. Let BRT$(t) \subseteq \mathbb{R}^n$ be the BRT at time $t$ (typically unsafe set):
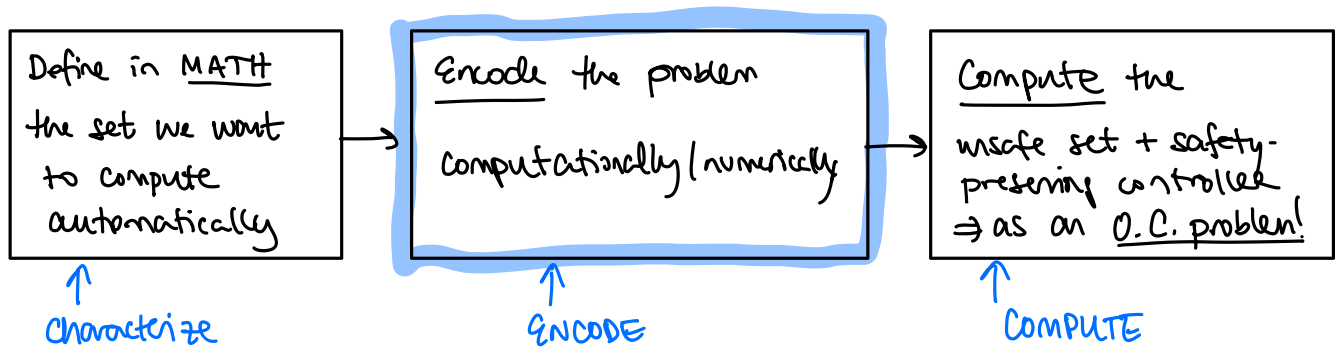
[ROBUST] BACKWARDS REACHABLE TUBE (BRT) of a set $\mathcal{L} \subseteq \mathbb{R}^n$ and a dynamical system $\dot{x} = f(x, u, d)$ is:

$$\text{BRT}(t) = \{ x \in \mathbb{R}^n : \forall u(\cdot) \in \mathbb{U}_t^T, \exists d(\cdot) \in \mathbb{D}_t^T, \ X_{x,t}^{u,d}(\tau) \in \mathcal{L} \text{ for some } \tau \in [t, T] \}$$

this is set of initial states @ time $t$ — initial states — for all ctrl signals — exists a disturbance — s.t. state traj. enters $\mathcal{L}$... — ...@ same time $\tau$

Intuitively, BRT(t) computes the set of all starting states from which no matter what the controller does, there exists a disturbance that drives the system into $\mathcal{L}$ (e.g., our failure set).

✳ TO compute unsafe set, you need to compute BRT of $\color{red}{F}$!

Failure set = constraint
unsafe set = BRT

| Define in MATH the set we want to compute automatically | Encode the problem computationally / numerically | Compute the unsafe set + safety-preserving controller ⟹ as an O.C. problem! |
|---|---|---|

↑ Characterize          ↑ ENCODE          ↑ COMPUTE

## HJ Reachability

know connection (mathematically) btwn. BRT & the failure/target set. let's talk about computing! HJ formulates this computation as an optimal control problem! lets us use all the O.C. tools (and in fact, SOTA algorithms) to compute BRT automatically.

HJ Reachability uses    level set methods to convert BRT characterization
to an O.C. problem.
⟶ numerical analysis on arbitrarily-shaped $\color{red}{F}$
⟶ propagate the influence of ctrl./dist. on the "growth" of BRT

**PROCEDURE:**

1) We have a failure set $F \subseteq \mathbb{R}$

2) Define a function $\ell(x): \mathbb{R}^n \to \mathbb{R}$ to implicitly encode this failure set:

$$\ell(x) \leq 0 \iff x \in F$$

One such function is Signed distance to $F$

→ S.D. $> 0$ when $x$ outside $F$
   S.D. $< 0$ —"— $x$ inside $F$
   S.D. $= 0$ —"— @ boundary

$$\ell(x) = \sqrt{x^2 + y^2} - R$$

3) we want to optimize $u(\cdot)$ w.r.t. $\ell(x)$ since $\ell(x)$ is our optimal control cost function!

$$J(x, u(\cdot), d(\cdot), t) = \min_{\tau \in [t, T]} \ell\left(\mathbf{x}_{x,t}^{u(\cdot), d(\cdot)}(\tau)\right)$$

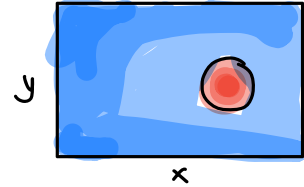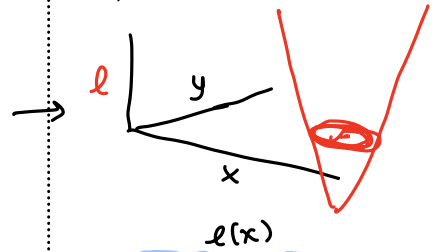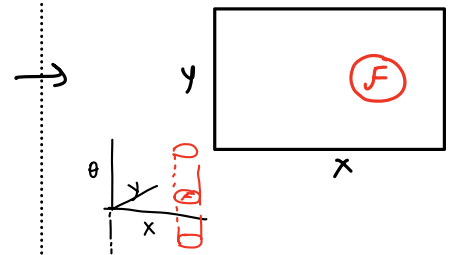⇒ "closest our system ever gets to $F$ when applying $(u(\cdot), d(\cdot))$ starting from $x$ at time $t$"

⇒ By looking @ the <u>sign of $J(\cdot, \cdot)$</u> we can tell if our traj. <u>ever entered</u> $F$!

$$V(x, t) = \max_{u(\cdot) \in \Gamma_u} \min_{d(\cdot) \in \Gamma_d} J(x, u(\cdot), d(\cdot), t)$$

$$= \max_{u(\cdot) \in \Gamma_u} \min_{d(\cdot) \in \Gamma_d} \left( \min_{\tau \in [t, T]} \ell(x(\tau)) \right)$$

⇒ $V(x^{init}, t) \leq 0$ for some state $x^{init}$ ⇒ controller $(u(\cdot))$ tried hardest but couldn't do anything to prevent $x_\tau \in F$
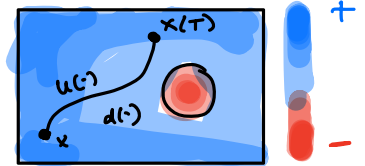
⇒ $x^{init} \in BRT(t)$
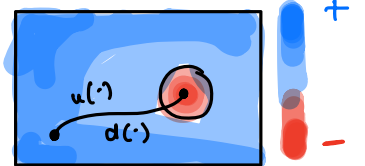
⇒ $\boxed{BRT(t) = \{x : V(x, t) \leq 0\}}$



$J(x, u, d, t) > 0$

$J(x, u, d, t) \leq 0$

$J(x, u, d, t) < 0$

Key diff. btwn what we had before is $\underline{\text{min over time!}}$

$\underline{\text{Good NEWS}}$: we can still use the principle of optimality & D.P. to compute V!

$$V(x,t) = \max_{u(\cdot)\in\Gamma_u} \min_{d(\cdot)\in\Gamma_d} \left( \min_{\tau\in[t,T]} \ell(x(\tau)) \right)$$

"... or, this failure happens in future"

$$= \max_{u(\cdot)\in\Gamma_u} \min_{d(\cdot)\in\Gamma_d} \boxed{\min} \left( \min_{\tau\in[t,t+\delta]} \ell(x(\tau)), \overbrace{\min_{s\in[t+\delta,T]} \ell(x(s))} \right)$$

"either my F violation happens right now"

$$= \max_{u(\cdot)\in\Gamma_u} \min_{d(\cdot)\in\Gamma_d} \min \left( \min_{\tau\in[t,t+\delta]} \ell(x(\tau)), \underbrace{\min_{s\in[t+\delta,T]} \ell(x(s))} \right)$$

$$:= J(x(t+\delta), u(\cdot), d(\cdot), t+\delta)$$

$$= \max_{u(\cdot)\in\Gamma_u} \min_{d(\cdot)\in\Gamma_d} \min \left( \min_{\tau\in[t,t+\delta]} \ell(x(\tau)), \underline{V(x(t+\delta), t+\delta)} \right)$$

by principle of opt

---

$\underline{\text{Hamilton-Jacobi Variational Inequality}}$ (HJI-VI)

"remember if we ever fail"

$$\min \left\{ \underline{\ell(x) - V(x,t)}, \underbrace{\frac{\partial V}{\partial t} + \max_u \min_d \frac{\partial V(x,t)}{\partial x} \cdot f(x,u,d)} \right\} = 0$$

HJB-PDE!

$$V(x,T) = \ell(x)$$

---

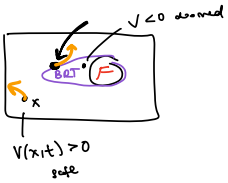| Define in $\underline{\text{MATH}}$ the set we want to compute automatically | → | Encode the problem computationally / numerically | → | Compute the unsafe set + safety-preserving controller $\Rightarrow$ as an $\underline{\text{O.C. problem!}}$ |

Characterize    ENCODE    COMPUTE

# Getting optimal control

Recall that system should stay away from $F$ encoded as $< 0$ in $\ell(x)$.

Then:

$$u^*_{safe}(x, t) = \operatorname*{argmax}_{u} \min_{d} \frac{\partial V^*(x, t)}{\partial x} \cdot f(x, u, d) \qquad (*)$$

optimal value function.



$V < 0$ doomed

BRT $F$

$x$

$V(x, t) > 0$
safe

# Least-restrictive Safety Filter

$$\underset{\substack{\uparrow \\ \text{executed ctrl.}}}{u^*(x)} = \begin{cases} \pi_{nom}(x) & \text{if } V^*(x) > 0 \\ \\ u^*_{safe}(x) & \text{if } V^*(x) = 0 \quad (\text{i.e. } x \in \partial BRT) \end{cases}$$

nominal policy (e.g. RL, MPC, ...)

"boundary of"

$\hookrightarrow$ from $(*)$

in practice $V^*(x) = \Delta$
where $\Delta > 0$ but
small (to account
for numerical error)