

Autonomous System Resilience and Guaranteed Performance in the Face of Unexpected Adversity

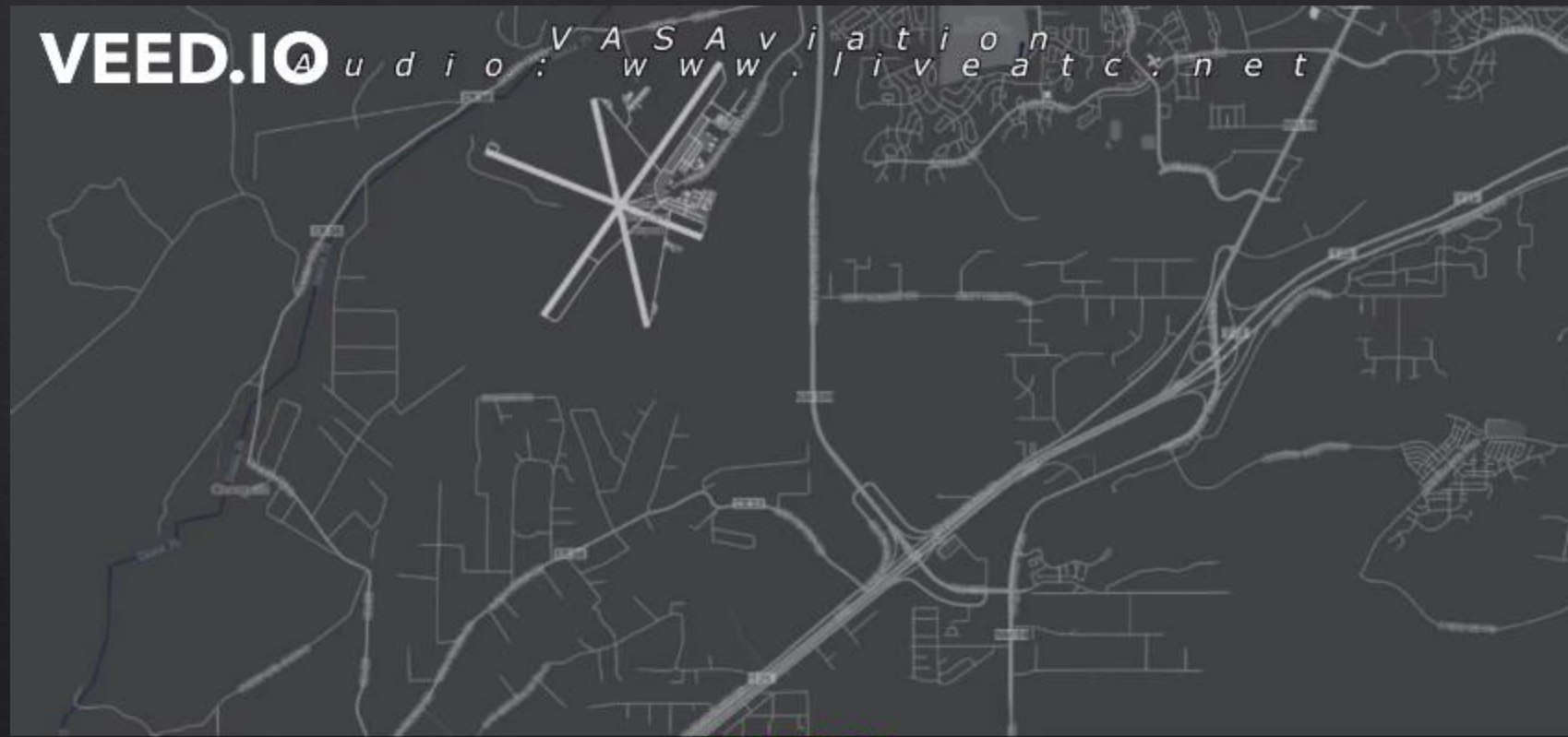
(hopefully different from all of my group's ULI-related SciTech papers)

Melkior Ornik



January 26, 2024

Motivating Example

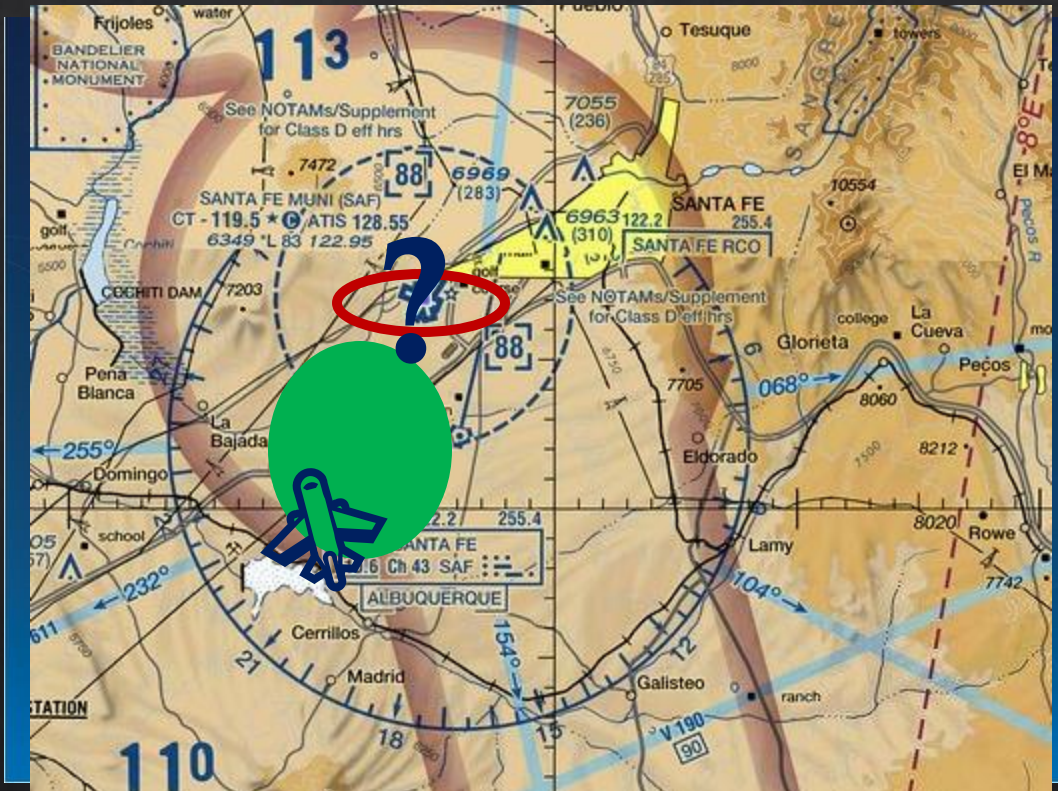


VEED.IO Audio: www.liveatc.net V A S Aviation

N1568Z
68Z, declaring an emergency! 68Z, declaring an emergency!
68Z, declaring an emergency!

The image is a screenshot of a VEED.IO audio player. The top left corner features the VEED.IO logo. The top center shows the audio source as 'www.liveatc.net' and 'V A S Aviation'. The main area is a map of an airport with a red starburst indicating an emergency location. The bottom of the screenshot shows a black bar with green text: 'N1568Z' followed by three lines of emergency declarations: '68Z, declaring an emergency!', '68Z, declaring an emergency!', and '68Z, declaring an emergency!'.

Maximal Understanding in the Face of Minimal Knowledge



Adaptive / Robust Control

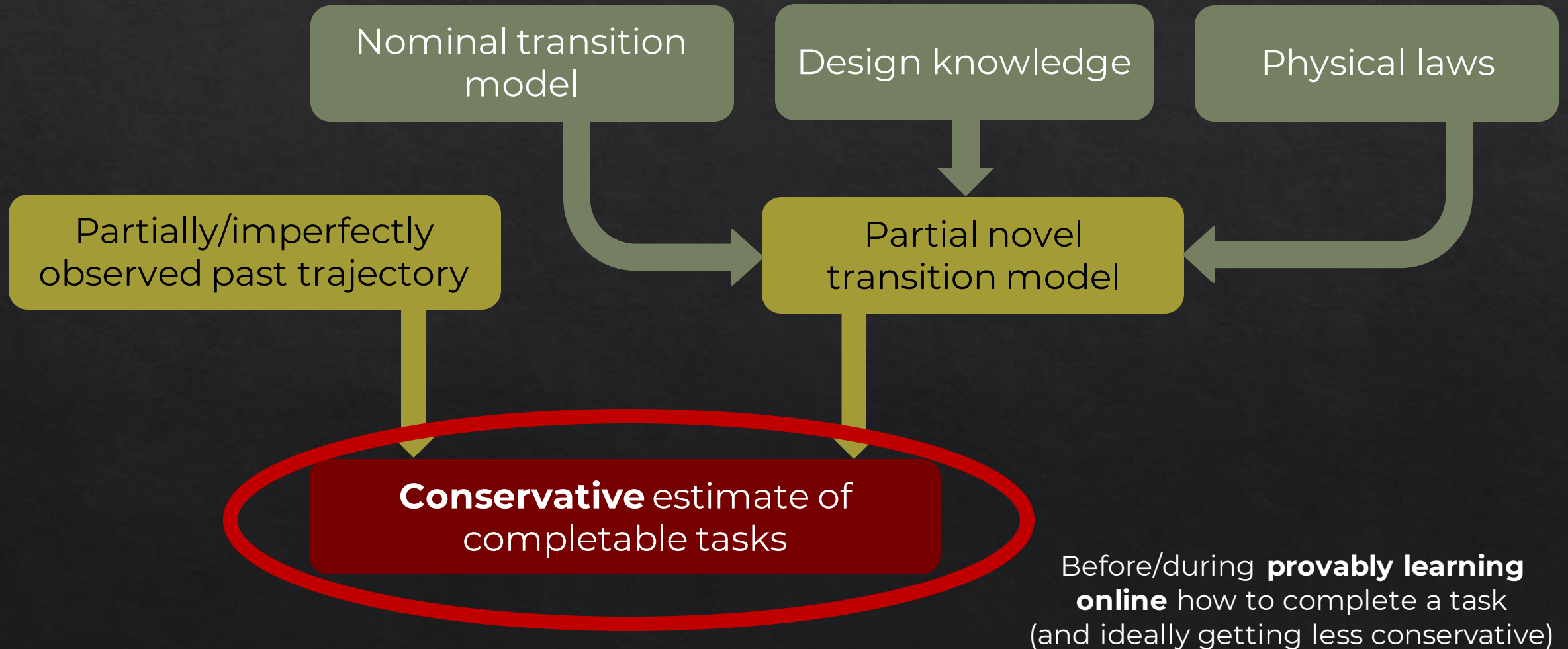
Reach the objective under (some) lack of knowledge of dynamics

Objective might not be reachable

First decide on a reachable target, *then* plan how to get there

Autopilot-in-command

Certifiable Capabilities



Disasters

$$\dot{x} = f(x, u), \quad u \in U$$

Change in dynamics
(e.g., physical damage)

Partial loss of control
(e.g., adversarial takeover)

Actuator
degradation

$$\dot{x} = \hat{f}(x, u, v), \quad (u, v) \in \hat{U}$$

For today: Nonlinear, continuous-time, *deterministic* (sort of – see later)

Actuator Degradation

$$\Leftrightarrow \dot{x} =$$

$$u \in \hat{W}$$

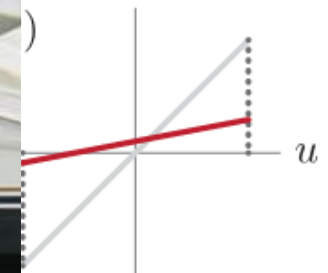
In theory **simple**

Dynamics
not exactly

reachable
set for the

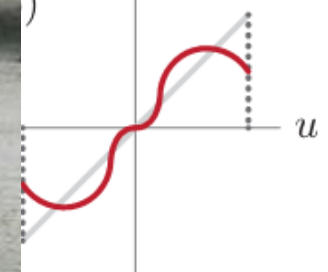


locally affine CDM



its are

reachable



chitz cont. CDM

Real-Time Reachability Analysis

In practice:

- ◇ Reachable set of an, even fully known, nonlinear system **is not computable** in real time (even meaningful **underapproximations** are difficult)
- ◇ *Operational envelope* for the nominal system is computed *beforehand*

We “know” the reachable set for the nominal system, and we know a bound on the degradation

- ◇ Our computation of the degraded reachable set **should not** start from scratch

Maximal Trajectory Difference

If an available control signal $\hat{u} : [0, T] \rightarrow \hat{U}$ differs from a desired nominal signal $u : [0, T] \rightarrow U$ by no more than ε (in some meaningful norm), then the state $\hat{x}(T)$ reached by \hat{u} cannot differ from $x(T)$ reached by u by more than some $h(\varepsilon)$

How to determine $h(\varepsilon)$?

- ◆ Looks like some version of **Grönwall's lemma**, but we can do better

Assumption 1. For all $x \in \mathbb{R}^n$, $u \in \mathcal{U}$, $t_0 \leq t < \infty$ $\|h(t, x, u)\| \leq a(t)w(\|x\|, \|u\|) + b(t)$.

Let $x(t)$ be a solution to the equation $\dot{x} = h(t, x, u)$, $0 \leq t_0 \leq t < \infty$, where $h(t, x, u) : [t_0, \infty) \times \mathbb{R}^n \times \mathcal{U} \rightarrow \mathbb{R}^n$ is continuous for $t_0 \leq t < \infty$, and $\mathcal{U} \subseteq \mathbb{R}^m$ is compact and satisfies $\max_{u \in \mathcal{U}} \|u\| = \delta$. Let Assumption 1 hold. Then,

$$\|x(t)\| \leq G^{-1} \left[G \left(\|x(t_0)\| + \int_{t_0}^t b(\tau) d\tau \right) + \int_{t_0}^t a(\tau) d\tau \right], \quad G(r) := \int_{r_0}^r \frac{ds}{w(s, \delta)}, \quad r > 0, r_0 > 0.$$

Wildness of system dynamics

Bihari Inequality

Intuitively, the result should depend on the “wildness” of the system dynamics

As in similar proofs derived from Grönwall’s lemma, bound the difference in trajectories by itself:

$$r(t) = f(\hat{x}(t), \hat{u}(t)) - f(x(t), u(t)),$$

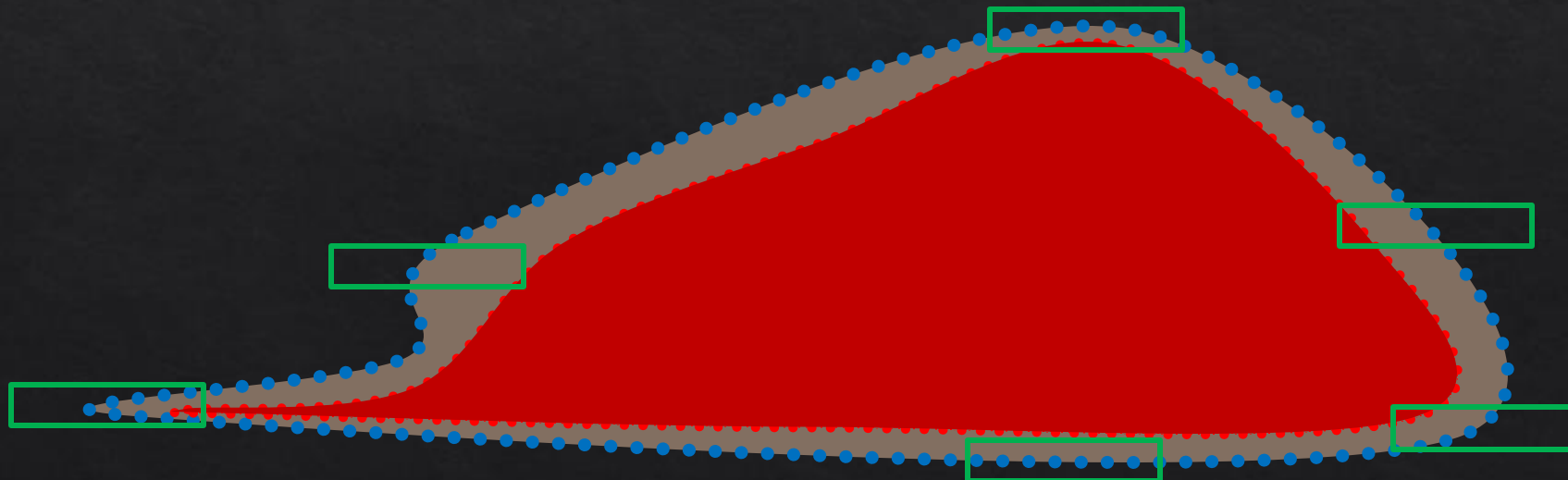
$$\|\hat{x}_i(t) - x_i(t)\| = \left\| \int_0^t r_i(\tau) d\tau \right\| \leq \int_0^t a_i(\tau) \omega_i(\hat{x}(\tau) - x(\tau), \hat{u}(\tau) - u(\tau)) + b_i(\tau) d\tau$$

Can use Lipschitz if needed, but can do better if more information is available

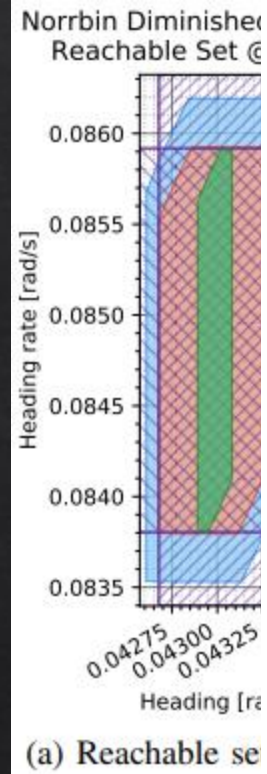
(Hyper)rectangular Tightening

Idea: if the Hausdorff distance of the nominal set of control inputs and the degraded set is bounded by M , then for each nominal control signal we can find an allowed “degraded” signal within an M -ball

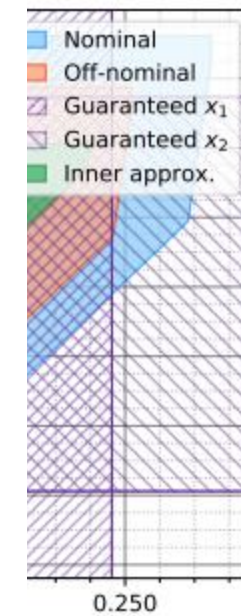
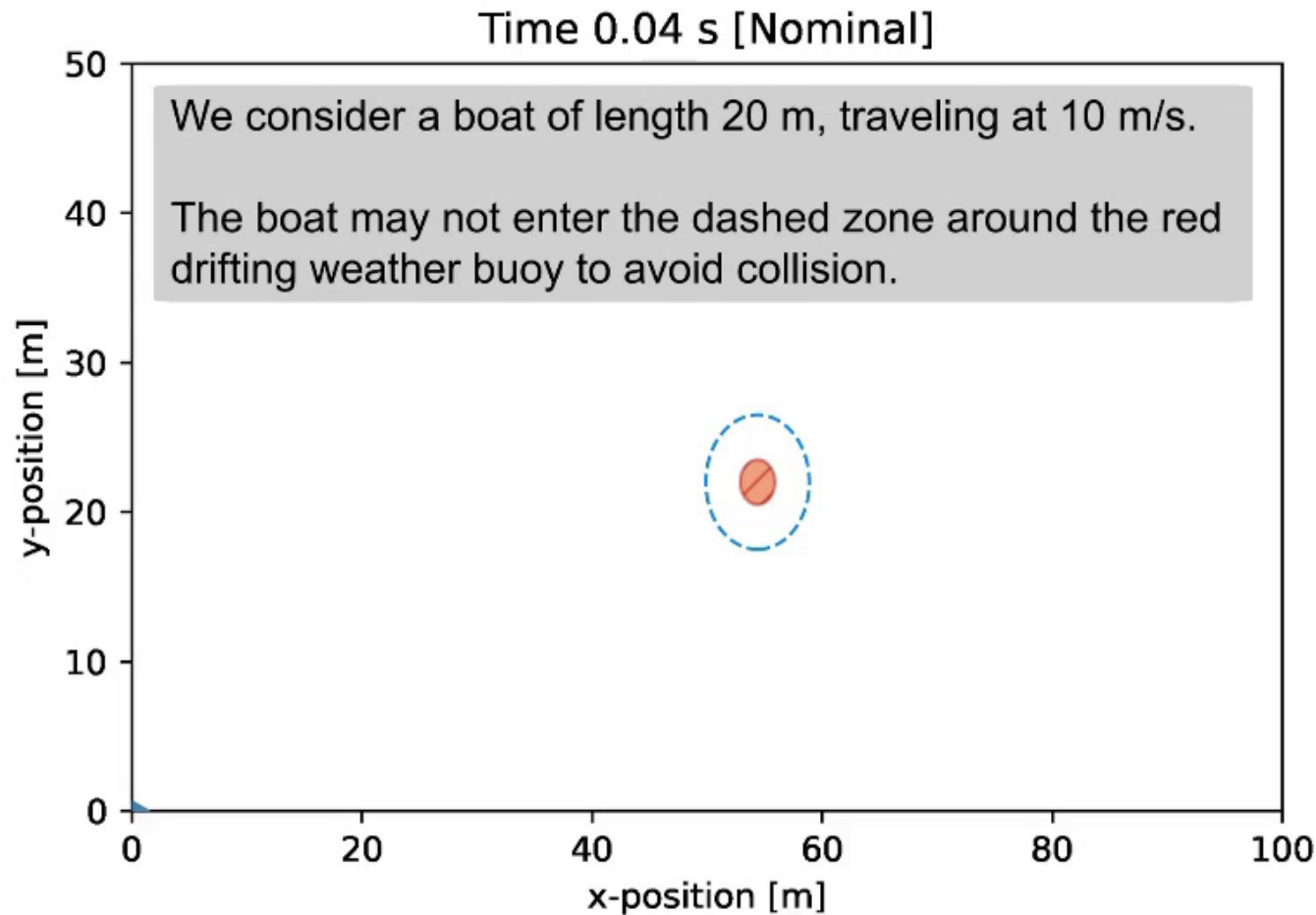
Any state in the boundary of the degraded reachable set maps to a “nearby” state on the boundary of the nominal reachable set



Norrbin's Ship Steering



(Better t



nates)

Disasters

$$\dot{x} = f(x, u), \quad u \in U$$

Change in dynamics
(e.g., physical damage)

Partial loss of control
(e.g., adversarial takeover)

Actuator
degradation

$$\dot{x} = \hat{f}(x, u, v), \quad (u, v) \in \hat{U}$$

Loss of Control Authority

$$\dot{x} =$$

$$\in U$$

Two-player

P1 (“cont

P2 (“envi

Players p

Can P1 win?



(time limit?)

Relative Actuator Strength

Intuitively, P1 wins for reaching **every** state if it can:

- ◇ “cancel out” the adversarial inputs and
- ◇ successfully fight (**or exploit!**) the drift at the same time

In linear systems, if the drift pushes the state “towards infinity” or “towards zero”, fighting the drift is impossible with bounded inputs: drift eventually overpowers the inputs – the only allowed drift is from $Re(\lambda(A)) = 0$

In bounded state spaces, bounded sets of states to reach, etc., fighting the drift may be possible. In all cases, the available input needs to be “stronger” than the adversary

Resilience to Loss of Authority

If the original system dynamics are $\dot{x} = Ax + \bar{B}\bar{u}$, loss of authority over several actuators means the system shifts to:

$$\dot{x} = Ax + Bu + Cv, \quad \bar{B} = [B \quad C], \quad u \in U, \quad v \in V$$

Available input stronger than the adversary: *almost the same as*

$$0 \in \text{Int} \left(\bigcap_{v \in V} BU + Cv \right)$$

If system is resilient to loss of authority, **how resilient is the system?**

Quantitative Resilience

If the target state is in the **guaranteed reachable set** of the initial state (*guaranteed* with respect to all possible adversarial control inputs):
it also matters how long it would take to reach the state, compared to the time it would take the system with nominal dynamics.

Idea: compare the quotient of the **nominal time-optimal reach time** and the **worst-case time-optimal reach time with adversarial input**

$$r_q = \frac{\inf_{\bar{u}} T_R^{\bar{u}}(x_0, x_{goal})}{\sup_v \inf_u T_R^{(u,v)}(x_0, x_{goal})}$$

The adversary chooses an input such that whatever the controller chooses, reach time will be large
Resilience quotient = 0: no resilience; **resilience quotient = 1:** perfect resilience

Driftless Systems

Computing resilience quotient: three (partly nested) optimal control problems

For **driftless systems and integrators**:

- ◇ Optimal control = optimization
- ◇ Optimal control inputs can be geometrically determined
- ◇ There is a state-invariant notion of “direction in which the adversary is the strongest”
- ◇ It is possible to determine the least resilient goal state (*“Minimax Maximax Quotient Theorem”*)

Lyapunov Bounds

Even for general linear systems,

- Optimal control signal is bang-bang, but impossible to analytically determine
- The set of velocities which the system can take depends on the state

Idea:

- A lower bound on the nominal optimal reach time: a lower bound on $\dot{V}(x(t))$
- An upper bound on the nominal optimal reach time: an upper bound on $\dot{V}(x(t))$
for a “good candidate for optimal control”
- A lower bound on the worst-case optimal reach time: a lower bound on $\dot{V}(x(t))$
for any response to the adversarial input
- An upper bound on the worst-case optimal reach time: an upper bound on $\dot{V}(x(t))$
for a “good candidate for optimal control” given the adversarial input

$$T_N^*(x_0) \geq 2 \frac{\lambda_{\min}^P}{\lambda_{\max}^Q} \ln \left(1 + \frac{\lambda_{\max}^Q \|x_0\|_P}{2\lambda_{\min}^P b_{\max}^P} \right)$$

$$T_N^*(x_0) \leq 2 \frac{\lambda_{\max}^P}{\lambda_{\min}^Q} \ln \left(1 + \frac{\lambda_{\min}^Q \|x_0\|_P}{2\lambda_{\max}^P b_{\min}^P} \right)$$

$$T_M^*(x_0) \geq 2 \frac{\lambda_{\min}^P}{\lambda_{\max}^Q} \ln \left(1 + \frac{\lambda_{\max}^Q \|x_0\|_P}{2\lambda_{\min}^P z_{\max}^P} \right)$$

$$T_M^*(x_0) \leq 2 \frac{\lambda_{\max}^P}{\lambda_{\min}^Q} \ln \left(1 + \frac{\lambda_{\min}^Q \|x_0\|_P}{2\lambda_{\max}^P z_{\min}^P} \right)$$

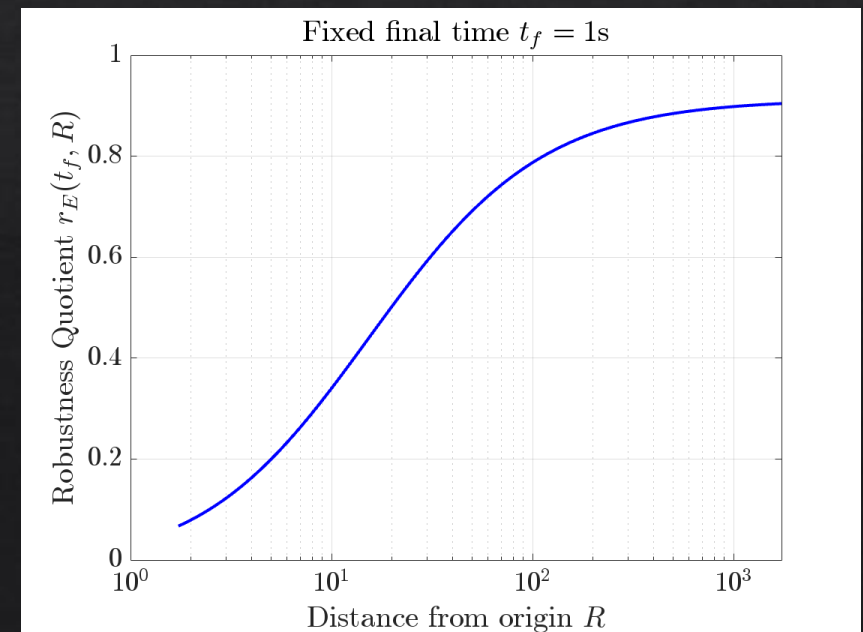
Energy Resilience

Time resilience is **hard**

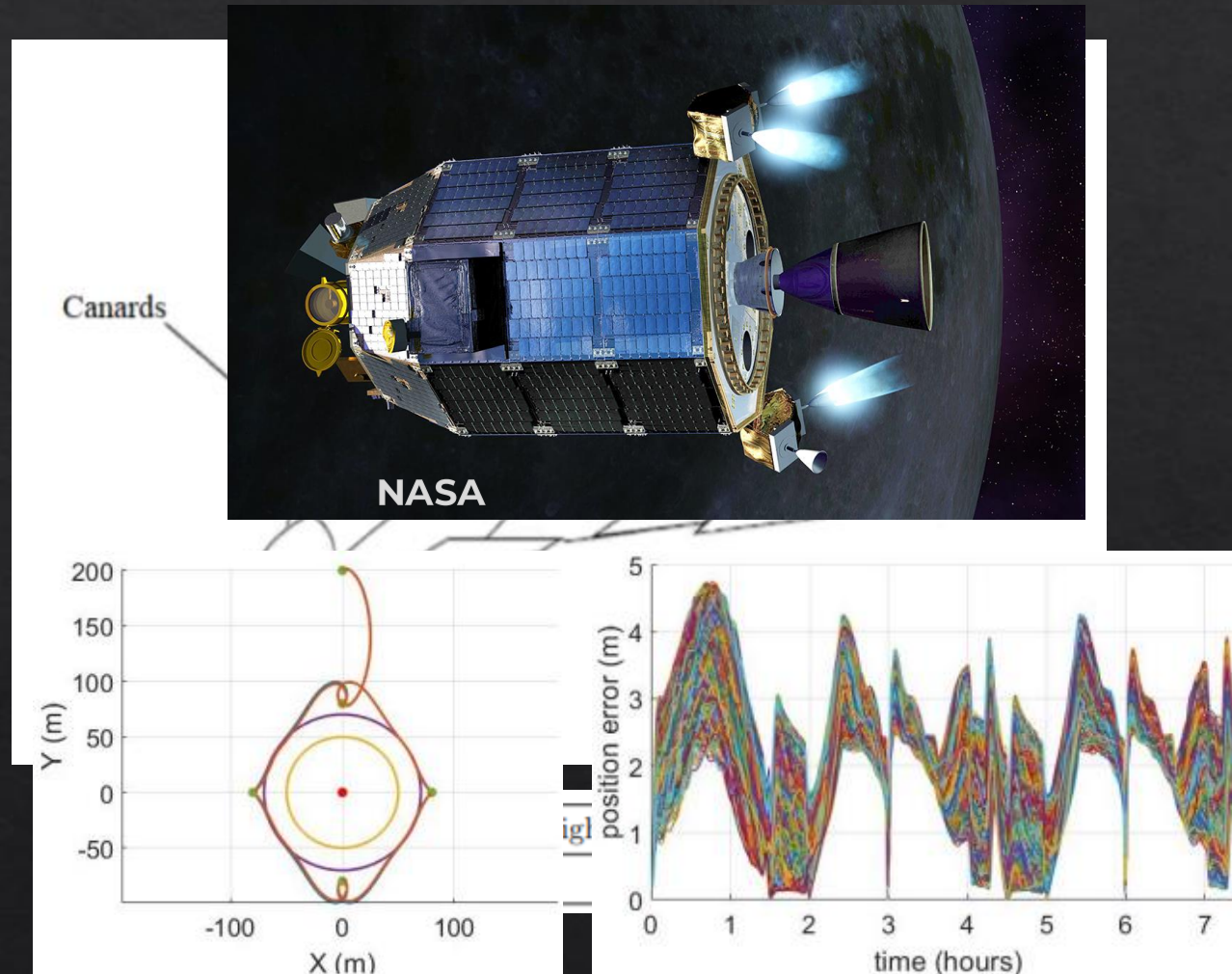
- ◇ Determining minimal time to reach a target for nominal dynamics with input constraints is not simple
- ◇ Determining minimal energy to reach a target (soft input constraints) for nominal dynamics **is simple**: controllability Gramian

for linear systems

- ◇ Pathway towards nonlinear systems: energy resilience – price of nonlinearity?
- ◇ Idea: Describe nonlinear system as a higher-dimensional linear system



Applications



Disasters

$$\dot{x} = f(x, u), \quad u \in U$$

Change in dynamics
(e.g., physical damage)

Partial loss of control
(e.g., adversarial takeover)

Actuator degradation

$$\dot{x} = \hat{f}(x, u, v), \quad (u, v) \in \hat{U}$$

Unknown Dynamics



Robust control

Classical adaptive control
(many) parameters

Both assume the

known (finitely)

Guaranteed Reachability

After a change in dynamics, the system can almost certainly no longer use the same control law to reach its target

- It may not even be able to reach its target using **any** control law

What is it **certifiably capable** of doing (even if we don't yet know how)?

- If **nothing** is known, nothing is certain
- Otherwise,

$$R^{\mathcal{G}}(x_0) = \bigcap_{\tilde{f} \in D_{con}} R^{\tilde{f}}(x_0)$$

All dynamics consistent with the current knowledge

Guaranteed Velocities

- ◇ In theory, guaranteed reachability set is well-defined
- ◇ In practice, how do compute it?

Idea: By finding all trajectories obtained by integrating **guaranteed velocities**

$V^G(t) = \bigcap_{\tilde{f} \in D_{con}} V^{\tilde{f}}(t)$, we will obtain at least some (if not all) guaranteed reachable states

Velocities guaranteed at time t = velocities guaranteed at time 0
 “modulo” maximal system wildness

Initial Velocities + Lipschitz

Set of initially available velocities **can** be computed in arbitrarily short time:
e.g., the learning part of **myopic control**

Bound on the change in set of available velocities: known system parameters + Lipschitz bound on the change in system dynamics (from physical laws and/or design knowledge)

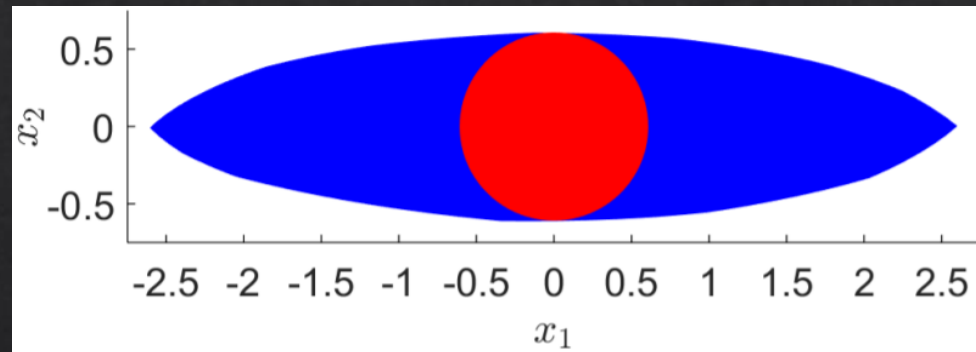
For control-affine systems $\dot{x} = f(x) + G(x)u, \quad u \in U,$

$$V^{\mathcal{G}}(x) = \bigcap_{(\hat{f}, \hat{G}) \in M_x^0(f(0)) \times M_x^1(G(0))} \hat{f} + \hat{G}U$$

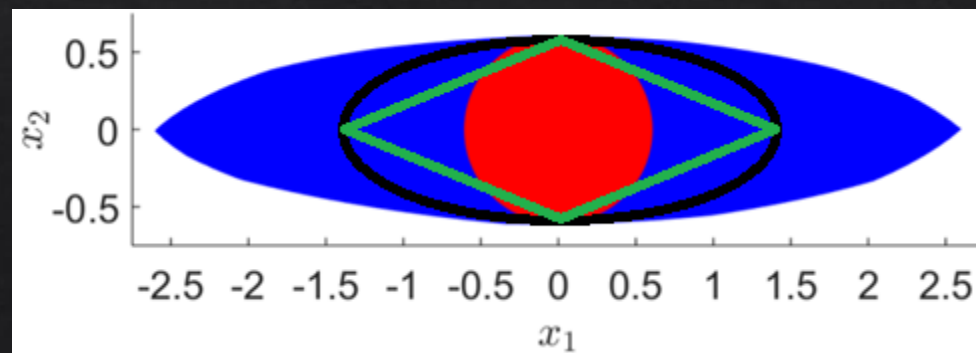
**Infinitely many
(particularly parametrized)
translated/dilated/rotated
copies of the same object**

Simple Guaranteed Velocities

Initial idea: fit a **maximal ball** within this infinite intersection

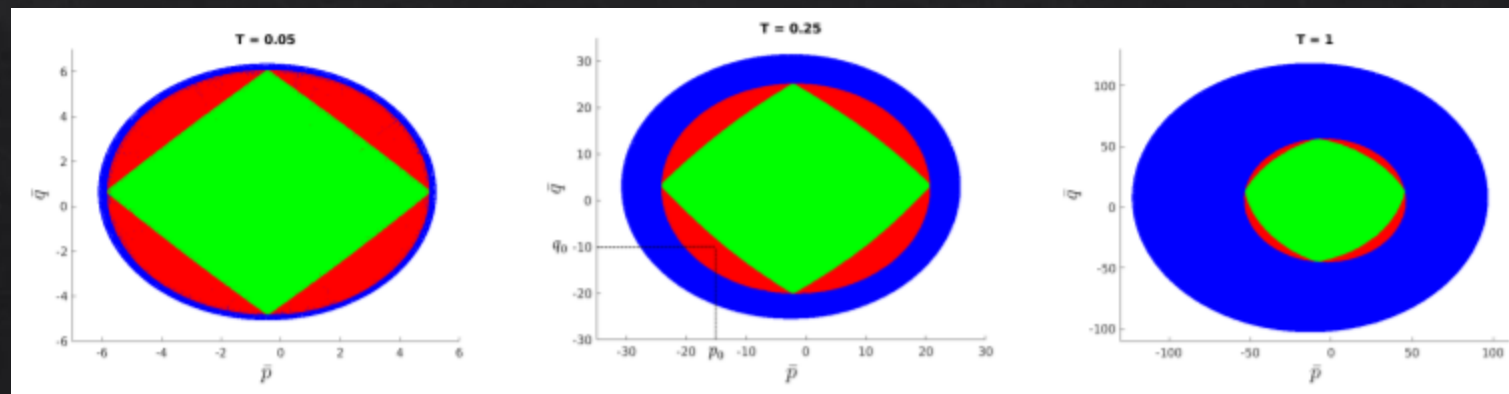
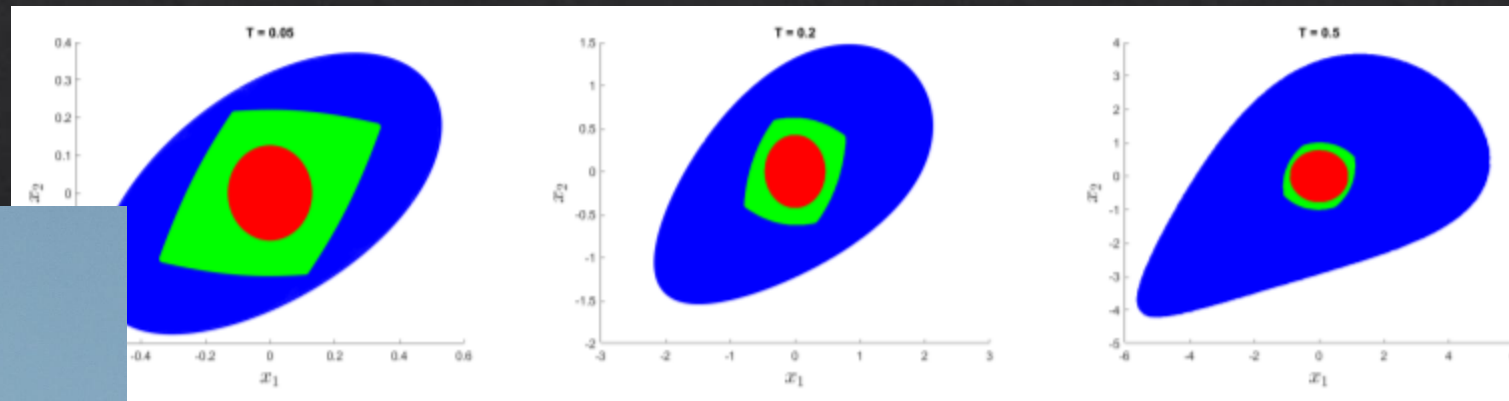


Better idea: fit a more complicated object, but still geometrically simple



Reachable Set Underapproximations

A **guaranteed** reachable set: reachable set of a “simple” control system



Learning-Control Pipeline

Once we have established what the system is capable of doing, we still do not know *how* to do it

◇ **End-to-end planning, learning and control:**

◇ **Task assignment:** what task can be provably completed

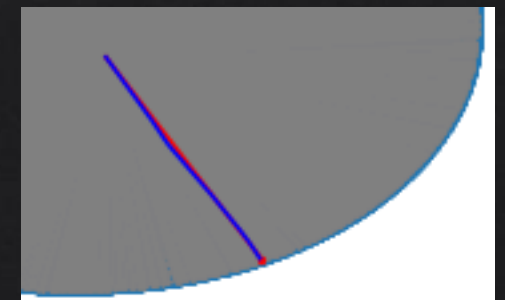
◇ **Real-time learning:** what do we need to know in order to be able to complete it

◇ persistent excitation allows us to learn local dynamics – we “need” a little time, but not much

◇ **Assured control:** complete the task

Lemma 4.5: Let $n \geq 1$ be fixed. Let $\mathbf{u}(t) = u_{n,j}$ for $t \in [\tau_n + j\delta t, \tau_n + (j+1)\delta t]$ and for $j \in \mathcal{I}$. Let \mathbf{x}_n be given and consider $z_n = \theta_n y$ with $|z_n - \mathbf{x}_n| = r$ for some $\theta_n \in (0, 1)$. Suppose that $\dot{d}_{z_n}(\phi_{\mathbf{u}}(t, \mathbf{x}_n)) < 0$ for all $t \in [\tau_n, \tau_{n+1}]$. Then, there exists a $z_{n+1} = \theta_{n+1} y$ such that $\theta_{n+1} > \theta_n$ and $|z_{n+1} - \mathbf{x}_{n+1}| = r$. Particularly, $r - |\mathbf{x}_{n+1} - z_n| \leq |z_{n+1} - z_n| < 2r$.

Corollary 4.12: Let $r = r(k, \tau)$ be the maximal radius of $\hat{\mathcal{R}}(k\tau, \mathbf{x}_n)$ for some k . Suppose that $2(M_0(m+1)^2\delta t + rL_{\max})M_0(m+1)^2\delta t + \mu(\delta t, \epsilon) \leq (r - M_0(m+1)^2\delta t)(b - c|\mathbf{x}_n|)$. Let $\mathbf{u}(\tau_n) = u_{n,0} := \operatorname{argmin}_{\mathbf{u} \in \mathcal{U}_\lambda} \dot{d}_{z_n}(\phi_{\mathbf{u}}(\tau_n, x_0))$ be the controller at τ_n . Then $\dot{d}_{z_n}(\phi_{\mathbf{u}}(t, \mathbf{x}_n)) < 0$ for all $t \in [\tau_n, \tau_n + \tau]$.



Disasters

$$\dot{x} = f(x, u), \quad u \in U$$

Change in dynamics
(e.g., physical damage)

Partial loss of control
(e.g., adversarial takeover)

Actuator
degradation

$$\dot{x} = \hat{f}(x, u, v), \quad (u, v) \in \hat{U}$$

Vision

Complexity



For what systems can we estimate capabilities? What capabilities can we estimate? What tasks can we complete?

Validation



But really...
what can we do *onboard*
during a disaster?

Short Term: Complexity

Much remains to be done: we only just started

- ◇ Technical assumptions (linear systems in loss of control authority, full-ish actuation for unknown dynamics, ...)
- ◇ Combination of scenarios: e.g., partially unknown dynamics with partial loss of control
- ◇ **We can be smarter**
 - ◇ Structural knowledge of “undamaged” part of the unknown dynamics
 - ◇ Not just reachability, but safe reachability (**reach-avoid**) and output reachability
- ◇ **Time-delayed systems, system operating under state constraints**

Teaser: Time Delays

Sensing of state and control inputs takes time:

$$\dot{x}(t) = f(x(t), u(t, x(t - \tau), v(t - \tau)), v(t))$$

No perfect response

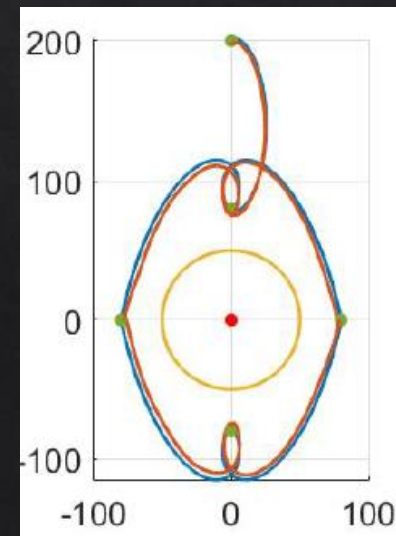
- ◇ Reachability is unlikely, but *approximate reachability* may be attainable
- ◇ How much damage can an adversary do by exploiting delay?
- ◇ How large is the delay?
- ◇ How wild are the system dynamics?

Adversary's strength

Delay size

$$\rho = \frac{\max\{\|Cv\| \mid v \in V\}}{\lambda_{\max}(A + A^T)/2} \left(e^{T_c \lambda_{\max}(A + A^T)/2} - 1 \right), \quad T_c \geq \tau$$

Wildness in dynamics



Teaser: Guaranteed Reachability on a **Manifold**

$$\dot{x} = f^?(x, u), \quad u \in U$$

$\dot{x} \in T_x \mathcal{M}$ $x \in \mathcal{M}$

Another kind of
side knowledge:
no “full actuation”

Finding the reachable set of a *known* system on a manifold is a challenge

- ◇ Analytical methods are hard enough on a Euclidean space
- ◇ How to approximate solving ODEs on a manifold?
- ◇ No good tools

Idea for guaranteed reachability: exploit “maximal wildness” of dynamics

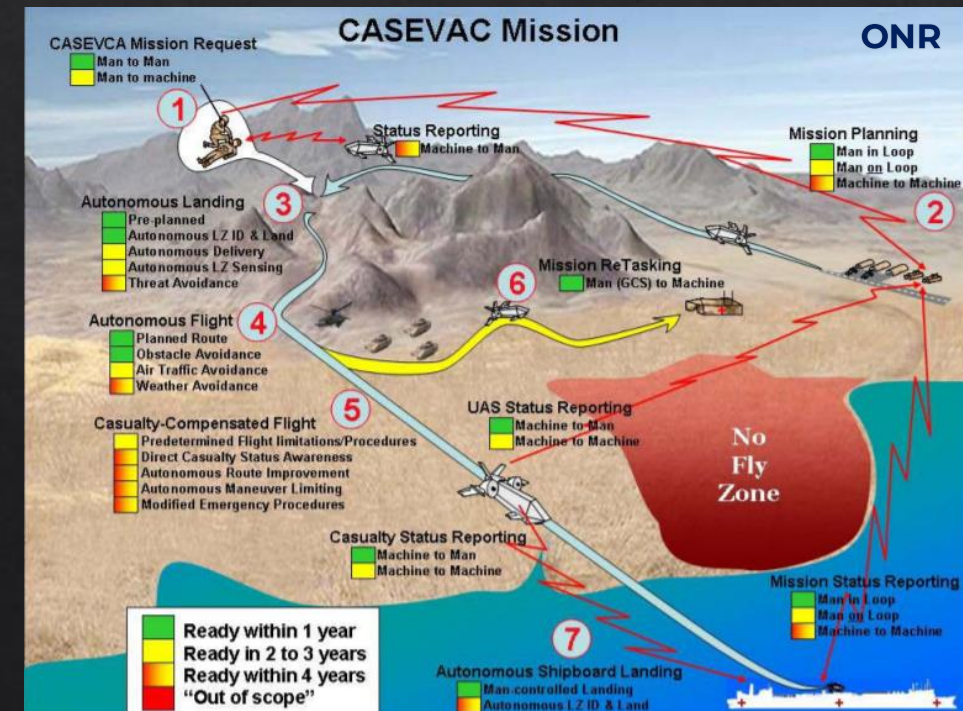
- ◇ On a Euclidean space: “maximal wildness” = Lipschitz constant
- ◇ Lipschitz on surfaces ← Riemannian metric

Validation and More

- ◇ **Using sensors and perception** to recognize fault type and gather information
 - ◇ Fault detection; sensor fusion; state estimation

- ◇ **Complex missions** in high-fidelity simulation
 - ◇ Concurrent sensing and actuation faults
 - ◇ Noise+hostile action

Ongoing joint work and demonstrations
with TC3 team (e.g., GUAM)



Acknowledgments

Extensive support from NASA (incl. current ULI), US Department of Defense – Air Force Office of Scientific Research, and Discovery Partners Institute

Real research (among others):

- ◆ **Actuator degradation:** Hamza El-Kebir, Ani Pirozmanishvili
collaboration with Joseph Bentsman, Changrak Choi et al. (JPL)
- ◆ **Partial loss of control:** Jean-Baptiste Bouvier, Kathleen Xu (UIUC/MIT),
Ram Padmanabhan
collaboration with Sai Pushpak Nandanoori et al. (PNNL),
Robyn Woollands, Himmat Panag
- ◆ **Unknown dynamics:** Taha Shafa, Yiming Meng

mornik@illinois.edu