# Security Analysis of Safe & Seldonian Reinforcement Learning Algorithms

A. Pinar Ozisik[1] and Philip S. Thomas[1]
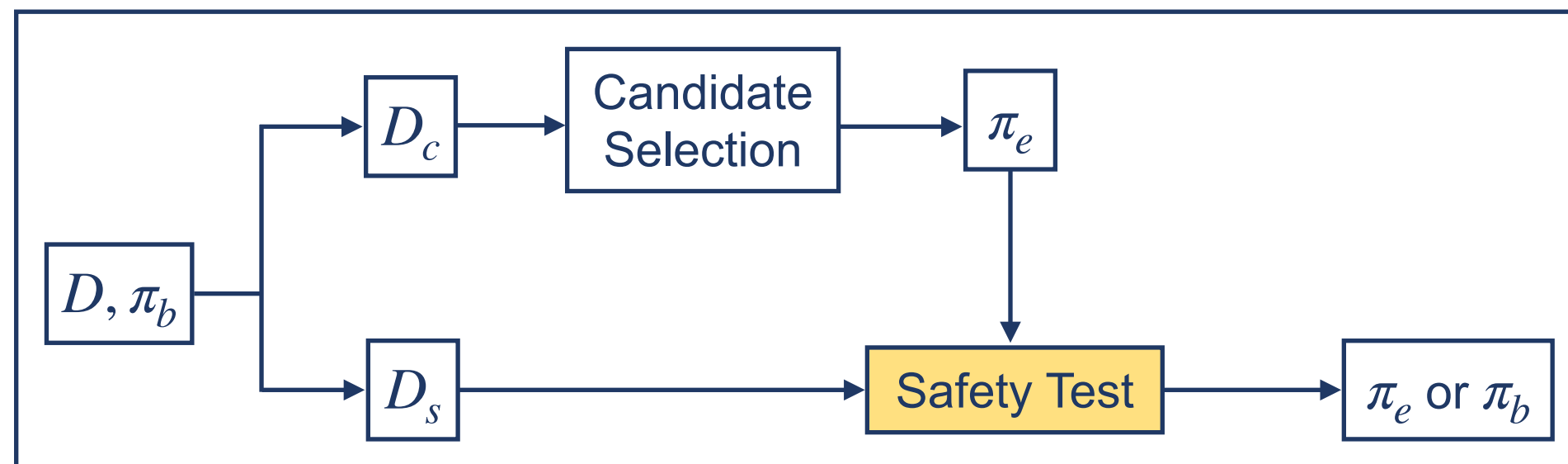University of Massachusetts[1]

## Problem Statement

- RL is proposed for high-risk applications, such as improving type 1 diabetes and sepsis treatments
- Safe and/or Seldonian RL provides high-confidence guarantees that the application will not cause undesirable behavior
- **Limitation:** Assumes that training data is free from anomalies such as errors, missing entries, and malicious attacks
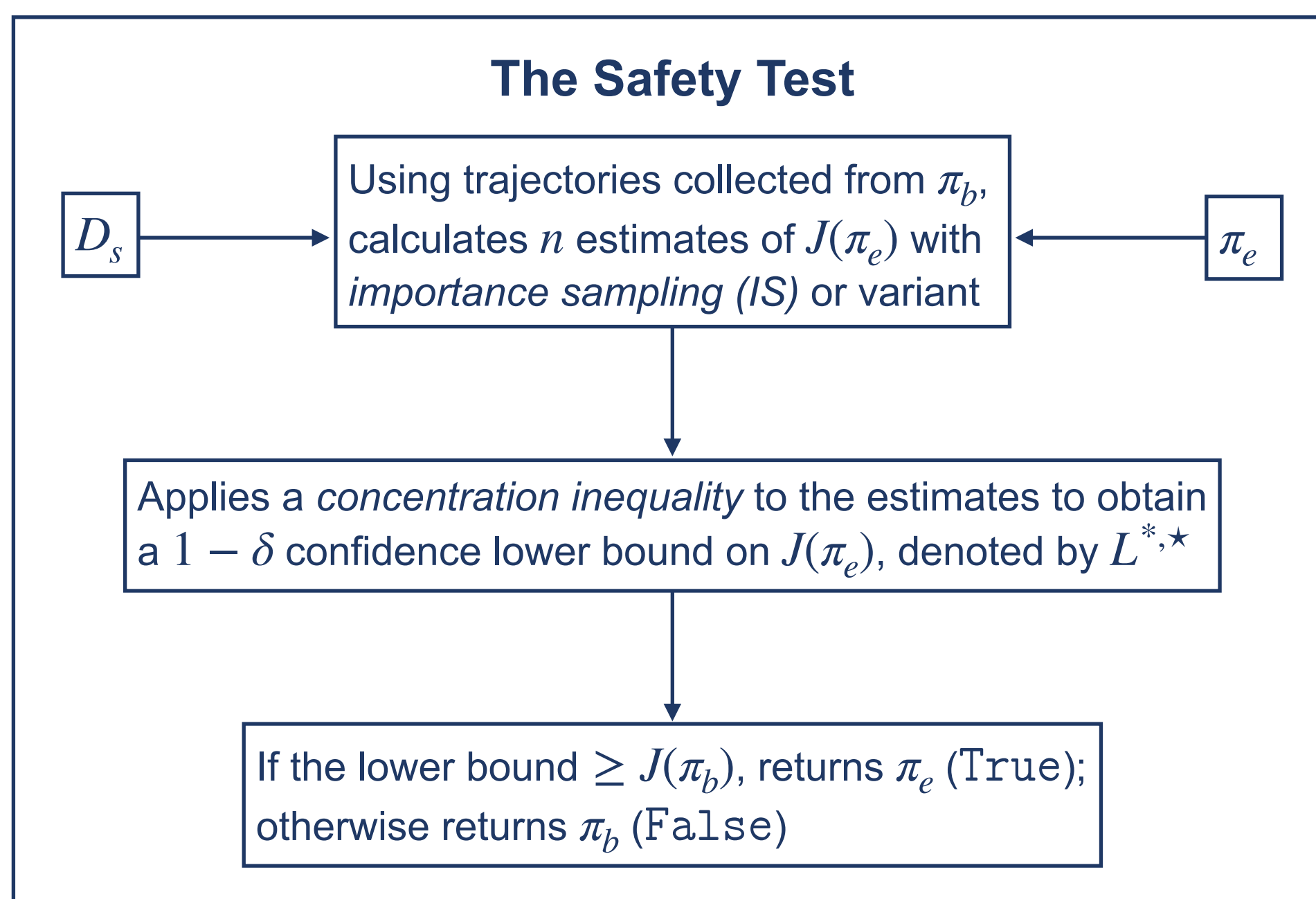- <u>How robust are these algorithms to perturbations in data?</u>

## Background

- **Goal:** Find a policy with larger performance than deployed policy $\pi_b$
- **Assumption:** Can collect data from policy $\pi_b$, but not from others
- Algorithm guarantees **safety:**

$$\Pr(J(a(D)) \geq J(\pi_b)) \geq 1 - \delta$$



### The Safety Test



$D_s$ → Using trajectories collected from $\pi_b$, calculates $n$ estimates of $J(\pi_e)$ with *importance sampling (IS)* or variant ← $\pi_e$

Applies a *concentration inequality* to the estimates to obtain a $1 - \delta$ confidence lower bound on $J(\pi_e)$, denoted by $L^{*,\star}$

If the lower bound $\geq J(\pi_b)$, returns $\pi_e$ (True); otherwise returns $\pi_b$ (False)

## Problem Formulation



- Robustness to perturbations is analyzed by robustness to adversarial attacks
- Algorithms robust to adversarial attacks will also be robust to non-adversarial anomalies in data
- **Attacker model:** Attacker adds fabricated trajectories to dataset

## Panacea: A New Algorithm



- Named after the Greek goddess of universal health
- Provides $\alpha$-security, with a user-specified $\alpha$, if the number of corrupt trajectories in $D$ is upper bounded
  - Takes as input number of corrupt trajectories
- Caps the importance weights using some clipping weight, $c$

## $\alpha$-security: A New Measure

- A measure that ensures safety even if data is corrupted by an adversary
- Larger $\alpha$ implies greater susceptibility to adversarial attacks
- **Assumption 1 (Inferior $\pi_e$).**

$$J(\pi_e) < J(\pi_b)$$

- **Assumption 2 (Absolute continuity).**

$$\big(\pi_b(s,a) = 0\big) \implies \big(\pi_e(s,a) = 0\big)$$

- **Assumption 3 ($\varphi$ safety).** Given Assumption 1,

$$\Pr(\varphi(\pi_e, D, J(\pi_b)) = \text{True}) < \delta$$

- Under Assumptions 1, 2, and 3, a safety test, $\varphi$, is secure with constant $\alpha$ for $\pi_e$, $\pi_b$, $k$ and $D$ collected from $\pi_b$, where $|D| = n$, if and only if,
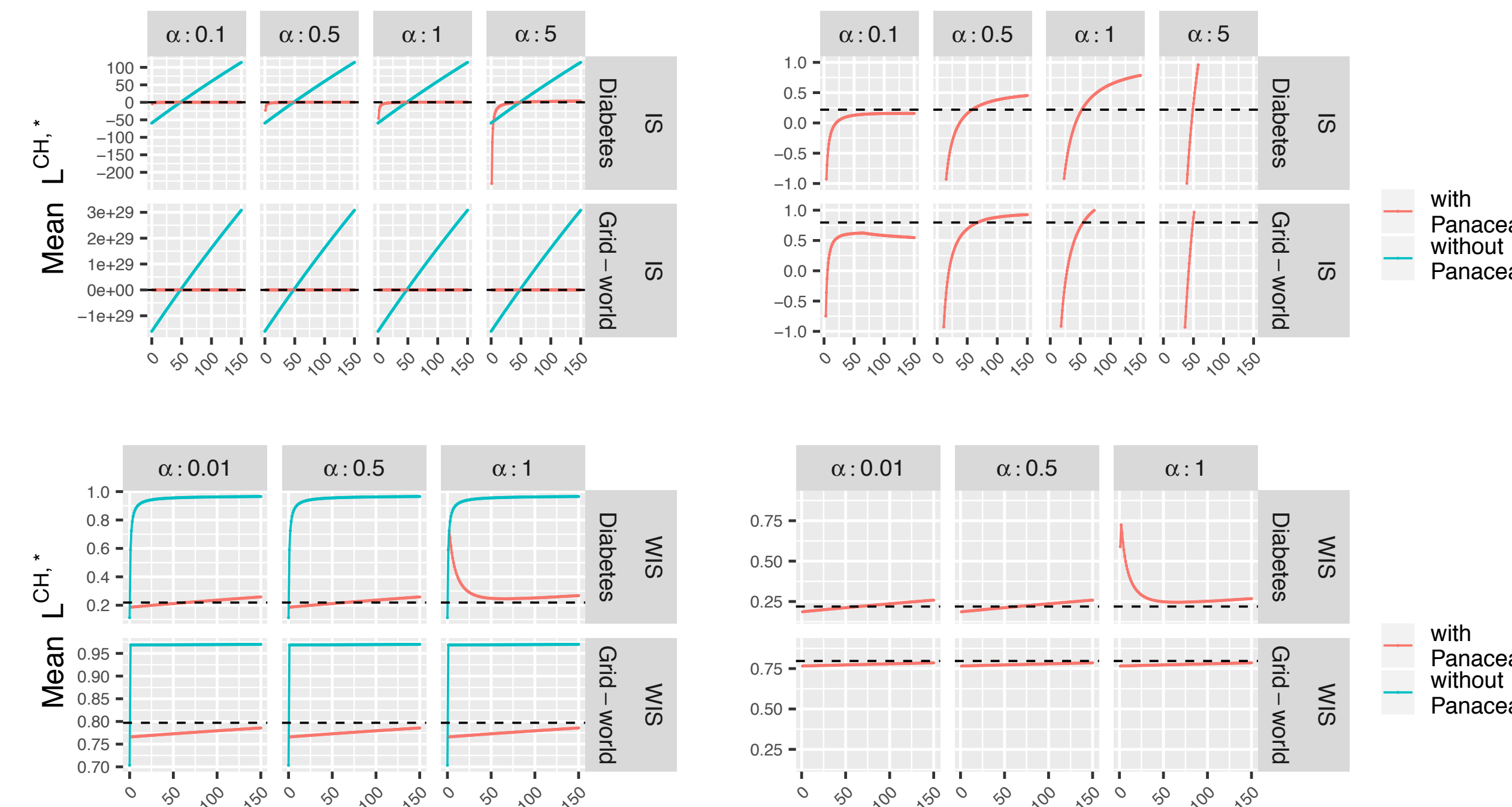
$$\forall m \in \mathcal{M}, \Pr\Big(\varphi\big(\pi_e, m(D,k), J(\pi_b) + \alpha\big) = \text{True}\Big) < \delta,$$

where $m$ is an attack function and $\mathcal{M}$ is the set of all attack functions

- $\alpha$ is computed based on **best attacker strategy**
  - Causes maximum artificial increase in $1 - \delta$ confidence lower bound on $J(\pi_e)$

## Empirical Evaluation on Grid-world & Type 1 Diabetes Treatment



Number of adversarial trajectories added to dataset of size 1,500

- Evaluated two safety tests
  - *Chernoff-Hoeffding* (CH) & IS
  - CH & *weighted importance sampling* (WIS)
- Randomly picked $\pi_b$ and $\pi_e$
- **Point where blue line crosses the black-dotted line represents the number of trajectories needed to break existing safety tests**
  - It takes 49 trajectories for IS and 1 trajectory for WIS for both domains
- **Point where red line crosses the black-dotted line represents the number of trajectories needed to break Panacea**
  - *IS with $\alpha = 0.5$*: 59 trajectories in grid-world and 68 trajectories in diabetes domain
  - *WIS with $\alpha = 0.5$*: does not cross black-dotted line in grid-world and 65 trajectories in diabetes domain
- **Conclusion:** Safety tests can be extremely fragile, but Panacea provides user-specified robustness