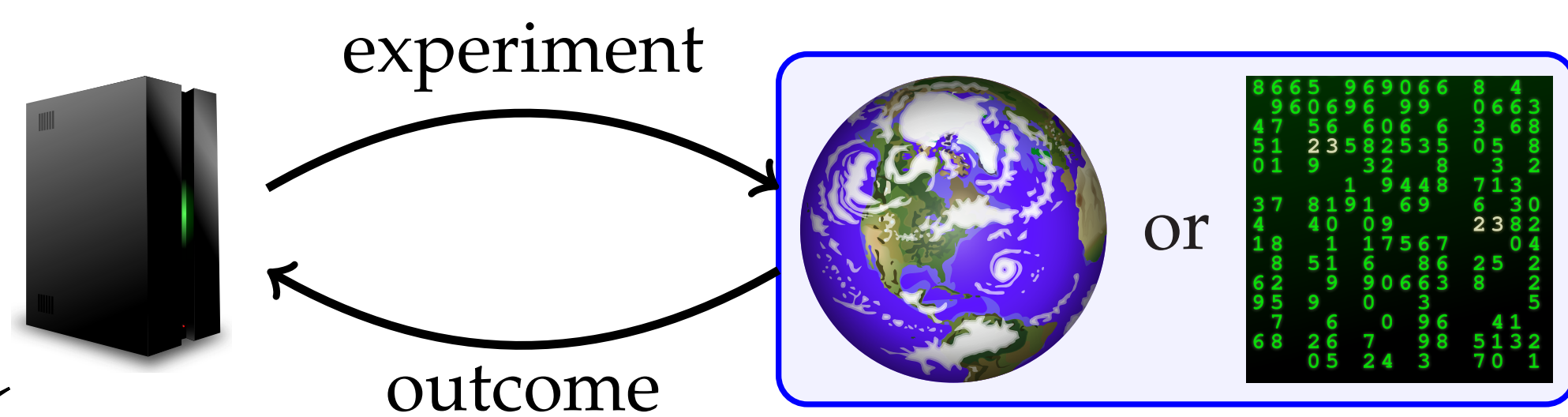




Topic: Pure Exploration

Task: most effective drug dose?
answer query most appealing website layout?
safest next robot action?

Setting: interactive learning



Main scientific questions

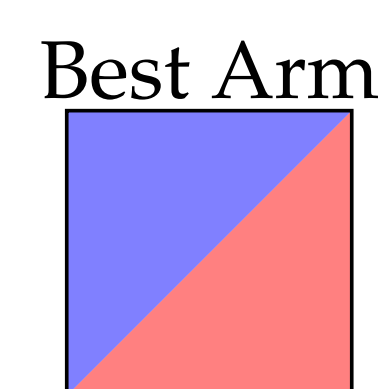
- Efficient systems
- Sample complexity as function of **query** and **environment**

Slogan

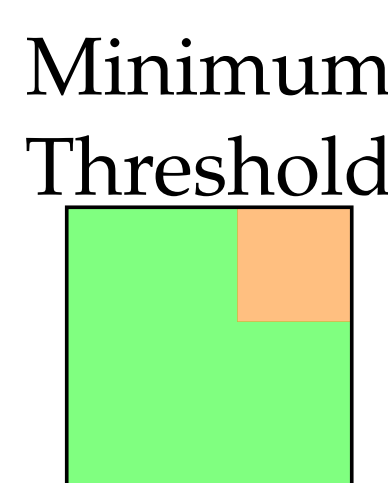
Pure exploration with **multiple correct answers** requires **stabilised methods** to pacify **discontinuity**.

Model

K -armed bandit, parameterised by arm means $\mu = (\mu_1, \dots, \mu_K)$.
Set \mathcal{M} of possible environments.



Set \mathcal{I} of possible answers. **Correct answer function** $i^* : \mathcal{M} \rightarrow \mathcal{I}$.



Strategy:

- Stopping rule** $\tau \in \mathbb{N}$
- In round $t \leq \tau$ **sampling rule** picks arm $A_t \in [K]$ and observes $X_t \sim \mu_{A_t}$.
- Recommendation rule** $\hat{I} \in [K]$.

Definition. A strategy is δ -PAC if $\mathbb{P}_\mu(\hat{I} \neq i^*(\mu)) \leq \delta$ for every bandit model $\mu \in \mathcal{M}$.

Goal: minimise **sample complexity** $\mathbb{E}_\mu[\tau]$ among δ -PAC strategies.

State of the Art: Lower Bound

Answer $i \in \mathcal{I}$ has **altern.** $\neg i := \{\lambda \in \mathcal{M} \mid i^*(\lambda) \neq i\}$

Theorem (Castro 2014, Garivier and Kaufmann 2016). Fix a δ -PAC strategy. Then for every bandit model $\mu \in \mathcal{M}$

$$\mathbb{E}_\mu[\tau] \geq \frac{\ln(1/\delta)}{\max_{w \in \Delta_K} \inf_{\lambda \in \neg i^*(\mu)} \sum_{i=1}^K w_i \text{KL}(\mu_i \parallel \lambda_i)}$$

State of the Art: Algorithm

Good algorithm **must** \rightarrow oracle proportions

$$w^*(\mu) = \operatorname{argmax}_{w \in \Delta_K} \inf_{\lambda \in \neg i^*(\mu)} \sum_{i=1}^K w_i \text{KL}(\mu_i \parallel \lambda_i)$$

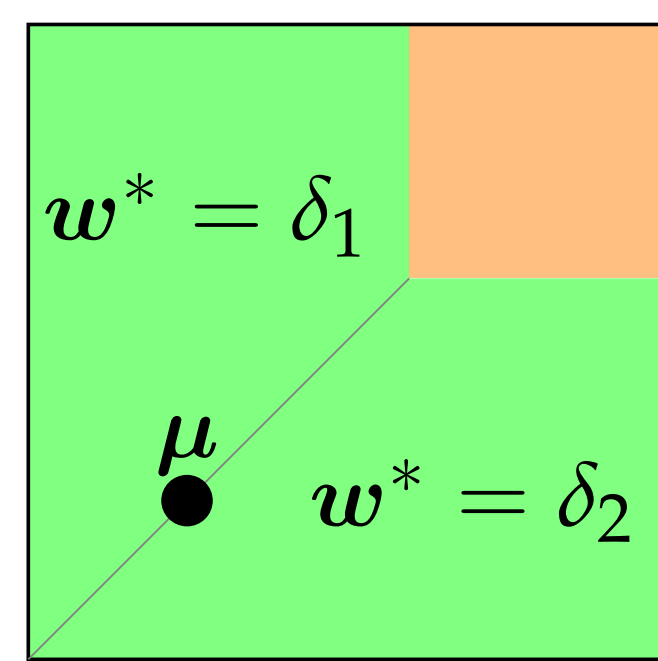
Track-and-Stop [Garivier and Kaufmann, 2016]

Crux: draw $A_t \sim w^*(\hat{\mu}(t))$.

- Ensure $\hat{\mu}(t) \rightarrow \mu$ by **forced exploration**
- this ensures $w^*(\hat{\mu}_t) \rightarrow w^*(\mu)$ **assuming w^* is continuous**
- Draw arm i with $N_i(t)/t$ below w_i^* (**tracking**)

Discontinuity with Single Answer

Can w^* really be discontinuous? At an instance μ where the lower bound does not diverge?



Example problem:

Minimum Threshold

$$w^*(\mu) = \operatorname{conv}(\{\delta_1, \delta_2\}).$$

First Result: Continuity Salvaged

Theorem. The oracle allocation w^* , when viewed as a **set-valued function**, is **upper hemicontinuous**. Moreover, its output is always a **convex set**.

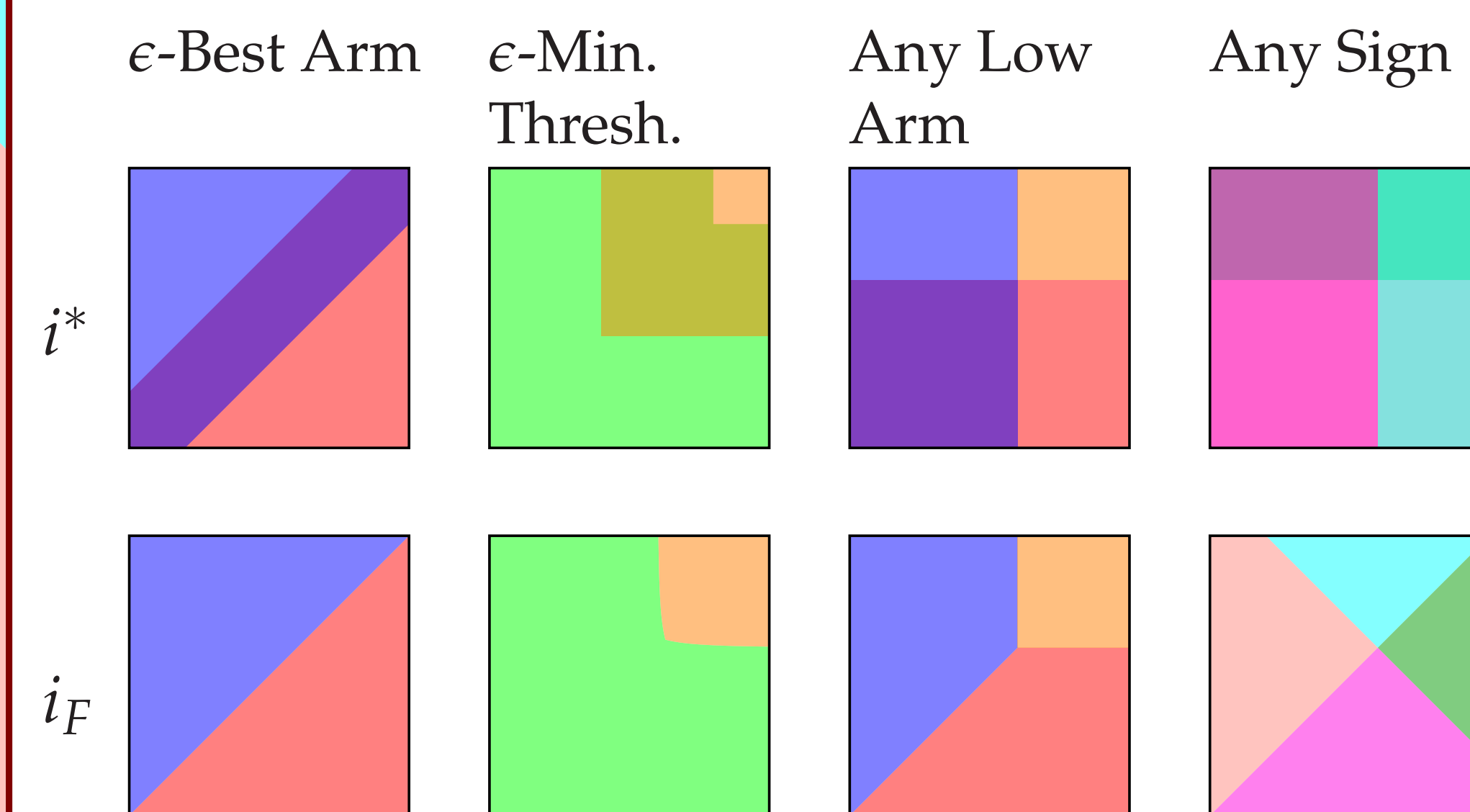
Intuition On bandit model μ , our empirical distribution will be a convex combination of δ_1 and δ_2 .

\Rightarrow **need to rethink Tracking!**

Theorem. Track-and-Stop with **C-tracking** is δ -PAC with **asymptotically optimal** sample complexity. Track-and-Stop with **D-tracking** may **fail** to converge.

Multiple-answer Problems

Set-valued correct answer function $i^* : \mathcal{M} \rightarrow 2^{\mathcal{I}}$.



Rethinking the Lower Bound

Single-answer lower bound is based on **KL contraction**. With **multiple correct answers**, this gives the **wrong leading constant**. \Rightarrow we do direct proof.

Theorem. Any δ -PAC algorithm verifies

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\tau_\delta]}{\ln(1/\delta)} \geq D(\mu)^{-1},$$

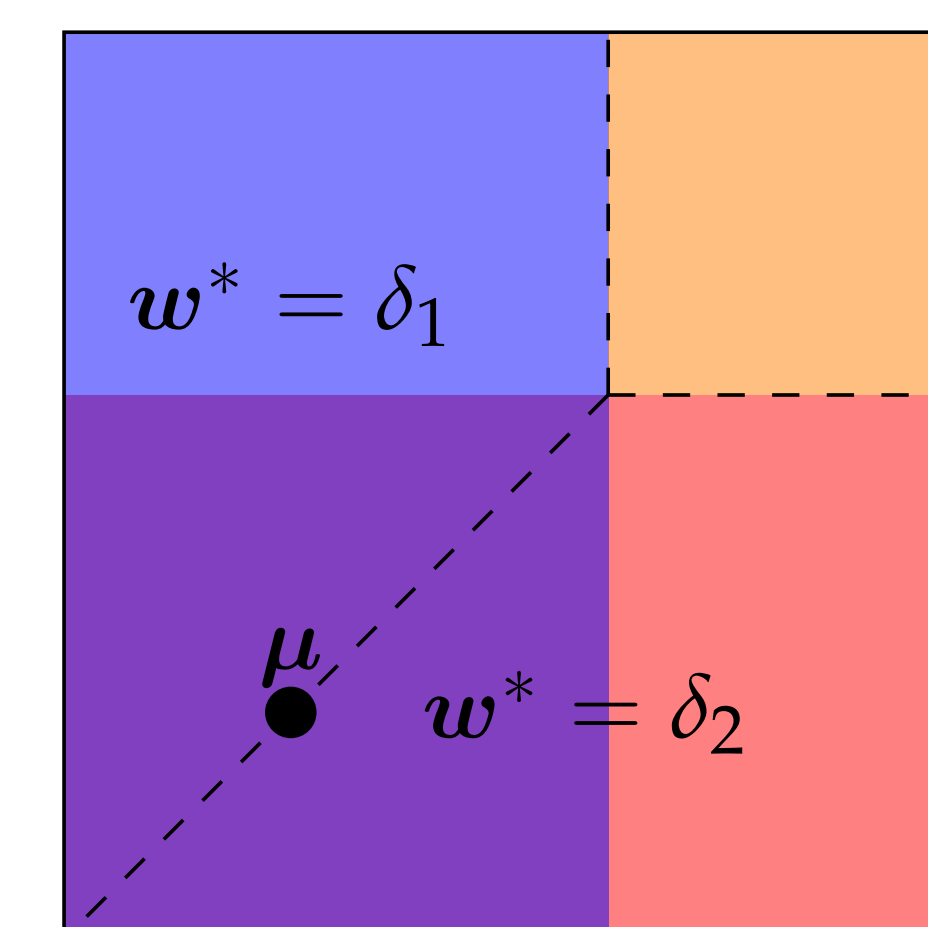
$$\text{where } D(\mu) = \max_{i \in i^*(\mu)} \max_{w \in \Delta_K} \inf_{\lambda \in \neg i} \sum_{k=1}^K w_k d(\mu_k, \lambda_k)$$

for any multiple answer instance μ with sub-Gaussian arm distributions.

On Matching the Lower Bound

Complication: $w^*(\mu)$, the set of maximisers of $D(\mu)$, is **unsalvageably discontinuous**.

Example: **Any Low Arm**

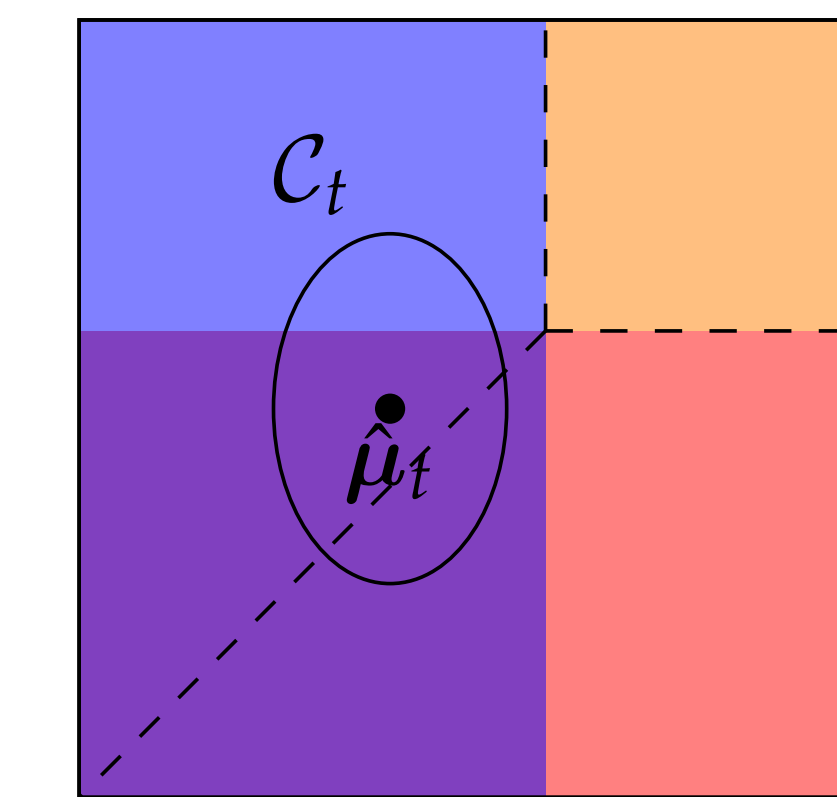


At μ , the oracle weights are $w^*(\mu) = \{\delta_1, \delta_2\}$. Tracking $w^*(\hat{\mu}_t)$ will play from convex hull.

Solution: Make it Sticky

New Sticky-Track-and-Stop Sampling Rule:

Find least (in "sticky order") oracle answer in confidence region \mathcal{C}_t . Track its oracle weights at $\hat{\mu}_t$.



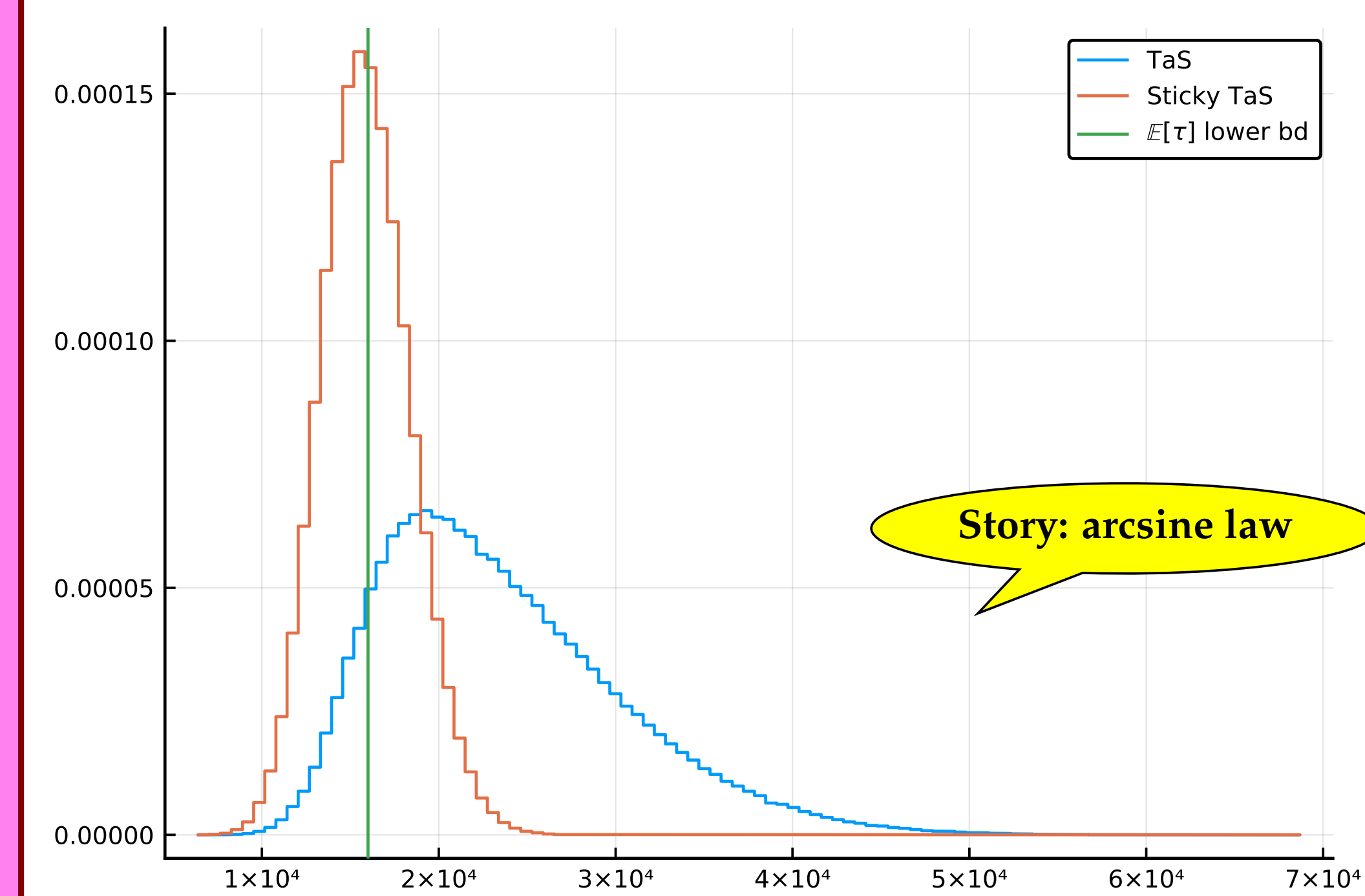
Main Result

When coupled with a good stopping rule,

Theorem. Sticky Track-and-Stop is **asymptotically optimal**, i.e. it verifies for all $\mu \in \mathcal{M}$,

$$\lim_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\tau_\delta]}{\ln(1/\delta)} = D(\mu)^{-1}.$$

Non-Sticky is Actually Dangerous



Histogram of stopping time, $\delta = e^{-80}$.

Where to go from here

- Practical efficiency
- Avoid forced exploration
- Regret (WIP), RL ...
- Moderate confidence $\delta \not\rightarrow 0$ regime; bounds
- Understand sparsity patterns
- Dynamically expanding planning horizon