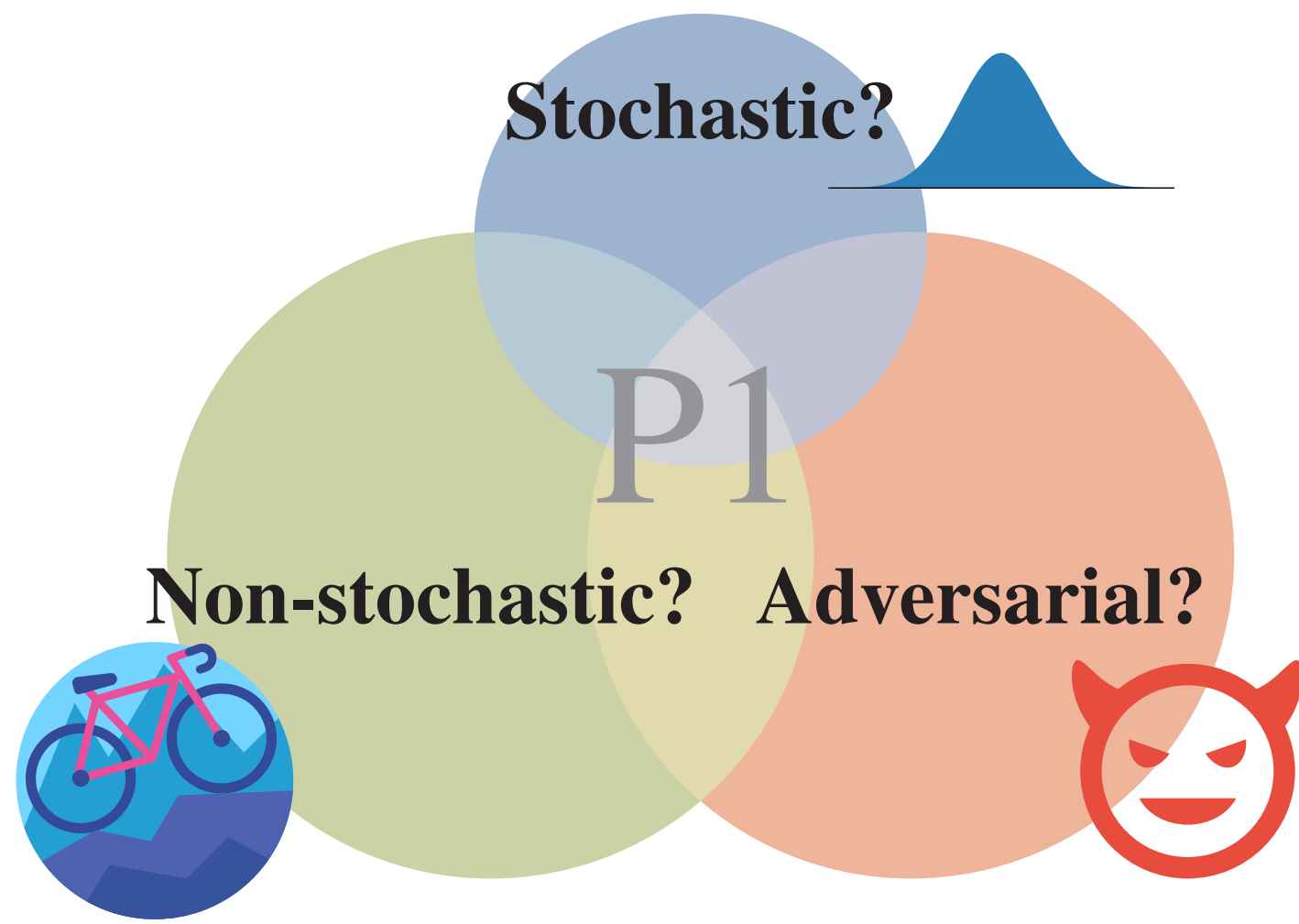


# BEST OF BOTH WORLDS: STOCHASTIC & ADVERSARIAL BEST-ARM IDENTIFICATION



YASIN ABBASI-YADKORI, PETER BARTLETT, VICTOR GABILLON, ALAN MALEK & MICHAL VALKO

## WHAT GIVES?



Find your best option when the data is potentially non-stochastic or adversarial!

## THE GAME: LEARNER VS ADVERSARY

For  $t = 1, 2, \dots, n$ ,

- simultaneously, **Learner** picks arm  $I_t \in [K]$ ,
  - / picks gain  $g_t \in [0, 1]^K$ .
  - Then, the learner observes  $g_{t, I_t}$ .
- Recommend arm  $J_n$  hoping  $J_n = k^*$ .

**Adversarial**

**Stochastic**

arbitrary  $g_{k,t}$

$g_{k,t}$  sampled i.i.d. from  $\nu_k$

$k_g^* = \arg \max_{k \in [K]} G_k$

$k_{STO}^* = \arg \max_{k \in [K]} \mu_k$

$G_k = \sum_{t=1}^n g_{k,t}$

$e_{ADV}(n) \triangleq \mathbb{P}(J_n \neq k_g^*)$

$e_{STO}(n) \triangleq \mathbb{P}(J_n \neq k_{STO}^*)$

**Gaps:**  $n\Delta_k^g \triangleq \begin{cases} G_{(1)} - G_k & \text{if } k \neq k_g^* \\ G_{(1)} - G_{(2)} & \text{if } k = k_g^* \end{cases}$

**Notions of complexity:**

$$H_{SR} \triangleq \max_{k \in [K]} \frac{k}{\Delta_{(k)}^2} \quad \text{and} \quad H_{UNIF} \triangleq \frac{K}{\Delta_{(1)}^2}$$

## OPTIMAL UNIFORM LEARNER

Rule:  $I_t$  uniformly at random.

**Th 1** (Rule vs. ). For all  $n$ , adversarial  $g$ ,

$$e_{ADV}(g)(n) = \mathcal{O}\left(\exp\left(-\frac{n}{H_{UNIF}(g)}\right)\right).$$

**Th 2** ( Lower bound). For any learner, a  $g^1$  of complexity  $H_{UNIF}$ ,

$$e_{g^1}(n) = \Omega\left(\exp\left(-\frac{n}{H_{UNIF}}\right)\right).$$

Rule: optimal gap-dependent rates against .

## ¿BEST OF BOTH WORLDS? (BOB)

Existing robust solutions?

	$e_{STO}(n)$	$e_{ADV}(g)(n)$
SR [1]	$e^{-\frac{n}{H_{SR} \log K}}$	1
Rule	$e^{-\frac{n}{H_{UNIF}}}$	$e^{-\frac{n}{H_{UNIF}}}$

**BOB question:** A learner performing optimally in both the *stochastic* and *adversarial* cases while not being aware of the nature of the rewards?

Why is the BOB question challenging?

- Bias of estimator  $\hat{G}_{k,t} = \frac{t \sum_{t'=1}^t \mathbf{1}\{I_{t'}=k\} g_{k,t'}}{\sum_{t'=1}^t \mathbf{1}\{I_{t'}=k\}}$
- Variance of  $\tilde{G}_{k,t} = \sum_{t'=1}^t \frac{g_{k,t'}}{p_{k,t'}} \mathbf{1}\{I_{t'}=k\}$

Pull uniformly for too long and incur a large variance of order  $K$  in  $\tilde{G}_{k,t}$ .

## ¡IMPOSSIBLE BOB!

New notion of complexity

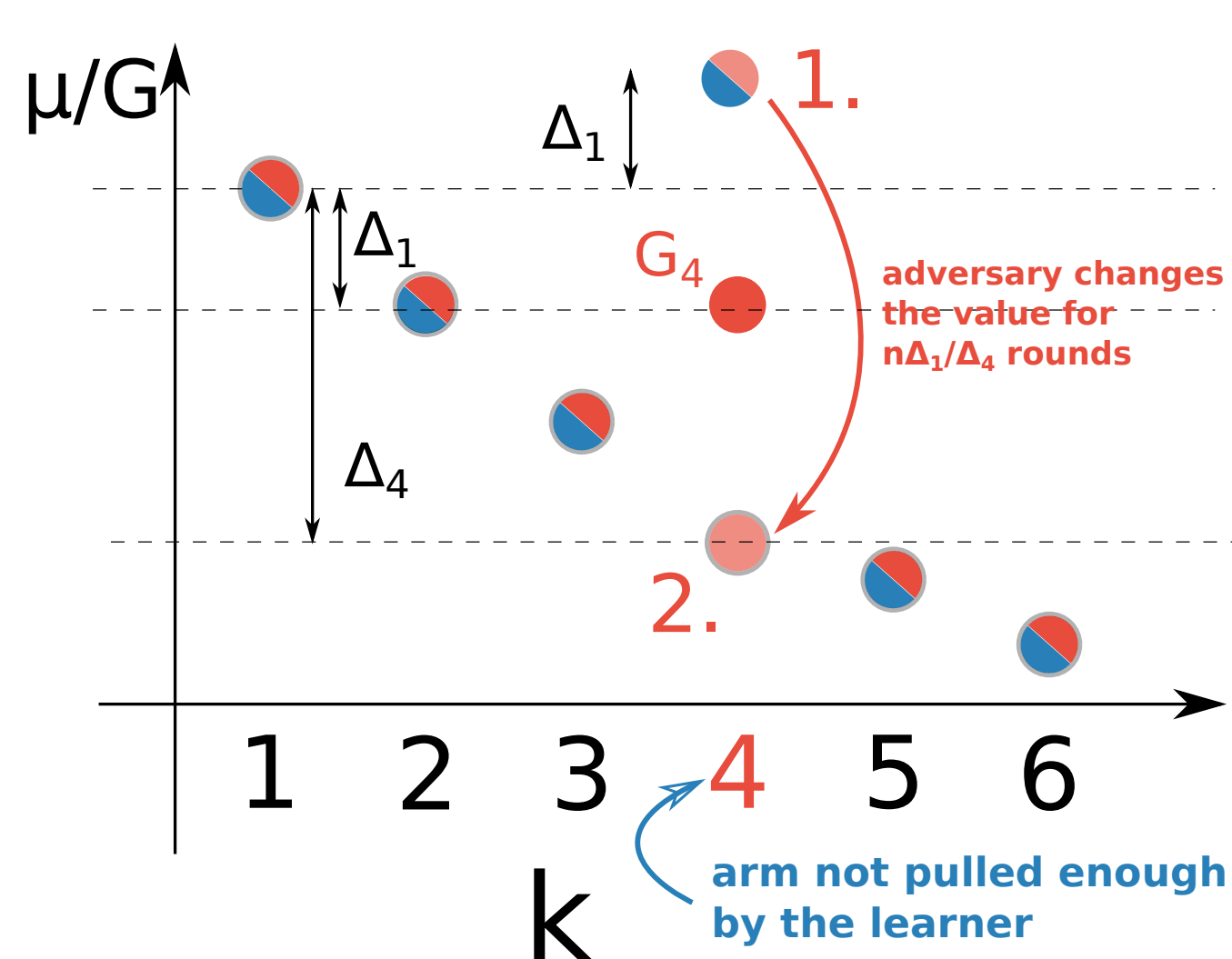
$$H_{BOB} \triangleq \frac{1}{\Delta_{(1)}} \max_{k \in [K]} \frac{k}{\Delta_{(k)}}$$

**Th 3** (Lower bound for the BOB challenge). For any learner, for any  $H_{BOB}$  there exists a stochastic problem with complexity  $H_{BOB}$  such that

$$\text{if } e_{STO}(n) \leq \frac{1}{64} \exp\left(-\frac{2048n}{H_{BOB}}\right),$$

then there exists an adversarial problem where

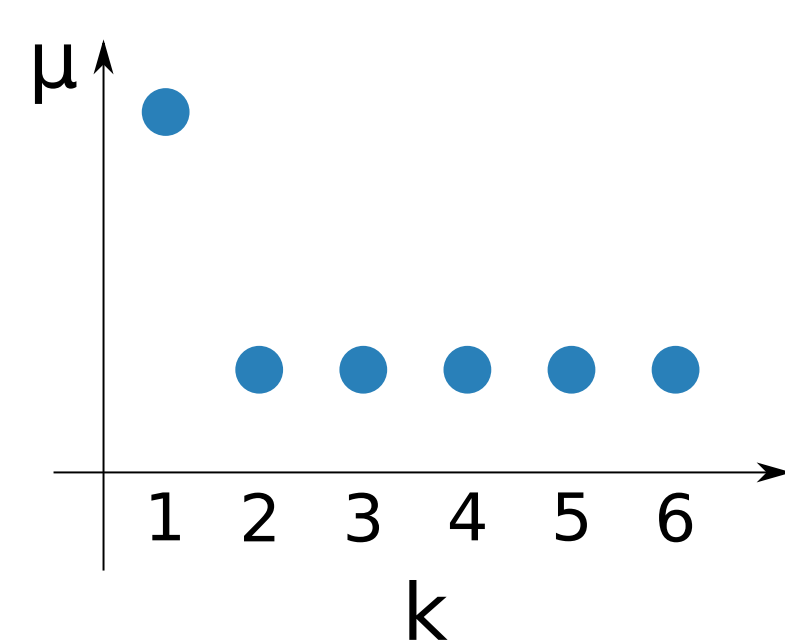
$$e_{ADV}(g)(n) \geq \frac{1}{16}.$$



## DIFFERENT REGIMES

$$H_{SR} \leq H_{BOB} \leq H_{UNIF}.$$

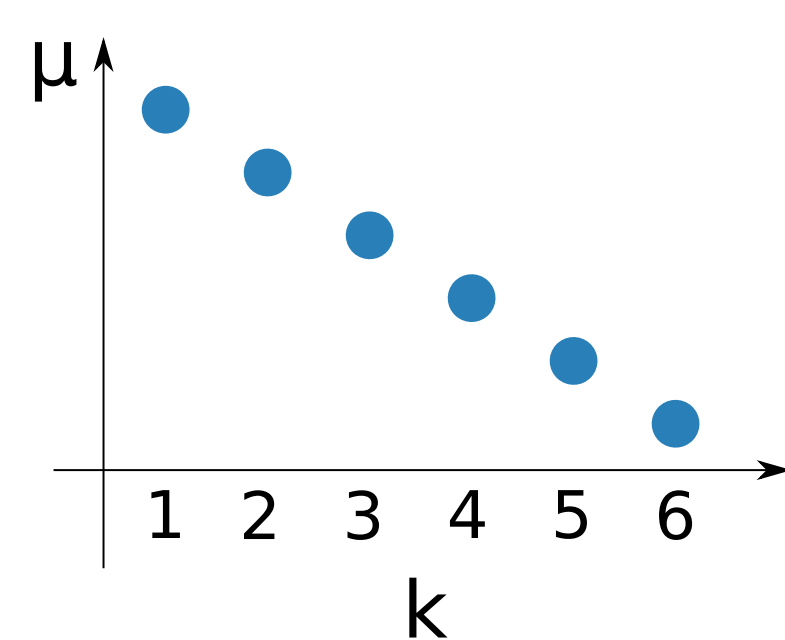
► **Flat regime**



$$H_{SR} = H_{BOB} = H_{UNIF}$$

**BOB is achieved by Rule.**

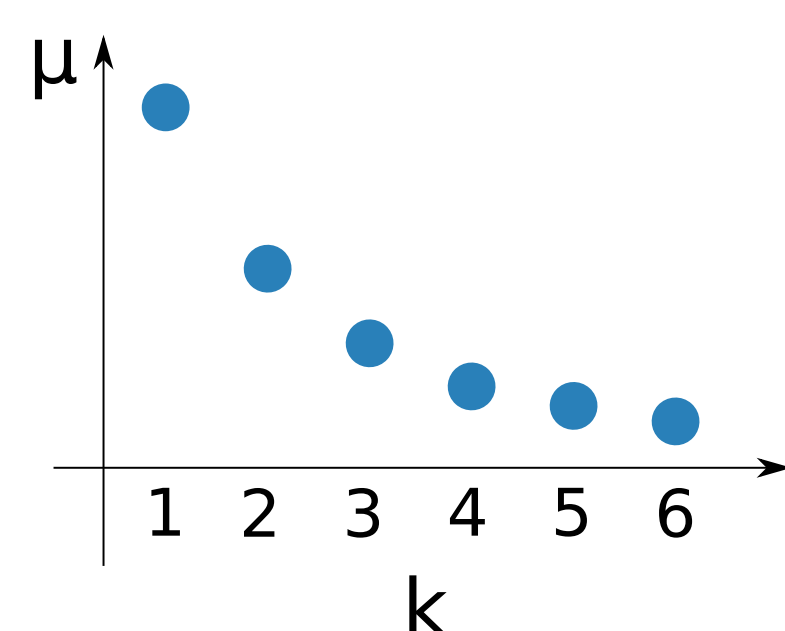
► **Linear regime**



$$H_{SR} = H_{BOB} = \frac{H_{UNIF}}{K}$$

BOB can be achieved but not by Rule. **We need a new learner!**

► **Square-root regime**



$$H_{SR} = \frac{H_{BOB}}{\sqrt{2K}} = \frac{H_{UNIF}}{K}$$

**No learner can do BOB!**

## THE BYCWF (THE BEST YOU CAN WISH FOR)

**Theorem 1** (Upper bounds for P1). For any problems:

- $e_{STO}(n) = \mathcal{O}\left(\exp\left(-\frac{n}{H_{BOB} \log^2(K)}\right)\right)$
- $e_{ADV}(g)(n) = \mathcal{O}\left(\exp\left(-\frac{n}{\log(K) H_{UNIF}(g)}\right)\right)$

[1] J.-Y. Audibert, S. Bubeck, and R. Munos. Best-arm identification in multi-armed bandits. In *Conference on Learning Theory*, 2010.

## THE P1 ALGORITHM

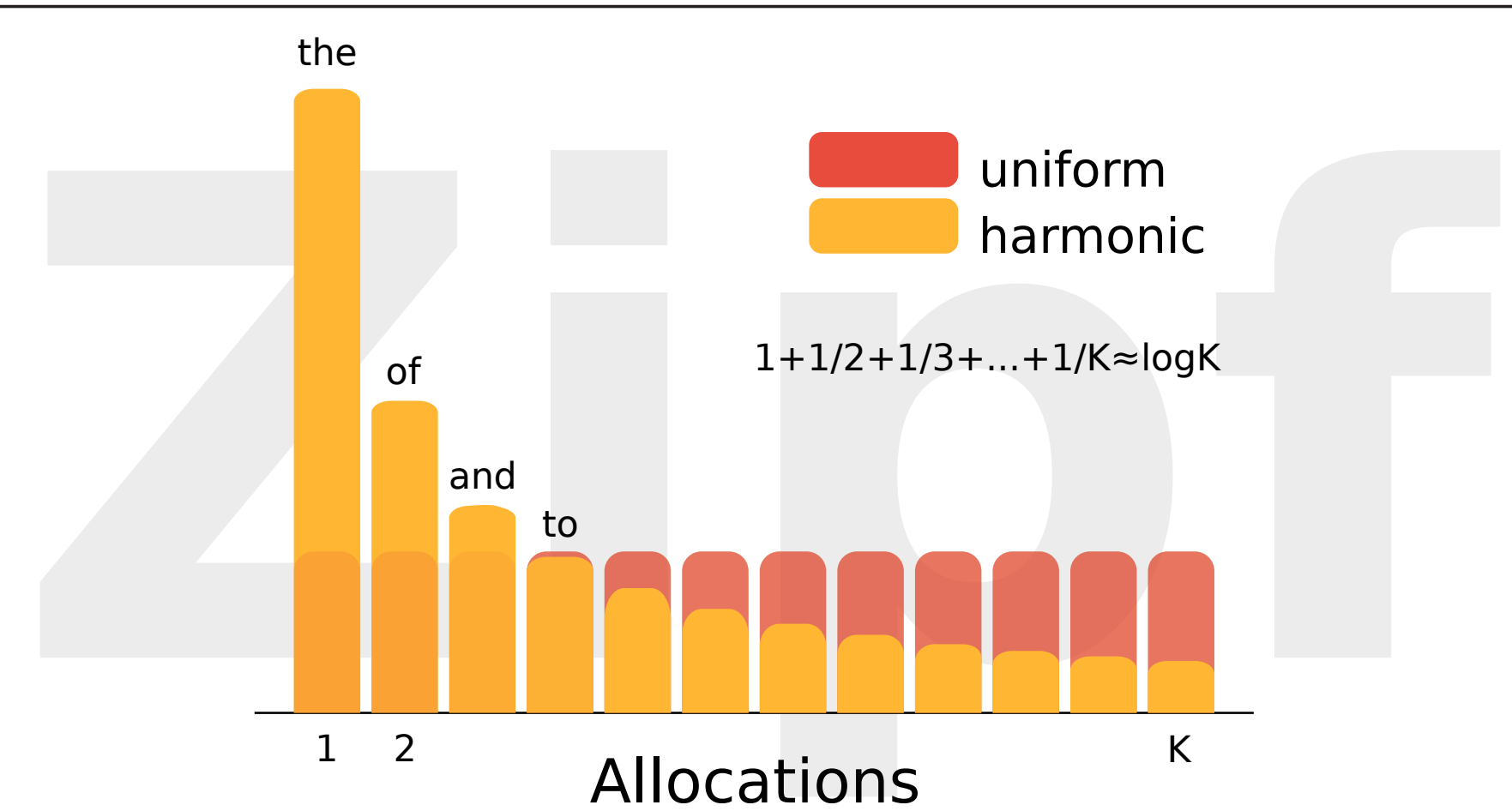
- P1 pulls
- the *best* arm with probability  $\frac{1}{2}$
  - the second *best* arm with proba  $\frac{1}{2}$
  - the third *best* arm with probability  $\frac{1}{3}$
  - and so on ... (and normalize)

For  $t = 1, 2, \dots$

- Sort & rank arms by decreasing  $\tilde{G}_{\cdot, t-1}$ : Rank arm  $k$  as  $\langle k \rangle_t \in [K]^a$ .
- Select  $I_t$  with  $\mathbb{P}(I_t = k) \triangleq \frac{1}{\langle k \rangle_t \log K}$ .

Recommend,  $J_t \triangleq \arg \max_{k \in [K]} \tilde{G}_{k,t}$ .

<sup>a</sup>Brake arbitrarily any problematic comparisons.

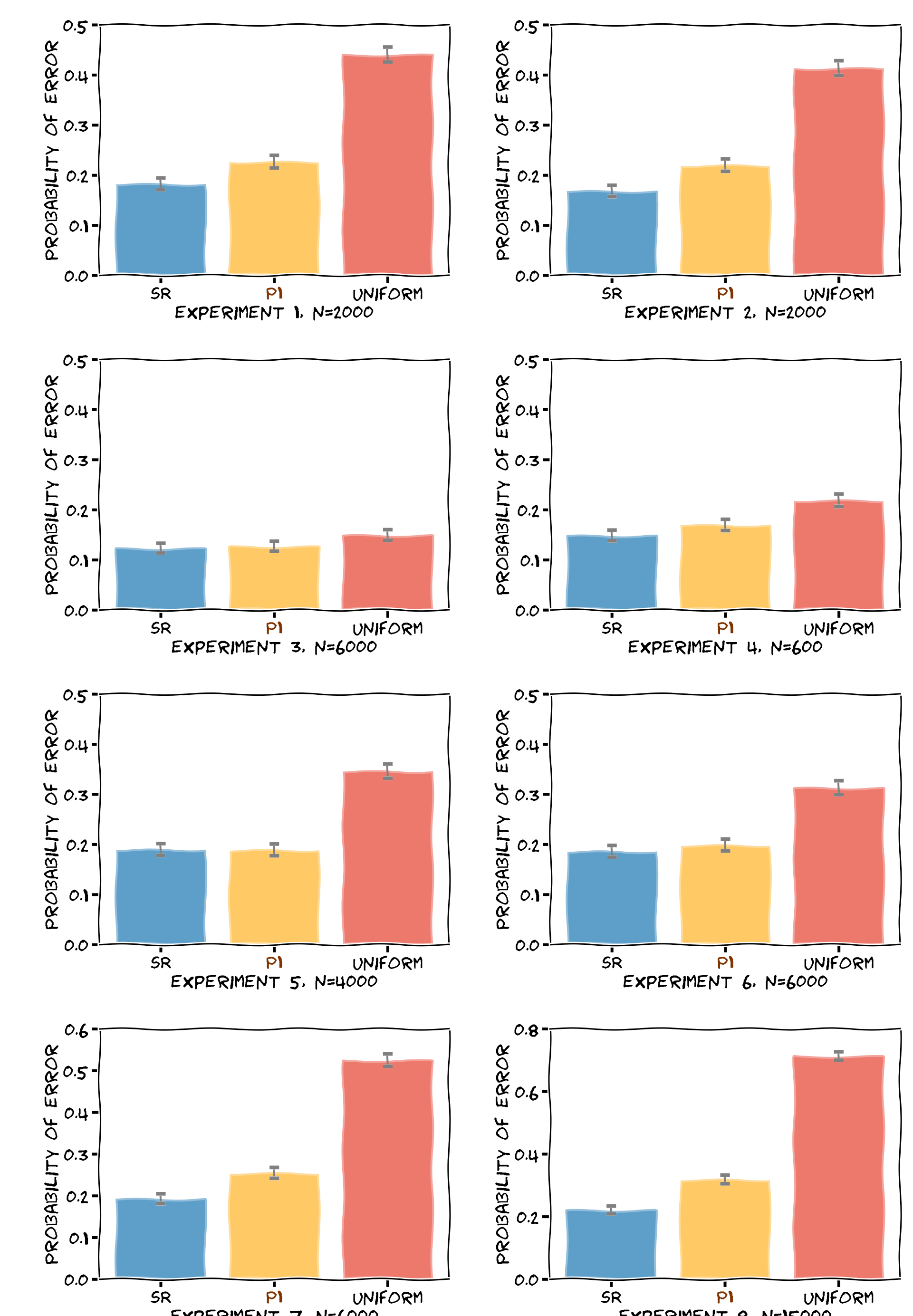


W.r.t. Rule, P1 early bets are almost costless!

P1 follows the allocation proportions of SR[1]

## STOCHASTIC CASE EXPERIMENTS

Experimental setup	$H_{SR}$	$H_{BOB}$	$H_{UNIF}$
1. 1 group of bad arms	2000	2000	2000
2. 2 groups of bad arms	1389	2083	3125
3. Geometric prog	5540	5540	11080
4. 3 groups of bad arms	400	500	938
5. Arithmetic prog	3200	3200	24000
6. 2 good, many bad	5000	7692	50000
7. 3 groups of bad arms	4082	5714	12000
8. Square-root gaps	3200	22M	160M



Empirical behavior in the figures mimics the behavior of the complexities in the table.