

## Problem and setting

- Learner interacts with  $K$  unknown arms denoted  $\nu_1, \dots, \nu_K$ .
- $X_{k,t}$  obtained by pulling arm  $k$  at time  $t$ .
- The learner seeks to collect the **largest possible reward**.
- Minimizing the extreme regret** which for a policy  $\pi$  that selects arm  $I_t$  at time  $t$  is defined by:

$$\mathcal{R}_T^\pi = \max_{k \leq K} \mathbb{E} \left[ \max_{t \leq T} X_{k,t} \right] - \mathbb{E}_\pi \left[ \max_{t \leq T} X_{I_t,t} \right].$$

Two notions of convergence:

- Weakly Vanishing Regret:**  $\mathcal{R}_T^\pi = o_{T \rightarrow \infty} (\max_{k \leq K} \mathbb{E}[\max_{t \leq T} X_{k,t}])$ .
- Strongly Vanishing Regret:**  $\lim_{T \rightarrow \infty} \mathcal{R}_T^\pi = 0$ .

## Challenge

- Relaxing parametric assumption on the distributions while obtaining strong theoretical guarantees**
  - Some works assume that the distribution are known (Frechet, Gumbel)
  - Other works have a semi-parametric assumption (second-order Pareto)
  - Some works with weaker assumptions but hard to obtain guarantees.
- Reducing computational and storage cost compared to existing approaches.**

Algorithm	Memory	Time
Extreme Hunter	$T$	$\mathcal{O}(T^2)$
MaxMedian	$T$	$\mathcal{O}(KT \log T)$
<b>QoMax-SDA</b>	$\mathcal{O}((\log T)^2)$	$\mathcal{O}(KT \log T)$
Extreme ETC	$\mathcal{O}(K(\log T)^3)$	$\mathcal{O}(K(\log T)^6)$
<b>QoMax-ETC</b>	$\mathcal{O}(K(\log T)^2)$	$\mathcal{O}(K(\log T)^3)$

Table: Average time and storage complexities of Extreme Bandit algorithms for a time horizon  $T$ .

## Dominating Tail

**Definition 1** (Exponential or polynomial tails). Let  $\nu$  be a distribution of survival function  $G$ . (1) If there exists  $C > 0$  and  $\lambda > 1$  such that  $G(x) \sim Cx^{-\lambda}$  we say that  $\nu$  has a **polynomial tail**. (2) If there exists  $C > 0, \lambda \in \mathbb{R}^+$  such that  $G(x) \sim C \exp(-\lambda x)$  we say that  $\nu$  has an **exponential tail**.

**Definition 2** (Dominating tail). Let  $G_1$  and  $G_2$  be the survival functions of two distributions  $\nu_1$  and  $\nu_2$ . We say that the tail of  $\nu_1$  **dominates** the tail of  $\nu_2$  (we write  $\nu_1 \succ \nu_2$ ) if there exists  $C > 1$  and  $x \in \mathbb{R}$  such that for all  $y > x$ ,  $G_1(y) > CG_2(y)$ .

## Quantile of Maxima (QoMax) estimator

→ Inspired by **Median of Means estimator**.  
 → Learner separates the data into **batches of equal sizes** and compute the quantile of order  $q$  of the maxima over the different batches. With  $N = b \times n$  data points, the learner allocates the data in  $b$  batches of size  $n$  and:

- find the maximum of each batch
- compute the quantile  $q$  over the  $b$  maxima.

→  $\bar{X}_{k,n,b}^q$  is the QoMax of order  $q$  computed from  $b$  batches of size  $n$  of i.i.d. replications from arm  $k$ .

## QoMax-ETC

For  $k \leq K$ :

Pull arm  $k$ ,  $b_T \times n_T$  times.

Allocate the data in  $b_T$  batches of size  $n_T$ . Compute their QoMax,  $\bar{X}_{k,n_T,b_T}^q$

For  $t = K \times n_T \times b_T + 1, \dots, T$ : Pull arm  $I_T = \arg \max_k \bar{X}_{k,n_T,b_T}^q$

## QoMax-SDA

A **round-based** algorithm based on three ingredients Beginning of round  $r$ :

- Selection of a leader: arm that has been pulled the most:  $\ell(r) = \arg \max_{k \leq K} n_k(r)$ .
  - Duels between the leader and the  $K - 1$  remaining arms: comparison of the QoMax of the challenger using its entire history and the **QoMax of the leader on a subsample of its history**.
  - Data collection procedure.
- $\mathcal{X}_k^r$  the history of arm  $k$  with  $b_k(r)$  batches of size  $n_k(r)$ .  $f(r)$  represents the *sampling obligation* at round  $r$ .
  - Duel:** Arm  $k \in \mathcal{A}_{r+1}$  (pulled arms) if (1) it wins its duel OR (2) undersampled i.e.  $n_k(r) \leq f(r)$ .
  - We assume that  $b_k(r)$  depends only on its number of queries  $n_k(r)$  i.e.  $b_k(r) = B(n_k(r))$  for some function  $B$ .

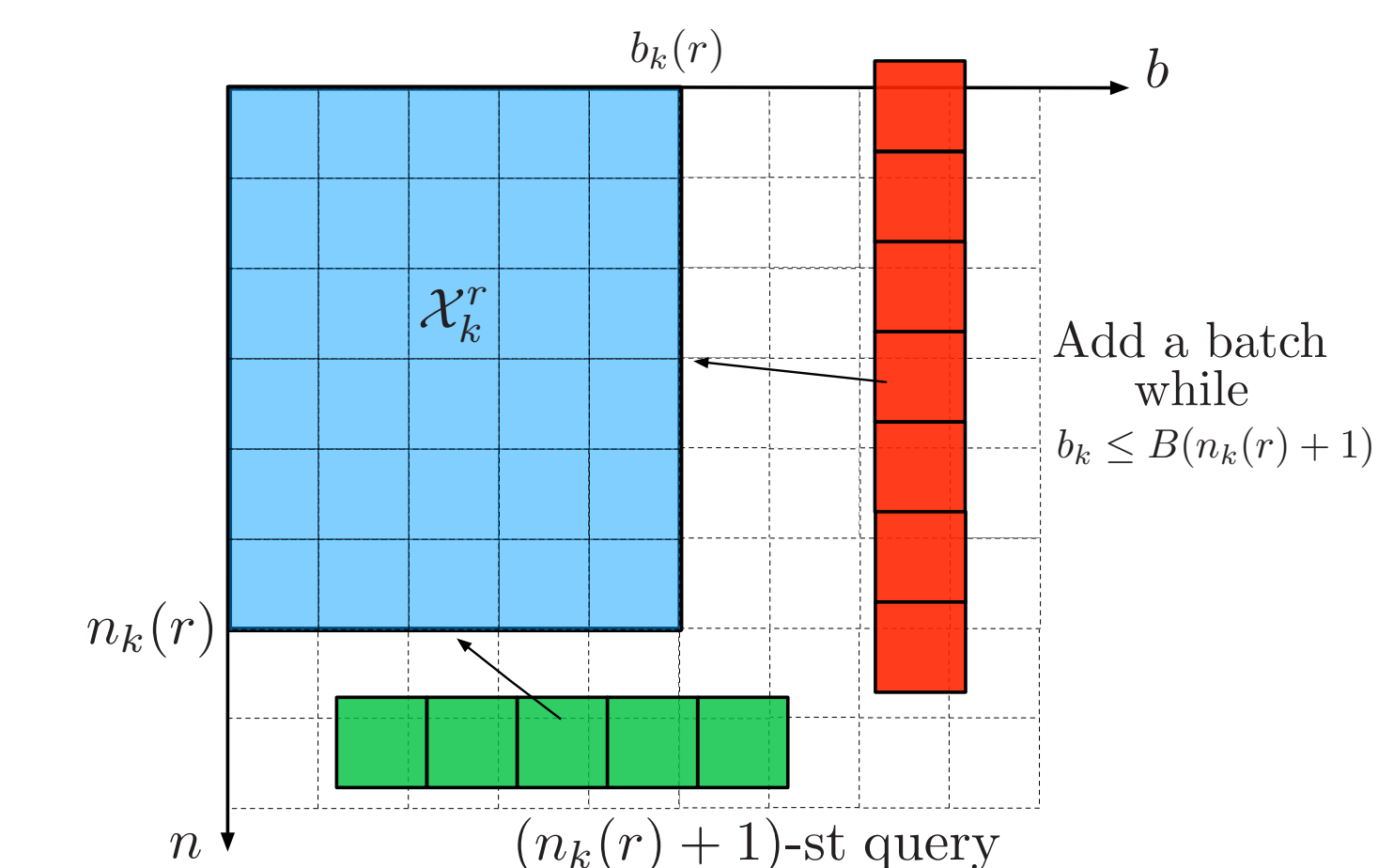


Figure 1: Illustration of the CollectData procedure at round  $r$  for a challenger  $k \in \mathcal{A}_{r+1}$ .

## Theoretical Guarantees

**Theorem 3** (Upper bound on the regret of QoMax-SDA). For any quantile  $q$ , any  $\gamma > 0$ , defining the parameters of QoMax-SDA as  $B(n) = n^\gamma$  and  $f(r) = (\log r)^{\frac{1}{\gamma}}$ .

The regret of QoMax-SDA is:

- Vanishing in the **strong sense** for exponential tails.
- Vanishing in the **weak sense** for polynomial tails.

## Numerical Results

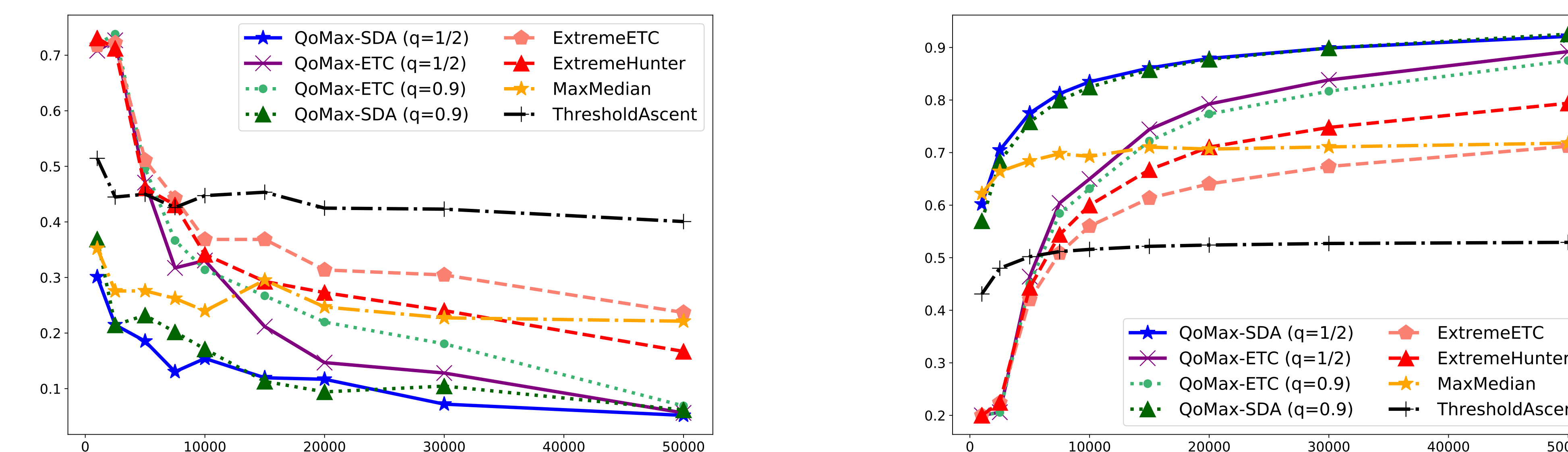


Figure 2: Proxy Empirical Regret (I) and Percentage of best arm pulls (II) averaged over  $10^4$  independent trajectories for  $T \in \{10^3, 2.5 \times 10^3, 5 \times 10^3, 7.5 \times 10^3, 10^4, 1.5 \times 10^4, 2 \times 10^4, 3 \times 10^4, 5 \times 10^4\}$ .