# TOWARDS FULLY ADAPTIVE REGRET MINIMIZATION IN HEAVY-TAILED BANDITS

GIANMARCO GENALTI, LUPO MARSIGLI, NICOLA GATTI AND ALBERTO MARIA METELLI

{gianmarco.genalti, nicola.gatti, albertomaria.metelli}@polimi.it, lupo.marsigli@mail.polimi.it

POLITECNICO MILANO 1863

- In stochastic multi-armed bandits, at each round $t \in \{1, \dots, T\}$ an agent choose among $K$ (unknowdistributions $\{v_i\}_{i \in \{1,\dots,K\}}$ and observes a sample reward $X_t$. The goal is to **minimize the expected cumulative regret** w.r.t. the best action.

- **Stochastic heavy-tailed bandits** gained popularity over the last years, extending the framework from sub-gaussian distributions to scenarios with (possibly) infinite variance, i.e.

$$\mathbb{E}_{v_i}[|X|^{1+\epsilon}] \leq u, \quad \epsilon \in (0,1], \quad \forall i \in \{1, \dots, K\}, \quad \textbf{all moments of order} > \mathbf{1 + \epsilon} \textbf{ are non-finite}.$$

- **Most of the literature assumes both $\epsilon$ and $u$ to be known** to the agent, but in practice this is usually a hard requirement to satisfy. We study the *adaptive heavy-tailed bandit problem*, in which the learner has no knowledge on these quantities.

- We show that, in general, **attaining optimal performance while being adaptive w.r.t. to either $\epsilon$ or $u$ is impossible**. However, under a specific assumption, **our algorithm is capable of matching the best possible performance of the setting**.