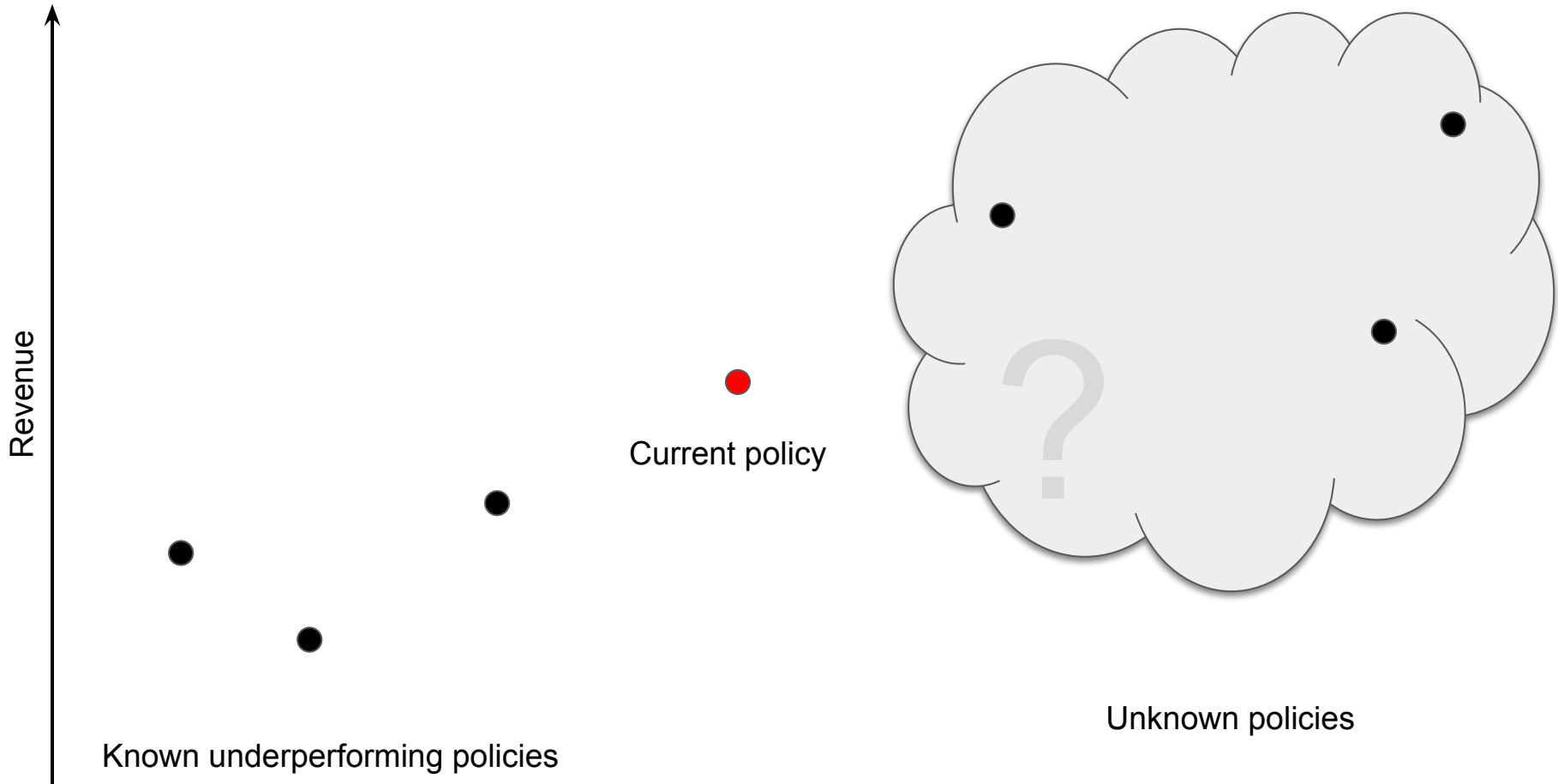
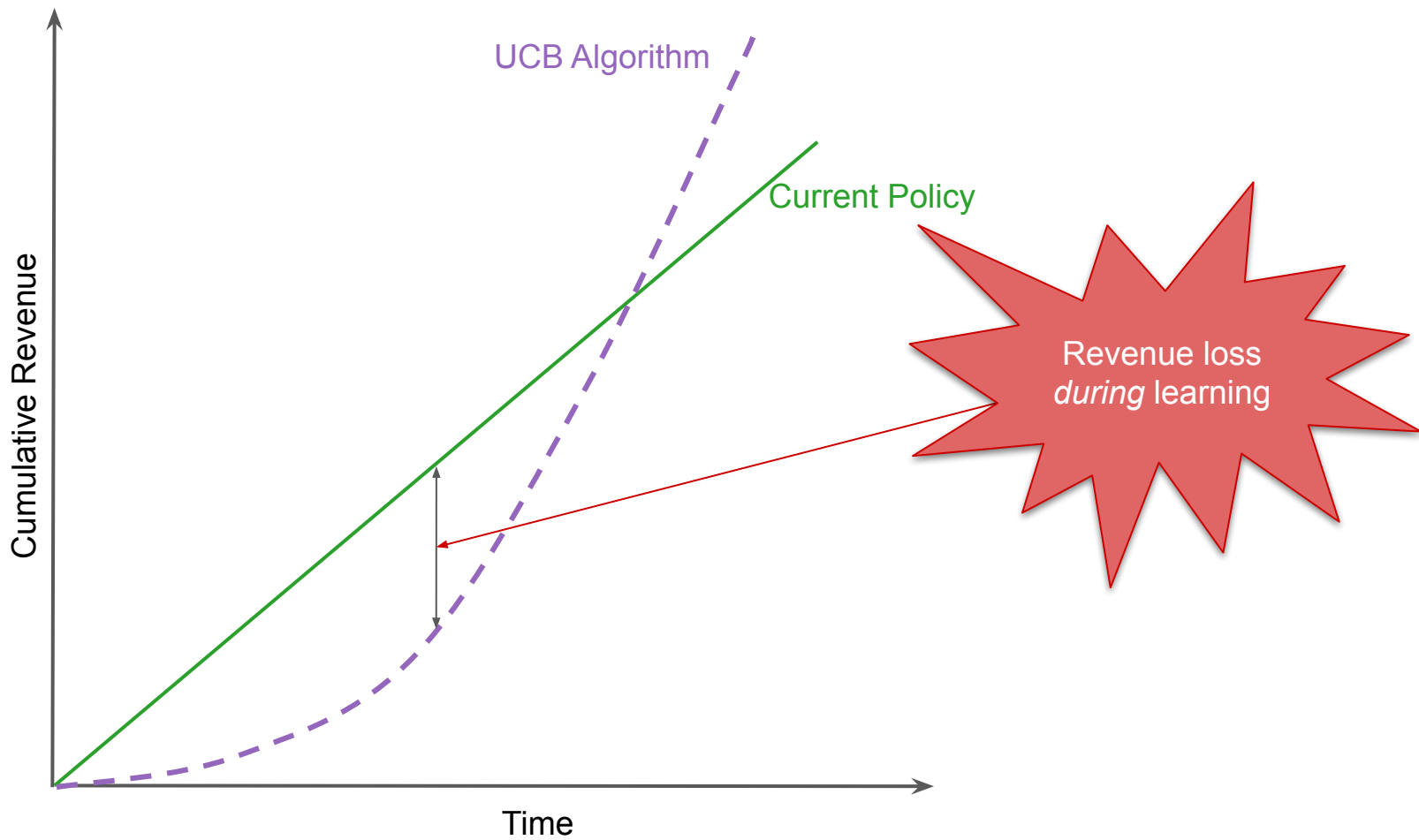


# Improved Algorithms for Conservative Exploration in Bandits

**Evrard Garcelon**, Mohammad Ghavamzadeh, Alessandro Lazaric and Matteo Pirodda

**Facebook AI Research**





*Problem:* How to learn an optimal policy without sacrificing much revenue?

(aka: how to perform exploration in a **conservative** way?)

# Conservative Condition

Should hold *uniformly*  
in time

$$\forall t > 0,$$

$$\mathbb{E} \left( \sum_{l=1}^t r_{l, a_l} \right)$$

$\geq$

$$(1 - \alpha) t \mu_b$$

Mean revenue of  
current policy

Mean *revenue* of the  
learning algorithm

Controls maximum  
revenue lost during  
learning

# Previous Work:

- Theoretically optimal algorithms for conservative exploration (CUCB) (Wu et al. 2016, Kazerouni et al. 2017)

# Contributions:

- *Improved empirical performance* in multi-armed and linear bandit (CUCB2)
- *Novel relaxed* conservative condition

## CUCB *(previous algorithm)*

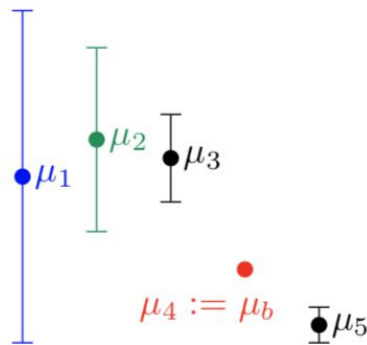
- Two phase algorithm
  - a. Computes optimistic arm
  - b. Checks a lower bound on the total revenue

=> impacts empirical performance!

## CUCB2 *(our algorithm)*

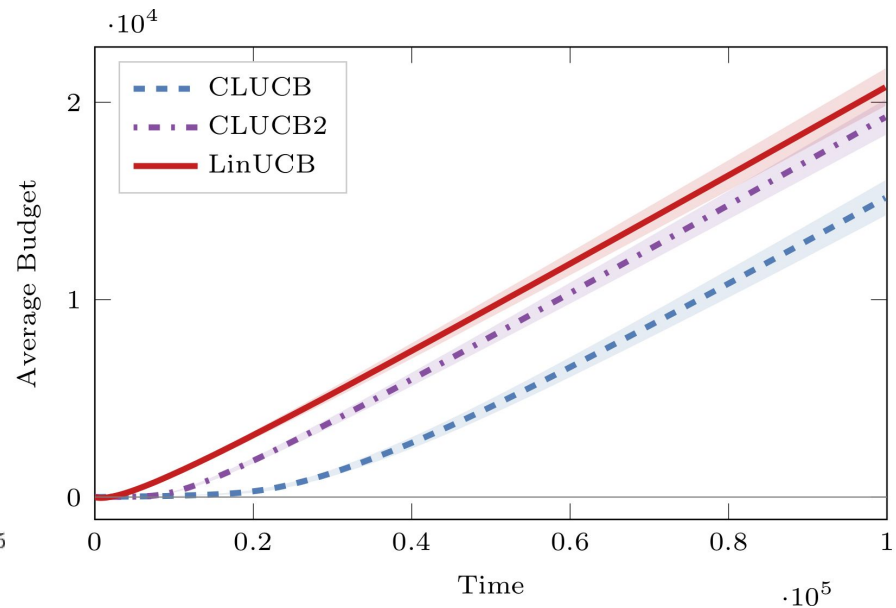
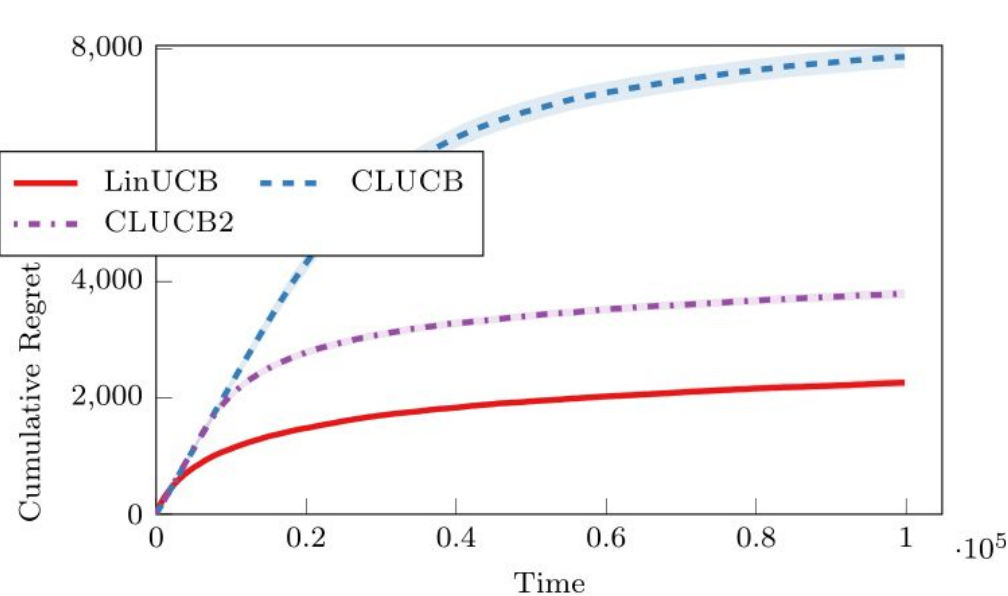
- Computes set of safe arms
- Plays the optimistic arm among safe arms

=> same regret but **better** performance!



*Example:* CUCB approach is suboptimal

# Jester Jokes Dataset (Goldberg et al. 2001)



- Cold start problem
- Linear features