

# Bandits: Best Arms Identification

Classical Settings and methods

Mengyan Zhang

Australian National University; Data61, CSIRO

Currently @ MSRA Social Computing

# Outline

- Definitions and Basic Settings
- Motivations and Applications
- Classical algorithms
  - Multi-armed bandits
    - UCB-based
    - Successive Rejects Type
    - Gap based – a unified algorithm
  - Black-box function optimization for many/continuous arms
  - Contextual bandits: Linear rewards

# Multi-armed Bandits: Sequential decision making



In each round  $t \in \{1, \dots, N\}$ ,

1. an agent selects an arm  $A_t = i \in \{1, \dots, K\}$  according policy  $\pi$
2. then receive a reward  $X_{i,T_i(t)}$  sampled from unknown distribution  $F_i$
3. update estimations over distribution  $F_i$  based on historical observations

# Multi-armed Bandits

## Sequential decision making



In each round  $t \in \{1, \dots, N\}$ ,

For a given user

1. an agent selects an arm  $A_t = i \in \{1, \dots, K\}$  according policy  $\pi$   
Recommender system      Item (e.g. news)      Recommendation strategy
2. then receive a reward  $X_{i,T_i(t)}$  sampled from unknown distribution  $F_i$   
CTR/click; non-click
3. update estimations over distribution  $F_i$  based on historical observations

# Multi-Armed Bandits

Simple regret  $r_t = \mu^* - \mu_{A_t}$  where  $\mu^* = \max_{k \in \{1, \dots, K\}} \mu_k$

## Best Arm Identification

to recommend best arm(s)  
at the end of exploration stage

### Fixed Budget:

the number of round for exploration phase is fixed and known

### Fixed Confidence:

the confidence level of quality of returned arms is fixed

**Regret minimization:** maximize the cumulative reward  
(i.e. minimize cumulative regret)

Cumulative regret  $R_T = \sum_{t=1}^T r_t$

Best choice with  
the current information

Other possibilities  
have not been tried

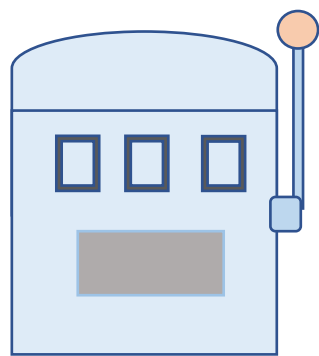
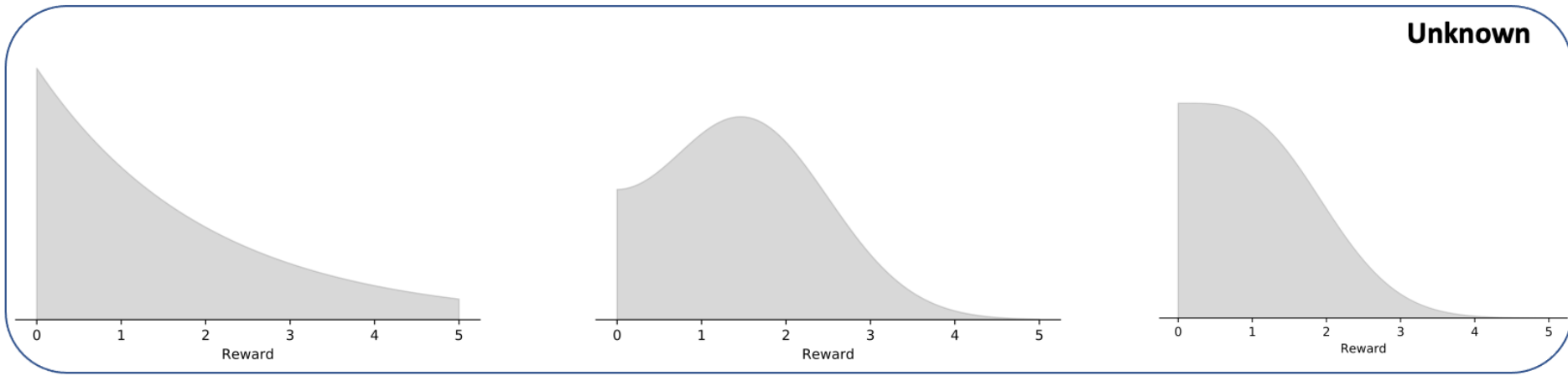


Regret Minimization:  
Exploitation vs. Exploration ?

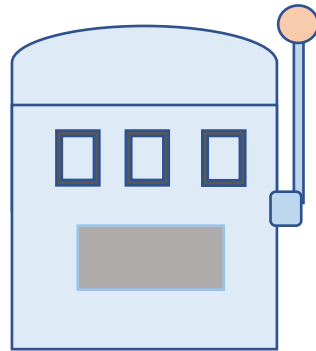
---

Best Arm Identification:  
How to allocate samples adaptively ?

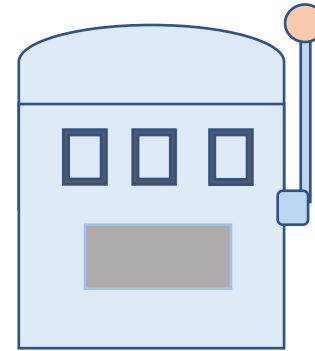
# Exploration Stage



**Arm 1**



**Arm 2**



**Arm 3**



Best Arms?

# Recommendation Stage



**Agent**



Fixed Budget

or

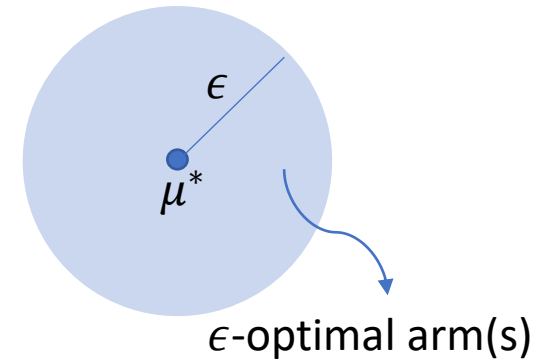
**Fixed Confidence**  
Recommend optimal arms  
with high probability



# Best Arm Identification

to recommend best arm(s) at the end of exploration stage

$$\text{Simple regret } r_t = \mu^* - \mu_{A_t} \quad \text{where } \mu^* = \max_{k \in \{1, \dots, K\}} \mu_k$$



**Fixed Budget:** the number of round for exploration phase is fixed and known

-> To maximize the probability of returning the  $\epsilon$ -optimal arm(s)

i.e.  $T$  is given, minimize Probability of error  $\delta_T = P[r_T \geq \epsilon]$

When  $\epsilon = 0, \delta_T = P[A_T \neq i^*]$

$$E[r_T] = \sum_{i \neq i^*} P[A_T = i] \Delta_i$$

Suboptimality gap  $\Delta_i = \mu^* - \mu_i$

$$\Delta_{\min} \delta_T \leq E[r_T] \leq \Delta_{\max} \delta_T$$

$\delta_T$  and  $E[r_T]$  behave similarly

**Fixed Confidence:** the confidence level of quality of returned arms is fixed

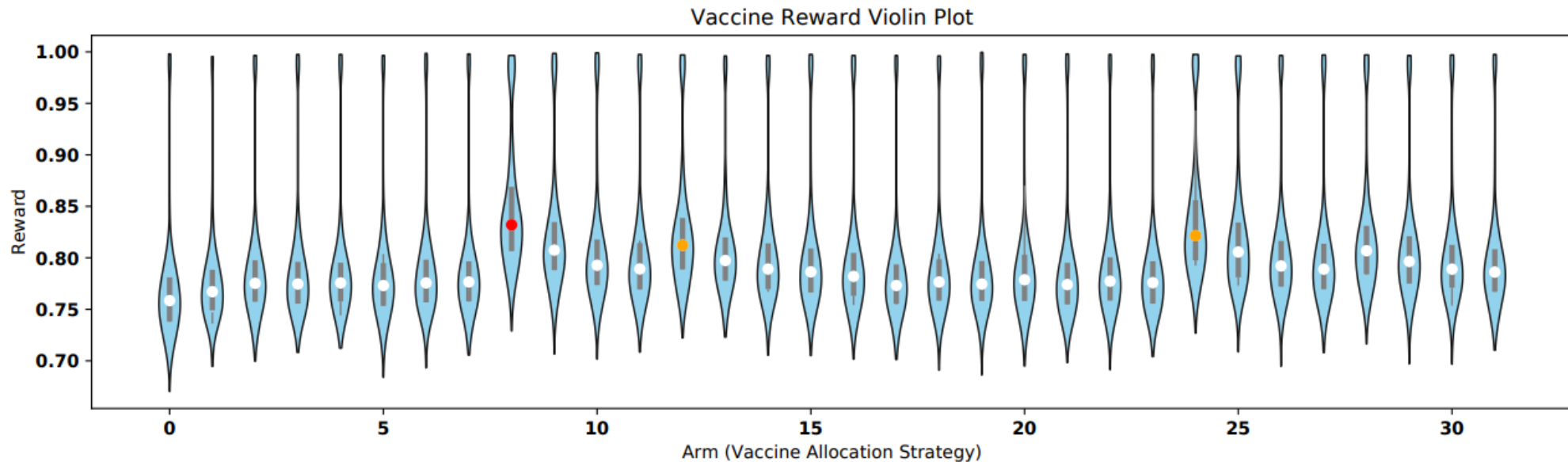
-> To minimize the number of rounds needed

i.e.  $\delta$  is given, minimize budget  $T$  s.t.  $P[r_T \geq \epsilon] \leq \delta$



# Applications: Vaccine testing

- Identify optimal strategies (highest mean/median reward) for allocation vaccines
- Arm: vaccine allocation strategy
- Reward: proportion of individuals that did not experience symptomatic infection



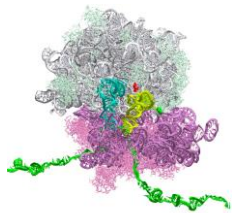
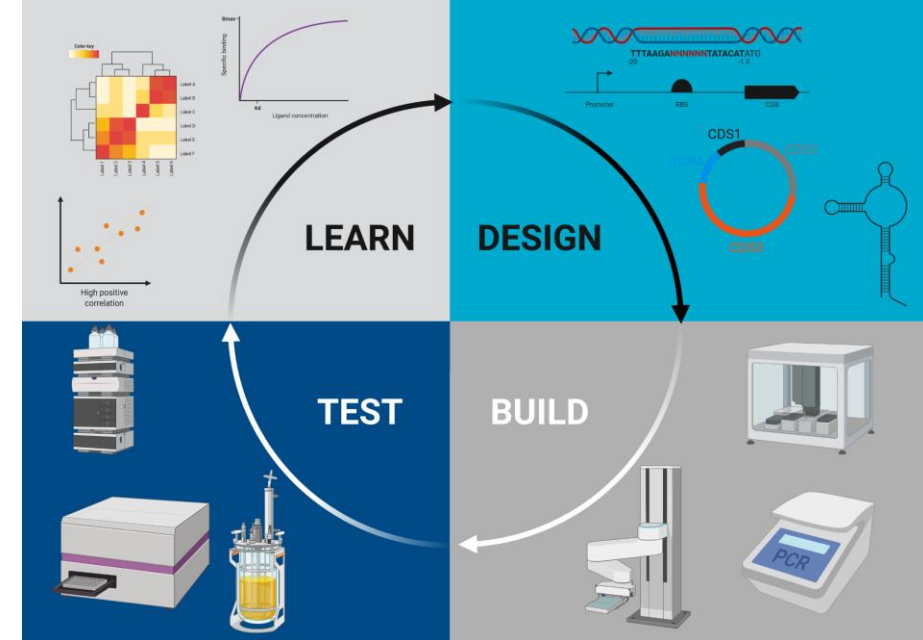
# Applications: Biological design

With fixed budget, design Ribosome Binding Site (RBS) sequences



Optimize the protein expression level

Identify the DNA sequences with highest possible protein expression level



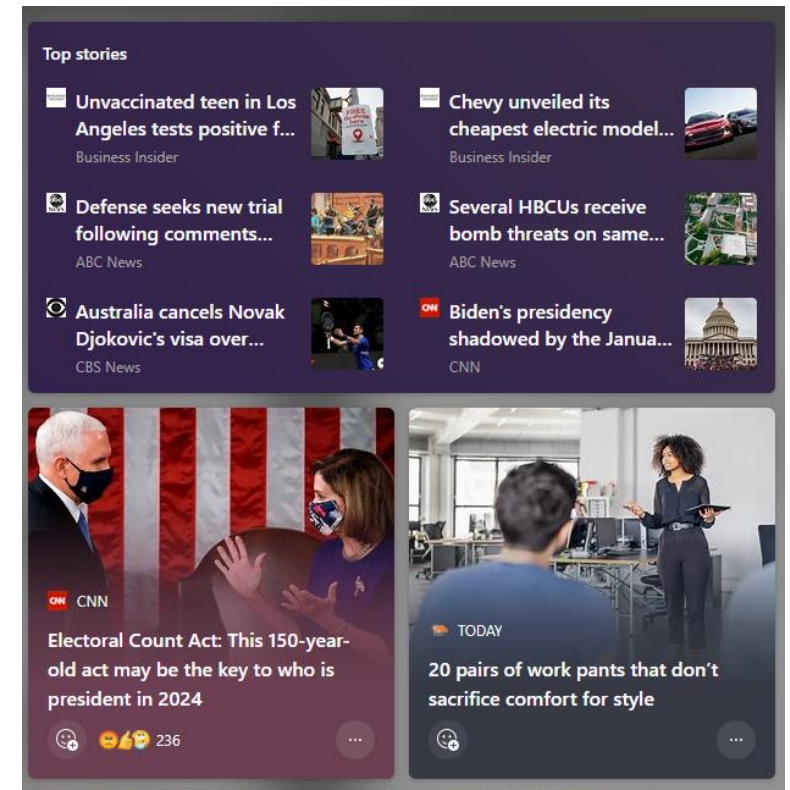
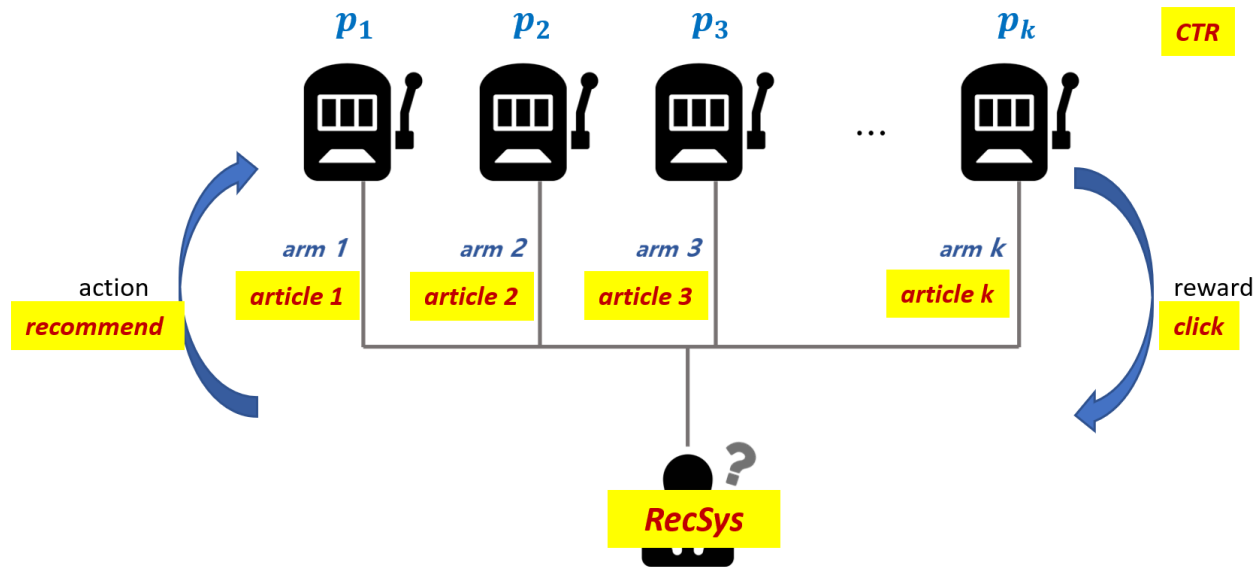
Arm: RBS sequence	Reward: Normalized* Protein Expression Level
TTTAAGAG <b>GTT</b> ATATACAT	1.58
TTTAAGA <b>ATATG</b> CTATACAT	1.42
TTTAAGAC <b>TCGG</b> ATATACAT	0.14
TTTAAGAG <b>GTTTT</b> TATACAT	2.88



\* zero mean and unit variance normalization  $z = \frac{x-\mu}{\sigma}$

# Applications: Recommendation System

- identify the most popular items (with potential high CTR) above some level of confidence using fewest possible samples
- Arm: item (e.g. news)
- Reward: click/ preference



# Outline

- Definitions and Basic Settings
- Motivations and Applications
- Classical algorithms
  - Multi-armed bandits
    - UCB-based
    - Successive Rejects Type
    - Gap based – a unified algorithm
  - Black-box function optimization for many/continuous arms
  - Contextual bandits: Linear rewards

# Best Arm Identification in Multi-Armed Bandits

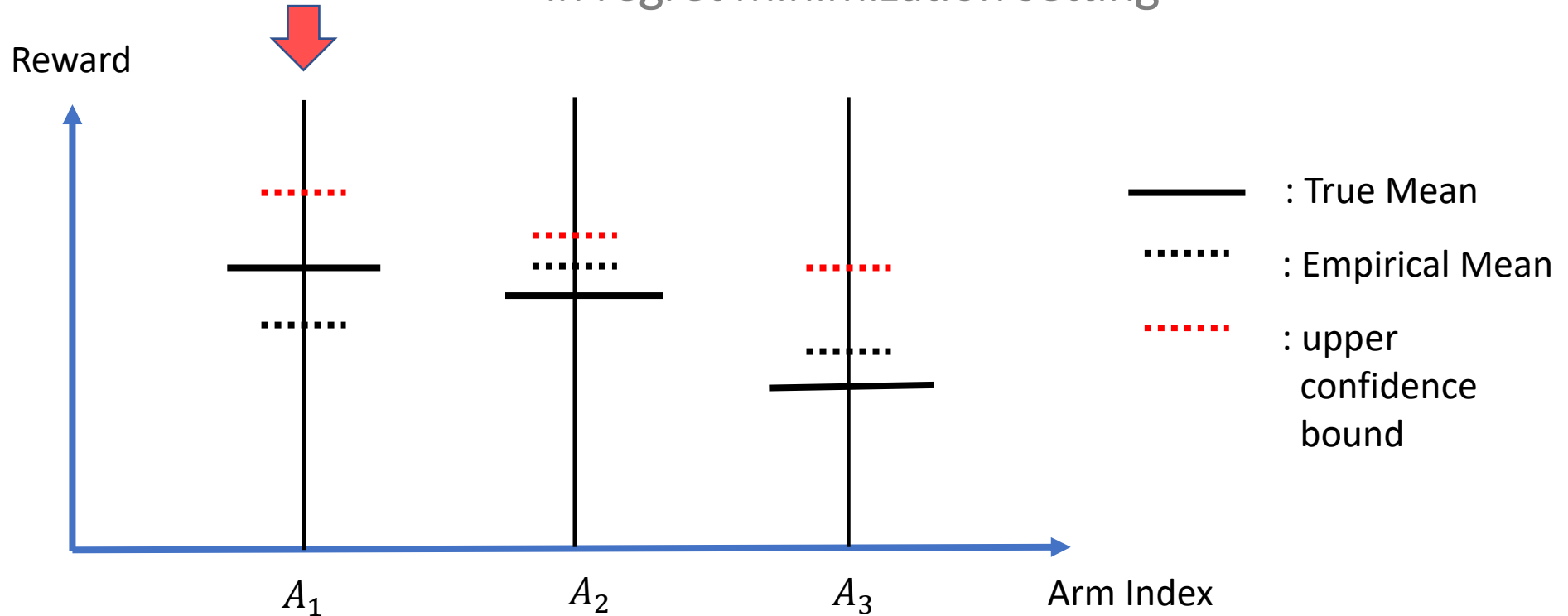
- K arms associated with K unknown reward distributions respectively (assume rewards are in  $[0,1]$  and there is a unique optimal arm),
- T rounds/budget,
- The agent select one arm  $A_t$  according to policy  $\pi$   $\rightarrow$  observe reward drawn from  $\nu_{A_t}$  independently from the past (actions and observations)

Simple regret  $r_T = \mu^* - \mu_{A_T}$  where  $\mu^* = \max_{k \in \{1, \dots, K\}} \mu_k$

Probability of error  $\delta_T = P[A_T \neq i^*]$

# Upper Confidence Bound (UCB)

in regret minimization setting



$$\operatorname{argmax}_{i \in \mathcal{K}} \text{Central Tendency} + \beta \text{Confidence Width}$$

Exploitation

Exploration

e.g. Empirical Mean

# UCB-E (Upper Confidence Bound Exploration) algorithm

$$\operatorname{argmax}_{i \in \{1, \dots, K\}} B_{i,s} = \hat{X}_{i,s} + \sqrt{\frac{a}{s}}$$

Suboptimality gap  $\Delta_i = \mu^* - \mu_i$

Hardness of the task

$$H_1 = \sum_{i=1}^K \frac{1}{\Delta_i^2} \quad \text{and} \quad H_2 = \max_{i \in \{1, \dots, K\}} i \Delta_{(i)}^{-2}.$$

Probability or error  $\delta_T \sim \exp(-c T/H_1)$

Number of draws

Regret minimisation  $\rightarrow a \sim O(\log T)$   
 BAI, need more exploration  $\rightarrow a \sim O(T)$

$$\text{UCB-E: } a = \frac{25}{36} (T - K)/H_1$$

$H_1$  is unknown and hard to be online estimated!



# Successive Rejects

$$\overline{\log}(K) = \frac{1}{2} + \sum_{i=2}^K \frac{1}{i}$$

$$n_k = \left\lceil \frac{1}{\overline{\log}(K)} \frac{T - K}{K + 1 - k} \right\rceil, n_0 = 0$$

- Devide budget  $T$  into  $K-1$  phases
- In each phase  $k$ , pull equally often each arm which has one been rejected yet
- At the end of each phase, reject the arm with lowest empirical mean
- Recommend the last surviving arm

Suboptimality gap  $\Delta_i = \mu^* - \mu_i$

Hardness of the task

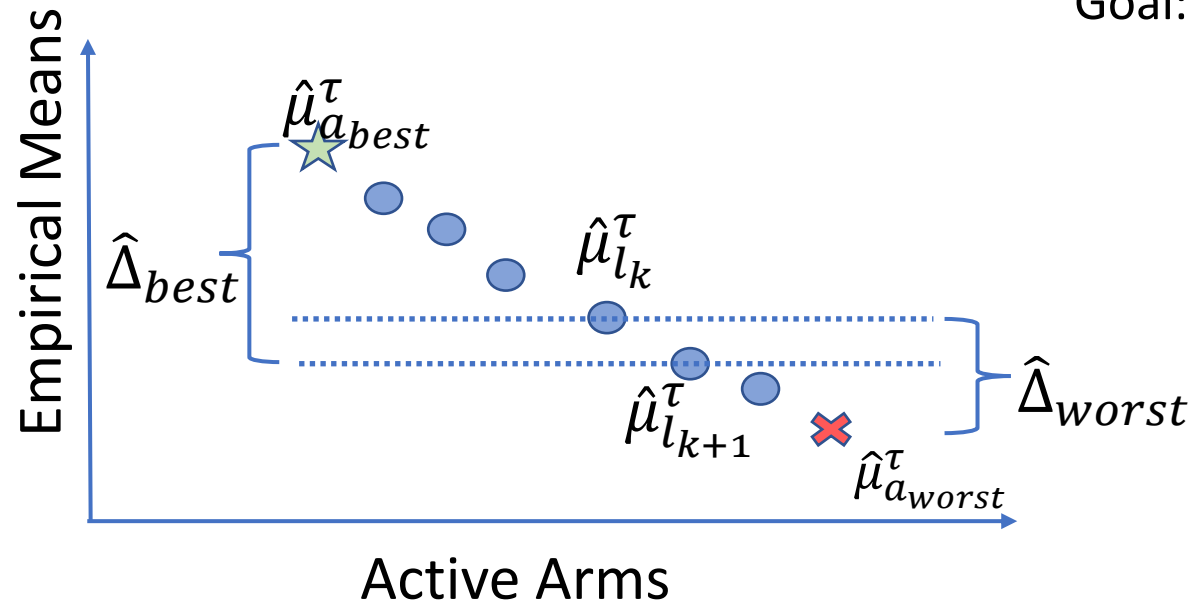
$$H_1 = \sum_{i=1}^K \frac{1}{\Delta_i^2} \quad \text{and} \quad H_2 = \max_{i \in \{1, \dots, K\}} i \Delta_{(i)}^{-2}.$$

$$H_2 \leq H_1 \leq \log(2K) H_2.$$

Probability of error  $\delta_T \sim \exp(-c T/H_2)$

# Multi-arm Identification – Successive Accepts and Rejects

Goal: recommend  $m \geq 1$  arms with fixed budget  $T$



$l_k$ : remaining number of arms to be found in current phase  $k$

$\hat{\Delta}_{best} > \hat{\Delta}_{worst}$  : Accept  $\star$

$\hat{\Delta}_{best} \leq \hat{\Delta}_{worst}$  : Reject  $\times$

$$\Delta_i^{(m)} = \begin{cases} \mu_i - \mu_{m+1} & \text{if } i \leq m \\ \mu_m - \mu_i & \text{if } i > m \end{cases},$$

$$H_1^{(m)} = \sum_{i=1}^K \frac{1}{\left(\Delta_i^{(m)}\right)^2},$$

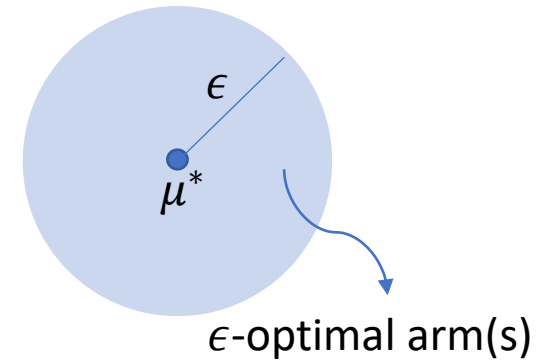
$$H_2^{(m)} = \max_{i \in \{1, \dots, K\}} i \left(\Delta_{(i)}^{(m)}\right)^{-2},$$

Probability of error  $\delta_T \sim \exp(-c T / H_2^{(m)})$

# Best Arm Identification

to recommend best arm(s) at the end of exploration stage

$$\text{Simple regret } r_t = \mu^* - \mu_{A_t} \quad \text{where } \mu^* = \max_{k \in \{1, \dots, K\}} \mu_k$$



**Fixed Budget:** the number of round for exploration phase is fixed and known

-> To maximize the probability of returning the  $\epsilon$ -optimal arm(s)

i.e.  $T$  is given, minimize Probability of error  $\delta_T = P[r_T \geq \epsilon]$

When  $\epsilon = 0$ ,  $\delta_T = P[A_T \neq i^*]$

$$E[r_T] = \sum_{i \neq i^*} P[A_T = i] \Delta_i$$

Suboptimality gap  $\Delta_i = \mu^* - \mu_i$

$$\Delta_{\min} \delta_T \leq E[r_T] \leq \Delta_{\max} \delta_T$$

$\delta_T$  and  $E[r_T]$  behave similarly

**Fixed Confidence:** the confidence level of quality of returned arms is fixed

-> To minimize the number of rounds needed

i.e.  $\delta$  is given, minimize budget  $T$  s.t.  $P[r_T \geq \epsilon] \leq \delta$

# Unified Methods: UGapE

## SELECT-ARM ( $t$ )

Compute  $B_k(t)$  for each arm  $k \in A$

Identify the set of  $m$  arms  $J(t) \in \arg \min_{k \in A}^{1..m} B_k(t)$

Pull the arm  $I(t) = \arg \max_{k \in \{l_t, u_t\}} \beta_k(t-1)$

Observe  $X_{I(t)}(T_{I(t)}(t-1) + 1) \sim \nu_{I(t)}$

Update  $\hat{\mu}_{I(t)}(t)$  and  $T_{I(t)}(t)$

High probability upper and lower bounds on the mean of arm  $k$

$$\begin{cases} U_k(t) = \hat{\mu}_k(t-1) + \beta_k(t-1) \\ L_k(t) = \hat{\mu}_k(t-1) - \beta_k(t-1). \end{cases}$$

Upper bound on simple regret  $B_k(t) = \max_{i \neq k}^m U_i(t) - L_k(t)$

Best possible arm left outside of  $J(t)$ ; worst possible arm among those in  $J(t)$ ;  
 $u_t = \arg \max_{j \notin J(t)} U_j(t)$  and  $l_t = \arg \min_{i \in J(t)} L_i(t)$ .

Represents: how bad the choice of  $J(t)$  could be

Fixed budget

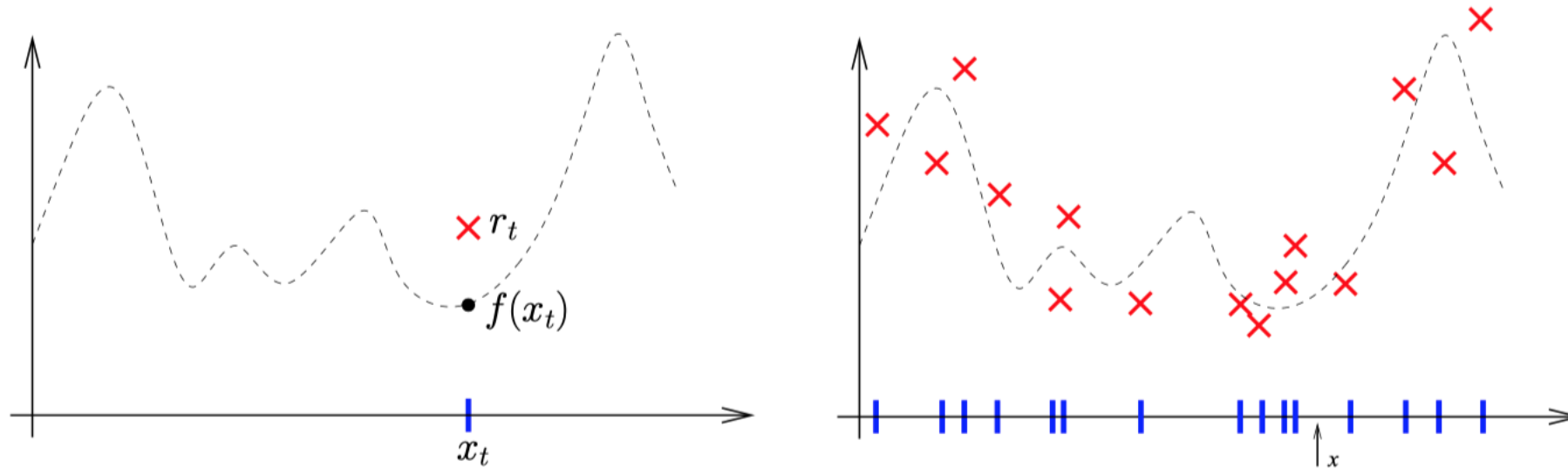
$$\text{UGapEb: } \beta_k(t-1) = b \sqrt{\frac{a}{T_k(t-1)}}$$

Fixed confidence

$$\text{UGapEc: } \beta_k(t-1) = b \sqrt{\frac{c \log \frac{4K(t-1)^3}{\delta}}{T_k(t-1)}}$$

# Black-box function optimisation

Optimism in the face of uncertainty for black-box optimisation of a function  $f$  given stochastic evaluation of the function.

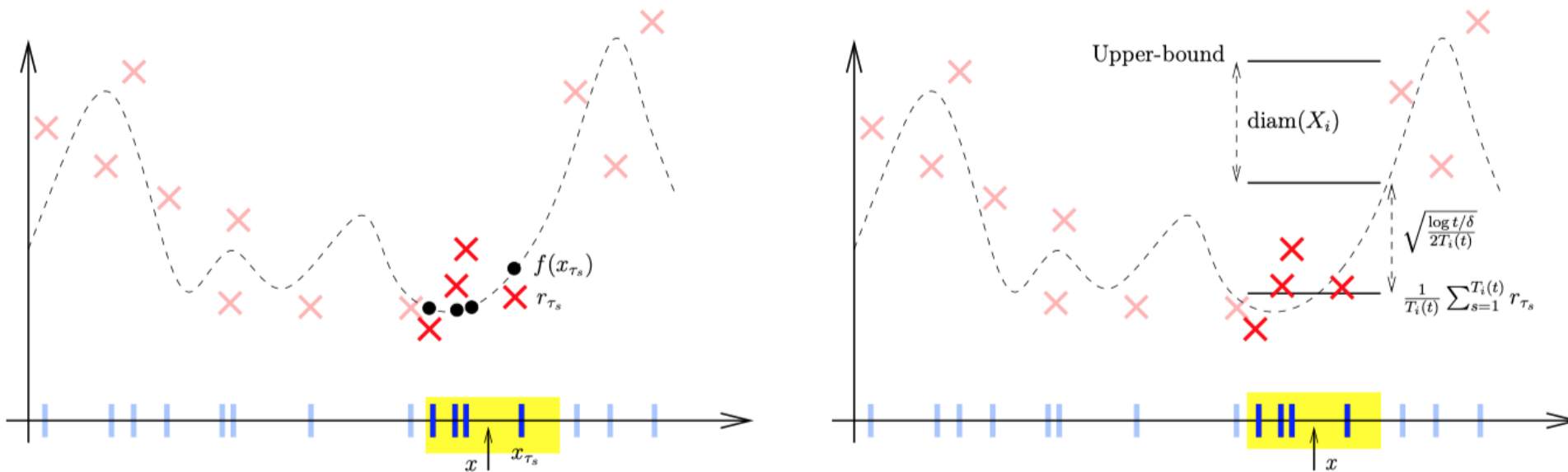


$$y_t = f(x_t) + \epsilon_t$$

$$E[\epsilon_t | x_t] = 0$$

Goal: Search for a good approximation of the maximum of a function  $\chi \rightarrow \mathbb{R}$  with a fixed number of function evaluations

# Optimism in the face of uncertainty



Given a subset  $X_i \in \mathcal{X}$

$$B_{t, T_i(t)}(X_i) \stackrel{\text{def}}{=} \frac{1}{T_i(t)} \sum_{s=1}^{T_i(t)} r_{\tau_s} + \sqrt{\frac{\log t / \eta}{2T_i(t)}} + \text{diam}(X_i) \geq \max_{x \in X_i} f(x).$$

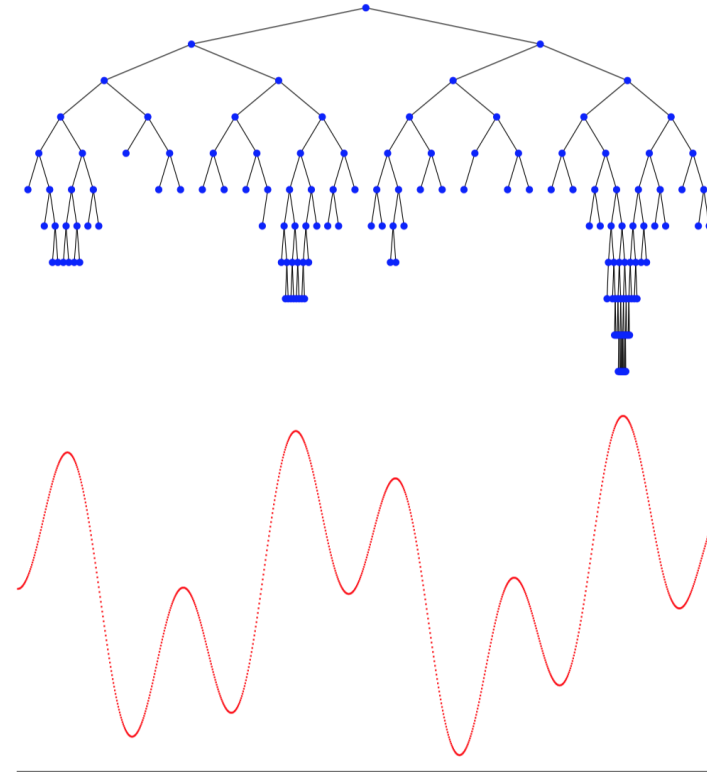
A trade-off in the choice of size of  $X_i$

# For many arms or continuous arms: Hierarchical Partition

For each round: 
$$b_{h,j}(t) \stackrel{\text{def}}{=} \hat{\mu}_{h,j}(t) + \sqrt{\frac{\log(n^2/\eta)}{2T_{h,j}(t)}} + \delta(h),$$

- Select the leaf node with highest b-value
- sample the center point and collect rewards
- Expand node: when # draws > some threshold

(depends on T, smoothness of f)



Extend: model f as a sample from Gaussian process, e.g. GPOO



# Contextual BAI with linear rewards – Generalized Successive Elimination

$$y_i = f(x_i) + \epsilon$$

Linear model  $\hat{\mu}_{i,t} = x_i^T \hat{\theta}_t$

Generalized linear model  $\hat{\mu}_{i,t} = g(x_i^T \hat{\theta}_t)$

e.g. logistic regression  $g(x) = (1 + \exp(-x))^{-1}$

**Repeat**

**Explore:** split budget evenly to each stage, then allocate the budget in each stage according to some **allocation policy**

**Estimate:** estimate each arm's mean  $\hat{\mu}_{i,t}$

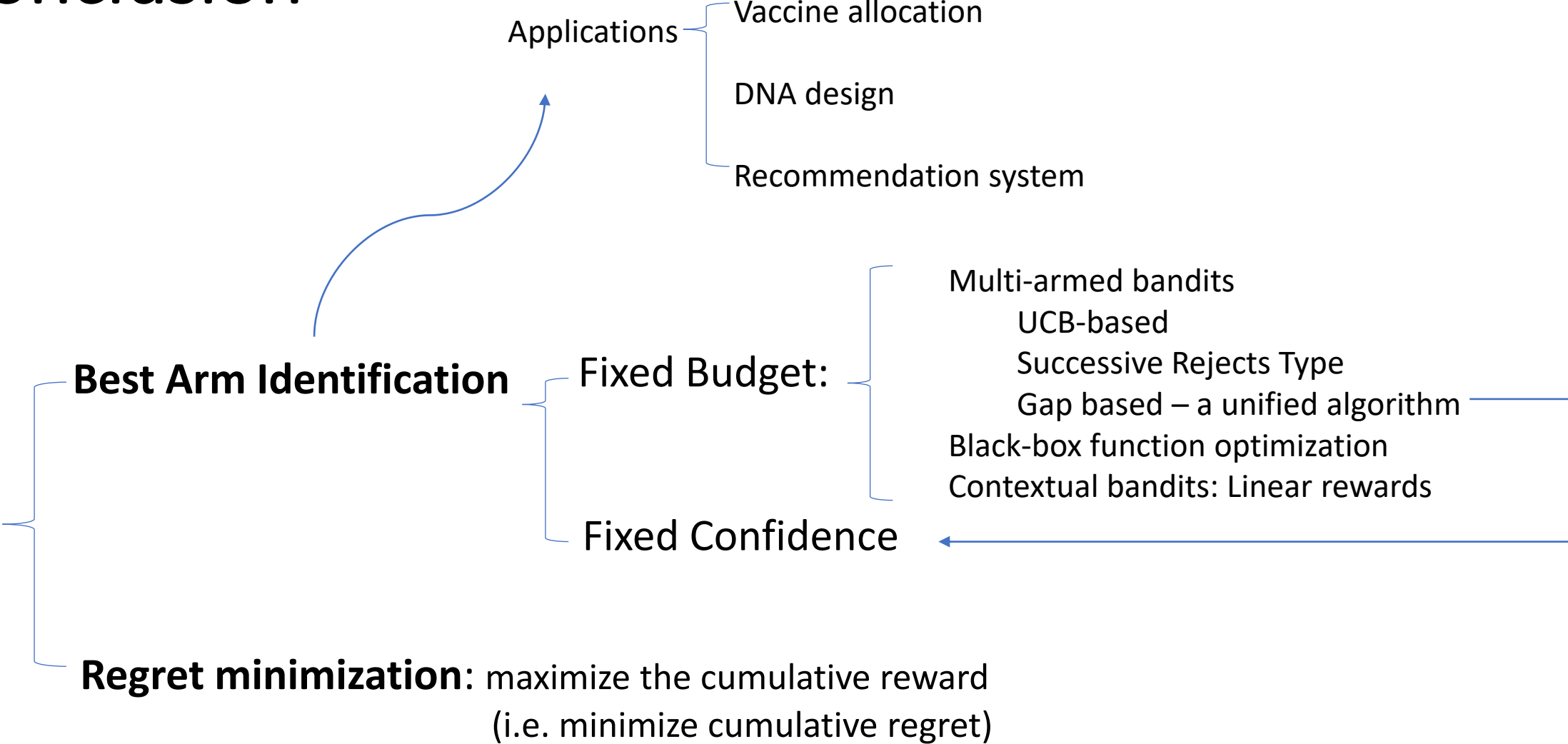
**Elimination:** sort arms in descending order of  $\hat{\mu}_{i,t}$ , keep the first  $1/\eta$  of them while eliminating the rest

**Until** only one arm is remained

Uniform allocation  
Optimal allocation

G-optimal design: minimize the maximum variance uniformly along all directions

# Conclusion



# References

- Even-Dar, Eyal, Shie Mannor, and Yishay Mansour. Action Elimination and Stopping Conditions for the Multi-Armed Bandit and Reinforcement Learning Problems. 2006.
- Audibert, Jean-Yves, and Sébastien Bubeck. Best Arm Identification in Multi-Armed Bandits. COLT2010.
- Gabillon, Victor, Mohammad Ghavamzadeh, and Alessandro Lazaric. Best Arm Identification: A Unified Approach to Fixed Budget and Fixed Confidence. In *Advances in Neural Information Processing Systems 25*. 2012.
- Bubeck, Sébastien, Tengyao Wang, and Nitin Viswanathan. Multiple Identifications in Multi-Armed Bandits. ICML2013.
- Munos, Rémi. “From Bandits to Monte-Carlo Tree Search: The Optimistic Principle Applied to Optimization and Planning,” 2014.
- **Mengyan Zhang**, Russell Tsuchida, and Cheng Soon Ong.. Opportunities and Challenges in Designing Genomic Sequences. ICML Workshop on Computational Biology 2021.
- **Mengyan Zhang**, and Cheng Soon Ong. Quantile Bandits for Best Arms Identification. ICML2021.
- Azizi, MohammadJavad, Branislav Kveton, and Mohammad Ghavamzadeh. “Fixed-Budget Best-Arm Identification in Contextual Bandits: A Static-Adaptive Algorithm.” *ArXiv:2106.04763 [Cs]*, June 9, 2021.
- **Mengyan Zhang**, Russell Tsuchida, and Cheng Soon Ong. Gaussian Process Bandits with Aggregated Feedback. AAI2022.