# Adaptive Recommendations with Bandit Feedback

Mengyan Zhang

Computational Media Lab, Australian National University
Machine Learning Research Group, Data61,CSIRO

Supervisors: Cheng Soon Ong, Lexing Xie, Eduardo Eyras

**Multi-armed Bandits:**
         **Sequential decision-making**



In each round t $\epsilon$ {1, .., N},

    1. an agent selects an arm $A_t = i \, \epsilon \, \mathcal{K}$ 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
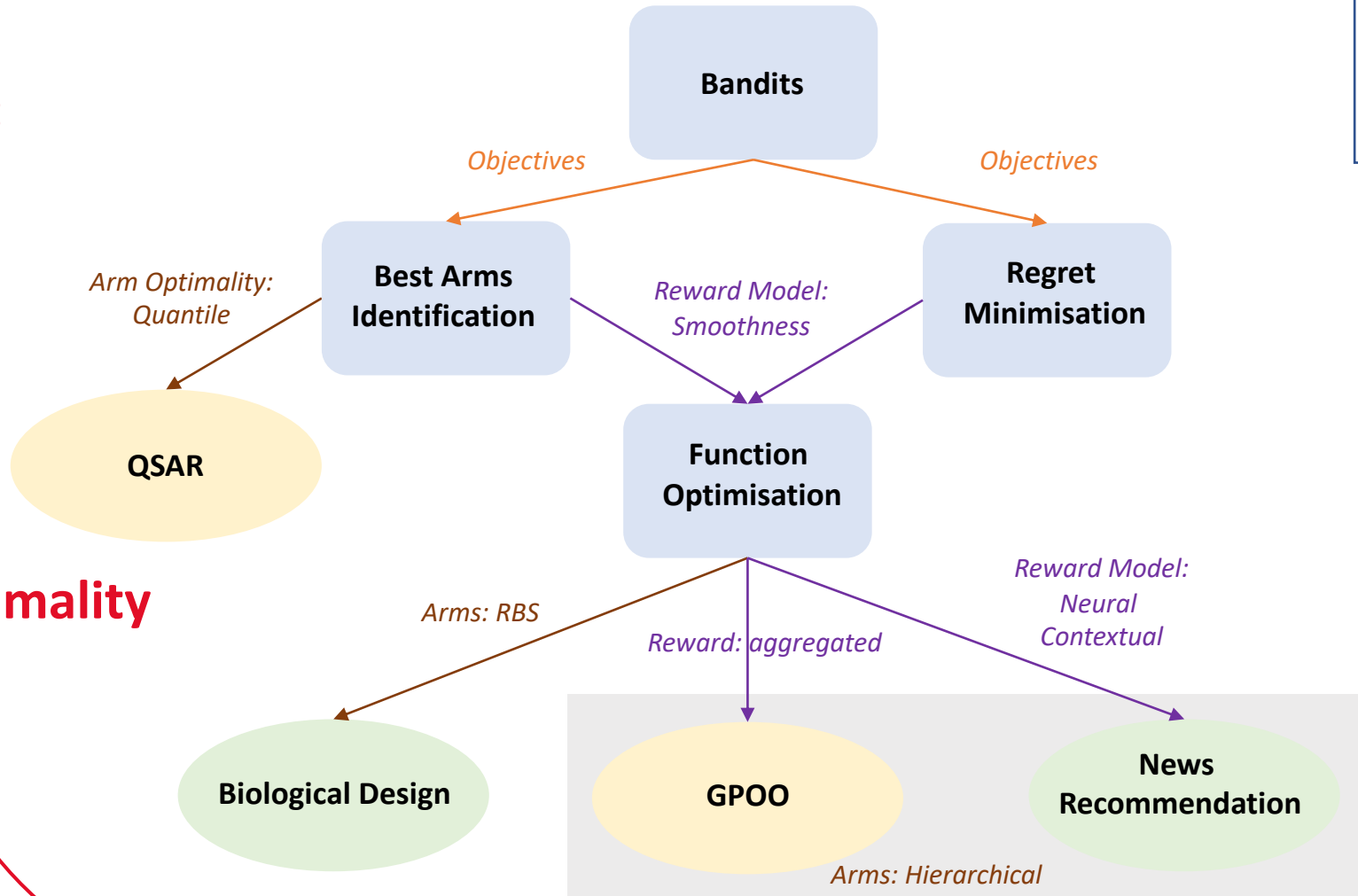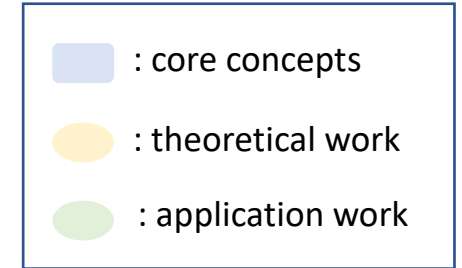
# Why Bandits?

- Challenge: **sequential decision-making with uncertainty**

- Provide model for **E & E dilemma**

- **Applications**:
  - Adaptive experimental design: clinical, drug, food
  - Configure web interfaces: item recommendations, dynamic pricing, ad placement
  - Plays a role in algorithms like Monte Carlo Tree Search

- Rich structure connecting to other branches of **math**: concentration analysis, information theory, etc.
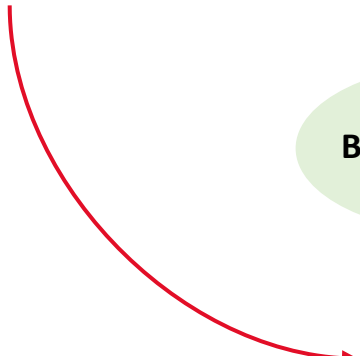
# Three Design Choices of bandit tasks

- Objectives:
  - What's the goal of designing a policy?
  - How to evaluate the performance?

- Rewards
  - How to model the rewards? – Smoothness, context

- Arms
  - How to define an arm?
  - Can we form new arms based on single arms?
  - How to define the optimality of arms?
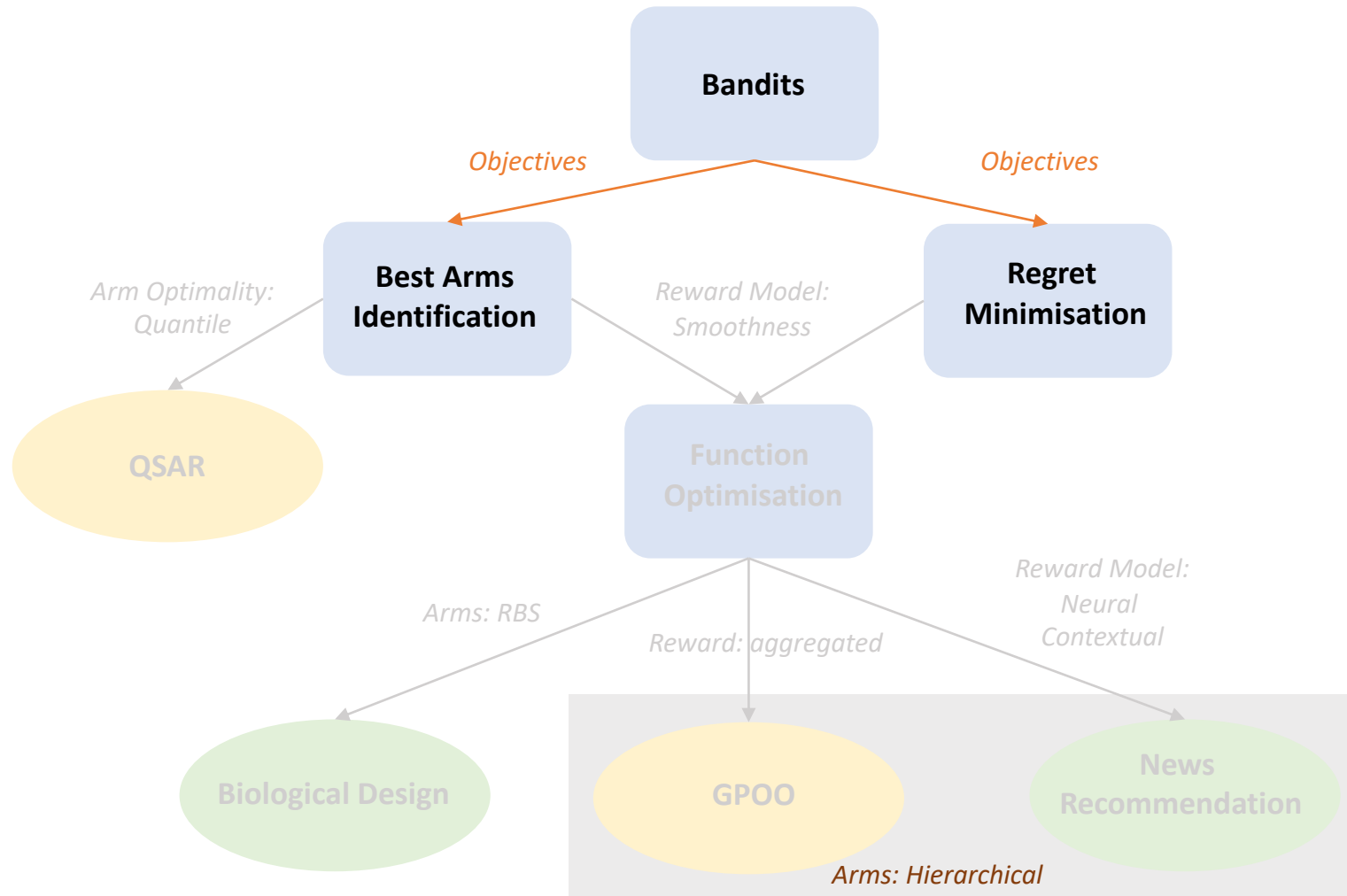
# Outline

# Objectives

- **Best arms identification (BAI):** have separate exploration stage

  identify the best m items when the exploration stage ends, e.g. with fixed budget N

  Simple regret $\quad r_N = \sum_{i=1}^{m}(\mu_{o_i} - \mathbb{E}[\mu_{A_N^i}]) \quad$ where $\quad \mu_{o_1} \geq \cdots \geq \mu_{o_K}$
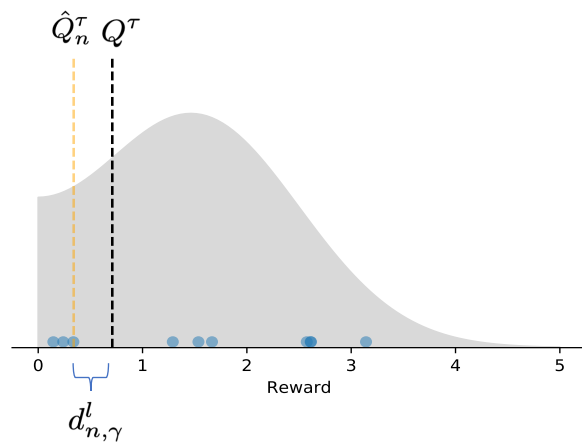
  Probability of error $\qquad e_N := \mathbb{P}\left(\mathcal{S}_m^N \neq \mathcal{S}_m^*\right)$

- **Regret minimization:** have no separate exploration stage

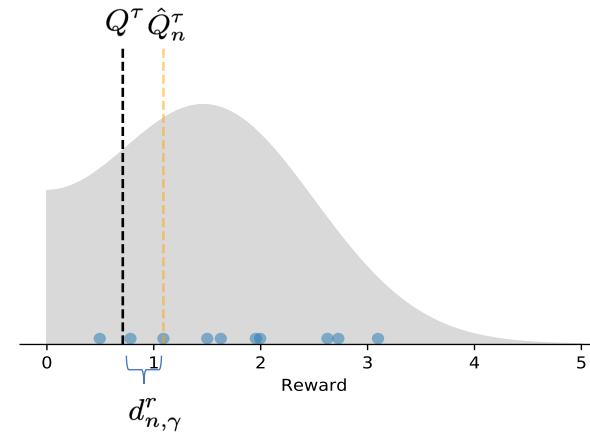  recommend items sequentially to users with the goal of minimising cumulative regret

  Cumulative regret $\qquad R_N = N\mu_* - \mathbb{E}\left[\sum_{t=1}^{N} X_{A_t}\right]$

# What matters and how to achieve? – In Theory

- Regret bounds (in expectation, or in high probability)
  - Sublinear regret e.g. $\lim_{N\to\infty} \frac{R_N}{N} = 0$
  - Probability or error e.g. decrease exponentially wrt budget, $O(\exp(-N))$
- How: utilise concentration inequalities:

$$\mathbb{P}\left(Q^\tau - \hat{Q}_n^\tau \geq d_{n,\gamma}^l\right) \leq \exp(-\gamma)$$

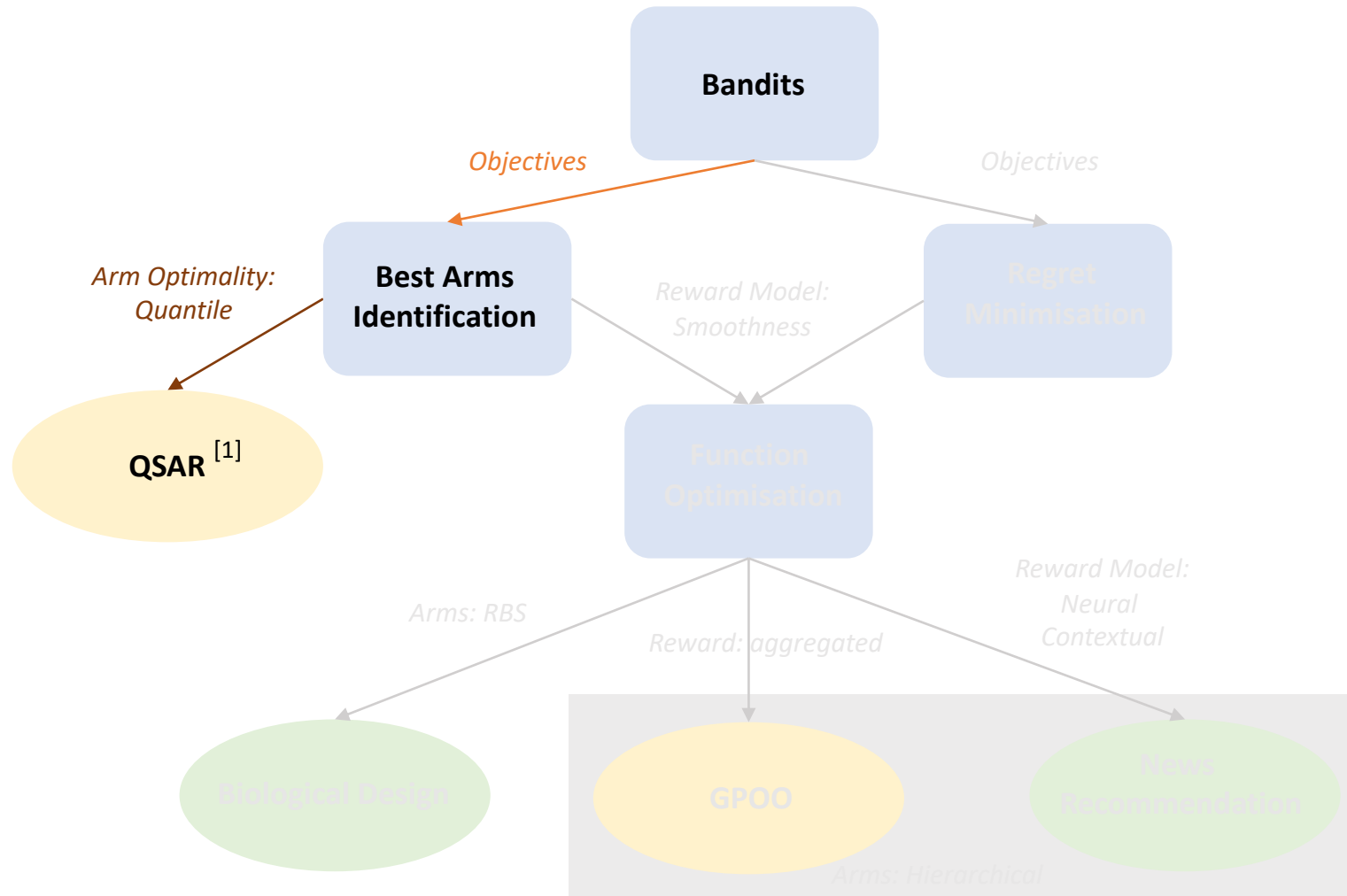$$\mathbb{P}\left(\hat{Q}_n^\tau - Q^\tau \geq d_{n,\gamma}^r\right) \leq \exp(-\gamma)$$

# What matters and how to achieve? – In Practice

- Performance
  - regret, probability of error
  - The improvement over random/baselines

- How:
  - Model assumption – fits the real applications
  - Quality of predictions of labels and uncertainty – representations, neural, Bayesian methods
  - Large design space – hierarchical design

# Outline



[1] Quantile Bandits for Best Arms Identification. **Mengyan Zhang**, Cheng Soon Ong. International Conference on Machine Learning 2021.

# Best Arms Identification with Fixed Budget

# BAI with Quantiles

# Applications: vaccine allocation

- Identify optimal strategies (highest **median** reward) for vaccine allocation
- **Arm**: vaccine allocation strategy (Allocate 100 vaccine doses to 5 age groups -- all combinations as arms )
- **Reward**: proportion of individuals that did not experience symptomatic infection



Vaccine Reward Violin Plot

# Contributions on Quantile BAI

**New Algorithm**: Quantile-based Successive Accepts and Rejects (Q-SAR)
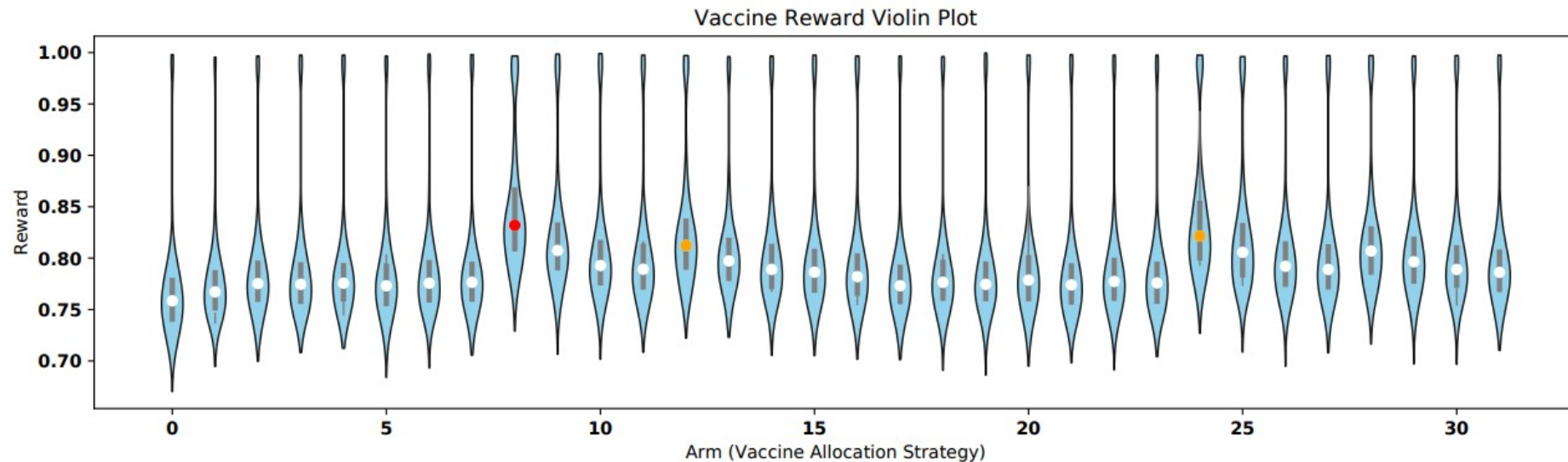
*Extends Bubuck et al. 2013:* Mean-based SAR



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

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

**New Concentration inequalities**

$$\mathbb{P}\left(Q^\tau - \hat{Q}_n^\tau \geq d_{n,\gamma}^l\right) \leq \exp(-\gamma)$$
$$\mathbb{P}\left(\hat{Q}_n^\tau - Q^\tau \geq d_{n,\gamma}^r\right) \leq \exp(-\gamma)$$

**Probability of error**

$$e_N := \mathbb{P}\left(\mathcal{S}_m^N \neq \mathcal{S}_m^*\right) \leq 2K^2 \exp\left(-\frac{N-K}{\overline{\log(K)}H^\tau}\right)$$

**Experiments** on vaccine allocation



Vaccine with 0.5-quantile (m=3)

Bubeck, S.; Wang, T.; and Viswanathan, N. Multiple identifications in multi-armed bandits. ICML 2013.

# Outline



Bandits

*Objectives*          *Objectives*

Best Arms Identification

*Reward Model: Smoothness*

Regret Minimisation

*Reward Summary: Quantile*

QSAR

Function Optimisation

$$f : \mathcal{X} \to \mathbb{R}$$

*Arms: RBS*

*Reward: aggregated*

*Reward Model: Neural Contextual*

Bio Design

GPOO

News Recommendation

*Arms: Hierarchical*

# Reward Smoothness – e.g. Gaussian Process

$$f(\mathbf{x}) \sim \mathcal{GP}\left(\mu(\mathbf{x}), k\left(\mathbf{x}, \mathbf{x}'\right)\right)$$

$$\mu(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] \quad \text{and} \quad k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - \mu(\mathbf{x}))(f(\mathbf{x}') - \mu(\mathbf{x}'))].$$



**Regression** is used to find a function (line) that represents a set of data points as closely as possible

y = 0

A **Gaussian process** is a probabilistic method that gives a confidence (shaded) for the predicted function

https://distill.pub/2019/visual-exploration-gaussian-processes/

# Acquisition Function: e.g. GP-UCB

Gaussian Process Upper Confidence Bound (GP-UCB) [1]
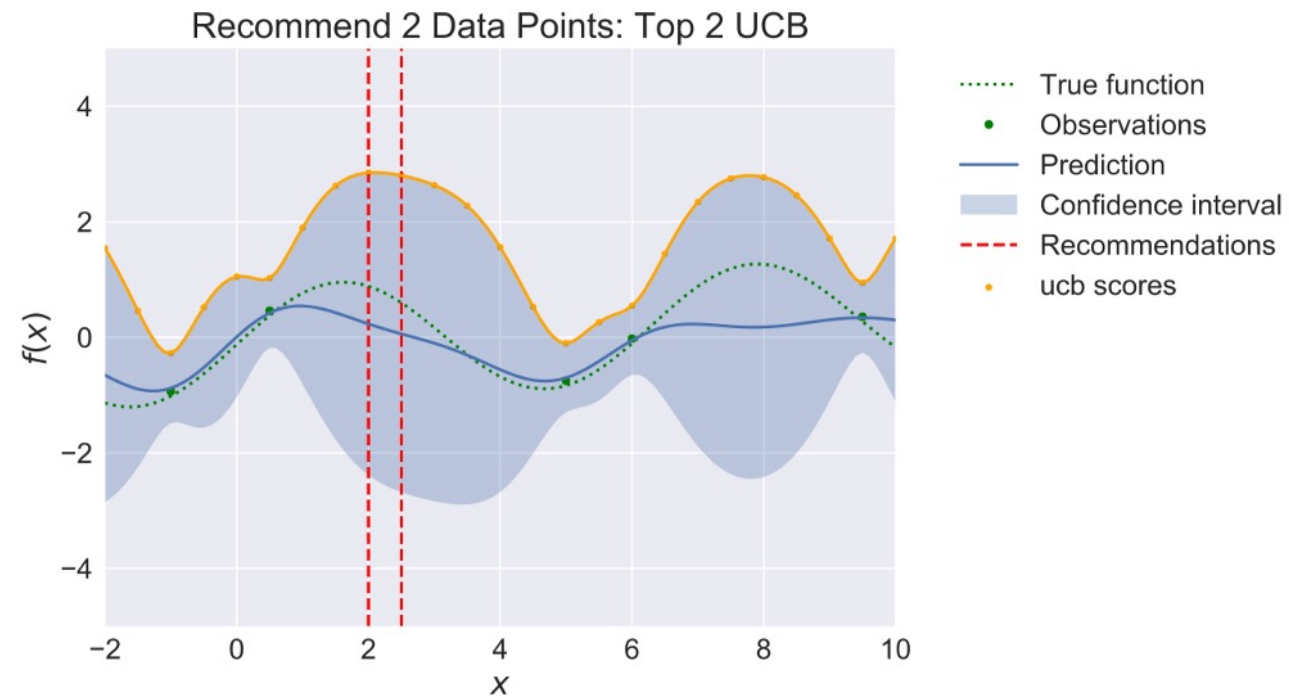
Posterior mean: exploitation

Posterior standard deviation: exploration

$$\text{argmax}_{\mathbf{x}_i \in \mathcal{K}} \left( \mu_t(\mathbf{x}_i) + \beta_t \sigma_{t-1}(\mathbf{x}_i) \right)$$

Balancing term



Recommend 2 Data Points: Top 2 UCB

Legend:
- ······ True function
- • Observations
- —— Prediction
- Confidence interval
- - - - Recommendations
- • ucb scores

[1] Srinivas, N.; Krause, A.; Kakade, S. M.; and Seeger, M., Gaussian process optimization in the bandit setting: No regret and experimental design. ICML 2009.

# Outline



[1] Machine learning guided batched design of a bacterial Ribosome Binding Site.

**Mengyan Zhang**, Maciej Bartosz Holowko, Huw Hayman Zumpe, Cheng Soon Ong. ACS Synthetic Biology Journal 2022.

[2] Opportunities and Challenges in Designing Genomic Sequences. **Mengyan Zhang**, Cheng Soon Ong. ICML Workshop on Computational Biology 2021.

# Bandits for Synthetic Biology

With fixed budget (450), design **Ribosome Binding Site (RBS)** sequences in batches (4)
(300 for bandit groups)

Optimize the protein expression level (translation initiation rate)
Identify the RBS sequences with highest possible protein expression level

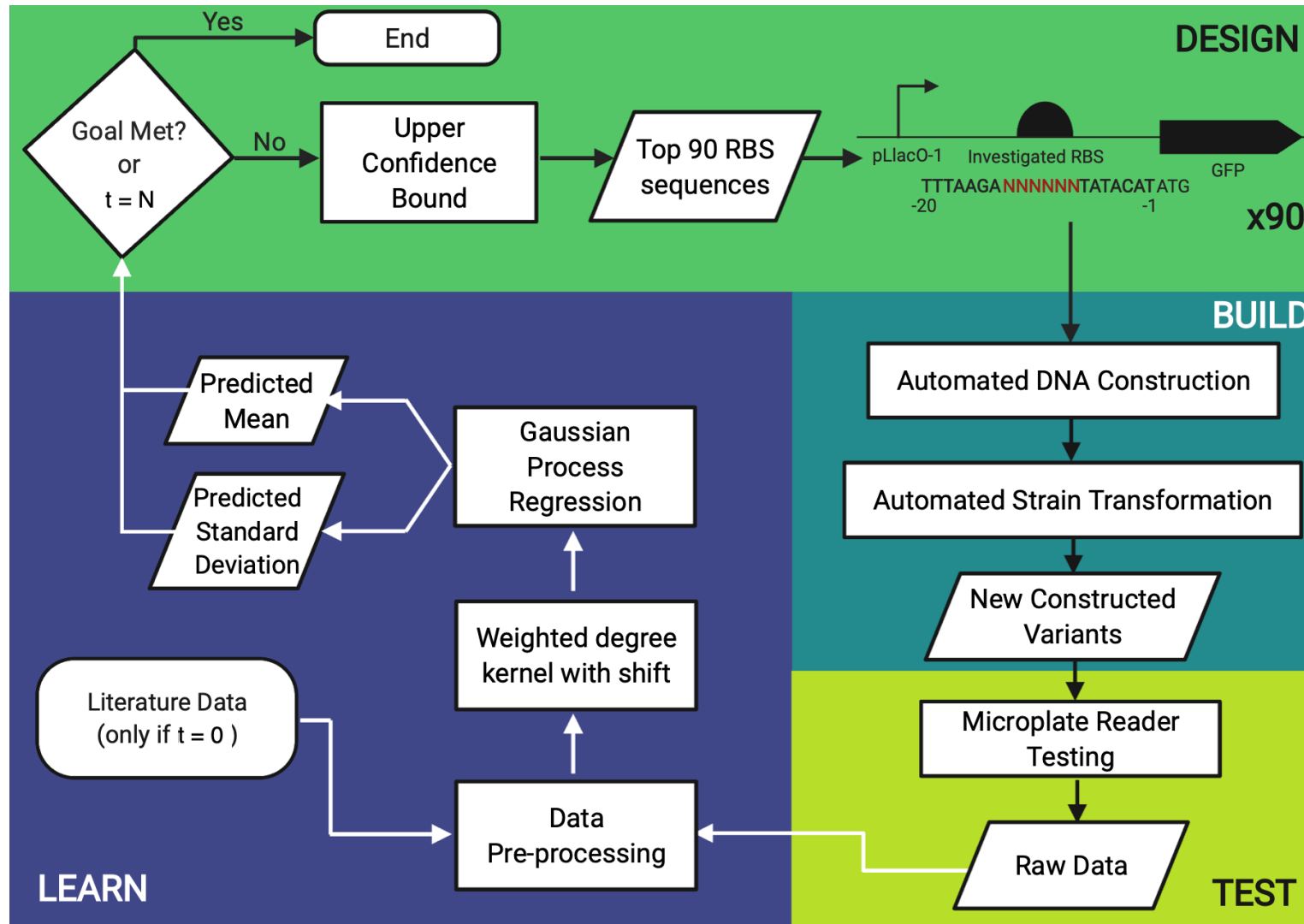| Arm: RBS sequence | Reward: Normalized* Translation Initiation Rate |
|---|---|
| TTTAAGAGTTATATATACAT | 1.58 |
| TTTAAGAATATGCTATACAT | 1.42 |
| TTTAAGACTCGGATATACAT | 0.14 |
| TTTAAGAGTTTTTTATACAT | 2.88 |

Green Fluorescent Protein (GFP)

- Design space: 4096 sequences

* zero mean and unit variance normalization $z = \frac{x-\mu}{\sigma}$
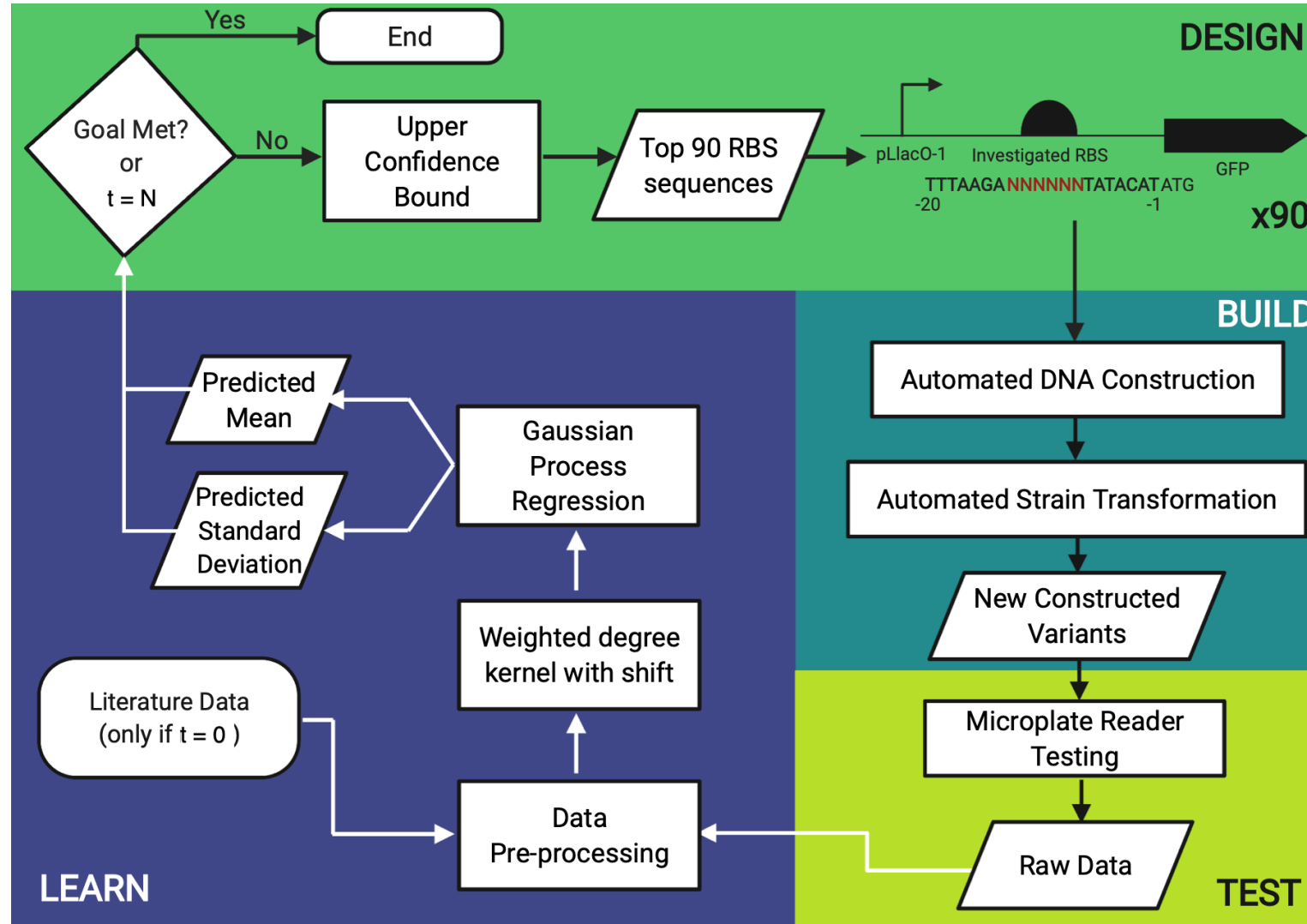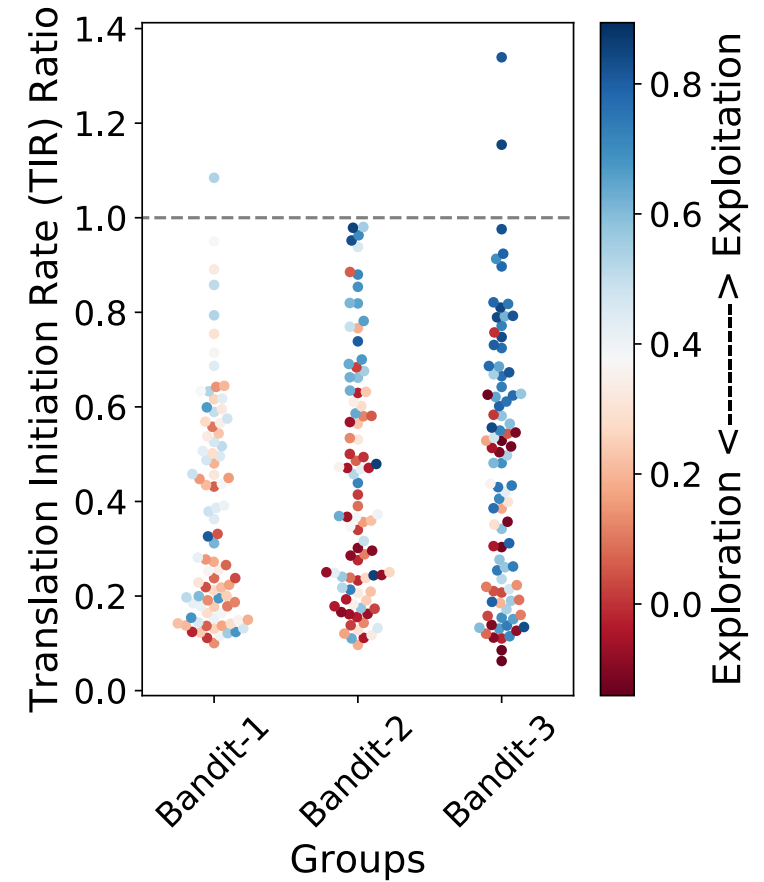
# Design-Build-Test-Learn (DBTL) Cycle



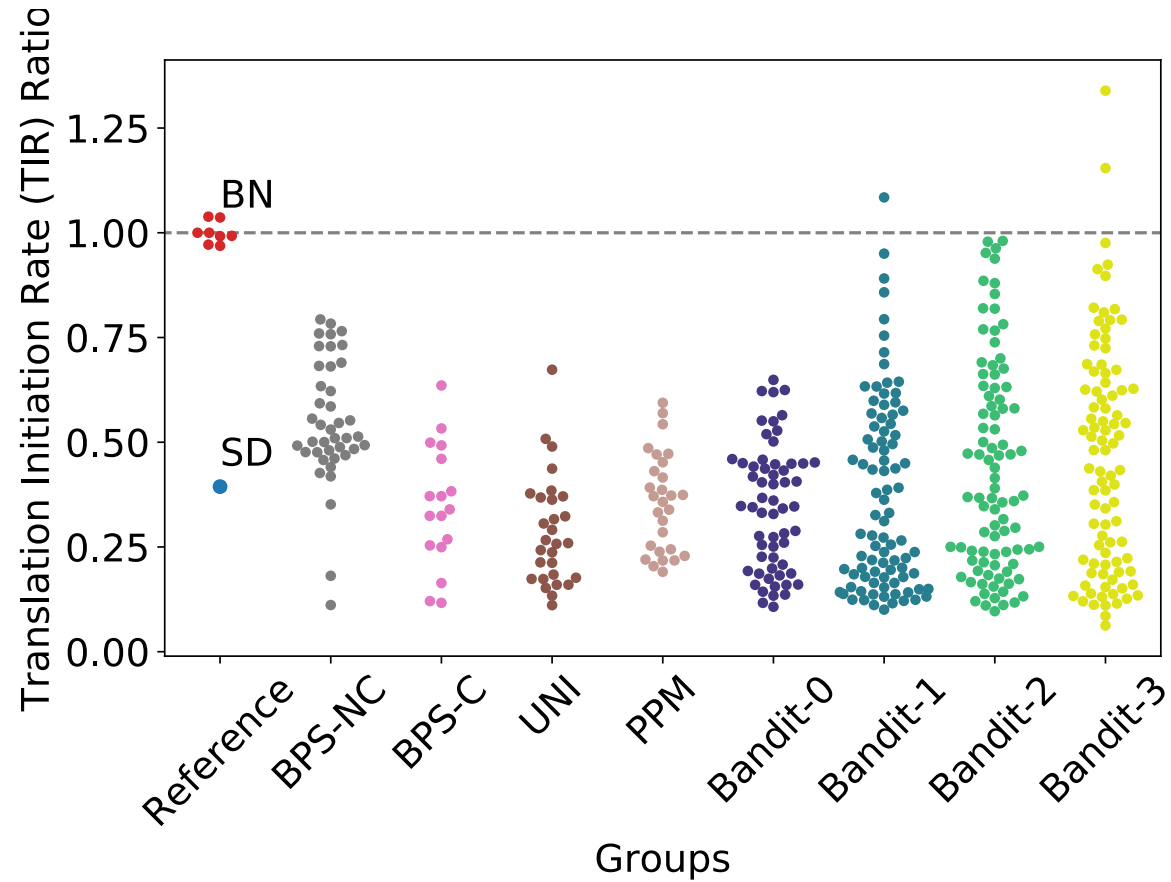CSIRO BioFoundry Lab

# Design-Build-Test-Learn (DBTL) Cycle

# Results: swarmplot

# Lessons learned and future opportunities

- ML (Bandits) guided DBTL cycle -- increase expression of our target protein by up to <span style="color:red">35%</span>, compared to a strong benchmark RBS

- Generalisation of our workflow: target on larger design space, more complicated genetic elements, e.g. promoters

# Outline



[1] Gaussian Process Bandits with Aggregated Feedback. **Mengyan Zhang**, Russell Tsuchida, Cheng Soon Ong. AAAI 2022.

# Computational cost for continuous space

$$\text{argmax}_{\mathbf{x}_i \in \mathcal{K}} \left( \mu_t(\mathbf{x}_i) + \beta_t \sigma_{t-1}(\mathbf{x}_i) \right)$$



Recommend 2 Data Points: Top 2 UCB

Legend:
- ······ True function
- ● Observations
- —— Prediction
- Confidence interval
- – – – Recommendations
- · ucb scores

# Problem Setting



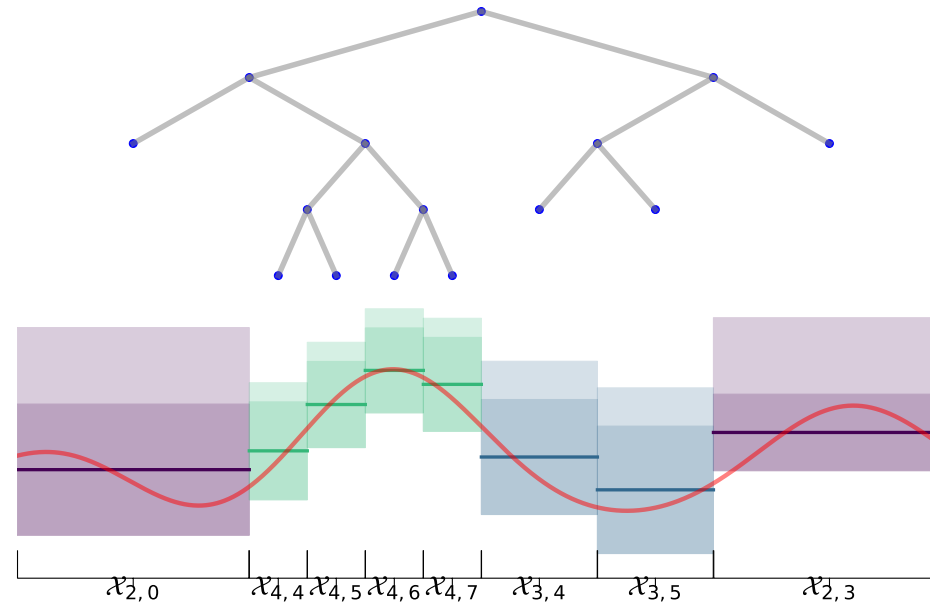- Arms: a leaf node, corresponding to a subset of continuous $[0,1]^d$
- Rewards: sampled from GP, only average reward for a node

$$r_t = \bar{F}(X_{h_t, i_t}) + \epsilon_t, \quad \bar{F}(X_{h_t, i_t}) := \frac{\sum_{\boldsymbol{x} \in \mathcal{C}_{h_t, i_t}} f(\boldsymbol{x})}{|\mathcal{C}_{h_t, i_t}|}$$

with   **Representative points** $\mathcal{C}_{h,i} = \{\boldsymbol{x}_{h,i^s}\}_{1 \leq s \leq S}$, where $\boldsymbol{x}_{h,i^s} \in \mathcal{X}_{h,i}$.

# Why Aggregated Feedback?

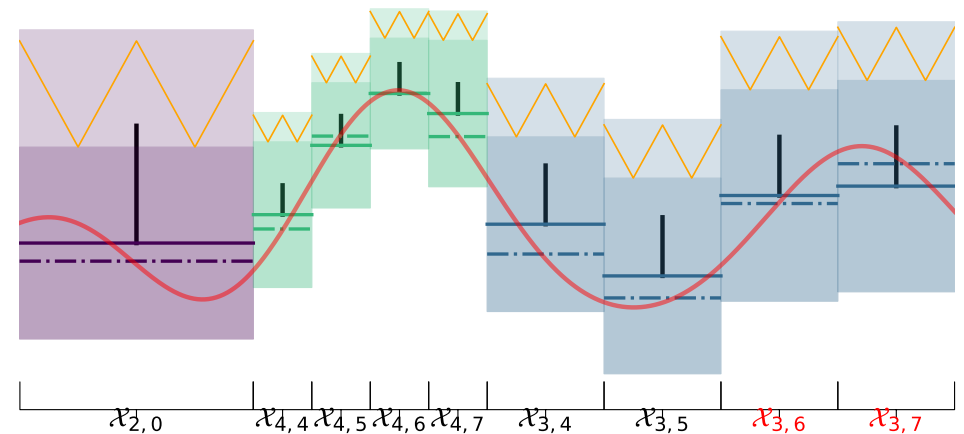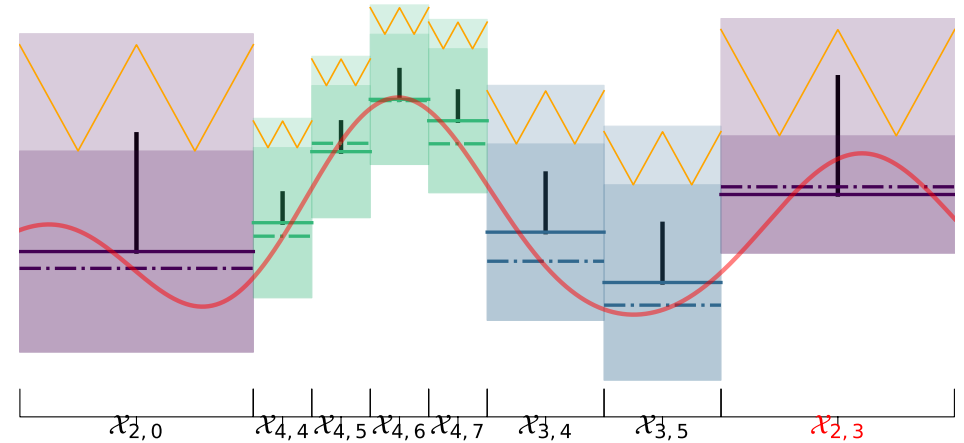| Application | Arm | Reward | Goal: to design a policy such that... | Why Aggregated? |
|---|---|---|---|---|
| **DNA Design** | DNA sequences | average protein expression level in a mixed culture | identify DNA sequences with the highest protein expression level with a given budget | expensive; search space is large |
| **Census Querying** | Respondent | average age of respondents inside queried area | identify the region with the highest average age with a fixed amount of querying | privacy concerns |
| **Radio Telescope** | spatial-frequency coordinates of objects in the sky | average radio wave energy from the queried area | identify the region with the highest average radio energy with a fixed amount of querying | hardware constraint |

# Gaussian Process Optimistic Optimisation (GPOO)
## How to choose node and when to split?

**Assumption:** **Decreasing Diameters:** $\sup_{\boldsymbol{x} \in \mathcal{X}_{h,i}} L\ell(\boldsymbol{x}_{h,i}, \boldsymbol{x}) \leq \delta(h)$ some decreasing sequence $\delta(h) > 0$.

- **Select leaf node with largest b-value:**

$$\boxed{b_{h,i}(t)} = \text{posterior mean} + \text{confidence interval} + \text{diameter } \delta(h_t)$$
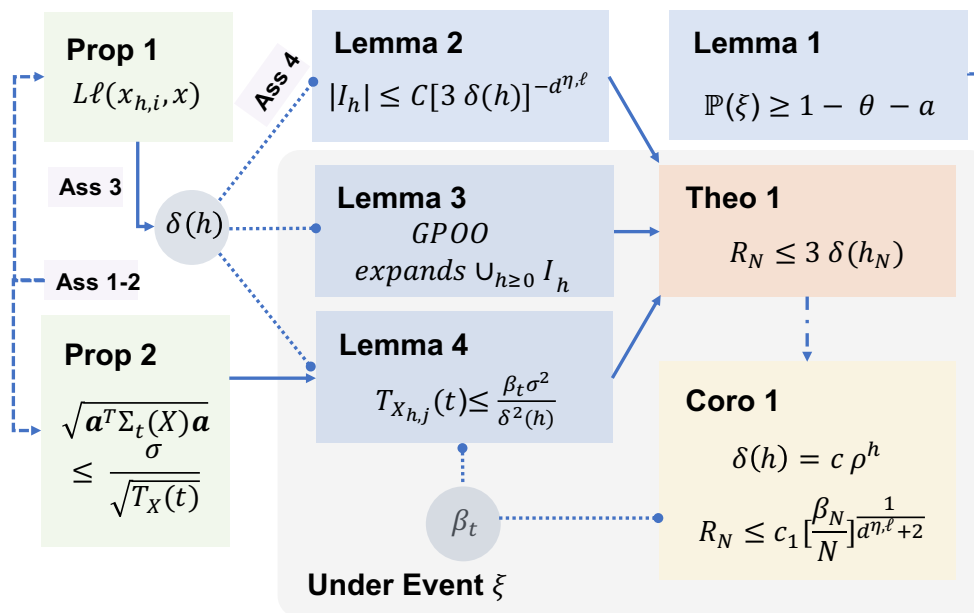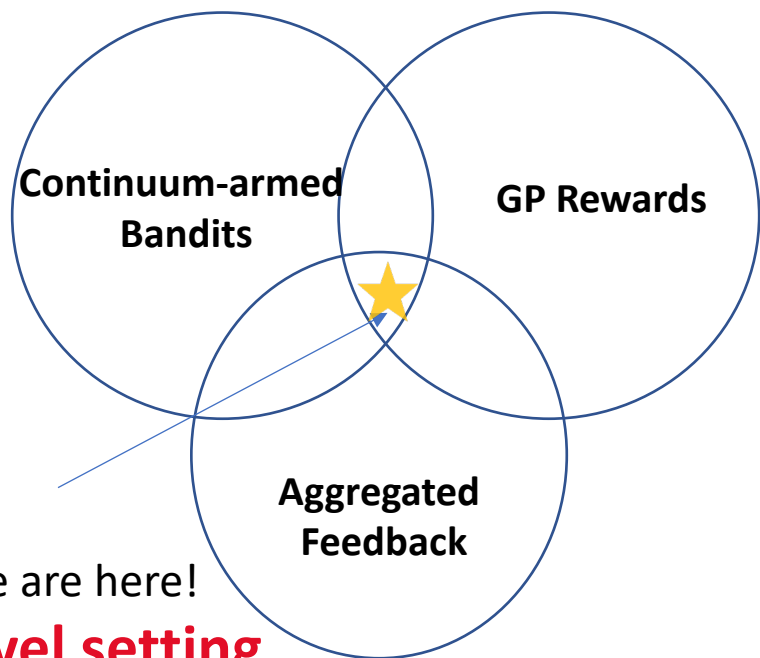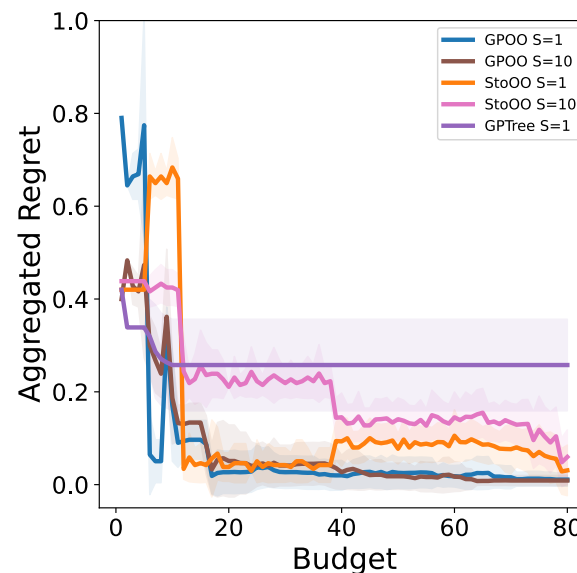
- **Expand:** if $\delta(h_t) >$ confidence interval

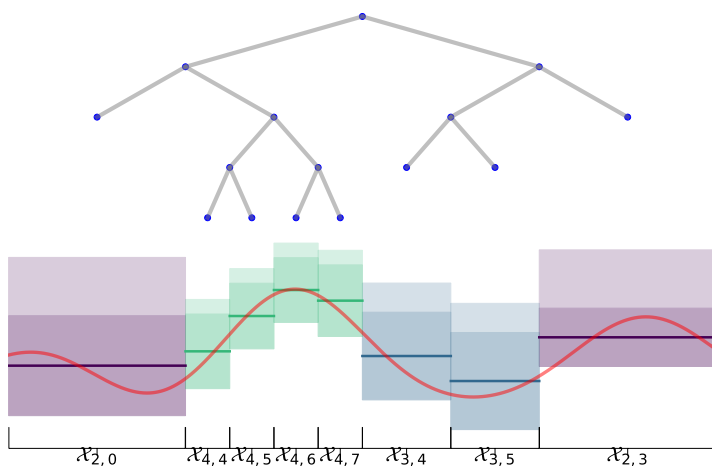# Contributions of GPOO

Upper bound on
(aggregated) regret



**Continuum-armed Bandits**

**GP Rewards**

**Aggregated Feedback**

We are here!

**A novel setting**

**New Algorithm**: **GPOO**

**Prop 1**
$L\ell(x_{h,i}, x)$

**Ass 4**

**Lemma 2**
$|I_h| \leq C[3\,\delta(h)]^{-d^{\eta,\ell}}$

**Lemma 1**
$\mathbb{P}(\xi) \geq 1 - \theta - a$

**Ass 3**

$\delta(h)$

**Ass 1-2**

**Lemma 3**
$GPOO$
$expands \cup_{h \geq 0} I_h$

**Theo 1**
$R_N \leq 3\,\delta(h_N)$

**Prop 2**
$\sqrt{\boldsymbol{a}^T \Sigma_t(X)\boldsymbol{a}} \leq \dfrac{\sigma}{\sqrt{T_X(t)}}$

**Lemma 4**
$T_{X_{h,j}}(t) \leq \dfrac{\beta_t \sigma^2}{\delta^2(h)}$

$\beta_t$

**Coro 1**
$\delta(h) = c\,\rho^h$
$R_N \leq c_1 [\dfrac{\beta_N}{N}]^{\frac{1}{d^{\eta,\ell}+2}}$

**Under Event $\xi$**

**Simulation results:**
Outperform baselines



$x_{2,0}$  $x_{4,4}$ $x_{4,5}$ $x_{4,6}$ $x_{4,7}$  $x_{3,4}$  $x_{3,5}$  $x_{2,3}$

# Lessons learned
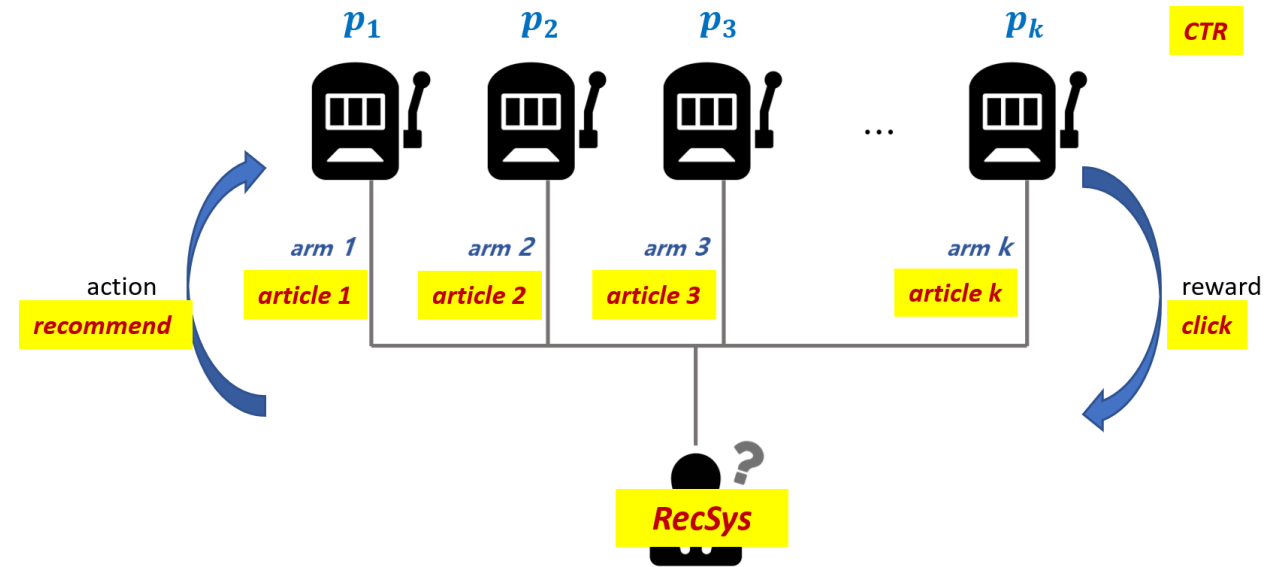
- Hierarchical design of arms -> computational efficient for large/continuous design space

- Average reward feedback -> same regret upper bound rate as single arm feedback (in our setting)

# Outline

[1] Two-Stage Neural Contextual Bandits for Personalised News Recommendation.
**Mengyan Zhang**, Thanh Nguyen-Tang, Fangzhao Wu, Zhenyu He, Xing Xie, Cheng Soon Ong. Under Review 2022
(work conducted during the internship in Microsoft Research Asia)

# Contextual Bandits:

## in recommender system



In each round t $\epsilon$ {1, .., N},
*Given a user*

1. an <u>agent</u> selects an <u>arm</u> $A_t = i \, \epsilon \, \{1, .., K\}$ according <u>policy</u> $\pi$
*Recommender system*      *Item (e.g. news)*                              *Recommendation strategy*

2. then receive a <u>reward</u> $X_{A_t, T_{A_t}(t)}$ sampled from unknown reward distribution $F_{A_t}$
*click; non-click*

3. update estimations over reward distributions based on historical observations

Two-Stage Deep Contextual Bandits News Recommendation

**Two-Stage Deep Contextual Bandits News Recommendation**

Hundreds of Candidate Topics — Stage One — Thousands of Candidate News — Stage Two — Recommended News

Millions of News Items

Topic Recommendation Policy

Item Recommendation Policy

**Neural Contextual Bandits Policies**
Generalised Additive Linear UCB, Generalised Bilinear UCB

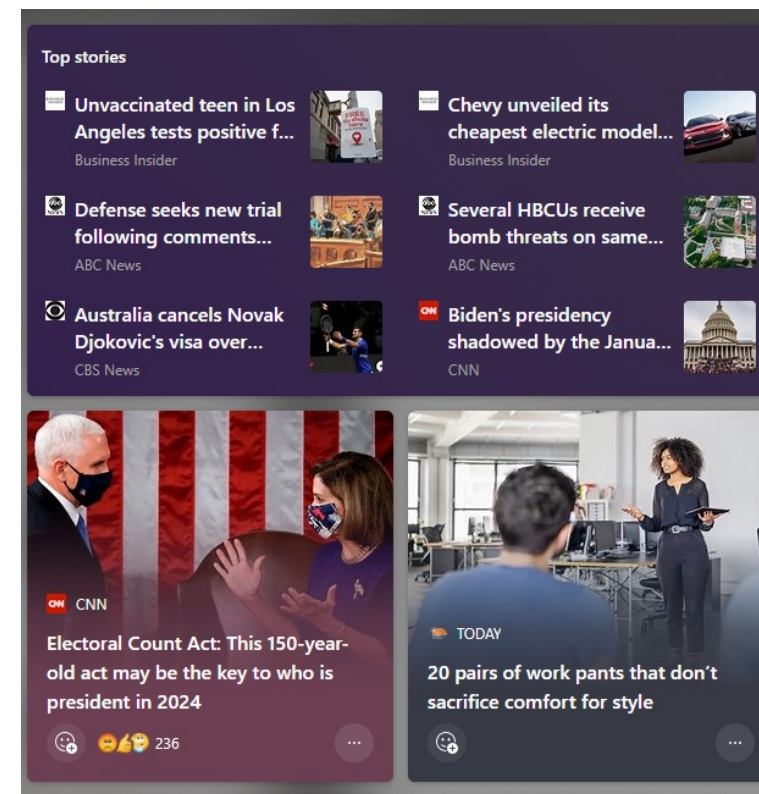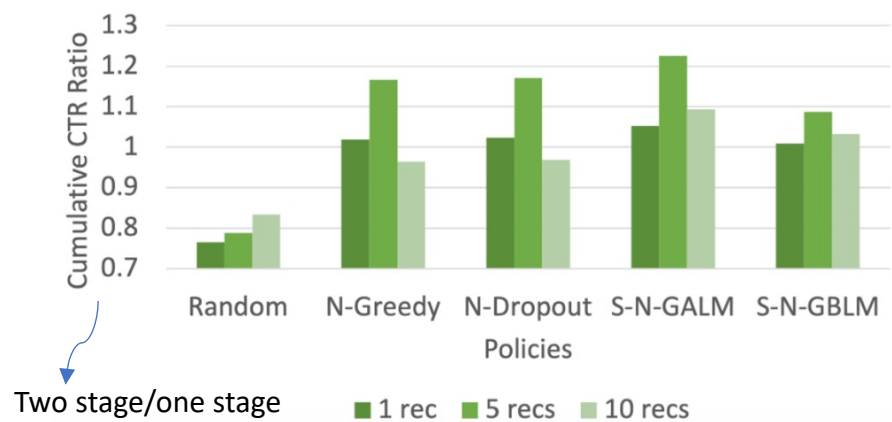S-N-GALM-UCB                S-N-GBLM-UCB

# Large-scale experiments

| Users | News | Topics | Samples |
|---|---|---|---|
| 1,000,000 | 161,013 | 285 | 24,155,470 |

Click-through-rate (CTR) $\quad \mathrm{CTR}_t^\tau = \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}\{y_t^\tau = 1\}$

Cumulative CTR $\quad \sum_{t=1}^{N} \mathrm{CTR}_t^\tau$



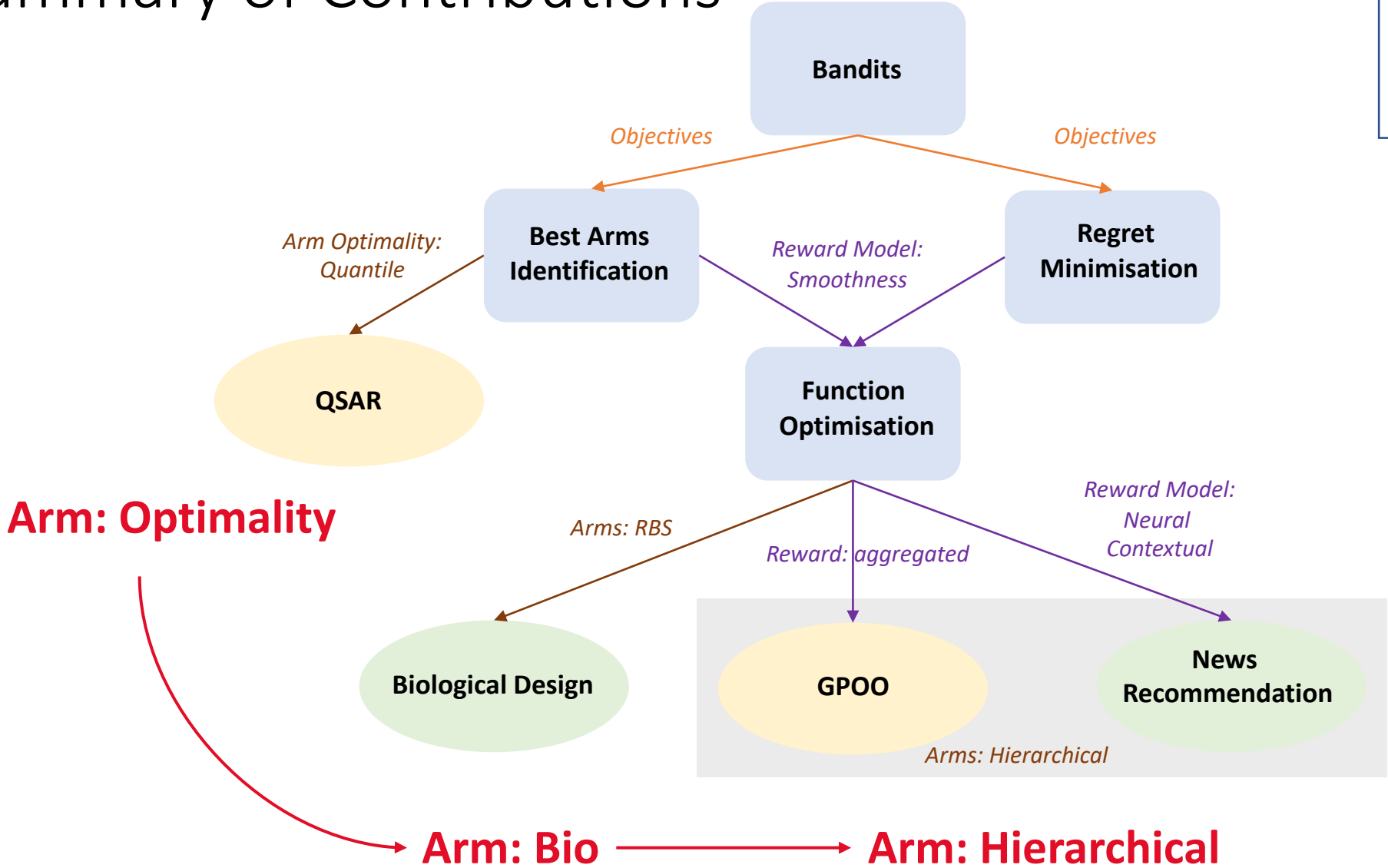Two stage/one stage


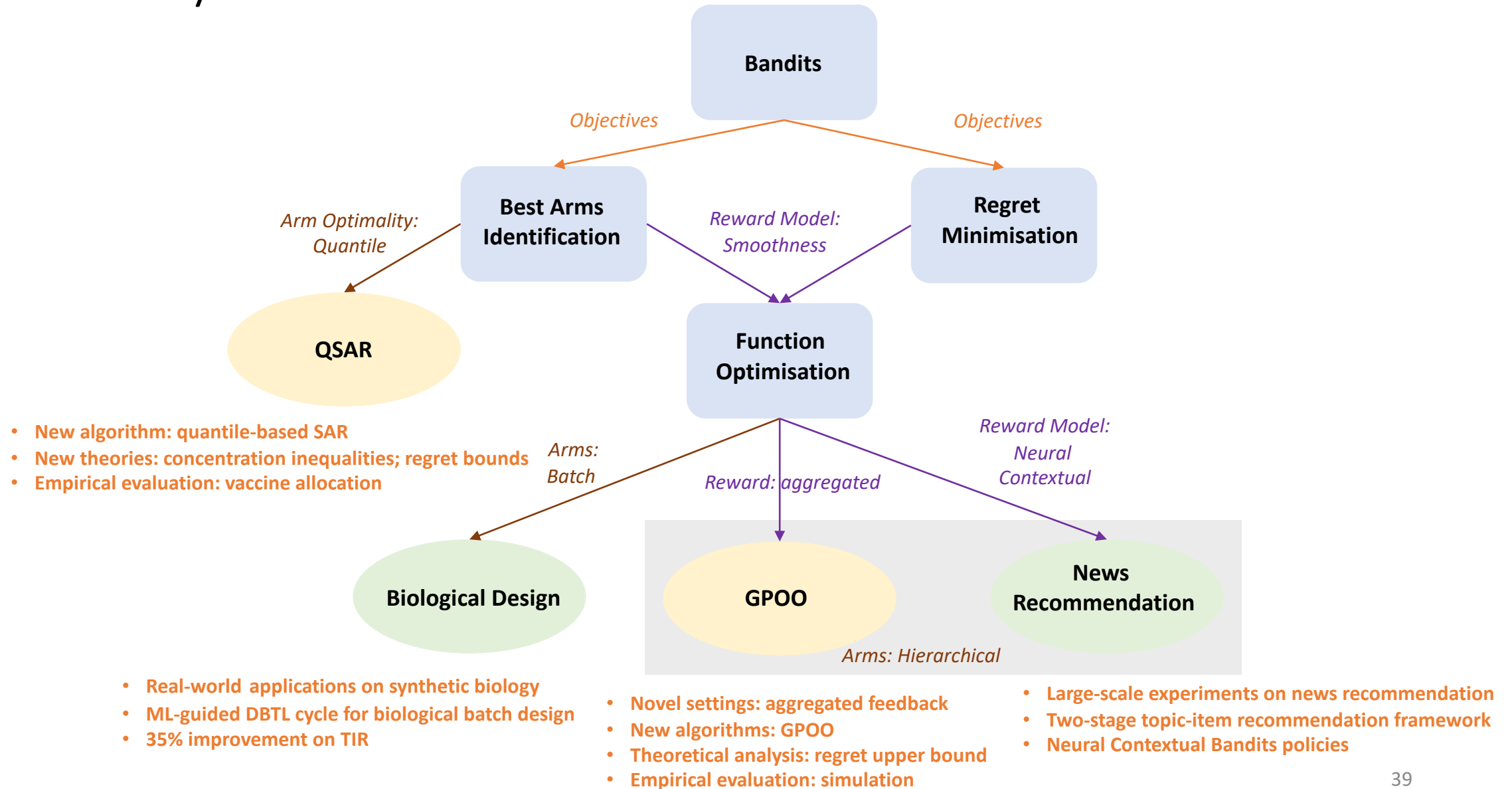
36

[1] https://msnews.github.io/

# Lessons learned

- Two-stage recommendation is useful: computational efficient

- Two-tower neural representation improves the performance

- Gaps between theory and practice: neural bandits, off-policy evaluation

# Summary of Contributions

# Summary of Contributions

# Publications

- Quantile Bandits for Best Arms Identification.

  **Mengyan Zhang**, Cheng Soon Ong. International Conference on Machine Learning 2021.

- Machine learning guided batched design of a bacterial Ribosome Binding Site.

  **Mengyan Zhang**, Maciej Bartosz Holowko, Huw Hayman Zumpe, Cheng Soon Ong. ACS Synthetic Biology Journal 2022.

- Opportunities and Challenges in Designing Genomic Sequences.

  **Mengyan Zhang**, Cheng Soon Ong. ICML Workshop on Computational Biology 2021.
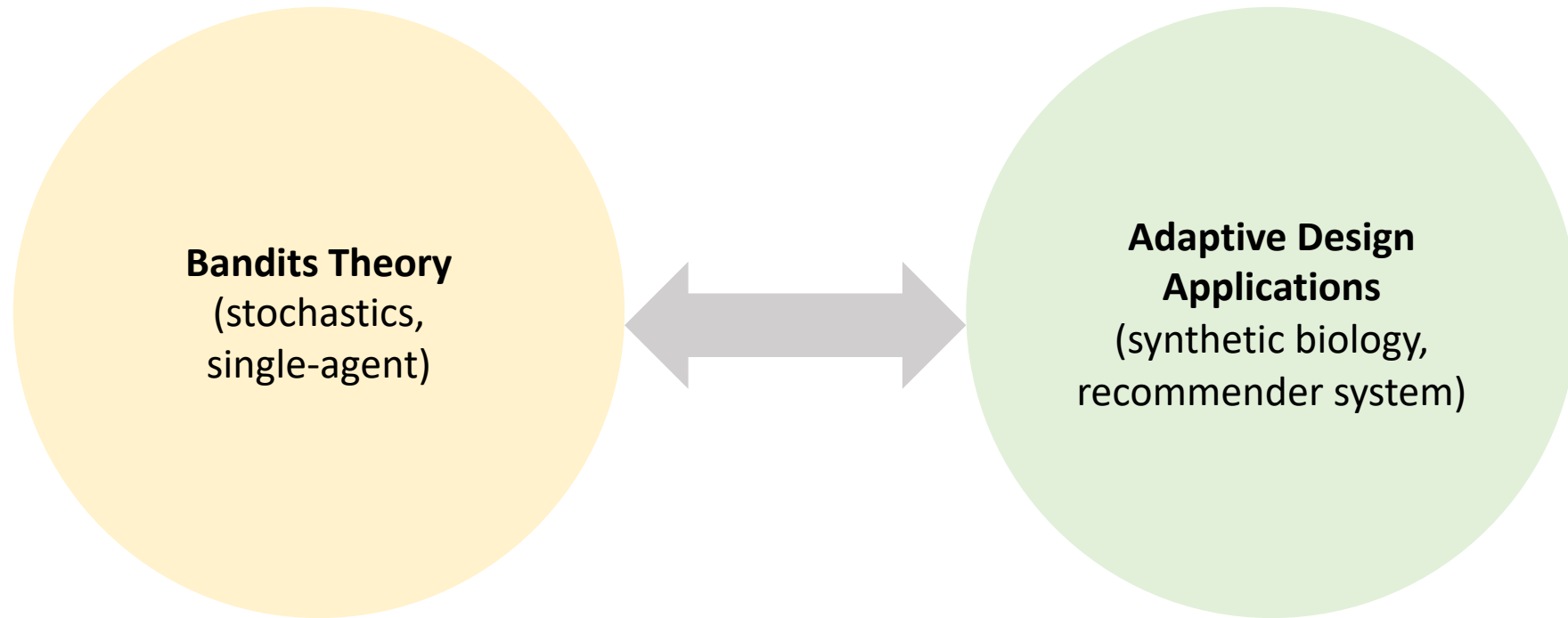
- Gaussian Process Bandits with Aggregated Feedback.

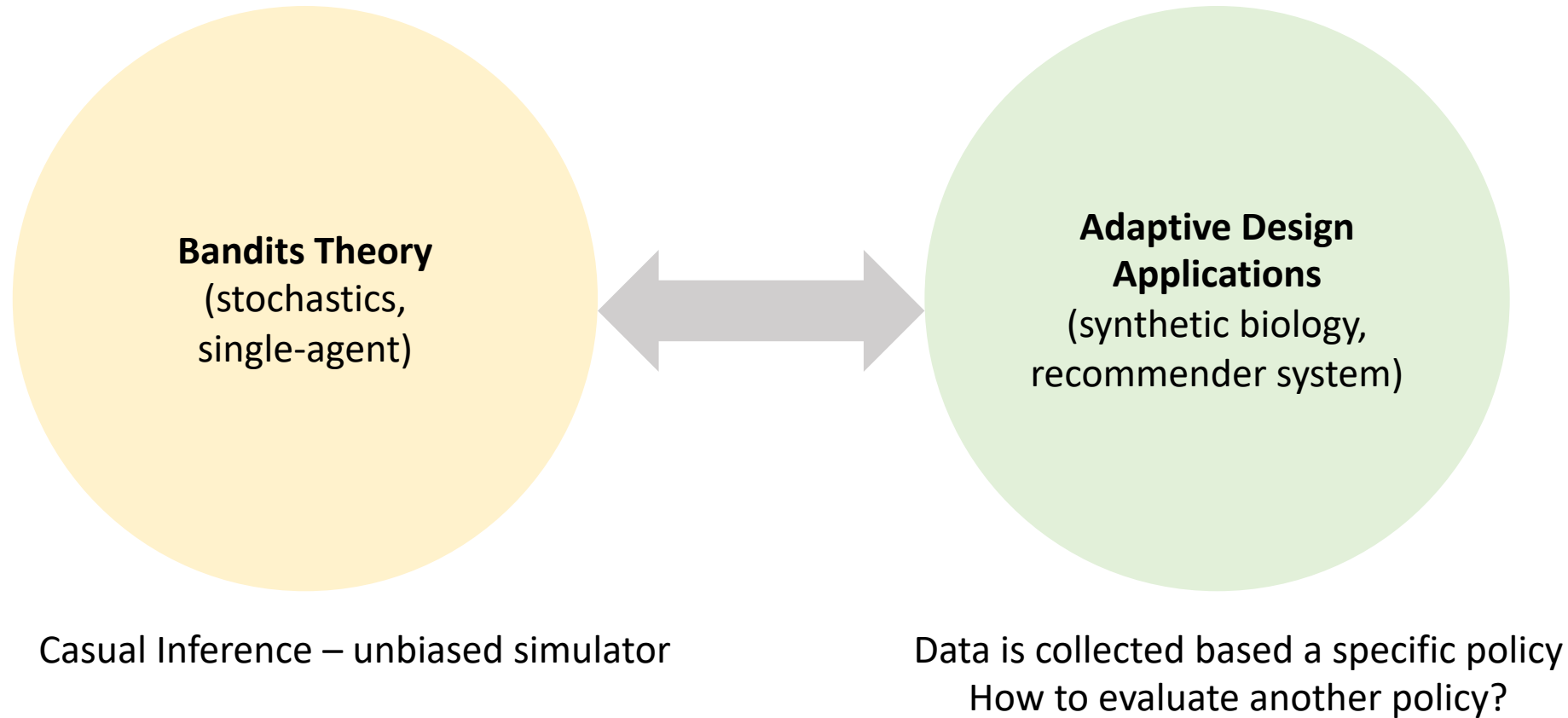  **Mengyan Zhang**, Russell Tsuchida, Cheng Soon Ong. AAAI 2022.

- Two-Stage Neural Contextual Bandits for Personalised News Recommendation.

  **Mengyan Zhang**, Thanh Nguyen-Tang, Fangzhao Wu, Zhenyu He, Xing Xie, Cheng Soon Ong. Under Review 2022.

# Future Work – Gaps between Theory and Practice

**Bandits Theory**
(stochastics,
single-agent)

**Adaptive Design
Applications**
(synthetic biology,
recommender system)

# Future Work – e.g. Off Policy Evaluation

**Bandits Theory**
(stochastics,
single-agent)

**Adaptive Design
Applications**
(synthetic biology,
recommender system)

Casual Inference – unbiased simulator

Data is collected based a specific policy
How to evaluate another policy?

# Future Work – Scientific Discovery

**Bandits Theory**
(stochastics,
single-agent)

**Adaptive Design
Applications**
(synthetic biology,
recommender system)

Call to action:

"scientific revolutions occur when there is cross pollination of ideas"

---- Thomas Kuhn

# Acknowledgements

- **Supervisory panel**:



Cheng Soon Ong        Lexing Xie        Eduardo Eyras

- **Mentors and Collaborators**:

  Sebastien Bubeck (Microsoft Research), Xing Xie (Microsoft Research Asian), Fangzhao Wu (Microsoft Research Asian), Russell Tsuchida (CSIRO), Maciej Bartosz Holowko (CSIRO), Thanh Nguyen-Tang (Deakin University), Huw Hayman Zumpe (CSIRO), Zhenyu He (UECSTC) and many others.

- **CMLab and ML research group folks**: for interesting discussions, useful suggestions and feedback.

- **Family and friends**: for unconditional support and care.

# Thanks for listening!

Mengyan Zhang

mengyan.zhang@anu.edu.au

https://mengyanz.github.io/