

## Introduction

$$\min_{x \in \mathbb{R}^{d_x}} F(x) := \mathbb{E}_{\xi \sim \mathbb{P}_\xi} [f(x, y^*(x; \xi); \xi)] \quad (\text{upper-level})$$

$$\text{where } y^*(x; \xi) := \operatorname{argmin}_{y \in \mathbb{R}^{d_y}} \mathbb{E}_{\eta \sim \mathbb{P}_{\eta|\xi}} [g(x, y; \eta, \xi)] \quad \forall \xi \quad (\text{lower-level})$$

- $\xi \sim \mathbb{P}_\xi$ : contextual information;  $\eta \sim \mathbb{P}_{\eta|\xi}$ : conditional distributions
- lower-level: contextual stochastic optimization

## Motivation

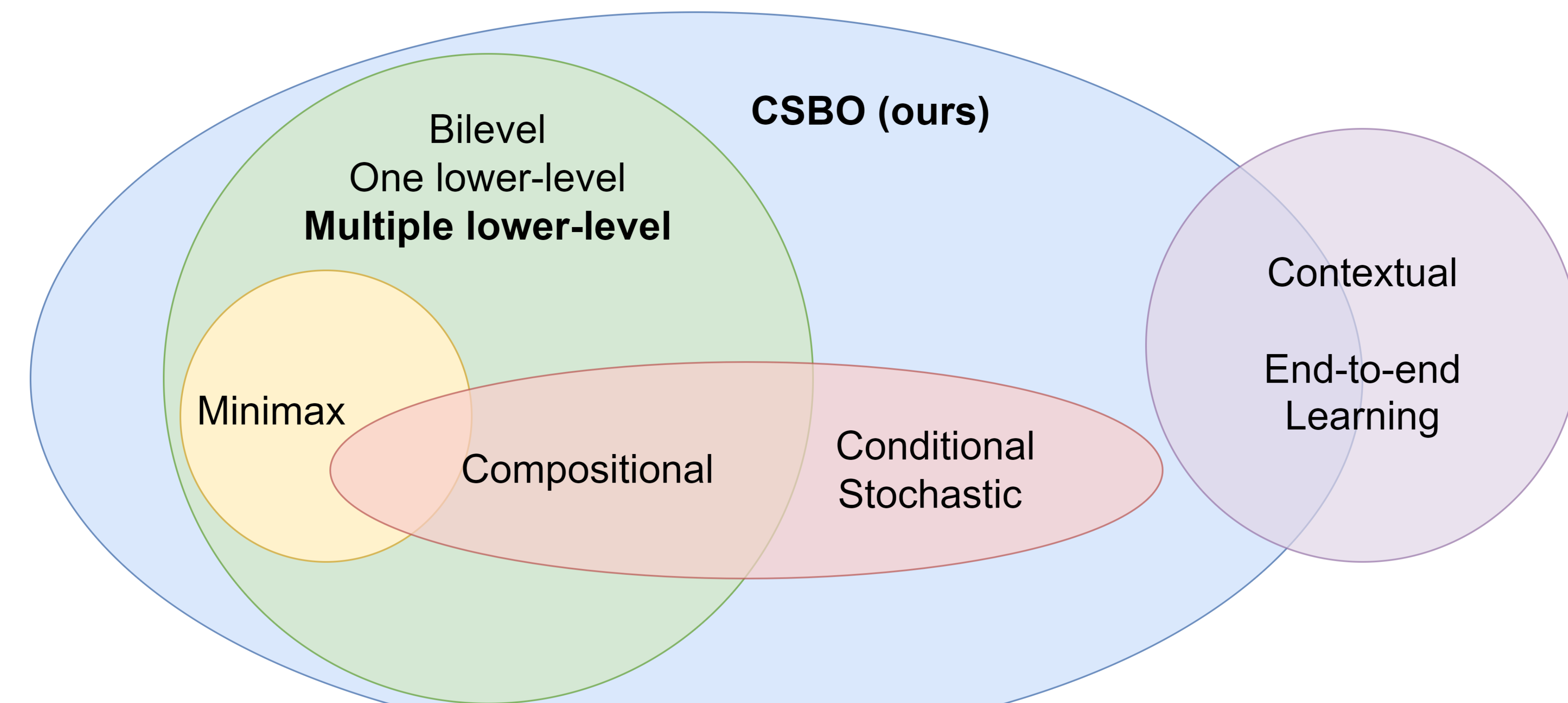
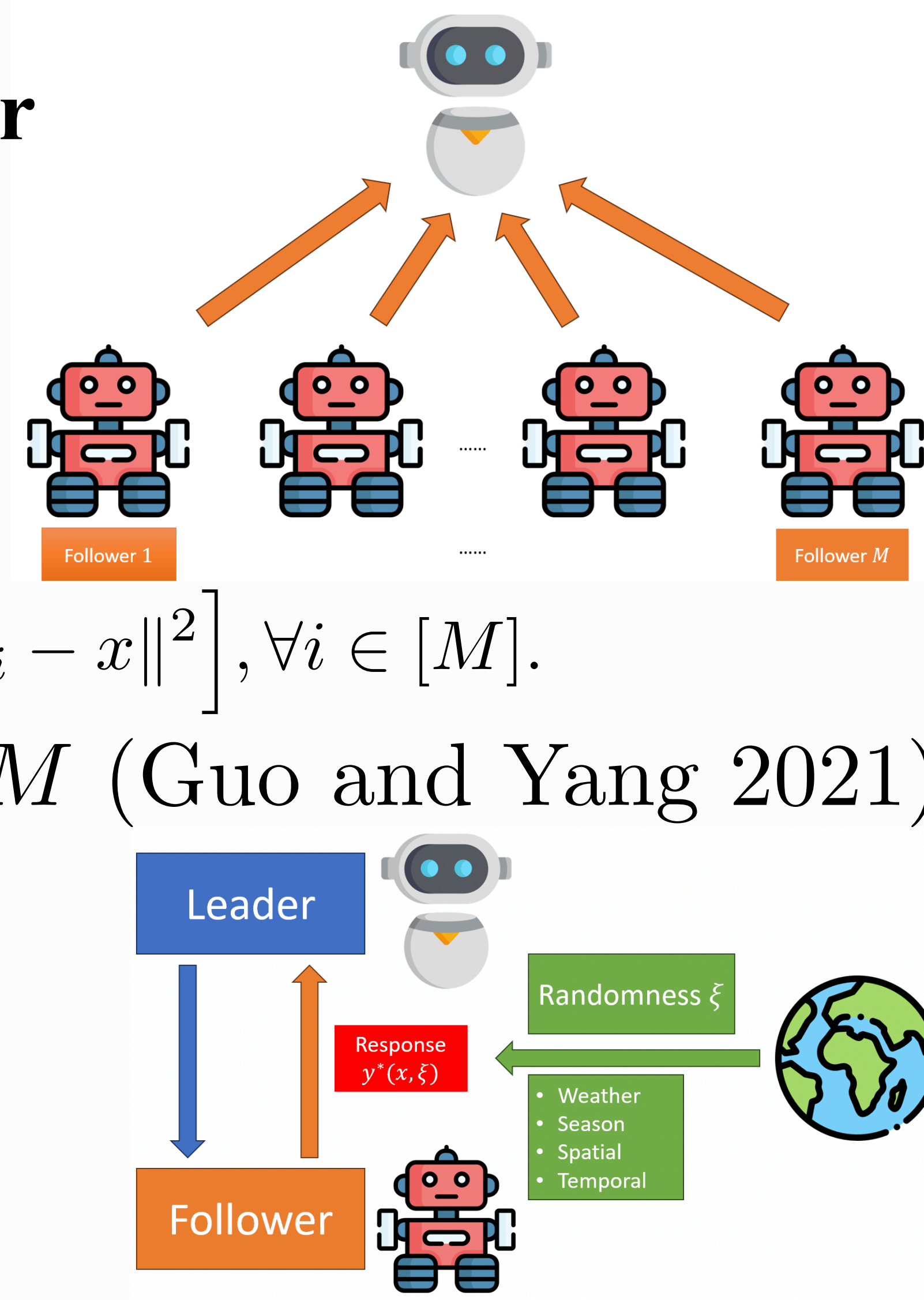
**I: Find a shared model parameter for multiple similar tasks/individuals.**

$$\min_x \frac{1}{M} \sum_{i=1}^M \mathbb{E}_{\eta \sim \rho_i} [l(y_i^*(x); \eta)]$$

where  $y_i^*(x) = \operatorname{argmin}_{y_i} \mathbb{E}_{\eta \sim \rho_i} [l(y_i; \eta) + \frac{\lambda}{2} \|y_i - x\|^2]$ ,  $\forall i \in [M]$ .

- Complexity depends linearly on  $M$  (Guo and Yang 2021)

**II: Optimal Response to Side Information.**



## Contextual Stochastic Bilevel Optimization

Yifan Hu<sup>1,2</sup>, Jie Wang<sup>3</sup>, Yao Xie<sup>3</sup>, Andreas Krause<sup>2</sup>, Daniel Kuhn<sup>1</sup>  
<sup>1</sup> EPFL, <sup>2</sup> ETH Zurich, <sup>3</sup> Gatech

## Algorithm Design

$$\nabla F(x) = \mathbb{E} \left[ \nabla_1 f(x, y^*(x; \xi); \eta, \xi) - \left( \mathbb{E}_{\eta' \sim \mathbb{P}_{\eta|\xi}} \nabla_{12}^2 g(x, y^*(x; \xi); \eta', \xi) \right) \times \left[ \mathbb{E}_{\eta \sim \mathbb{P}_{\eta|\xi}} \nabla_{22}^2 g(x, y^*(x; \xi); \eta, \xi) \right]^{-1} \times \nabla_2 f(x, y^*(x; \xi); \eta, \xi) \right]$$

## Challenges:

- Estimate Hessian inverse
- Estimate the optimal response  $y^*(x; \xi)$

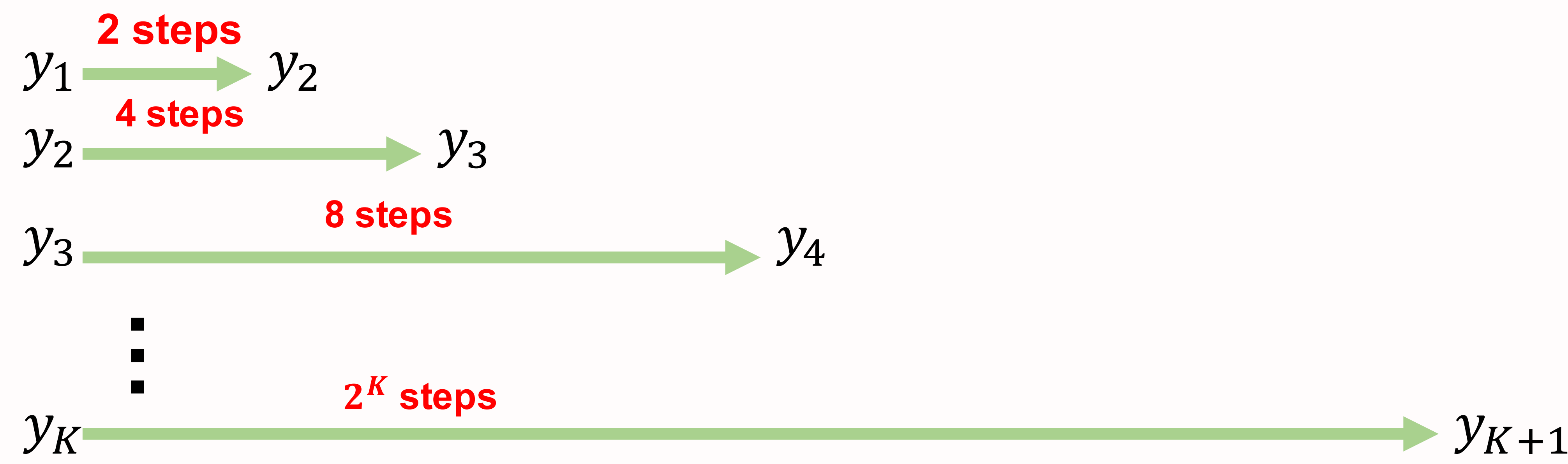
## Matrix Inverse Estimation: Neumann Series

For random matrix  $A$  such that  $0 \prec A \prec I$ :

$$[\mathbb{E}A]^{-1} = \sum_{i=0}^{\infty} (I - \mathbb{E}A)^i = \sum_{i=0}^{\infty} \prod_{n=1}^i \mathbb{E}(I - A_n) \approx \sum_{i=0}^N \prod_{n=1}^i \mathbb{E}(I - A_n).$$

- Bias: exponentially decreasing in  $N$ , i.e.,  $N = \mathcal{O}(\log(\epsilon^{-1}))$

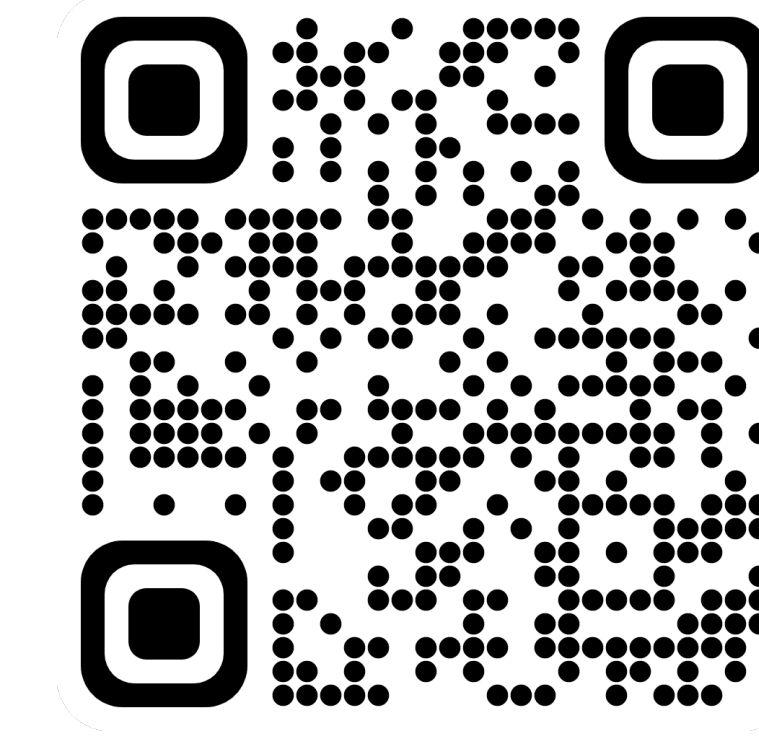
## Optimal Response Estimation: Epoch SGD



- Gradient estimator:  $\hat{v}(x; y_{K+1})$ .
- **Con:** Need  $\mathcal{O}(\epsilon^{-2})$  operations to get  $y_{K+1}$ .



Paper



Yifan's Homepage

## Random Sampling Gradient Estimator (Hu et al. 2021)

$$\hat{v}(x; y_{K+1}) = \hat{v}(x; y_1) + \sum_{k=1}^K [\hat{v}(x; y_{k+1}) - \hat{v}(x; y_k)]$$

$$= \hat{v}(x; y_1) + \sum_{k=1}^K p_k \frac{\hat{v}(x; y_{k+1}) - \hat{v}(x; y_k)}{p_k} = \mathbb{E}_{k \sim \mathbb{P}_k} \left[ \hat{v}(x; y_1) + \frac{\hat{v}(x; y_{k+1}) - \hat{v}(x; y_k)}{p_k} \right].$$

- Sample  $k$  according to pmf  $p_k \propto 2^{-k}$ ,  $\sum_{k=1}^K p_k = 1$ . Construct estimator  $\hat{v}(x) = \hat{v}(x; y_1) + \frac{\hat{v}(x; y_{k+1}) - \hat{v}(x; y_k)}{p_k}$ .
- **High probability:** generate small  $k$ , **Low probability:** generate large  $k$ .
- **Per-iteration cost reduction:** from  $\mathcal{O}(2^K) = \mathcal{O}(\epsilon^{-2})$  to  $\mathcal{O}(K) = \tilde{\mathcal{O}}(1)$ .
- **Variance reduction effect** as  $\hat{v}(x; y_{k+1}) - \hat{v}(x; y_k) \rightarrow 0$  for large  $k$ .

## Takeaway

To find an  $\epsilon$ -stationary point, the sample complexity is

- For vanilla SGD, it is  $\tilde{\mathcal{O}}(\epsilon^{-6})$ .

- For random sampling method, it is  $\tilde{\mathcal{O}}(\epsilon^{-4})$ .

**Remark:** No dependence on number of tasks or individuals  $M$ .

## Numerical Study: Meta-learning on Mini-ImageNet

