

# Towards Efficient Machine Learning: Algorithms, Theoretical Foundations, and Applications

---

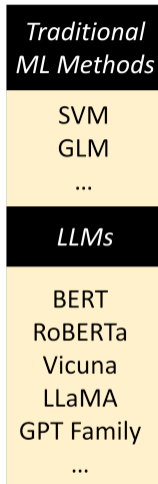
Shaocong Ma

February, 2026

Computer Science, University of Maryland, College Park







## Data



| Date collected | Plot | Species | Sex | Weight |
|----------------|------|---------|-----|--------|
| 1/9/78         | 1    | DM      | M   | 40     |
| 1/9/78         | 1    | DM      | F   | 36     |
| 1/9/78         | 1    | DS      | F   | 133    |
| 1/20/78        | 1    | DM      | F   | 39     |
| 1/20/78        | 2    | DM      | M   | 43     |
| 1/20/78        | 2    | DS      | F   | 144    |
| 3/13/78        | 2    | DM      | F   | 51     |
| 3/13/78        | 2    | DM      | F   | 64     |
| 3/13/78        | 2    | DS      | F   | 146    |

Train



VLMs

CLIP  
BLIP  
FLAVA  
Sora  
DALL-E  
Gemini  
...

Adaptation



## Applications

Image/Vedio generation



Visual search



Data generation

| Plot | Species | Individual | Population |
|------|---------|------------|------------|
| 1    | DM      | 145        | 1000       |
| 1    | DM      | 146        | 1000       |
| 1    | DM      | 147        | 1000       |
| 1    | DM      | 148        | 1000       |
| 1    | DM      | 149        | 1000       |
| 1    | DM      | 150        | 1000       |
| 1    | DM      | 151        | 1000       |
| 1    | DM      | 152        | 1000       |
| 1    | DM      | 153        | 1000       |
| 1    | DM      | 154        | 1000       |
| 1    | DM      | 155        | 1000       |
| 1    | DM      | 156        | 1000       |
| 1    | DM      | 157        | 1000       |
| 1    | DM      | 158        | 1000       |
| 1    | DM      | 159        | 1000       |
| 1    | DM      | 160        | 1000       |
| 1    | DM      | 161        | 1000       |
| 1    | DM      | 162        | 1000       |
| 1    | DM      | 163        | 1000       |
| 1    | DM      | 164        | 1000       |
| 1    | DM      | 165        | 1000       |
| 1    | DM      | 166        | 1000       |
| 1    | DM      | 167        | 1000       |
| 1    | DM      | 168        | 1000       |
| 1    | DM      | 169        | 1000       |
| 1    | DM      | 170        | 1000       |
| 1    | DM      | 171        | 1000       |
| 1    | DM      | 172        | 1000       |
| 1    | DM      | 173        | 1000       |
| 1    | DM      | 174        | 1000       |
| 1    | DM      | 175        | 1000       |
| 1    | DM      | 176        | 1000       |
| 1    | DM      | 177        | 1000       |
| 1    | DM      | 178        | 1000       |
| 1    | DM      | 179        | 1000       |
| 1    | DM      | 180        | 1000       |
| 1    | DM      | 181        | 1000       |
| 1    | DM      | 182        | 1000       |
| 1    | DM      | 183        | 1000       |
| 1    | DM      | 184        | 1000       |
| 1    | DM      | 185        | 1000       |
| 1    | DM      | 186        | 1000       |
| 1    | DM      | 187        | 1000       |
| 1    | DM      | 188        | 1000       |
| 1    | DM      | 189        | 1000       |
| 1    | DM      | 190        | 1000       |
| 1    | DM      | 191        | 1000       |
| 1    | DM      | 192        | 1000       |
| 1    | DM      | 193        | 1000       |
| 1    | DM      | 194        | 1000       |
| 1    | DM      | 195        | 1000       |
| 1    | DM      | 196        | 1000       |
| 1    | DM      | 197        | 1000       |
| 1    | DM      | 198        | 1000       |
| 1    | DM      | 199        | 1000       |
| 1    | DM      | 200        | 1000       |

Image analysis



Photo editing



## VLA Models



OpenVLA



...

Noisy Data



Model Actions



## Applications

### Embodied AI



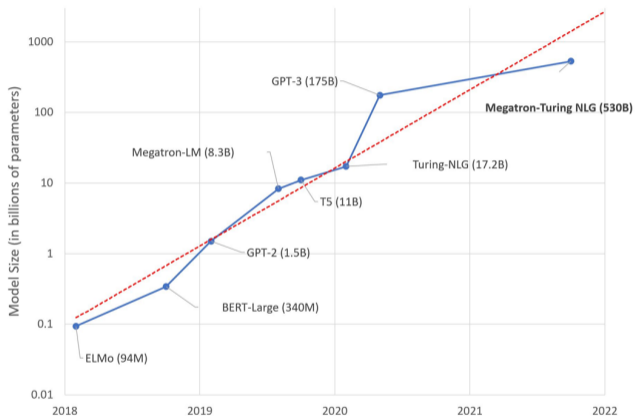
### Robots



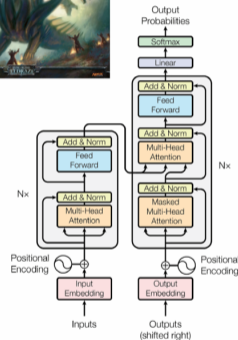
### Autonomous Driving



# Challenge 1: Increasing Model Sizes



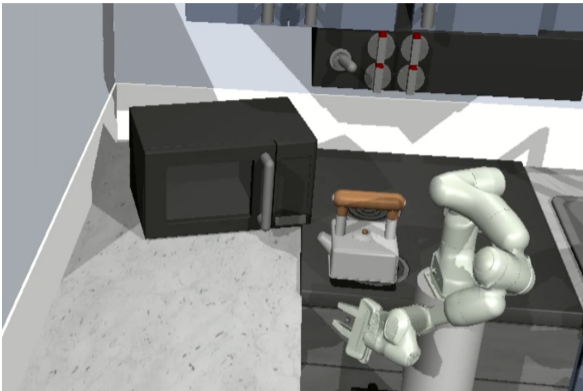
Source: <https://huggingface.co/blog/large-language-models>



**Question:**

How to save the computational resources?

## Challenge 2: Noisy or Adversarial Data

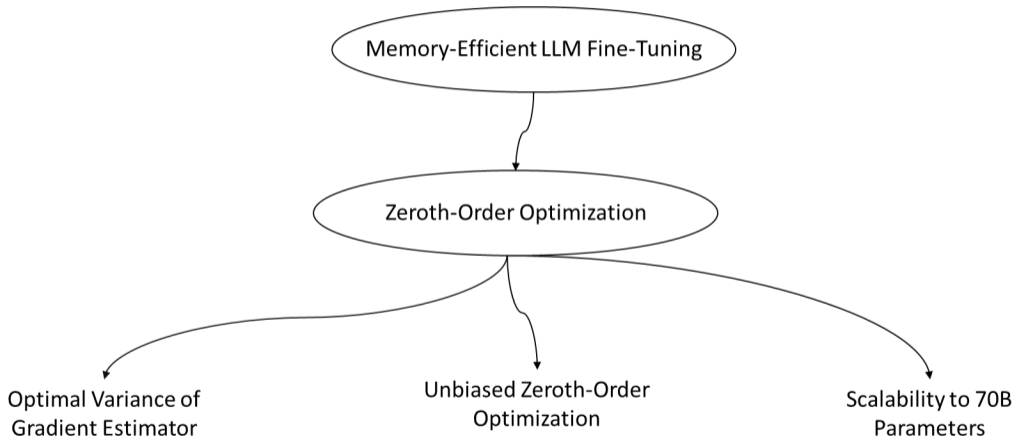


Source: Gu, Shangding, et al. "Robust Gymnasium: A Unified Modular Benchmark for Robust Reinforcement Learning" ICLR 2025.

**Question:**

How to efficiently train a robust ML model?

# Memory-Efficient LLM Fine-Tuning: Outline



# Motivation: Why Fine-Tuning A Local LLM?

## ■ Data privacy:

- Healthcare data
- Trading decision
- Computer control history
- ...



# Motivation: Why Fine-Tuning A Local LLM?

## ■ Data privacy:

- Healthcare data
- Trading decision
- Computer control history
- ...



## ■ Domain Adaptation:

- Medical Diagnosis
- Chip design
- Specific programming language
- ...



# Challenge: The “Memory Wall” Challenge in LLM Fine-tuning

- A single GPU cannot handle backpropagation for entire large models.

**Table 1:** VRAM Requirements and GPU Configuration

| Model Size | First-Order (Full FT) | Est. GPU Setup |
|------------|-----------------------|----------------|
| OPT-1.3B   | ≈ 27 GB               | 1 × A100       |
| OPT-6.7B   | ≈ 156 GB              | 2 × A100       |
| OPT-13B    | ≈ 356 GB              | 4 × A100       |
| OPT-30B    | ≈ 633 GB              | 8 × A100       |

Source: Malladi, Sadhika, et al. "Fine-tuning language models with just forward passes." NeurIPS 2023.

## Question:

How to save the computational resources?

# Zeroth-Order Optimization (ZOO)

Optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x).$$

Gradient Descent:

$$x' \leftarrow x - \eta \nabla f(x)$$

Notation:

- $f(x)$ : The loss function.
- $\eta$ : The learning rate.

# Zeroth-Order Optimization (ZOO)

Optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x).$$

Gradient Descent:

$$x' \leftarrow x - \eta \nabla f(x)$$

Notation:

- $f(x)$ : The loss function.
- $\eta$ : The learning rate.

**Memory-Consuming:** Deep Neural Network  $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial L_N} \frac{\partial L_N}{\partial L_{N-1}} \dots \frac{\partial L_1}{\partial x}$ .

Maintain all intermediate states for backpropagation.

# Zeroth-Order Optimization (ZOO)

Core Formula (Two-Point Estimator):

$$\nabla f(x) \approx \hat{\nabla} f(x) = \frac{f(x + \mu v) - f(x)}{\mu} v$$
$$x' \leftarrow x - \eta \hat{\nabla} f(x)$$

Notation:

- $f(x)$ : The loss function.
- $v$ : A random perturbation vector (e.g., drawn from a Gaussian distribution  $\mathcal{N}(0, I_d)$ ).
- $\mu$ : The perturbation stepsize.

# Zeroth-Order Optimization (ZOO)

Core Formula (Two-Point Estimator):

$$\nabla f(x) \approx \hat{\nabla} f(x) = \frac{f(x + \mu v) - f(x)}{\mu} v$$
$$x' \leftarrow x - \eta \hat{\nabla} f(x)$$

Notation:

- $f(x)$ : The loss function.
- $v$ : A random perturbation vector (e.g., drawn from a Gaussian distribution  $\mathcal{N}(0, I_d)$ ).
- $\mu$ : The perturbation stepsize.

A high-level explanation:

- $\frac{f(x+\mu v)-f(x)}{\mu} > 0 \implies$  Loss increases  $\implies$  Move to the opposite direction of  $v$ ;
- $\frac{f(x+\mu v)-f(x)}{\mu} < 0 \implies$  Loss decreases  $\implies$  Move to the direction of  $v$ .

# Advantages and Challenges of ZOO

## Core Advantage:

- Requires only the **Forward Pass**.
- Memory footprint is comparable to **Inference** only.

# Advantages and Challenges of ZOO

## Core Advantage:

- Requires only the **Forward Pass**.
- Memory footprint is comparable to **Inference** only.

## Key Challenges (Focus of this Talk):

1. **High Variance:** Gradient estimates are volatile/noisy.
2. **Biased:** Finite difference methods rely on approximations.

# Advantages and Challenges of ZOO

## Core Advantage:

- Requires only the **Forward Pass**.
- Memory footprint is comparable to **Inference** only.

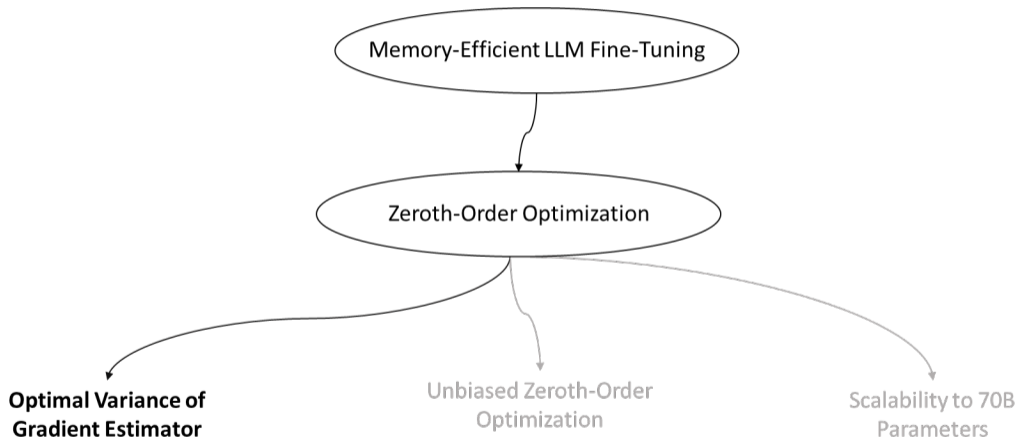
## Key Challenges (Focus of this Talk):

1. **High Variance:** Gradient estimates are volatile/noisy.
2. **Biased:** Finite difference methods rely on approximations.

## Roadmap: Two theoretical works improving ZOO

1. Derive the condition for achieving the optimal variance.
2. Propose a unbiased gradient estimator family.

# Memory-Efficient LLM Fine-Tuning: Outline



# Minimum Variance: Directionally Aligned Perturbation (DAP)

## Recap: Zero-Order Gradient Estimator

$$\nabla f(x) \approx \hat{\nabla} f(x) := \frac{f(x + \mu v) - f(x)}{\mu} v$$

where  $v$  is typically sampled from a Gaussian distribution.

**Q:** Is the standard choice of distribution actually optimal?

# Minimum Variance: Directionally Aligned Perturbation (DAP)

Recap: Zero-Order Gradient Estimator

$$\nabla f(x) \approx \hat{\nabla} f(x) := \frac{f(x + \mu v) - f(x)}{\mu} v$$

where  $v$  is typically sampled from a Gaussian distribution.

**Goal: Minimize Estimation Error**

Find the optimal distribution  $V$  for the perturbation vector  $v$ :

$$\begin{aligned} \min_V \quad & \mathbb{E}_{v \sim V} \left\| \frac{f(x + \mu v) - f(x)}{\mu} v - \nabla f(x) \right\|^2, \\ \text{s.t.} \quad & \mathbb{E}_{v \sim V} [v v^\top] = I_d. \end{aligned}$$

# Minimum Variance: Directionally Aligned Perturbation (DAP)

Recap: Zero-Order Gradient Estimator

$$\nabla f(x) \approx \hat{\nabla} f(x) := \frac{f(x + \mu v) - f(x)}{\mu} v$$

where  $v$  is typically sampled from a Gaussian distribution.

**Goal: Minimize Estimation Error**

Find the optimal distribution  $V$  for the perturbation vector  $v$ :

$$\begin{aligned} \min_V \quad & \mathbb{E}_{v \sim V} \left\| \frac{f(x + \mu v) - f(x)}{\mu} v - \nabla f(x) \right\|^2, \\ \text{s.t.} \quad & \mathbb{E}_{v \sim V} [v v^\top] = I_d. \end{aligned}$$

**Optimization Challenges:**

- **Functional space:** Optimization is taken over all probability distributions.
- **Constraints with an empty interior:** The empty interior precludes the use of Interior Point Methods.

# Minimize Estimation Error (Part I)

Find the optimal distribution  $V$  for the perturbation vector  $v$ :

$$\begin{aligned} \min_V \quad & \mathbb{E}_{V \sim V} \left\| \frac{f(x + \mu V) - f(x)}{\mu} \Big|_V - \nabla f(x) \right\|^2, \\ \text{s.t.} \quad & \mathbb{E}_{V \sim V} [VV^\top] = I_d. \end{aligned}$$

Taylor Expansion:

$$\begin{aligned} f(x + \mu v) - f(x) &\approx \mu \nabla f(x)^\top v + \mu^2 \cdot \text{Bias} \\ \frac{f(x + \mu v) - f(x)}{\mu} \Big|_v &\approx vv^\top \nabla f(x) + \mu \cdot \text{Bias} \\ \frac{f(x + \mu V) - f(x)}{\mu} \Big|_V - \nabla f(x) &\approx (VV^\top - I_d) \nabla f(x) + \underbrace{\mu \cdot \text{Bias}}_{\text{Ignored}} \end{aligned}$$

Definition of Bias:  $\mu \cdot \text{Bias} := \mathbb{E}[\hat{\nabla} f(x)] - \nabla f(x)$

## Minimize Estimation Error (Part II)

Find the optimal distribution  $V$  for the perturbation vector  $v$ :

$$\begin{aligned} \min_V \quad & \mathbb{E}_{v \sim V} \left\| \frac{f(x + \mu v) - f(x)}{\mu} - \nabla f(x) \right\|^2, \\ \text{s.t.} \quad & \mathbb{E}_{v \sim V} [v v^\top] = I_d. \end{aligned}$$

Plug in the objective function:

$$\mathbb{E}_{v \sim V} \left\| \frac{f(x + \mu v) - f(x)}{\mu} - \nabla f(x) \right\|^2 = \mathbb{E}_{v \sim V} \nabla f(x)^\top (v v^\top) \nabla f(x) - \|\nabla f(x)\|^2$$

## Minimize Estimation Error (Part II)

Find the optimal distribution  $V$  for the perturbation vector  $v$ :

$$\begin{aligned} \min_V \quad & \mathbb{E}_{v \sim V} \left\| \frac{f(x + \mu v) - f(x)}{\mu} - \nabla f(x) \right\|^2, \\ \text{s.t.} \quad & \mathbb{E}_{v \sim V} [v v^\top] = I_d. \end{aligned}$$

Plug in the objective function:

$$\mathbb{E}_{v \sim V} \left\| \frac{f(x + \mu v) - f(x)}{\mu} - \nabla f(x) \right\|^2 = \mathbb{E}_{v \sim V} \nabla f(x)^\top (v v^\top) \nabla f(x) - \|\nabla f(x)\|^2$$

Simplified Objective:

$$\begin{aligned} \min_V \quad & \mathbb{E}_{v \sim V} a^\top (v v^\top) a \\ \text{s.t.} \quad & \mathbb{E}_{v \sim V} [v v^\top] = I_d. \end{aligned}$$

## Minimize Estimation Error (Part II)

We analytically solve this functional optimization problem.

Derive the lower bound + Tell when the lower bound is achieved

### Theorem

Let  $v$  be a random vector following the distribution  $V$  with  $\mathbb{E}_{v \sim V} v v^T = I_d$  and  $a \in \mathbb{R}^d$  be a fixed vector. Then

$$d \|a\|^2 \leq \mathbb{E}_{v \sim V} a^T (v v^T)^2 a \leq d \|a\|^2 + \frac{\|a\|^2}{2} \left( \rho_V + \sqrt{\rho_V^2 + 4(d-1)\rho_V} \right)$$

where  $\rho_V := \mathbb{E} \|v\|^4 - d^2$ .

Ma, S. & Huang, H. Revisiting Zeroth-Order Optimization: Minimum-Variance Two-Point Estimators and Directionally Aligned Perturbations. ICLR, 2025. Spotlight

## Minimize Estimation Error (Part II)

We analytically solve this functional optimization problem.

Derive the lower bound + Tell when the lower bound is achieved

### Theorem

Let  $v$  be a random vector following the distribution  $V$  with  $\mathbb{E}_{v \sim V} vv^T = I_d$  and  $a \in \mathbb{R}^d$  be a fixed vector. Then

$$d\|a\|^2 \leq \mathbb{E}_{v \sim V} a^T (vv^T)^2 a \leq d\|a\|^2 + \frac{\|a\|^2}{2} \left( \rho_V + \sqrt{\rho_V^2 + 4(d-1)\rho_V} \right)$$

where  $\rho_V := \mathbb{E}\|v\|^4 - d^2$ .

Ma, S. & Huang, H. Revisiting Zeroth-Order Optimization: Minimum-Variance Two-Point Estimators and Directionally Aligned Perturbations. ICLR, 2025. Spotlight

What kind of distribution actually achieves this lower bound?

# DAPs: Directionally Aligned Perturbation

We analytically solve this functional optimization problem.

(Equality Condition)

■ Constant Magnitude Perturbations:

- $\mathbb{E}_{v \sim V}[v v^\top] = I_d$ .
- $\|v\|$  is fixed.

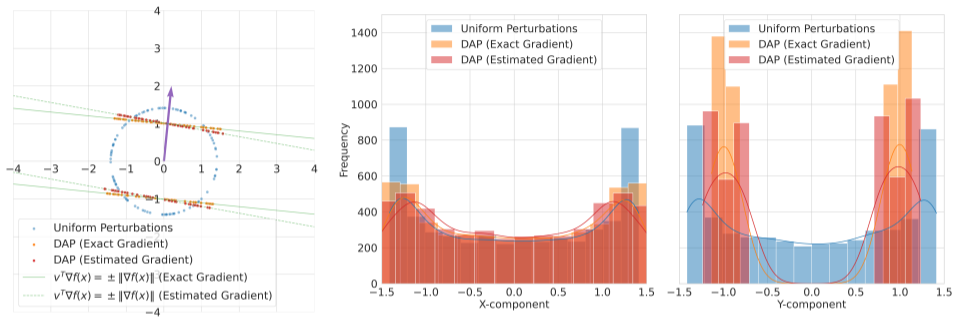
■ Directionally Aligned Perturbations (DAPs):

- $\mathbb{E}_{v \sim V}[v v^\top] = I_d$ .
- $\nabla f(x)^\top v$  is fixed.

⇒ Both estimators achieve the **minimum variance**.

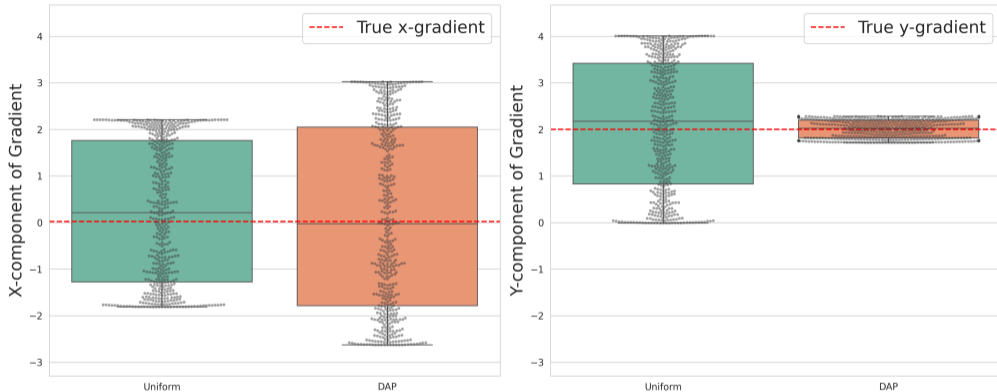
⇒ DAPs have some nice properties.

# Traditional Methods Cannot Identify the Important Directions



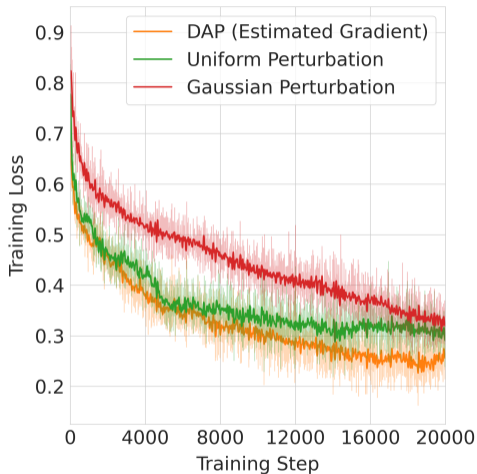
**Figure 1:** Illustration of the *directional alignment* property of DAP in  $d = 2$  with estimating the gradient of  $f(x) = x_1^2 + x_2^2$  at  $x = [0.1 \ 1]^T$ . Traditional estimator is **symmetric**, but we need a **non-symmetric** estimator.

# Traditional Methods Cannot Identify the Important Directions



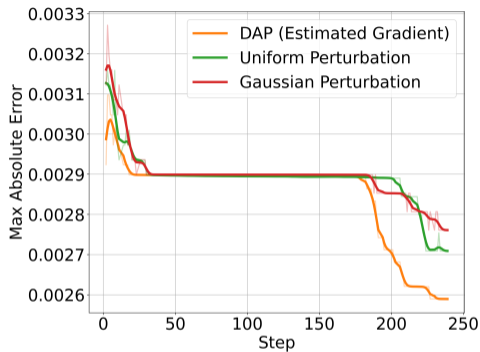
**Figure 2:** Comparison of gradient estimation performance with estimating the gradient of  $f(x) = x_1^2 + x_2^2$  at  $x = [0.1 \ 1]^T$  between uniform random perturbations and DAPs. The **non-symmetric** estimator is more accurate in the direction with larger gradient.

# Applications in LLM Fine-Tuning



**Figure 3:** Comparison of training loss curves among different random perturbations on **Large Language Model Fine Tuning**.

# Applications in Scientific Optimization



**Figure 4:** Comparison of training loss curves among different random perturbations on **Mesh Optimization for the Physical Numerical Solver**.

## Summary

Derived the optimal distribution of  $v$  to achieve the minimum variance.

$$\nabla f(x) \approx \hat{\nabla} f(x) := \frac{f(x + \mu v) - f(x)}{\mu} v$$

## Summary

Derived the optimal distribution of  $v$  to achieve the minimum variance.

$$\nabla f(x) \approx \hat{\nabla} f(x) := \frac{f(x + \mu v) - f(x)}{\mu} v$$

(Taylor approximation)

$$\frac{f(x + \mu v) - f(x)}{\mu} v - \nabla f(x) \approx (vv^T - I_d) \nabla f(x) + \underbrace{\mu \cdot \text{Bias}}_{\text{Ignored}}$$

## Summary

Derived the optimal distribution of  $v$  to achieve the minimum variance.

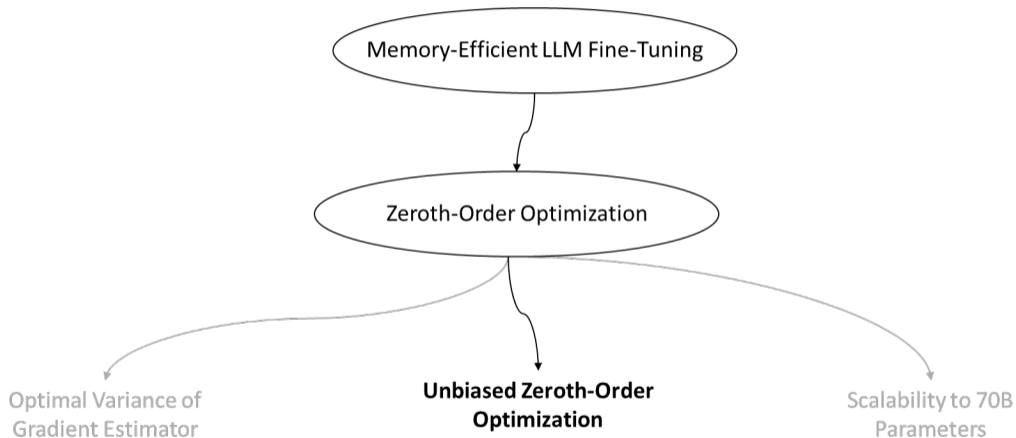
$$\nabla f(x) \approx \hat{\nabla} f(x) := \frac{f(x + \mu v) - f(x)}{\mu} v$$

(Taylor approximation)

$$\frac{f(x + \mu v) - f(x)}{\mu} v - \nabla f(x) \approx (vv^T - I_d) \nabla f(x) + \underbrace{\mu \cdot \text{Bias}}_{\text{Ignored}}$$

Is it possible to eliminate the bias completely?

# Memory-Efficient LLM Fine-Tuning: Outline



# Inherent Bias of Two-Point Estimator

## Recap: Zero-Order Gradient Estimator

$$\hat{\nabla}f(x) := \frac{f(x + \mu v) - f(x)}{\mu} v.$$

## Taylor Expansion:

$$f(x + \mu v) - f(x) \approx \mu \nabla f(x)^\top v + \mu^2 \cdot \text{Bias}$$

$$\frac{f(x + \mu v) - f(x)}{\mu} v \approx v v^\top \nabla f(x) + \mu \cdot \text{Bias}$$

$$\frac{f(x + \mu v) - f(x)}{\mu} v - \nabla f(x) \approx (v v^\top - I_d) \nabla f(x) + \underbrace{\mu \cdot \text{Bias}}_{\text{Not Ignored?}}$$

▷ When  $\mu$  is large, the two-point estimator exhibits significant bias.

Unbiased zeroth-order gradient estimator using only function evaluations.

- *Step 1.* Directional derivative along the direction  $v$ .

$$\nabla_v f(x) = \lim_{\mu \rightarrow 0} \frac{f(x + \mu v) - f(x)}{\mu}.$$

- *Step 2.* Telescoping series. Let  $\mu_n \rightarrow 0$ .

$$\begin{aligned} \nabla_v f(x) &= \frac{f(x + \mu_1 v) - f(x)}{\mu_1} \\ &+ \sum_{n=1}^{\infty} \left[ \frac{f(x + \mu_{n+1} v) - f(x)}{\mu_{n+1}} - \frac{f(x + \mu_n v) - f(x)}{\mu_n} \right]. \end{aligned}$$

# Unbiased Estimator based on Multi-Level Monte Carlo

- Step 3. Expectation representation.

Let  $\sum_n p_n = 1$  and  $0 < p_n < 1$ .

$$\begin{aligned} \nabla_v f(x) = & \sum_{n=1}^{\infty} p_n \left[ \frac{f(x + \mu_1 v) - f(x)}{\mu_1} \right. \\ & \left. + \frac{1}{p_n} \left( \frac{f(x + \mu_{n+1} v) - f(x)}{\mu_{n+1}} - \frac{f(x + \mu_n v) - f(x)}{\mu_n} \right) \right]. \end{aligned}$$

Then  $\nabla_v f(x)$  can be represented as

$$\mathbb{E}_{n \sim \{p_n\}_{n=1}^{\infty}} \left[ \frac{f(x + \mu_1 v) - f(x)}{\mu_1} + \frac{1}{p_n} \left( \frac{f(x + \mu_{n+1} v) - f(x)}{\mu_{n+1}} - \frac{f(x + \mu_n v) - f(x)}{\mu_n} \right) \right].$$

$\implies$  **Unbiased Estimator Family**

# Unbiased Estimator Family

- $P_4$ -Estimator:

$$P_4(n, v) := \frac{f(x + \mu_1 v) - f(x)}{\mu_1} + \frac{1}{\rho_n} \left( \frac{f(x + \mu_{n+1} v) - f(x)}{\mu_{n+1}} - \frac{f(x + \mu_n v) - f(x)}{\mu_n} \right)$$

- $P_3$ -Estimator:

$$P_3(n, v) := \frac{f(x + \mu_1 v) - f(x)}{\mu_1} U_2 + \frac{1}{\rho_n} \left( \frac{f(x + \mu_{n+1} v) - f(x)}{\mu_{n+1}} - \frac{f(x + \mu_n v) - f(x)}{\mu_n} \right) (1 - U_2)$$

where  $U_2 \sim \text{Uniform}(\{0, 1\})$ .

- We can also define  $P_2$ -Estimator and  $P_1$ -Estimator.

Ma, S. & Huang, H. On the Optimal Construction of Unbiased Gradient Estimators for Zeroth-Order Optimization. NeurIPS, 2025.

Spotlight

# Unbiased Estimator Family: Variance

## Theorem

Suppose that  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is second-order continuously differentiable and has  $L$ -Lipschitz continuous gradient.  $\sum_{n=1}^{\infty} \mu_n < \infty$  and  $V \sim \sqrt{d} \text{Uniform}(\mathbb{S}^{d-1})$ . Define

$$\mu := \mu_1, \quad \varrho := \sum_{n=1}^{\infty} \frac{(\mu_{n+1} - \mu_n)^2}{\rho_n}, \quad \text{and} \quad \varphi := \sum_{n=1}^{\infty} \frac{\mu_n^2}{\rho_n}.$$

Then

$$\text{Var}[P_4(n, v)v] \leq (d-1) \|\nabla f(x)\|^2 + \frac{3L^2}{4} d^3 \mu^2 + \frac{L^2 d^3}{2} \varrho,$$

$$\text{Var}[P_3(n, v)v] \leq \text{Var}[P_4(n, v)v] + \frac{L^2}{8} d^3 \mu^2 + \frac{L^2 d^3}{8} \varrho.$$

▷ This variance results in the optimal oracle complexity.

# Unbiased Estimator Family: Variance

## Theorem

Let  $\{\mu_n\}_{n=1}^{\infty}$  be a positive, decreasing sequence with  $\sum_{n=1}^{\infty} \mu_n < \infty$ , and let  $\{p_n\}_{n=1}^{\infty}$  be a Probability Mass Function. Then

$$\varrho \geq \mu^2.$$

The equality holds if and only if

$$p_n = \frac{\mu_n - \mu_{n-1}}{\mu}.$$

- ▷ We obtain a simple and elegant relation to derive the optimal sequence  $\{p_n\}$  and  $\{\mu_n\}$ .
- Geometric  $P_k$ -Estimator:  $n \sim \text{Geom}(c)$  and  $\mu_n = \mu_1 c^{n-1}$ .
  - Zipf's  $P_k$ -Estimator:  $n \sim \text{Zipf}(s)$  ( $s > 1$ ) and  $\mu_n = \mu_1 \left[ 1 - \left( \sum_{j=1}^{n-1} \frac{1}{j^s} \right) / \zeta(s) \right]$ .

# Unbiased Estimator Leads to Better Accuracy

The quadratic loss  $f_{\text{reg}} : \mathbb{R}^d \rightarrow \mathbb{R}$  and the logistic loss  $f_{\text{cls}} : \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$f_{\text{reg}}(x) = x^T A^T A x, \quad f_{\text{cls}}(x) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-b_i \cdot (a_i^T \cdot x))).$$

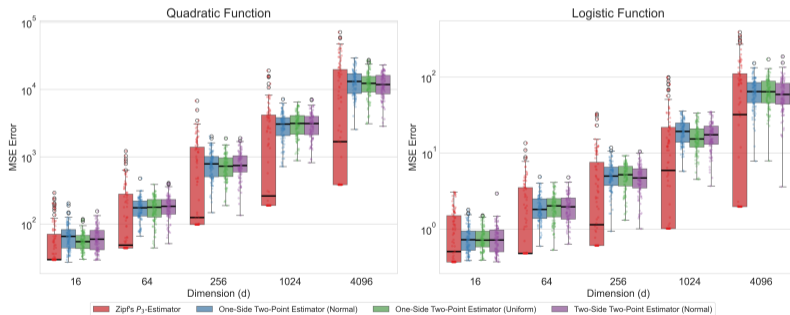


Figure 5: The MSE error (Left:  $f_{\text{reg}}$ , Right:  $f_{\text{cls}}$ ) of different estimators.

# Applications in LLM Fine-Tuning

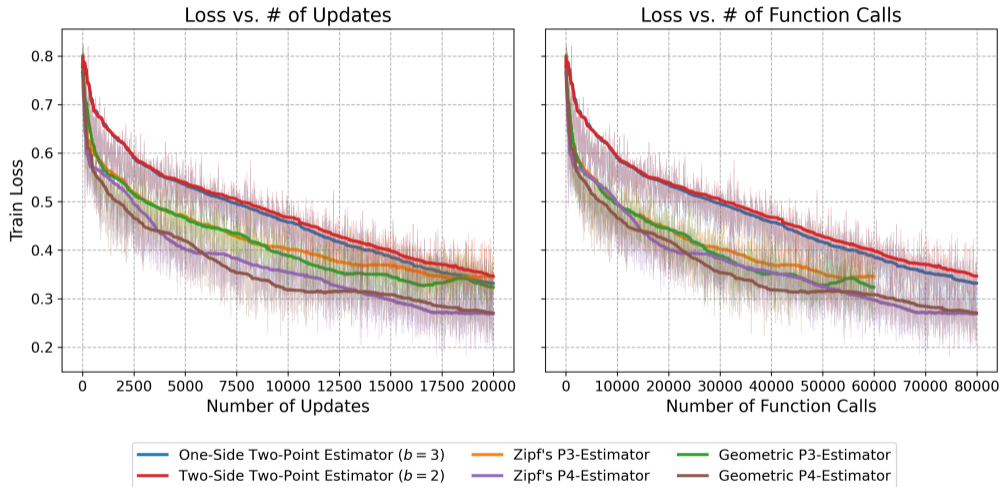


Figure 6: Fine-tuning the OPT-1.3B model on SST-2 using different gradient estimators.

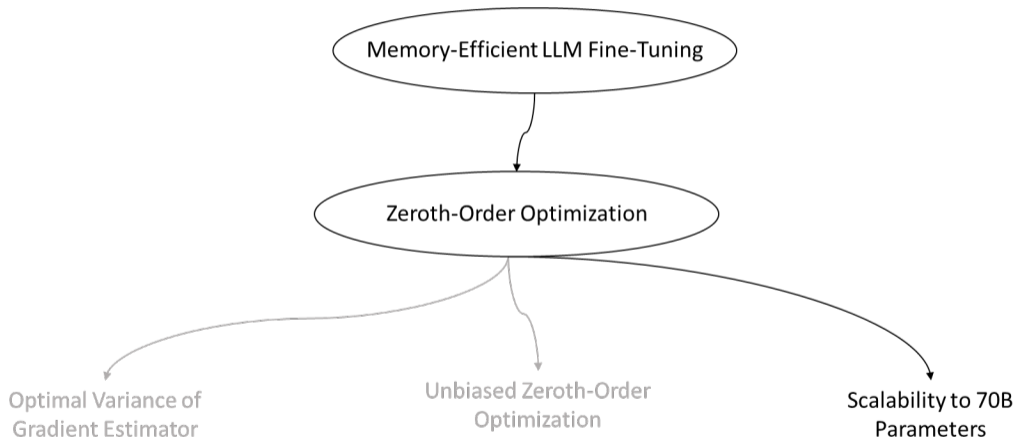
## Summary

- Constructed the family of unbiased zeroth-order gradient estimators.
- Provided the theoretical framework to minimize its variance.
- Validated its performance in synthetic and LLM experiments.

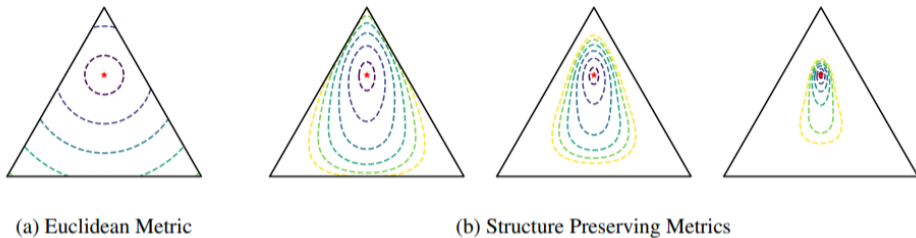
- Constructed the family of unbiased zeroth-order gradient estimators.
- Provided the theoretical framework to minimize its variance.
- Validated its performance in synthetic and LLM experiments.

Can we further scale up the Zeroth-Order Optimization method?

# Memory-Efficient LLM Fine-Tuning: Outline



# Geometric Constraints in Quantized LLMs

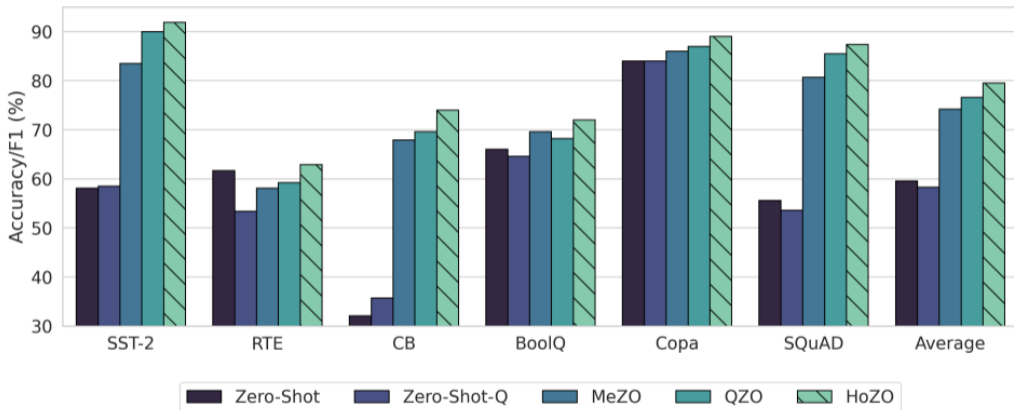


**Figure 7:** Visualization of different metrics on the probability simplex.

▷ Scale parameters in quantized LLMs form a geodesically incomplete Riemannian manifold. We propose structure preserving metrics to handle this issue.

Ma, S. & Huang, H. Riemannian Zeroth-Order Gradient Estimation with Structure-Preserving Metrics for Geodesically Incomplete Manifolds. ICLR, 2026.

# Memory-Efficient Fine-Tuning of Quantized LLMs



**Figure 8:** On the INT4 Llama-2-7B model across 6 downstream tasks, HoZO achieves consistently better performance than all baselines.

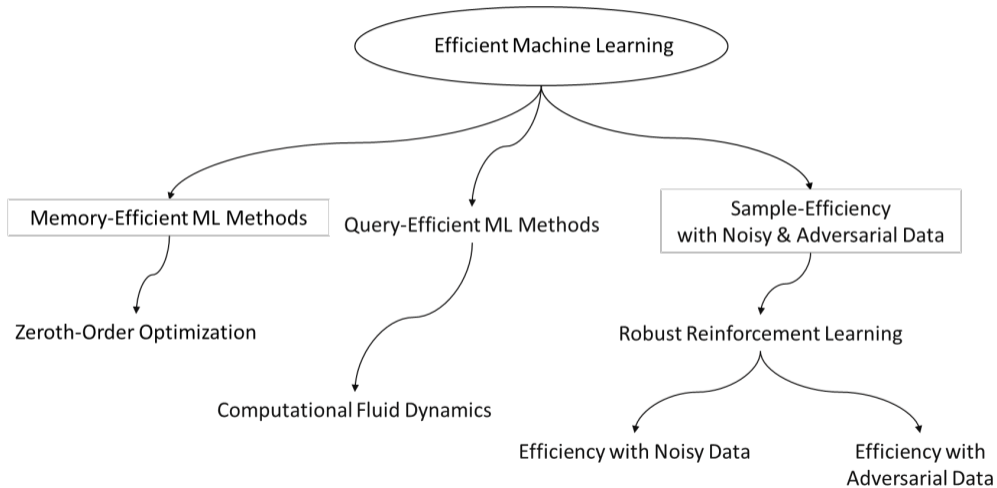
# Memory-Efficient Fine-Tuning of Quantized LLMs

| Model       | Method      | Model Precision | SST-2       | RTE         | CB          | BoolQ       | Copa        | SQuAD       | Average     |       |
|-------------|-------------|-----------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------|
| Llama-2-7b  | Zero-Shot   | 16 bits         | 58.1*       | 61.7*       | 32.1*       | 66.0*       | 84.0        | 55.6*       | 59.6        | -     |
|             | Zero-Shot-Q | 4 bits          | 58.5*       | 53.4*       | 35.7*       | 64.6*       | 84.0        | 53.6*       | 58.3        | ↓1.3  |
|             | MeZO        | 16 bits         | 83.5*       | 58.1*       | 67.9*       | 69.6*       | 86.0        | 80.7*       | 74.2        | ↑14.6 |
|             | QZO         | 4 bits          | 90.0*       | 59.2*       | 69.6*       | 68.2*       | 87.0        | 85.5*       | 76.6        | ↑17.0 |
|             | <b>HoZO</b> | 4 bits          | <b>91.9</b> | <b>62.9</b> | <b>74.0</b> | <b>72.0</b> | <b>89.0</b> | <b>87.4</b> | <b>79.5</b> | ↑19.9 |
| Llama-2-13b | Zero-Shot   | 16 bits         | 61.1        | 50.9        | 44.0        | 74.1        | 89.0        | 63.7        | 63.8        | -     |
|             | Zero-Shot-Q | 4 bits          | 60.0        | 47.3        | 47.0        | 74.7        | 87.0        | 63.1        | 63.2        | ↓0.6  |
|             | MeZO        | 16 bits         | 90.7        | 58.5        | <b>77.0</b> | 81.6        | <b>92.0</b> | 87.5        | 81.2        | ↑17.4 |
|             | QZO         | 4 bits          | 91.9        | 62.8        | <b>77.0</b> | <b>82.4</b> | <b>92.0</b> | 89.2        | 82.5        | ↑18.7 |
|             | <b>HoZO</b> | 4 bits          | <b>92.4</b> | <b>68.6</b> | 75.0        | 81.6        | <b>92.0</b> | <b>89.3</b> | <b>83.2</b> | ↑19.4 |
| Llama-2-70b | Zero-Shot-Q | 4 bits          | 56.4        | 60.6        | 47.0        | 74.7        | 92.0        | 71.4        | 67.0        | -     |
|             | QZO         | 4 bits          | 90.6        | <b>80.8</b> | 82.0        | <b>83.8</b> | 93.0        | 90.4        | 86.8        | ↑19.8 |
|             | <b>HoZO</b> | 4 bits          | <b>91.5</b> | 79.1        | <b>83.0</b> | <b>83.8</b> | <b>95.0</b> | <b>91.2</b> | <b>87.3</b> | ↑20.3 |

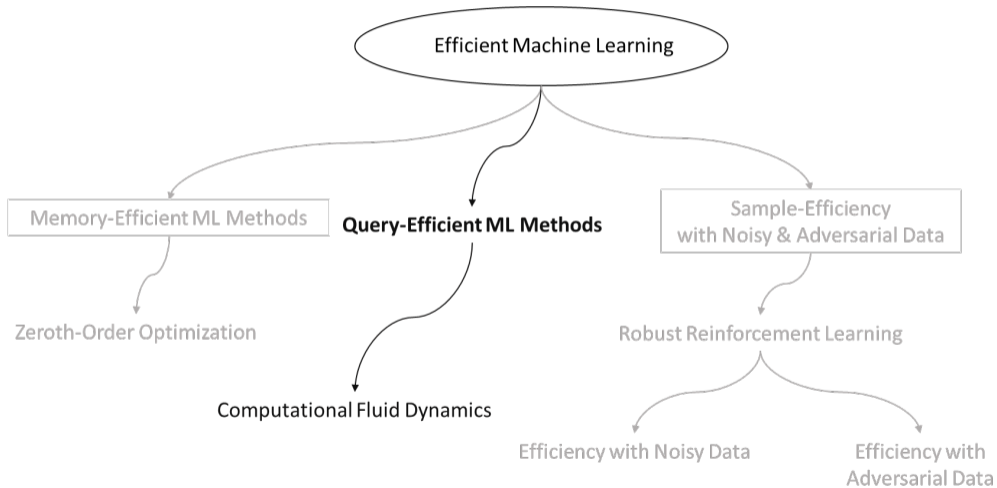
Ma, S. & Huang, H. Hyper-Octant Zeroth-Order Optimization: Fine-Tuning Quantized LLMs on the Positive Orthant Manifold.

Submitted to ICML, 2026.

# My Other Research Summary

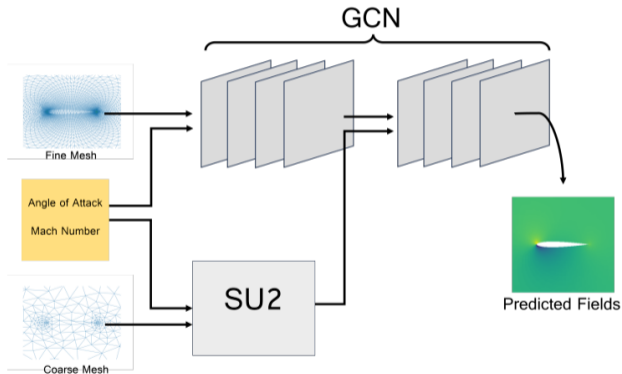


# Query-Efficient ML Methods



# Query-Efficient ML Methods

**Motivation:** A single query may take minutes or hours long.



**Figure 9:** The CFD-GCN model. The PDE solver can be slow. How to improve the query efficiency?

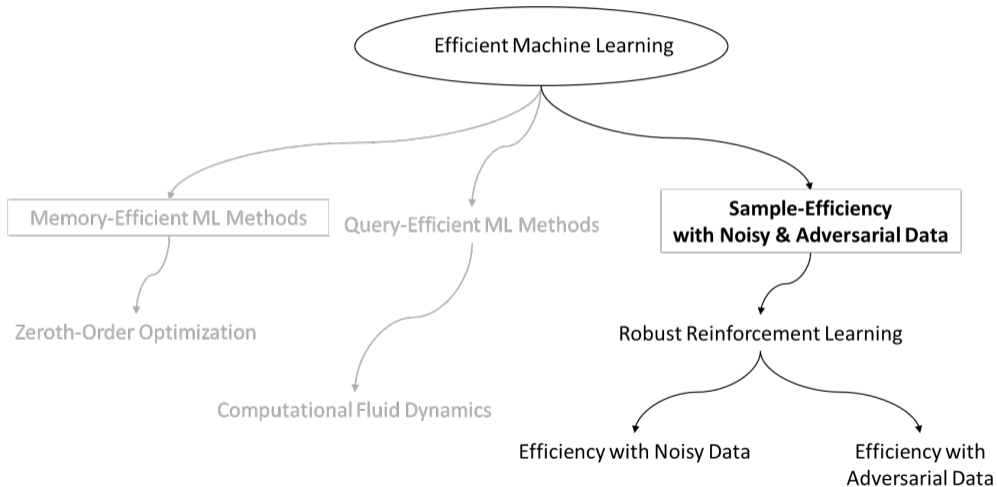
**Our Solution:** Estimate the non-differentiable part only.

$$\nabla_x f(g(x)) = \frac{\partial f}{\partial g} \times \underbrace{\frac{\widehat{\partial g}}{\partial x}}_{\text{Est. this}}$$

Ma, S., Diffenderfer, J., Kailkhura, B., & Zhou, Y. Deep learning of PDE correction and mesh adaption without automatic differentiation. Machine Learning, 2025.

Ma, S., Diffenderfer, J., Kailkhura, B., & Zhou, Y. End-to-End Mesh Optimization of a Hybrid Deep Learning Black-Box PDE Solver. NeurIPS, 2023 (ML4PS Workshop).

# Sample-Efficiency with Noisy & Adversarial Data



**Motivation:** Overcome sample-inefficiency caused by non-IID Markovian data.

**Examples:**

- Agent-environment interactions in RL.
- Continuous sensor streams in autonomous driving or IoT.
- Network traffic or queueing systems.

## Our Solution:

- Data Reshuffling.

Ma, S., & Zhou, Y. Understanding the impact of model incoherence on convergence of incremental sgd with random reshuffle. ICML, 2020.

- Larger Batch Size.

Ma, S., Chen, Z., Zhou, Y., Ji, K., & Liang, Y. Data sampling affects the complexity of online SGD over dependent data. UAI, 2022.

- Variance Reduction.

Ma, S., & Zhou, Y. Variance-Reduced Off-Policy TDC Learning: Non-Asymptotic Convergence Analysis. NeurIPS, 2020.

Ma, S., Chen, Z., & Zhou, Y. Greedy-GQ with Variance Reduction: Finite-time Analysis and Improved Complexity. ICLR, 2021.

# Sample-Efficiency with Adversarial Data

## Motivation:

- Environment mismatch between the simulation and training environments.
- Multi-agent systems.
- Adversarial attacking.

## Question:

How to efficiently train a robust ML model?

## Our Solution:

- Environment mismatch between the simulation and training environments.  
Ma, S., & Huang, H. Robust Reinforcement Learning in Finance: Modeling Market Impact with Elliptic Uncertainty Sets. NeurIPS, 2025.
- Multi-agent systems.  
Ma, S., Chen, Z., Zou, S. & Zhou, Y. Decentralized Robust V-Learning for Solving Markov Games with Model Uncertainty. JMLR, 2023.
- Adversarial attacking.  
Chen, Z., Ma, S., & Zhou, Y. Accelerated Proximal Alternating Gradient-Descent-Ascent for Nonconvex Minimax Machine Learning. IEEE ISIT, 2022.

Future Work

## ■ Motivation:

- **Data Privacy & Security:**  
Sensitive user data (e.g., personal messages, health records) never leaves the device.
- **Reduced Latency:**  
Real-time analysis of wearable data (e.g., heart attack detection, EEG translation).
- **Personalization:**  
Adapts the model to the specific user's habits and local context.

## ■ Motivation:

- **Data Privacy & Security:**

Sensitive user data (e.g., personal messages, health records) never leaves the device.

- **Reduced Latency:**

Real-time analysis of wearable data (e.g., heart attack detection, EEG translation).

- **Personalization:**

Adapts the model to the specific user's habits and local context.

## ■ Challenges:

- Extreme memory constraints  $\implies$  Memory Efficiency.

- Extreme computational constraints  $\implies$  Computation Efficiency.

## ■ Motivation:

- **High-cost of LLM agents evaluation & inference:**  
A complete LLM agent run can take hours.
- **Non-differentiable modules & structures:**  
Black-box models; topological routing structures.

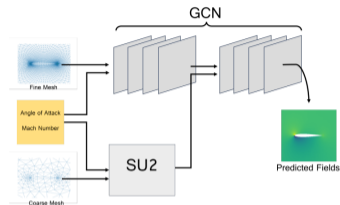
# Fine-Tuning Agentic Framework with Zeroth-Order Optimization

## ■ Motivation:

- **High-cost of LLM agents evaluation & inference:**  
A complete LLM agent run can take hours.
- **Non-differentiable modules & structures:**  
Black-box models; topological routing structures.

## ■ Challenges:

- Extreme long time cost  $\implies$  Query Efficiency.
- Extreme sparse reward signals  $\implies$  Sample Efficiency.



## Motivation:

- The scientific software can be slow and non-differentiable  $\implies$  Query Efficiency.

## Future directions:

- **Turbulence:** More challenging CFD problems.
- **Extension to protein prediction:** Molecular Dynamics + DNN.

# Future Research Directions and Strategies



- Continue collaborations with
  - UMD, Univ. of Utah, USC, TAMU, UTA, OSU, ASU, Buffalo, *etc.*
  - Lawrence Livermore National Lab, MD Anderson Cancer Center, *etc.*
- Seek more collaborations with
  - Texas Tech University CS (AI, ML, System, Data-Intensive Computing, Cyber Security, HCI), Physics, Chemistry & Biochemistry (CFD, AI4Science)
  - Texas Tech University Health Sciences Center (TTUHSC) (AI4Healthcare)
  - Local Industrials (UMC Health System, Covenant Health, Texas Instruments, *etc.*)
  - General Industrials (Google, Amazon, UnitedHealthcare, *etc.*)

# Successful Funding Experiences in Proposal Writing

## AI for Healthcare

- A Real-World Test Bed for Post-Market Surveillance and Stress Testing of AI-Enabled Imaging Tools  
2025–2027, FDA, **\$1.2M**.
- Ultrascale Machine Learning to Empower Discovery in Alzheimer's Disease Biobanks  
2026–2031 (Recommended), NIH center grant, **\$15M**.

## Robust Machine Learning

- Advanced AI Framework to Improve Understanding and Prediction of Wildland Fire  
2026–2028, NSF-RISE, **\$1,856,577**.

## Other Writing Experiences

- NSF MFAI, NSF GCR, NSF PCL, NSF SLES, NSF/NIH SCH, NIH R01s.

- First several proposals: NSF MFAI, NSF Early Career Development, NSF-CISE, *etc.*
- Collaborate with colleagues at TTU to seek: NSF Future CoRe (larger budget), NSF ACED, NSF AIMing, NIH-NIA R01, NIH-NIGMS R01, NIH-NIBIB R01, *etc.*

## Teaching Experiences

- Teaching Assistant at UCSB
  - PSTAT 5A: Statistics
  - PSTAT 5LS: Statistics for Life Science
  - PSTAT 109: Statistics for Economics
  - PSTAT 172A: Actuarial Statistics
  - PSTAT 175: Survival Analysis
- Teaching Assistant at Univ. of Utah
  - ECE 3500: Fundamentals of Signals and Systems
- Co-teach at UMD
  - CMSC422: Introduction to Machine Learning

## I can teach various courses:

- Lower-Level Undergraduate Courses:
  - Data Analysis & Data Science
  - Machine Learning & AI
  - Numerical Algorithms
  - Discrete Mathematics
  - Linear Algebra
- Upper-Level Undergraduate or Graduate Courses:
  - Advanced Machine Learning & Statistical Learning
  - Modern Machine Learning Models
  - Advanced Stochastic Algorithms
  - Reinforcement Learning

**Thank You!**