

Towards Efficient Machine Learning Algorithms: Theoretical Foundations and Applications

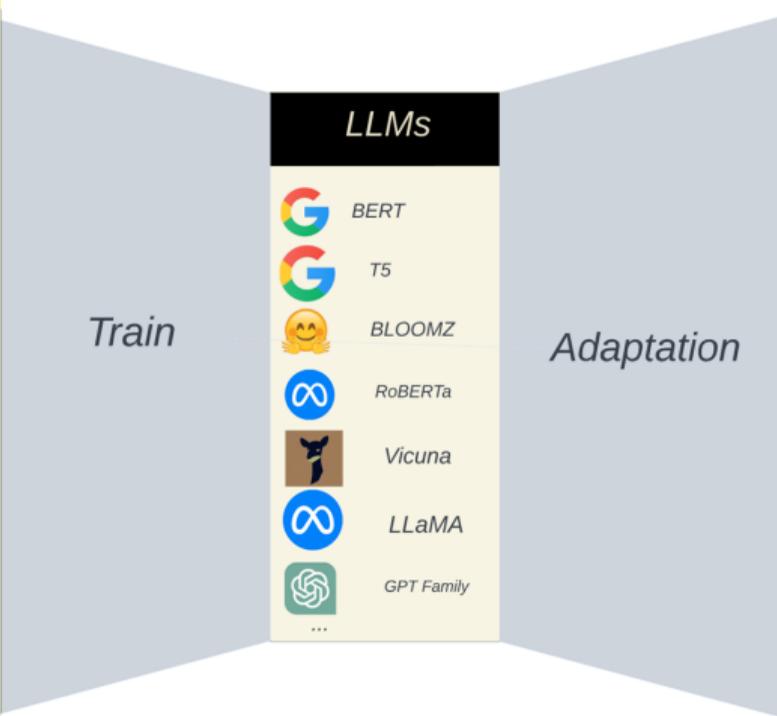
Shaocong Ma

February, 2026

University of Maryland, College Park



Data



Applications

Sentiment Analysis

Translation

Recommendation

Information extraction

Code development

```

for i in people.data.users:
    response = client.get_status(user_timeline.get(screen_name=i.screen_name)
    print "Got", len(response.data), "tweets from", i.screen_name
    if len(response.data) == 0:
        tweets = response.data[0]['created_at']
        tdate2 = datetime.strptime(tdate, "%a %b %d %H:%M:%S %Y")
        today = datetime.now()
        howlong = (today-tdate2).days
        if howlong >= daywindow:
            print i.screen_name, "has tweeted in the past", daywindow,
            totaltweets += len(response.data)
        for j in response.data:
            j.entities.urls
            for k in j.entities.urls:
                newurl = k.expanded_url
                urlset.add(newurl, j.user.screen_name)
    else:
        print i.screen_name, "has not tweeted in the past", daywindow
                    
```

Question answering

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	135
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	148

Train

Multimodal Models



CLIP



BLIP



FLAVA



ViLT



Sora



DALL-E



Gemini

..

Adaptation

Applications

Image/Vedio generation



Visual search



Data generation

Plot	Species	Sex	Weight
1/9/78	DM	M	40
1/9/78	DM	F	36
1/9/78	DS	F	135
1/20/78	DM	F	39
1/20/78	DM	M	43
1/20/78	DS	F	144
3/13/78	DM	F	51
3/13/78	DM	F	64
3/13/78	DS	F	148

Image analysis



Photo editing



...

VLA Models



OpenVLA



...

Noisy Data



Model Actions

Applications

Embodied AI



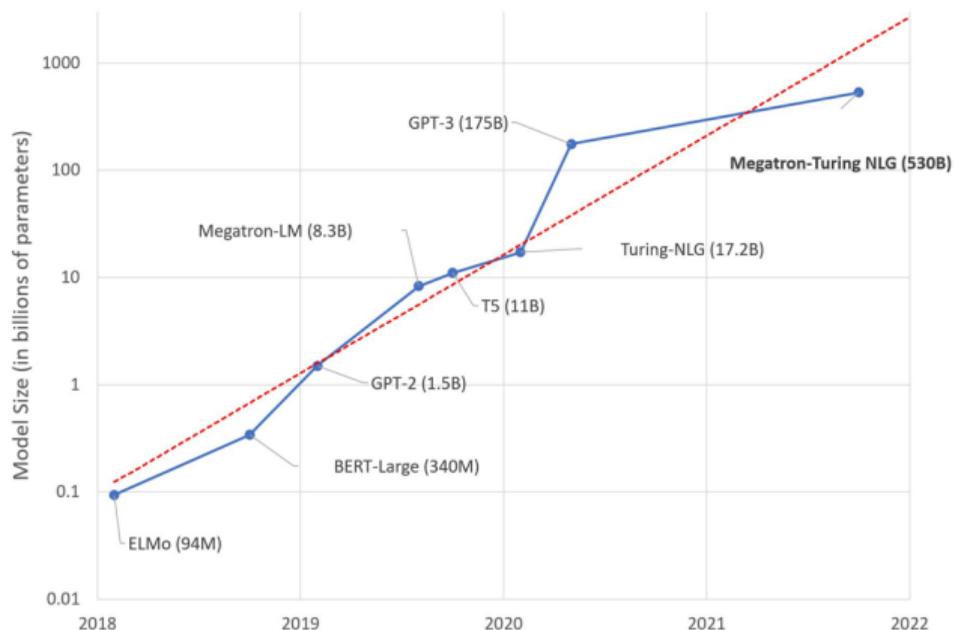
Robots



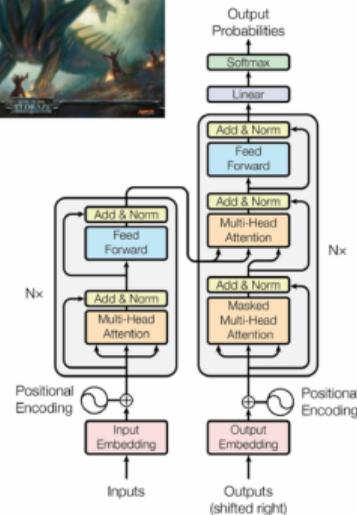
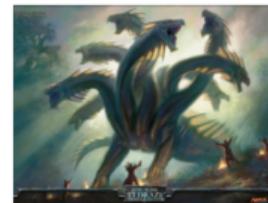
Autonomous Driving



Challenge 1: Increasing Model Sizes



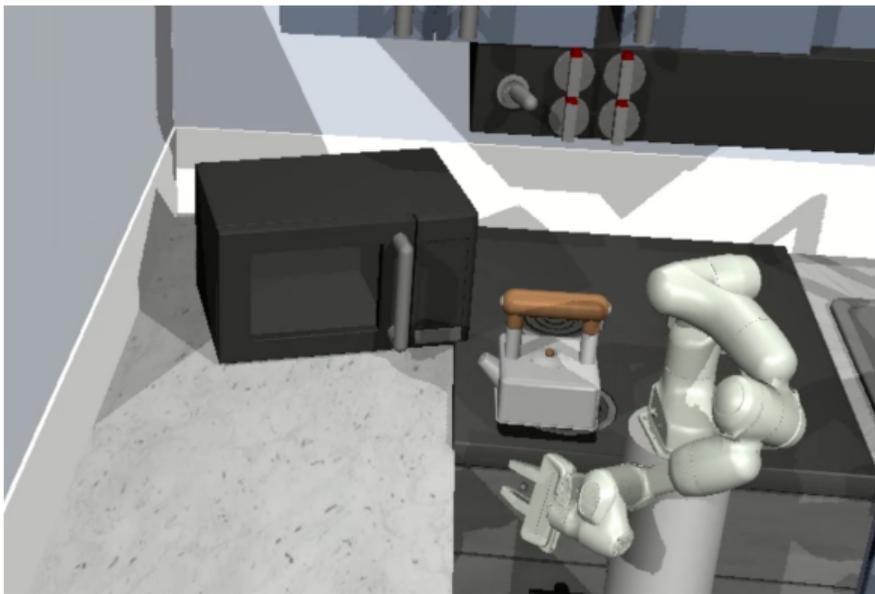
Source: <https://huggingface.co/blog/large-language-models>



Question:

How to save the computational resources?

Challenge 2: Noisy or Adversarial Data

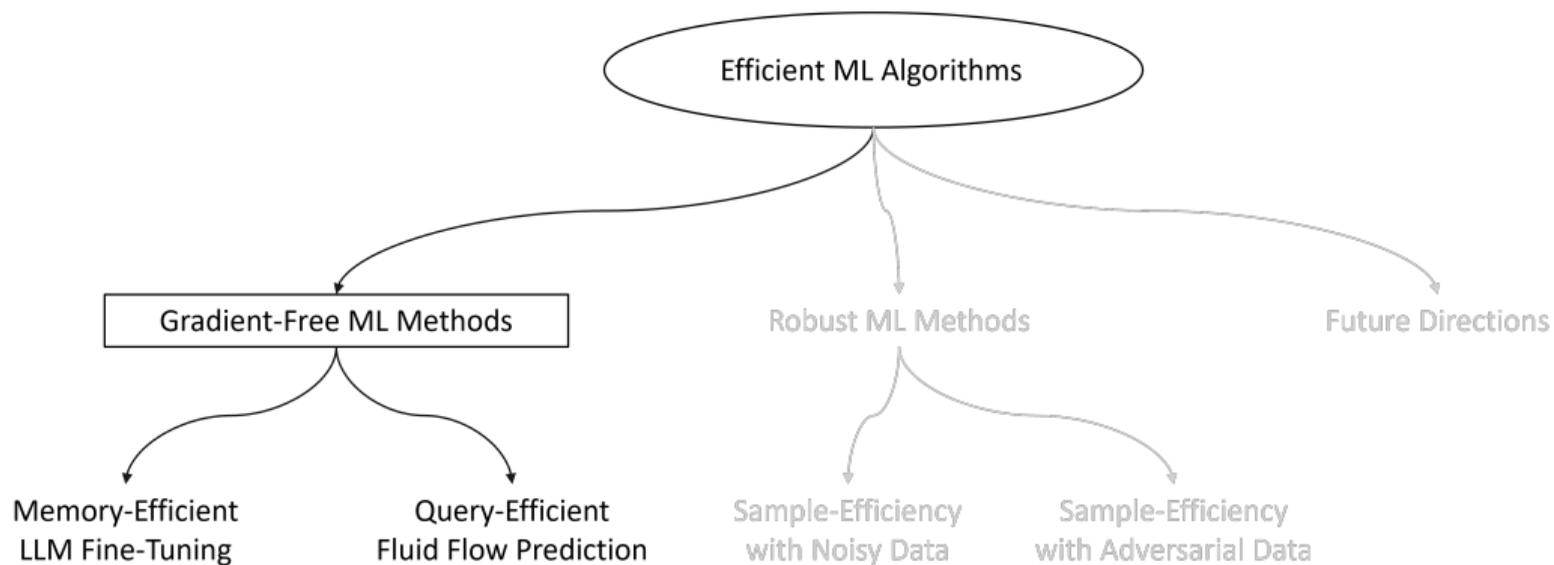


Source: Gu, Shangding, et al. "Robust Gymnasium: A Unified Modular Benchmark for Robust Reinforcement Learning." ICLR 2025.

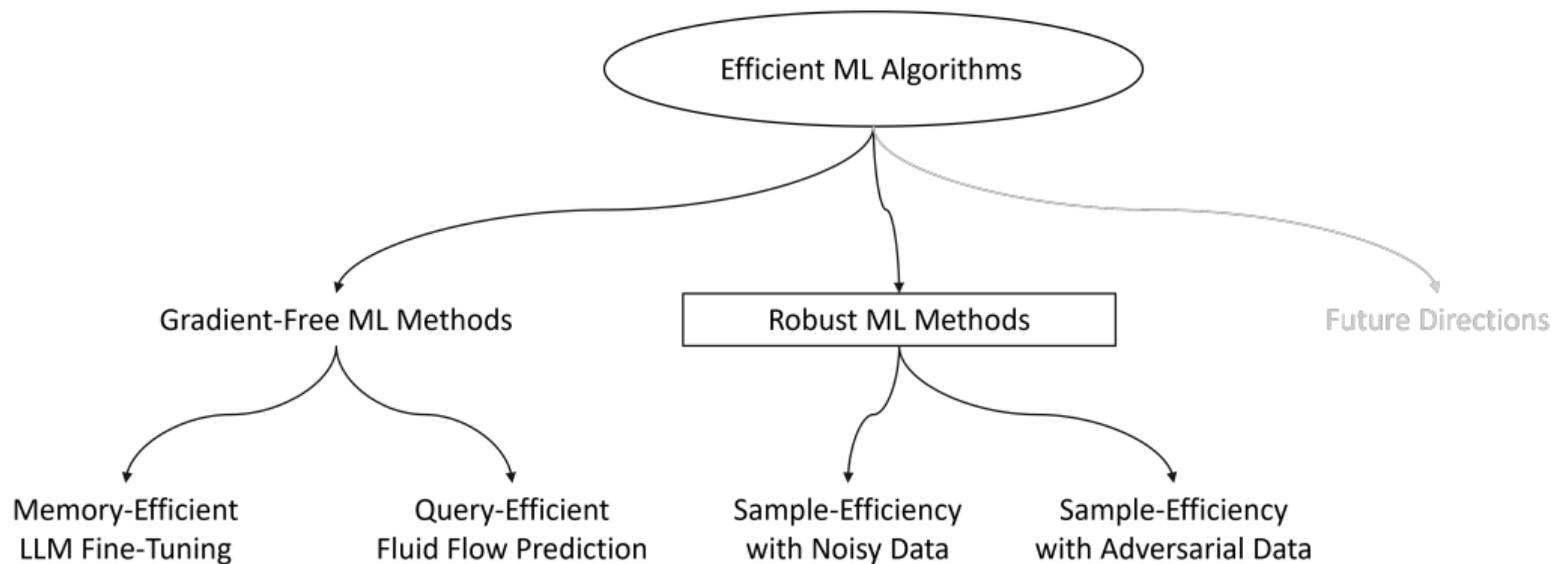
Question:

How to efficiently train a robust ML model?

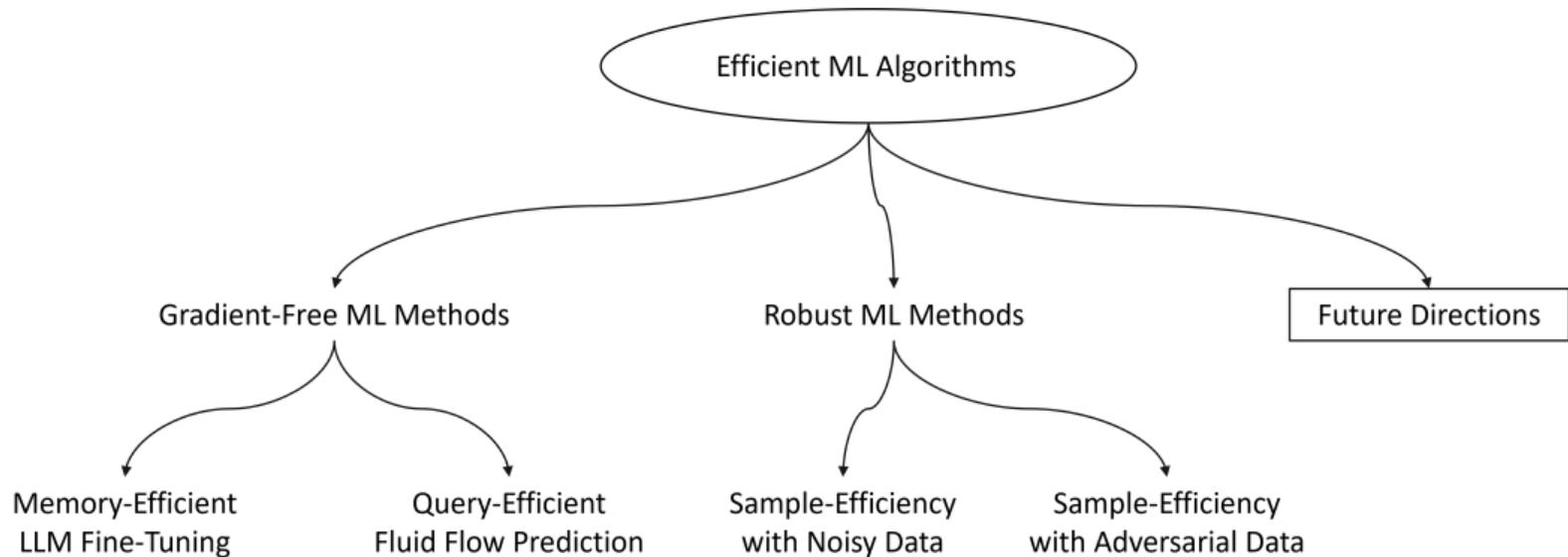
My Research Summary and Talk Outline



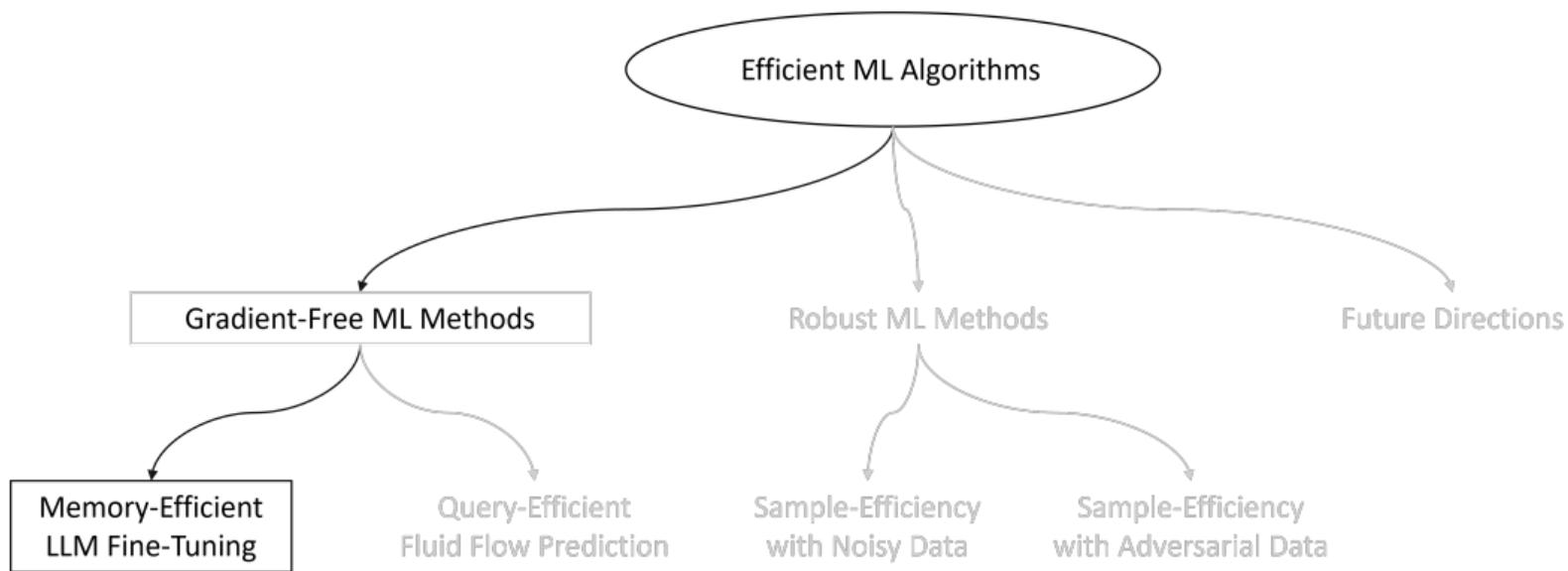
My Research Summary and Talk Outline



My Research Summary and Talk Outline



My Research Summary and Talk 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
- ...



Motivation: 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.
- η : 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.
- η : 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 $\mathcal{N}(0, I_d)$).
- μ : 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 $\mathcal{N}(0, I_d)$).
- μ : 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 minimum variance.
2. Propose a unbiased gradient estimator family.

Revisiting Zeroth-Order Optimization: Minimum-Variance Two-Point Estimators and Directionally Aligned Perturbations

Shaocong Ma, Heng Huang.

University of Maryland, College Park

ICLR 2025 Spotlight

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.

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:

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

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 :

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

Plug in the objective function:

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

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.

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

Minimize Estimation Error (Part II)

We analytically solve this functional optimization problem.

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

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.

\implies Both estimators achieve the **minimum variance**.

\implies DAPs have some nice properties.

Traditional Methods Cannot Identify the Important Directions

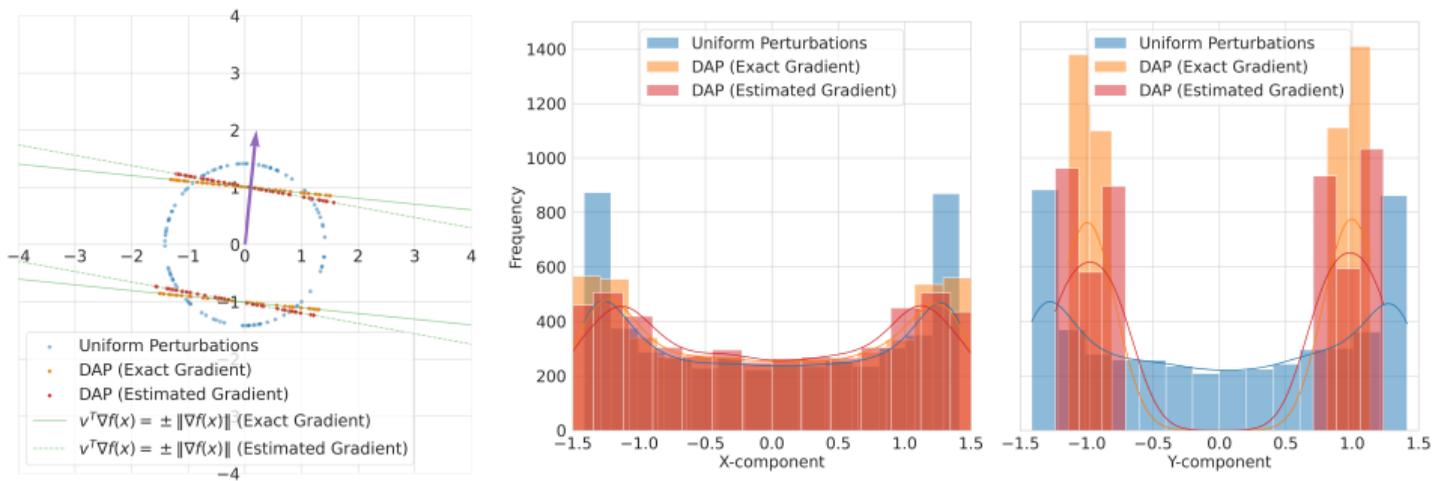


Figure 1: Illustration of the *directional alignment* property of DAP in $d = 2$ with estimating the gradient of $f(x) = x_1^2 + x_2^2$ at $x = [0.1 \ 1]^T$. Traditional estimator is **symmetric**, but we need a **non-symmetric** estimator.

Traditional Methods Cannot Identify the Important Directions

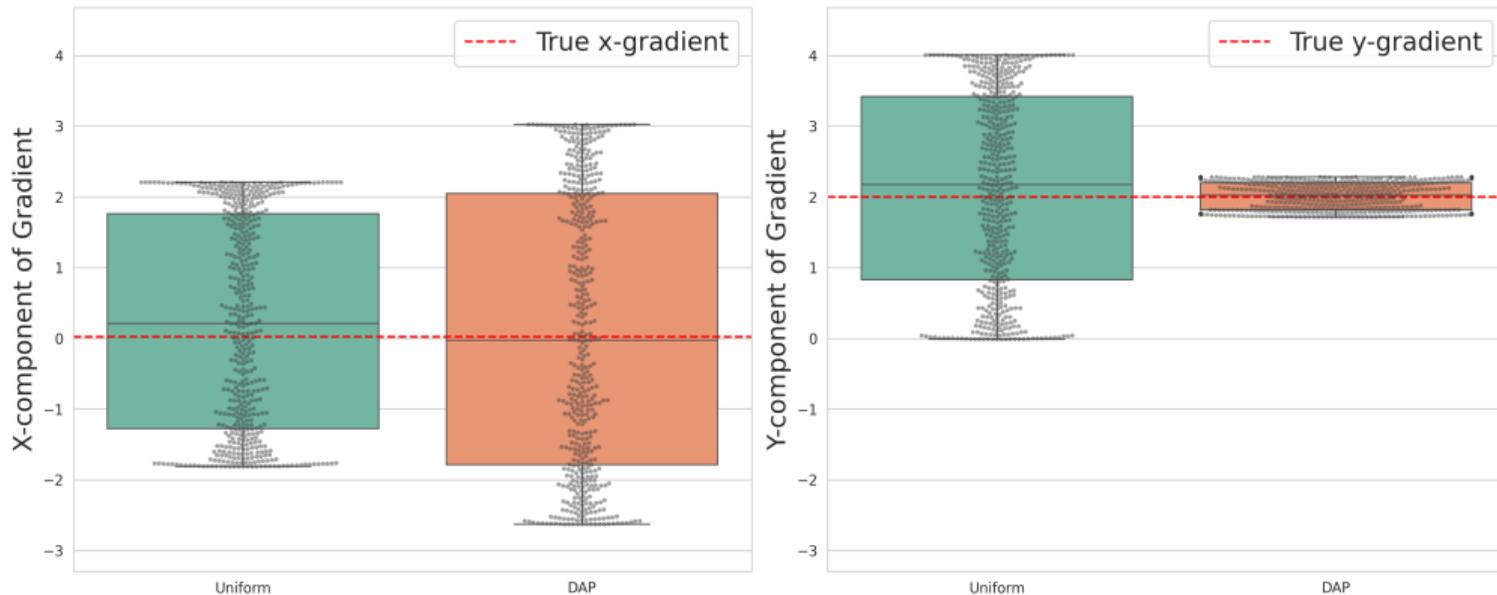


Figure 2: Comparison of gradient estimation performance with estimating the gradient of $f(x) = x_1^2 + x_2^2$ at $x = [0.1 \ 1]^T$ between uniform random perturbations and DAPs. The **non-symmetric** estimator is more accurate in the direction with larger gradient.

Applications in LLM Fine-Tuning

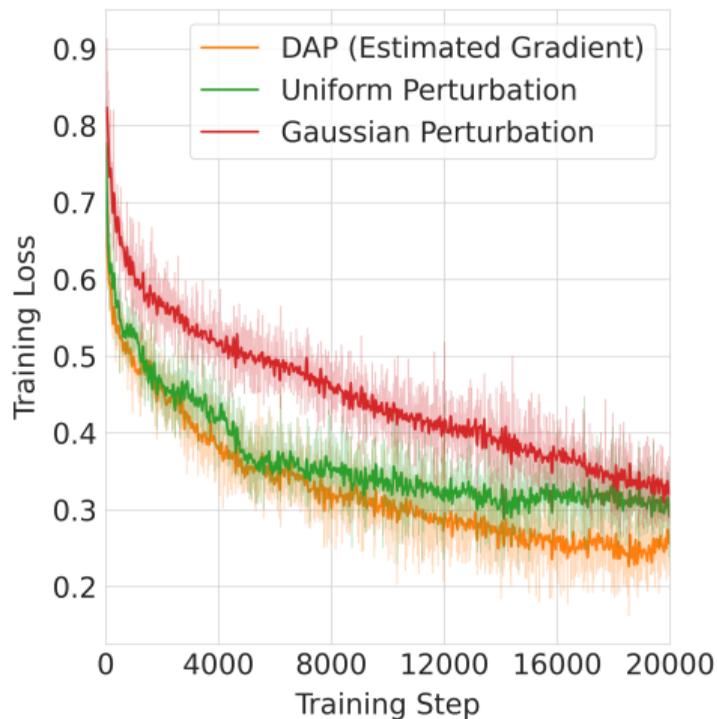


Figure 3: Comparison of training loss curves among different random perturbations on [Large Language Model Fine Tuning](#).

Applications in Scientific Optimization

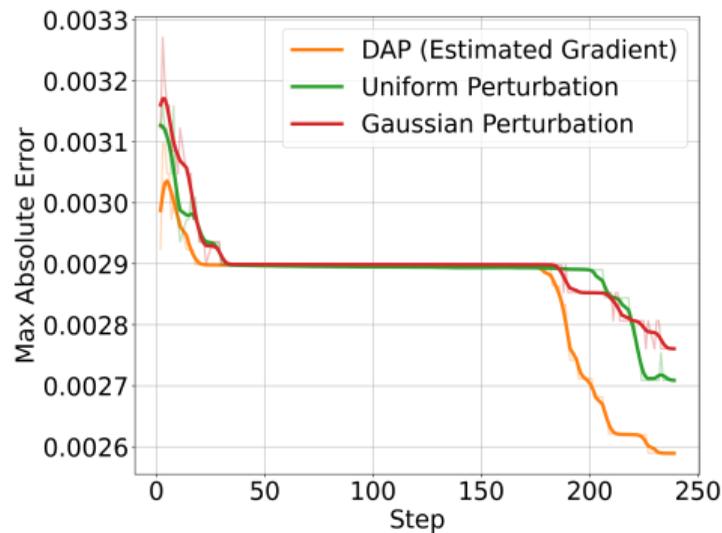


Figure 4: Comparison of training loss curves among different random perturbations on **Mesh Optimization for the Physical Numerical Solver**.

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?

On the Optimal Construction of Unbiased Gradient Estimators for Zeroth-Order Optimization

Shaocong Ma, Heng Huang.

University of Maryland, College Park

NeurIPS 2025 Spotlight

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 μ 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**

- 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.

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 $f_{\text{cls}} : \mathbb{R}^d \rightarrow \mathbb{R}$:

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

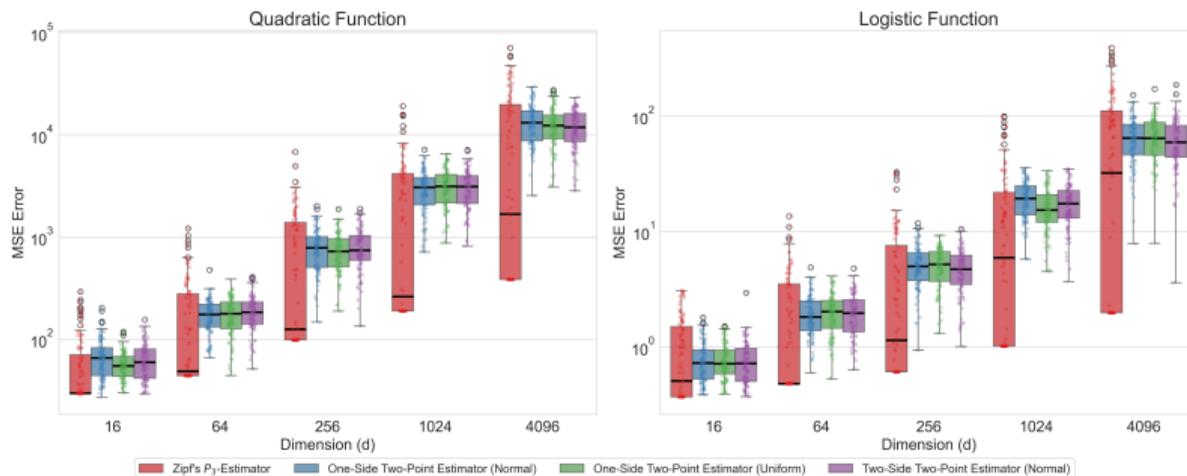


Figure 5: The MSE error (Left: f_{reg} , Right: f_{cls}) of different estimators.

Applications in LLM Fine-Tuning

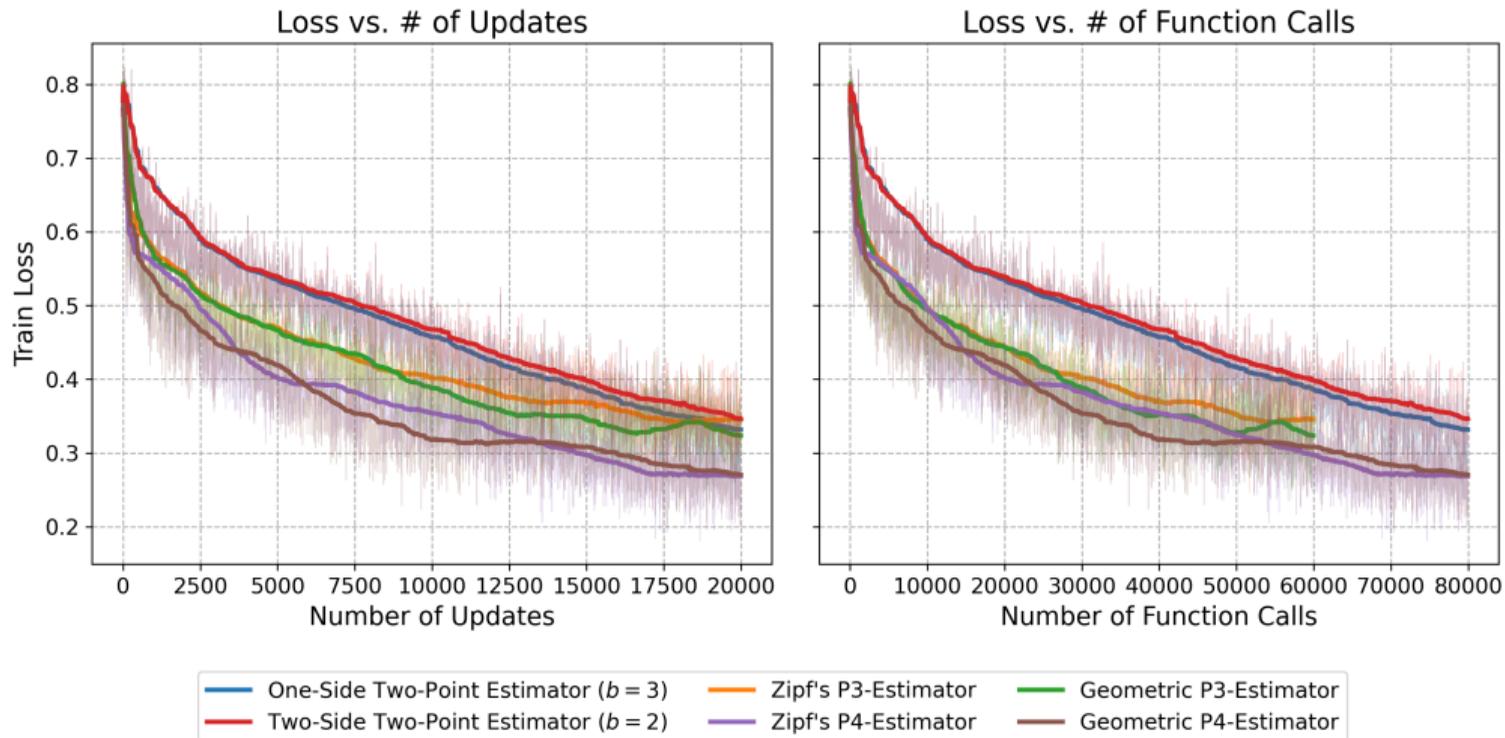


Figure 6: Fine-tuning the OPT-1.3B model on SST-2 using different gradient estimators.

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?

Riemannian Zeroth-Order Gradient Estimation with Structure-Preserving Metrics for Geodesically Incomplete Manifolds

Shaocong Ma, Heng Huang.

University of Maryland, College Park

ICLR 2026

Hyper-Octant Zeroth-Order Optimization: Fine-Tuning Quantized LLMs on the Positive Orthant Manifold

Shaocong Ma, Heng Huang.

University of Maryland, College Park

ICML 2026 Submission

Geometric Constraints in Quantized LLMs

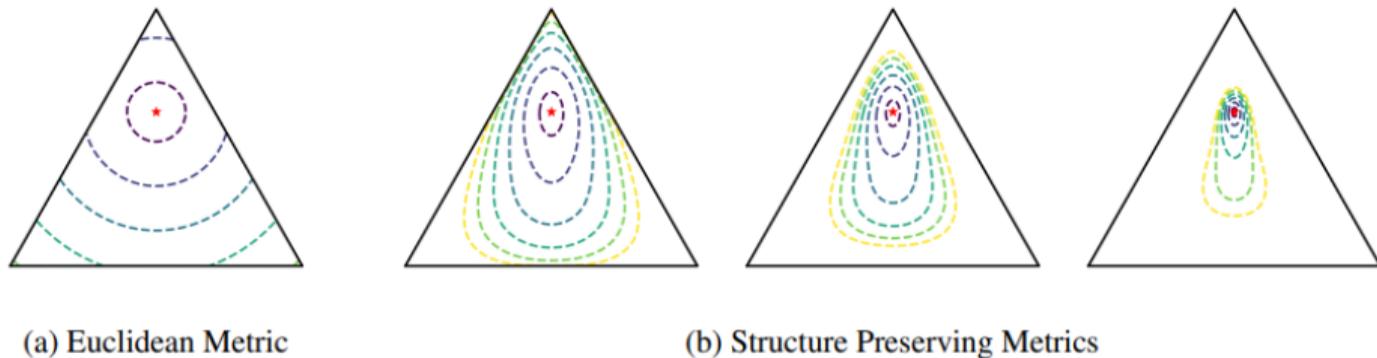


Figure 7: Visualization of different metrics on the probability simplex.

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

Memory-Efficient Fine-Tuning of Quantized LLMs

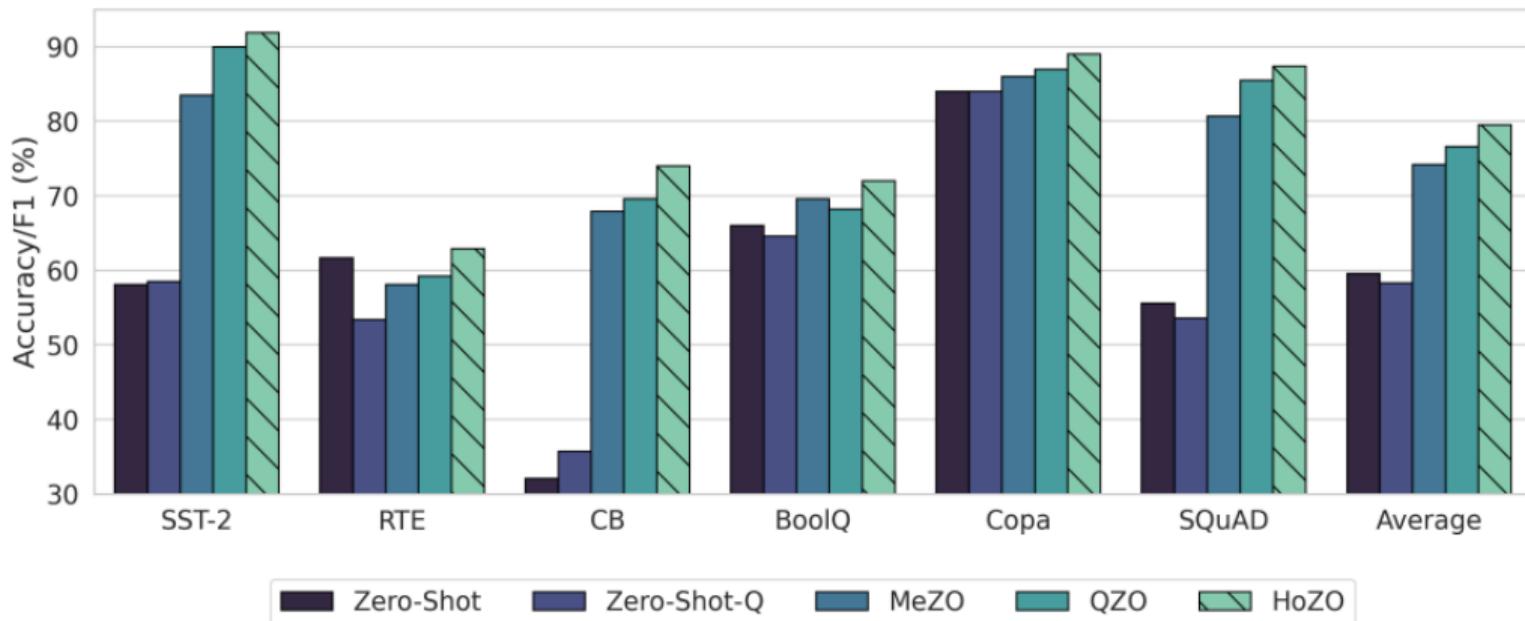


Figure 8: On the INT4 Llama-2-7B model across 6 downstream tasks, HoZO achieves consistently better performance than all baselines.

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
	HoZO	4 bits	91.9	62.9	74.0	72.0	89.0	87.4	79.5	↑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	77.0	81.6	92.0	87.5	81.2	↑17.4
	QZO	4 bits	91.9	62.8	77.0	82.4	92.0	89.2	82.5	↑18.7
	HoZO	4 bits	92.4	68.6	75.0	81.6	92.0	89.3	83.2	↑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	80.8	82.0	83.8	93.0	90.4	86.8	↑19.8
	HoZO	4 bits	91.5	79.1	83.0	83.8	95.0	91.2	87.3	↑20.3

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.

Efficient Machine Learning on Edge Devices

■ 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.

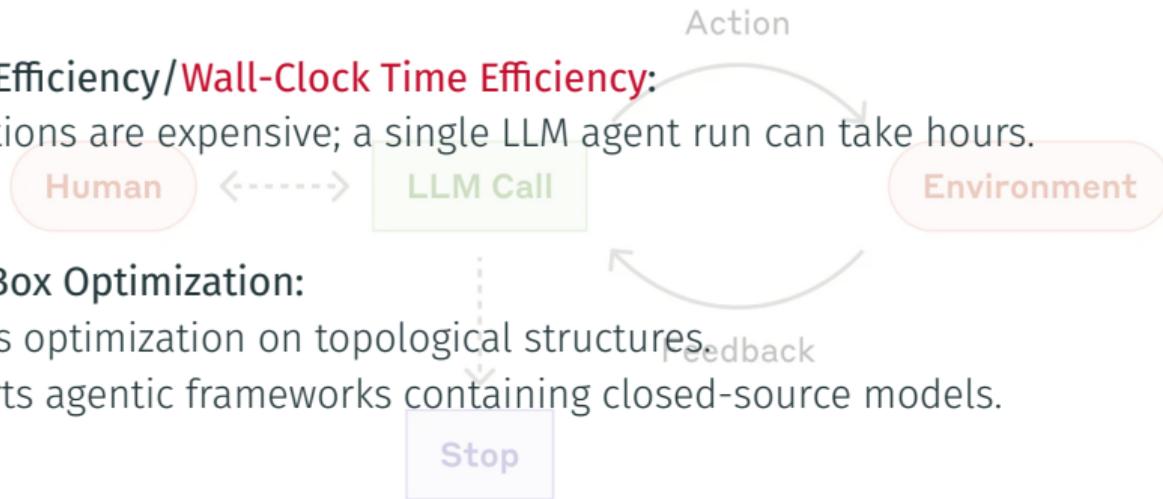
■ Challenges:

- Extreme memory constraints \implies Memory Efficiency.

- Extreme computational constraints \implies **Computation Efficiency.**

Fine-Tuning Agentic Framework with Zeroth-Order Optimization

- **Query Efficiency/Wall-Clock Time Efficiency:**
Evaluations are expensive; a single LLM agent run can take hours.



- **Black-Box Optimization:**
Enables optimization on topological structures.
Supports agentic frameworks containing closed-source models.

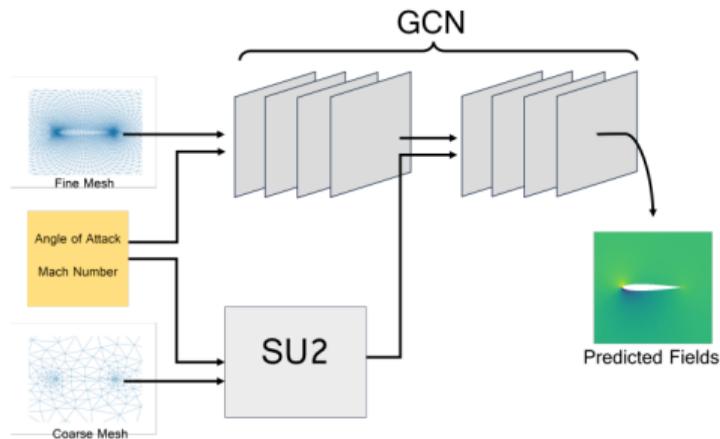


Figure 9: The CFD-GCN model. The PDE solver can be slow. How to improve the query efficiency?

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

Future Research Directions and Strategies



■ Continue collaborations with

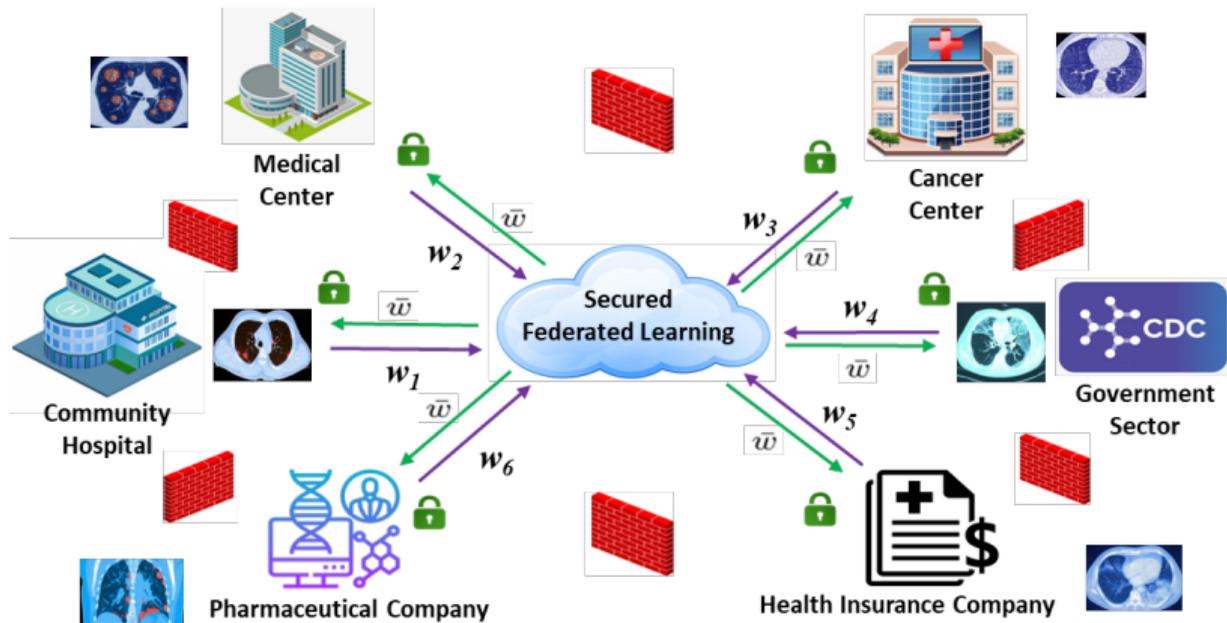
- UMD, Univ. of Utah, USC, TAMU, OSU, ASU, Buffalo, *etc.*
- Lawrence Livermore National Lab, MD Anderson Cancer Center, *etc.*

■ Seek more collaborations with

- Kent CS (Core Machine Learning, AI), Kent Physics, MS, A&E *etc.* (AI4Science, Robotics)
- Brain Health Research Institute (AI4Healthcare)
- Local Industrials (Cleveland Clinic, GE HealthCare, Progressive Insurance, *etc.*)
- General Industrials (Google, Amazon, UnitedHealthcare, *etc.*)

Collaboration: AI for Healthcare & Federated Learning

Data Privacy, Trustworthy, Communication Efficiency ...



Other Experiences

Successful Experience in Helping 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.

My Future Research Funding Plan

- First several proposals: NSF MFAI, NSF Early Career Development, NSF-CISE, *etc.*
- Collaborate with colleagues at Kent to seek: NSF Core Medium, 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

Sample Syllabus: Reinforcement Learning (1/2)

Target Audience & Prerequisites

- **Audience:** Senior Undergraduate / First-year Graduate Students
- **Prerequisites:** Linear Algebra, Probability, Proficiency in Python (PyTorch)

Part I: Foundations (Weeks 1-5)

- Markov Decision Processes (MDPs) & Bellman Equations
- Tabular Methods: Dynamic Programming, Monte Carlo
- Temporal-Difference Learning (TD-Learning, Q-Learning)

Part II: Deep Reinforcement Learning (Weeks 6-10)

- **Value-Based:** Deep Q-Networks (DQN) and variants (Double, Dueling)
- **Policy Gradient:** REINFORCE, TRPO, PPO
- **Actor-Critic**

Sample Syllabus: Reinforcement Learning (2/2)

Part III: Frontiers & Applications (Weeks 11-14)

- Offline Reinforcement Learning
- Multi-Agent Systems (MARL)
- RL for Large Language Models (RLHF & DPO)

Grading Scheme: Research-Oriented

- **Coding Assignments (30%):**
 - Implement algorithms from scratch (e.g., PPO) using PyTorch
- **Midterm Exam (30%):**
 - Theoretical concepts and paper reviews
- **Final Project (40%):**
 - Open-ended research project (Team of 2-3)
 - Encouraged to reproduce recent NeurIPS/ICLR papers or apply RL to new domains (e.g., LLM Optimization)

Thank You!