

# Stochastic **Auto-Conditioned** Fast Gradient Methods with Optimal Rates

---

Yao Ji & Guanghui (George) Lan

H. Milton Stewart School of Industrial and Systems Engineering, Georgia Tech

PART I

# | Problem Setup

## What This Talk Is About

---

We study composite stochastic optimization:

$$\Phi^* = \min_{x \in X} \{f(x) + h(x)\}, \quad X \subseteq \mathbb{R}^n \text{ closed convex}$$

- >  $f$ : convex, differentiable,  $L$ -smooth with **unknown**  $L$
- >  $h$ : closed, convex,  $\text{prox}_h$  efficiently computable
- > Gradients via stochastic oracle  $G(x, \xi)$

### Main Question

Can we achieve **optimal convergence rates** without knowing  $L$  and without line search?

# Accelerated Unbounded Stochastic

---

KREISLER ET AL., '24

OURS

**Distribution**

Restricted: subgaussian ( $\sigma_x$ )

Bounded variance

# Accelerated Unbounded Stochastic

---

	KREISLER ET AL., '24	OURS
<b>Distribution</b>	Restricted: subgaussian ( $\sigma_x$ )	Bounded variance
<b>Rate</b>	Sub-optimal $\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{\epsilon}} + \frac{\sigma^2}{\epsilon^2} + \frac{\sigma_*^2}{\epsilon}\right)$	Optimal $\mathcal{O}\left(\frac{1}{\sqrt{\epsilon}} + \frac{\sigma^2}{\epsilon^2}\right)$

# Accelerated Unbounded Stochastic

---

	KREISLER ET AL., '24	OURS
<b>Distribution</b>	Restricted: subgaussian ( $\sigma_x$ )	Bounded variance
<b>Rate</b>	Sub-optimal $\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{\varepsilon}} + \frac{\sigma^2}{\varepsilon^2} + \frac{\sigma_*^2}{\varepsilon}\right)$	Optimal $\mathcal{O}\left(\frac{1}{\sqrt{\varepsilon}} + \frac{\sigma^2}{\varepsilon^2}\right)$
<b>Parameter</b>	Need subgaussian $\sigma(x)$	Variance estimation

# Accelerated Unbounded Stochastic

	KREISLER ET AL., '24	OURS
<b>Distribution</b>	Restricted: subgaussian ( $\sigma_x$ )	Bounded variance
<b>Rate</b>	Sub-optimal $\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{\varepsilon}} + \frac{\sigma^2}{\varepsilon^2} + \frac{\sigma_*^2}{\varepsilon}\right)$	Optimal $\mathcal{O}\left(\frac{1}{\sqrt{\varepsilon}} + \frac{\sigma^2}{\varepsilon^2}\right)$
<b>Parameter</b>	Need subgaussian $\sigma(x)$	Variance estimation
<b>Gradient</b>	Bounded gradient, minibatch	No such assumption

# Accelerated Unbounded Stochastic

	KREISLER ET AL., '24	OURS
<b>Distribution</b>	Restricted: subgaussian ( $\sigma_x$ )	Bounded variance
<b>Rate</b>	Sub-optimal $\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{\varepsilon}} + \frac{\sigma^2}{\varepsilon^2} + \frac{\sigma_*^2}{\varepsilon}\right)$	Optimal $\mathcal{O}\left(\frac{1}{\sqrt{\varepsilon}} + \frac{\sigma^2}{\varepsilon^2}\right)$
<b>Parameter</b>	Need subgaussian $\sigma(x)$	Variance estimation
<b>Gradient</b>	Bounded gradient, minibatch	No such assumption
<b>Horizon</b>	Need total iteration $N$	$N$ free, auto-conditioned

PART II

# | The Algorithm

## Recap: Auto-Conditioned Fast Gradient Method (Li & Lan, 2024)

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle \nabla f(x_{k-1}), z \rangle + \frac{1}{2\eta_k} \|z - y_{k-1}\|^2 \right\}$$

$$x_k = \frac{\tau_k}{1 + \tau_k} x_{k-1} + \frac{1}{1 + \tau_k} z_k \quad y_k = (1 - \beta_k) y_{k-1} + \beta_k z_k$$

where

$$\eta_k \propto \min \left\{ \frac{k\eta_{k-1}}{k-1}, \frac{k}{\bar{L}_{k-1}} \right\}$$

$$\bar{L}_{k-1} = \frac{\|\nabla f(x_{k-1}) - \nabla f(x_{k-2})\|^2}{2[f(x_{k-2}) - f(x_{k-1}) - \langle \nabla f(x_{k-1}), x_{k-2} - x_{k-1} \rangle]}$$

# Stochastic Auto-Conditioned Fast Gradient Method

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle G(x_{k-1}, \xi), z \rangle + \frac{1}{2\eta_k} \|z - y_{k-1}\|^2 \right\}$$

$$x_k = \frac{\tau_k}{1 + \tau_k} x_{k-1} + \frac{1}{1 + \tau_k} z_k \quad y_k = (1 - \beta_k) y_{k-1} + \beta_k z_k$$

where

$$\eta_k \propto \min \left\{ \frac{k\eta_{k-1}}{k-1}, \frac{k}{\bar{L}_{k-1}} \right\}$$

$$\bar{L}_{k-1} = \frac{\|G(x_{k-1}, \bar{\xi}) - G(x_{k-2}, \bar{\xi})\|^2}{2 [F(x_{k-2}, \bar{\xi}) - F(x_{k-1}, \bar{\xi}) - \langle G(x_{k-1}, \bar{\xi}), x_{k-2} - x_{k-1} \rangle]}$$

# Stochastic Auto-Conditioned Fast Gradient Method

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle G(x_{k-1}, \xi), z \rangle + \frac{1}{2\eta_k} \|z - y_{k-1}\|^2 \right\}$$

$$x_k = \frac{\tau_k}{1 + \tau_k} x_{k-1} + \frac{1}{1 + \tau_k} z_k \quad y_k = (1 - \beta_k) y_{k-1} + \beta_k z_k$$

where

**ATTEMPT FAILS!**

$$\eta_k \propto \min \left\{ \frac{k\eta_{k-1}}{k-1}, \frac{1}{\bar{L}_{k-1}} \right\}$$

$$\bar{L}_{k-1} = \frac{\|G(x_{k-1}, \bar{\xi}) - G(x_{k-2}, \bar{\xi})\|^2}{2 [F(x_{k-2}, \bar{\xi}) - F(x_{k-1}, \bar{\xi}) - \langle G(x_{k-1}, \bar{\xi}), x_{k-2} - x_{k-1} \rangle]}$$

# Three Batch Types & Filtration

---

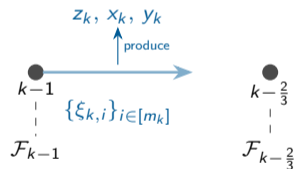
●  
 $k-1$   
⋮  
 $\mathcal{F}_{k-1}$

## Three Batch Types & Filtration

### Gradient Est.

Samples  $\{\xi_{k,i}\}_{i=1}^{m_k}$

Compute  $G_k, z_k, x_k, y_k$



## Type I: iterate update batch $\xi_{k,i}$

---

### Gradient estimate

$$G_k = \frac{1}{m_k} \sum_{i=1}^{m_k} G(x_{k-1}, \xi_{k,i})$$

## Type I: iterate update batch $\xi_{k,i}$

---

### Gradient estimate

$$G_k = \frac{1}{m_k} \sum_{i=1}^{m_k} G(x_{k-1}, \xi_{k,i})$$

### Proximal step

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle G_k, z \rangle + h(z) + \frac{1}{2\eta_k} \|y_{k-1} - z\|^2 \right\}$$

## Type I: iterate update batch $\xi_{k,i}$

### Gradient estimate

$$G_k = \frac{1}{m_k} \sum_{i=1}^{m_k} G(x_{k-1}, \xi_{k,i})$$

### Proximal step

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle G_k, z \rangle + h(z) + \frac{1}{2\eta_k} \|y_{k-1} - z\|^2 \right\}$$

### Output solution

$$x_k = \frac{z_k + \tau_k x_{k-1}}{1 + \tau_k}, \quad \tau_k = \frac{k}{2}$$

## Type I: iterate update batch $\xi_{k,i}$

### Gradient estimate

$$G_k = \frac{1}{m_k} \sum_{i=1}^{m_k} G(x_{k-1}, \xi_{k,i})$$

### Proximal step

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle G_k, z \rangle + h(z) + \frac{1}{2\eta_k} \|y_{k-1} - z\|^2 \right\}$$

### Output solution

$$x_k = \frac{z_k + \tau_k x_{k-1}}{1 + \tau_k}, \quad \tau_k = \frac{k}{2}$$

### Output center

$$y_k = (1 - \beta_k) y_{k-1} + \beta_k z_k, \quad \beta_1 = 0, \quad \beta_k \equiv \beta$$

## Type I: iterate update batch $\xi_{k,i}$

### Gradient estimate

$$G_k = \frac{1}{m_k} \sum_{i=1}^{m_k} G(x_{k-1}, \xi_{k,i})$$

### Proximal step

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle G_k, z \rangle + h(z) + \frac{1}{2\eta_k} \|y_{k-1} - z\|^2 \right\}$$

### Output solution

$$x_k = \frac{z_k + \tau_k x_{k-1}}{1 + \tau_k}, \quad \tau_k = \frac{k}{2}$$

### Output center

$$y_k = (1 - \beta_k) y_{k-1} + \beta_k z_k, \quad \beta_1 = 0, \quad \beta_k \equiv \beta$$

For the next iteration: compute  $\eta_{k+1}$ ,  $m_{k+1}$ ,  $n_{k+1}$

## Three Batch Types & Filtration

### Gradient Est.

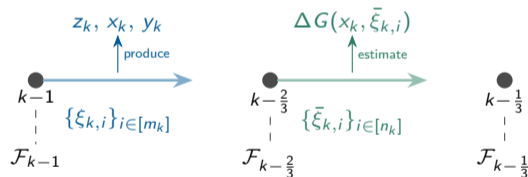
Samples  $\{\xi_{k,i}\}_{i=1}^{m_k}$

Compute  $G_k, z_k, x_k, y_k$

### Stepsize update. I

Samples  $\{\bar{\xi}_{k,i}\}$

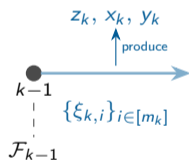
Compute  $\Delta G$



## Three Batch Types & Filtration

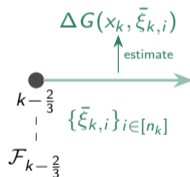
### Gradient Est.

Samples  $\{\xi_{k,i}\}_{i=1}^{m_k}$   
 Compute  $G_k, z_k, x_k, y_k$



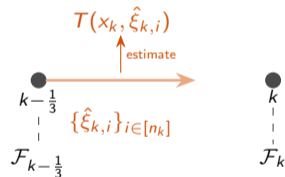
### Stepsize update. I

Samples  $\{\bar{\xi}_{k,i}\}$   
 Compute  $\Delta G$



### Stepsize update. II

Samples  $\{\hat{\xi}_{k,i}\}$   
 Compute  $T$



**Type II: stepsize  $\eta_{k+1}$  update batches  $\bar{\xi}_{k,i}$  ,  $\hat{\xi}_{k,i}$**

---

## Type II: stepsize $\eta_{k+1}$ update batches $\bar{\xi}_{k,i}$ , $\hat{\xi}_{k,i}$

---

### 1. Stochastic gradient difference:

$$\Delta G(x_k, \bar{\xi}_k) := \frac{1}{n_k} \sum_{i=1}^{n_k} [G(x_k, \bar{\xi}_{k,i}) - G(x_{k-1}, \bar{\xi}_{k,i})].$$

## Type II: stepsize $\eta_{k+1}$ update batches $\bar{\xi}_{k,i}$ , $\hat{\xi}_{k,i}$

---

### 1. Stochastic gradient difference:

$$\Delta G(x_k, \bar{\xi}_k) := \frac{1}{n_k} \sum_{i=1}^{n_k} [G(x_k, \bar{\xi}_{k,i}) - G(x_{k-1}, \bar{\xi}_{k,i})].$$

### 2. Empirical first-order Taylor remainder:

$$T(x_k, \hat{\xi}_k) := \frac{1}{n_k} \sum_{i=1}^{n_k} [F(x_{k-1}, \hat{\xi}_{k,i}) - F(x_k, \hat{\xi}_{k,i}) - \langle G(x_k, \hat{\xi}_{k,i}), x_{k-1} - x_k \rangle]$$

## Type II: stepsize $\eta_{k+1}$ update batches $\bar{\xi}_{k,i}$ , $\hat{\xi}_{k,i}$

### 1. Stochastic gradient difference:

$$\Delta G(x_k, \bar{\xi}_k) := \frac{1}{n_k} \sum_{i=1}^{n_k} [G(x_k, \bar{\xi}_{k,i}) - G(x_{k-1}, \bar{\xi}_{k,i})].$$

### 2. Empirical first-order Taylor remainder:

$$T(x_k, \hat{\xi}_k) := \frac{1}{n_k} \sum_{i=1}^{n_k} \left[ F(x_{k-1}, \hat{\xi}_{k,i}) - F(x_k, \hat{\xi}_{k,i}) - \langle G(x_k, \hat{\xi}_{k,i}), x_{k-1} - x_k \rangle \right]$$

### 3. The ratio: empirical local cocoercivity-based smoothness estimator

$$\bar{L}_k := \begin{cases} \frac{\|\Delta G(x_k, \bar{\xi}_{k,i})\|^2}{2T(x_k, \hat{\xi}_{k,i})}, & \text{if } T(x_k, \hat{\xi}_{k,i}) > 0, \\ 0, & \text{if } T(x_k, \hat{\xi}_{k,i}) = 0. \end{cases}$$

## Type II: stepsize $\eta_{k+1}$ update batches $\bar{\xi}_{k,i}$ , $\hat{\xi}_{k,i}$

### 1. Stochastic gradient difference:

$$\Delta G(x_k, \bar{\xi}_k) := \frac{1}{n_k} \sum_{i=1}^{n_k} [G(x_k, \bar{\xi}_{k,i}) - G(x_{k-1}, \bar{\xi}_{k,i})].$$

### 2. Empirical first-order Taylor remainder:

$$T(x_k, \hat{\xi}_k) := \frac{1}{n_k} \sum_{i=1}^{n_k} \left[ F(x_{k-1}, \hat{\xi}_{k,i}) - F(x_k, \hat{\xi}_{k,i}) - \langle G(x_k, \hat{\xi}_{k,i}), x_{k-1} - x_k \rangle \right]$$

### 3. The ratio: empirical local cocoercivity-based smoothness estimator

$$\bar{L}_k := \begin{cases} \frac{\|\Delta G(x_k, \bar{\xi}_k)\|^2}{2T(x_k, \hat{\xi}_k)}, & \text{if } T(x_k, \hat{\xi}_k) > 0, \\ 0, & \text{if } T(x_k, \hat{\xi}_k) = 0. \end{cases}$$

- > Computes the ratio of **stochastic gradient variation** to **stochastic local curvature**.

## Type II: stepsize $\eta_{k+1}$ update batches $\bar{\xi}_{k,i}$ , $\hat{\xi}_{k,i}$

### 1. Stochastic gradient difference:

$$\Delta G(x_k, \bar{\xi}_k) := \frac{1}{n_k} \sum_{i=1}^{n_k} [G(x_k, \bar{\xi}_{k,i}) - G(x_{k-1}, \bar{\xi}_{k,i})].$$

### 2. Empirical first-order Taylor remainder:

$$T(x_k, \hat{\xi}_k) := \frac{1}{n_k} \sum_{i=1}^{n_k} \left[ F(x_{k-1}, \hat{\xi}_{k,i}) - F(x_k, \hat{\xi}_{k,i}) - \langle G(x_k, \hat{\xi}_{k,i}), x_{k-1} - x_k \rangle \right]$$

### 3. The ratio: empirical local cocoercivity-based smoothness estimator

$$\bar{L}_k := \begin{cases} \frac{\|\Delta G(x_k, \bar{\xi}_{k,i})\|^2}{2T(x_k, \hat{\xi}_{k,i})}, & \text{if } T(x_k, \hat{\xi}_{k,i}) > 0, \\ 0, & \text{if } T(x_k, \hat{\xi}_{k,i}) = 0. \end{cases}$$

- > Computes the ratio of **stochastic gradient variation** to **stochastic local curvature**.
- >  $\bar{L}_k$  cleanly proxies the local Lipschitz constant near  $x_k$  to set the next stepsize  $\eta_{k+1}$ .

## Stochastic adaptive stepsize $\eta_{k+1}$ update

---

Using  $\bar{L}_k$ , we define the stepsize recursively as follows. Let  $\eta_1 > 0$  and define

$$\eta_2 = \min \left\{ \frac{1}{16\beta\bar{L}_1}, 2(1 - \beta)\eta_1 \right\}, \quad \eta_{k+1} = \min \left\{ \frac{k}{16\bar{L}_k}, \frac{(k+1)\eta_k}{k} \right\}, \quad \forall k \geq 2, \text{ a.s.}$$

- > Large  $\bar{L}_k \Rightarrow$  small  $\eta_{k+1}$ : **conservative step** near curved regions  $x_k$
- > Small  $\bar{L}_k \Rightarrow \eta_{k+1}$  may grow as  $\mathcal{O}(k)$ : **accelerated regime**
- > No knowledge of global  $L$  at any point in the algorithm

For the next iteration: compute  $m_{k+1}, n_{k+1}$

## Adaptive minibatch $m_{k+1}$ for iterate update

Minibatch in parameter free setting

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ v_k^{\max}, \frac{\sigma_k^2}{D_0^2} \right\} \right\rceil \right)$$

$$\text{Local } \bar{L}_k, \sigma_k^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

- >  $\sigma_k^2$  : local variance of the stochastic gradient error
- >  $v_k^{\max}$  : local variance of the local smoothness estimator

## Adaptive minibatch $m_{k+1}$ for iterate update

### Minibatch in parameter free setting

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ v_k^{\max}, \frac{\sigma_k^2}{D_0^2} \right\} \right\rceil \right)$$

Local  $\bar{L}_k, \sigma_k^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$

### Minibatch AS-SA Ghadimi & Lan 2016

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{L^2} \cdot \frac{\sigma^2}{D_0^2} \right\rceil \right)$$

Global smoothness  $L$ , variance  $\sigma^2$

- >  $\sigma_k^2$  : local variance of the stochastic gradient error
- >  $v_k^{\max}$  : local variance of the local smoothness estimator

## Adaptive minibatch $m_{k+1}$ for iterate update

### Minibatch in parameter free setting

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ v_k^{\max}, \frac{\sigma_k^2}{D_0^2} \right\} \right\rceil \right)$$

Local  $\bar{L}_k, \sigma_k^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$

### Minibatch AS-SA Ghadimi & Lan 2016

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{L^2} \cdot \frac{\sigma^2}{D_0^2} \right\rceil \right)$$

Global smoothness  $L$ , variance  $\sigma^2$

- >  $\sigma_k^2$  : local variance of the stochastic gradient error
- >  $v_k^{\max}$  : local variance of the local smoothness estimator
- >  $m_{k+1}$  is chosen to control the variability introduced by the previous stepsize  $\eta_{k+1} \propto \bar{L}_k^{-1}$

## Adaptive minibatch $m_{k+1}$ for iterate update

### Minibatch in parameter free setting

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ v_k^{\max}, \frac{\sigma_k^2}{D_0^2} \right\} \right\rceil \right)$$

$$\text{Local } \bar{L}_k, \sigma_k^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

### Minibatch AS-SA Ghadimi & Lan 2016

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{L^2} \cdot \frac{\sigma^2}{D_0^2} \right\rceil \right)$$

Global smoothness  $L$ , variance  $\sigma^2$

- >  $\sigma_k^2$  : local variance of the stochastic gradient error
- >  $v_k^{\max}$  : local variance of the local smoothness estimator
- >  $m_{k+1}$  is chosen to control the variability introduced by the previous stepsize  $\eta_{k+1} \propto \bar{L}_k^{-1}$
- >  $\bar{L}_k^{-1}$  is not unbiased!  $\implies v_k^{\max}$  controls the bias of  $\bar{L}_k^{-1}$ .

## Adaptive minibatch $m_{k+1}$ for iterate update

### Minibatch in parameter free setting

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ v_k^{\max}, \frac{\sigma_k^2}{D_0^2} \right\} \right\rceil \right)$$

Local  $\bar{L}_k, \sigma_k^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$

### Minibatch AS-SA Ghadimi & Lan 2016

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{L^2} \cdot \frac{\sigma^2}{D_0^2} \right\rceil \right)$$

Global smoothness  $L$ , variance  $\sigma^2$

- >  $\sigma_k^2$  : local variance of the stochastic gradient error
- >  $v_k^{\max}$  : local variance of the local smoothness estimator
- >  $m_{k+1}$  is chosen to control the variability introduced by the previous stepsize  $\eta_{k+1} \propto \bar{L}_k^{-1}$
- >  $\bar{L}_k^{-1}$  is not unbiased!  $\implies v_k^{\max}$  controls the bias of  $\bar{L}_k^{-1}$ .
- >  $m_{k+1}$  is determined by previously observed quantities:  $v_k^{\max}$ ,  $\bar{L}_k$ , and  $\sigma_k^2$  — **Adaptive!**

## Adaptive minibatch $m_{k+1}$ for iterate update

### Minibatch in parameter free setting

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ v_k^{\max}, \frac{\sigma_k^2}{D_0^2} \right\} \right\rceil \right)$$

Local  $\bar{L}_k, \sigma_k^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$

### Minibatch AS-SA Ghadimi & Lan 2016

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{L^2} \cdot \frac{\sigma^2}{D_0^2} \right\rceil \right)$$

Global smoothness  $L$ , variance  $\sigma^2$

- >  $\sigma_k^2$  : local variance of the stochastic gradient error
- >  $v_k^{\max}$  : local variance of the local smoothness estimator
- >  $m_{k+1}$  is chosen to control the variability introduced by the previous stepsize  $\eta_{k+1} \propto \bar{L}_k^{-1}$
- >  $\bar{L}_k^{-1}$  is not unbiased!  $\implies v_k^{\max}$  controls the bias of  $\bar{L}_k^{-1}$ .
- >  $m_{k+1}$  is determined by previously observed quantities:  $v_k^{\max}$ ,  $\bar{L}_k$ , and  $\sigma_k^2$  — **Adaptive!**

For the next iteration: compute  $n_{k+1}$

## Adaptive minibatch $n_{k+1}$ for stepsize update

### Stochastic AC-FGM

$$n_{k+1} = \mathcal{O}\left(\left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max\left\{v_k^{\max}, \frac{\tilde{\sigma}_k^2 + \tilde{\sigma}_{k+1}^2}{D_0^2}\right\}\right\rceil\right)$$

$$\text{Local } \bar{L}_k, \tilde{\sigma}_{k+1}^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

## Adaptive minibatch $n_{k+1}$ for stepsize update

### Stochastic AC-FGM

$$n_{k+1} = \mathcal{O}\left(\left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max\left\{v_k^{\max}, \frac{\tilde{\sigma}_k^2 + \tilde{\sigma}_{k+1}^2}{D_0^2}\right\}\right\rceil\right)$$

$$\text{Local } \bar{L}_k, \tilde{\sigma}_{k+1}^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

>  $\tilde{\sigma}_{k+1}^2$  : after  $m_{k+1}$ ,  $x_{k+1}$  can be computed, and thus  $\tilde{\sigma}_{k+1}^2$  is known for  $n_{k+1}$ .

$$\mathbb{E}_\xi[\|G(x_k, \xi) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \tilde{\sigma}_k^2, \quad \tilde{\sigma}_k^2 \in \mathcal{G}_k := \sigma(x_1, \dots, x_k).$$

## Adaptive minibatch $n_{k+1}$ for stepsize update

### Stochastic AC-FGM

$$n_{k+1} = \mathcal{O}\left(\left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max\left\{v_k^{\max}, \frac{\tilde{\sigma}_k^2 + \tilde{\sigma}_{k+1}^2}{D_0^2}\right\}\right\rceil\right)$$

$$\text{Local } \bar{L}_k, \tilde{\sigma}_{k+1}^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

>  $\tilde{\sigma}_{k+1}^2$  : after  $m_{k+1}$ ,  $x_{k+1}$  can be computed, and thus  $\tilde{\sigma}_{k+1}^2$  is known for  $n_{k+1}$ .

$$\mathbb{E}_\xi[\|G(x_k, \xi) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \tilde{\sigma}_k^2, \quad \tilde{\sigma}_k^2 \in \mathcal{G}_k := \sigma(x_1, \dots, x_k).$$

>  $n_{k+1}$  is determined by previously observed quantities:  $v_k^{\max}$ ,  $\bar{L}_k$ , and  $\sigma_k^2$  — **Adaptive!**

## Adaptive minibatch $n_{k+1}$ for stepsize update

### Stochastic AC-FGM

$$n_{k+1} = \mathcal{O}\left(\left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max\left\{v_k^{\max}, \frac{\tilde{\sigma}_k^2 + \tilde{\sigma}_{k+1}^2}{D_0^2}\right\}\right\rceil\right)$$

$$\text{Local } \bar{L}_k, \tilde{\sigma}_{k+1}^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

>  $\tilde{\sigma}_{k+1}^2$  : after  $m_{k+1}$ ,  $x_{k+1}$  can be computed, and thus  $\tilde{\sigma}_{k+1}^2$  is known for  $n_{k+1}$ .

$$\mathbb{E}_\xi[\|G(x_k, \xi) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \tilde{\sigma}_k^2, \quad \tilde{\sigma}_k^2 \in \mathcal{G}_k := \sigma(x_1, \dots, x_k).$$

>  $n_{k+1}$  is determined by previously observed quantities:  $v_k^{\max}$ ,  $\bar{L}_k$ , and  $\sigma_k^2$  — **Adaptive!**

>  $n_{k+1}$  is used to control the variability induced by the previous stepsize  $\eta_{k+1} \propto \bar{L}_k^{-1}$

## Adaptive minibatch $n_{k+1}$ for stepsize update

### Stochastic AC-FGM

$$n_{k+1} = \mathcal{O}\left(\left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max\left\{v_k^{\max}, \frac{\tilde{\sigma}_k^2 + \tilde{\sigma}_{k+1}^2}{D_0^2}\right\}\right\rceil\right)$$

$$\text{Local } \bar{L}_k, \tilde{\sigma}_{k+1}^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

>  $\tilde{\sigma}_{k+1}^2$  : after  $m_{k+1}$ ,  $x_{k+1}$  can be computed, and thus  $\tilde{\sigma}_{k+1}^2$  is known for  $n_{k+1}$ .

$$\mathbb{E}_\xi[\|G(x_k, \xi) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \tilde{\sigma}_k^2, \quad \tilde{\sigma}_k^2 \in \mathcal{G}_k := \sigma(x_1, \dots, x_k).$$

>  $n_{k+1}$  is determined by previously observed quantities:  $v_k^{\max}$ ,  $\bar{L}_k$ , and  $\sigma_k^2$  — **Adaptive!**

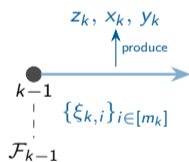
>  $n_{k+1}$  is used to control the variability induced by the previous stepsize  $\eta_{k+1} \propto \bar{L}_k^{-1}$

>  $\bar{L}_k^{-1}$  is not unbiased!  $\implies v_k^{\max}$  controls the bias of  $\bar{L}_k^{-1}$ .

## Three Batch Types & Filtration

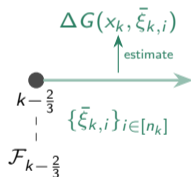
### Gradient Est.

Samples  $\{\xi_{k,i}\}_{i=1}^{m_k}$   
 Compute  $G_k, z_k, x_k, y_k$



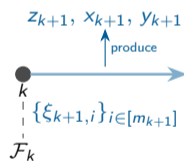
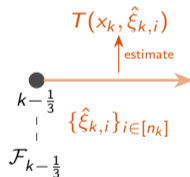
### Stepsize update. I

Samples  $\{\bar{\xi}_{k,i}\}$   
 Compute  $\Delta G$



### Stepsize update. II

Samples  $\{\hat{\xi}_{k,i}\}$   
 Compute  $T$



## Three Batch Types & Filtration

### Gradient Est.

Samples  $\{\xi_{k,i}\}_{i=1}^{m_k}$   
 Compute  $G_k, z_k, x_k, y_k$

### Stepsize update. I

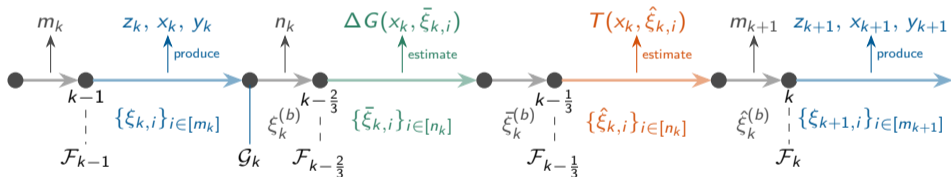
Samples  $\{\bar{\xi}_{k,i}\}$   
 Compute  $\Delta G$

### Stepsize update. II

Samples  $\{\hat{\xi}_{k,i}\}$   
 Compute  $T$

### Variance Est.

$r_k$  samples  
 Update  $m_k, n_k$



## Type III: minibatch update batches

Minibatch in parameter free setting (unknown var)

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ \hat{v}_k^{\max}, \frac{\hat{\sigma}_k^2}{D_0^2} \right\} \right\rceil \right)$$

$$\text{Local } \bar{L}_k, \hat{\sigma}_k^2, \hat{v}_k^{\max} := \max_{0 \leq i \leq k} \hat{v}_i$$

Minibatch in parameter free setting (known var)

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ v_k^{\max}, \frac{\sigma_k^2}{D_0^2} \right\} \right\rceil \right)$$

$$\text{Local } \bar{L}_k, \sigma_k^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

## Type III: minibatch update batches

Minibatch in parameter free setting (unknown var)

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ \hat{v}_k^{\max}, \frac{\hat{\sigma}_k^2}{D_0^2} \right\} \right\rceil \right)$$

$$\text{Local } \bar{L}_k, \hat{\sigma}_k^2, \hat{v}_k^{\max} := \max_{0 \leq i \leq k} \hat{v}_i$$

Minibatch in parameter free setting (known var)

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ v_k^{\max}, \frac{\sigma_k^2}{D_0^2} \right\} \right\rceil \right)$$

$$\text{Local } \bar{L}_k, \sigma_k^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

- > The third type  $\{\xi_{k,i}^b\}_{i=1}^{r_k} \Rightarrow \hat{\sigma}_k$ ,  $\{\bar{\xi}_{k,i}^b\}_{i=1}^{r_k} \Rightarrow \hat{\delta}_k$ , and  $\{\hat{\xi}_{k,i}^b\}_{i=1}^{r_k} \Rightarrow \hat{v}_k$  fresh batches estimate variance without sacrificing adaptivity or acceleration.

## Type III: minibatch update batches

### Minibatch in parameter free setting (unknown var)

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ \hat{v}_k^{\max}, \frac{\hat{\sigma}_k^2}{D_0^2} \right\} \right\rceil \right)$$

$$\text{Local } \bar{L}_k, \hat{\sigma}_k^2, \hat{v}_k^{\max} := \max_{0 \leq i \leq k} \hat{v}_i$$

### Minibatch in parameter free setting (known var)

$$m_{k+1} = \mathcal{O} \left( \left\lceil \frac{N(k+1)^2}{\bar{L}_k^2} \max \left\{ v_k^{\max}, \frac{\sigma_k^2}{D_0^2} \right\} \right\rceil \right)$$

$$\text{Local } \bar{L}_k, \sigma_k^2, v_k^{\max} := \max_{0 \leq i \leq k} v_i$$

- > The third type  $\{\xi_{k,i}^b\}_{i=1}^{r_k} \Rightarrow \hat{\sigma}_k$ ,  $\{\bar{\xi}_{k,i}^b\}_{i=1}^{r_k} \Rightarrow \hat{\delta}_k$ , and  $\{\hat{\xi}_{k,i}^b\}_{i=1}^{r_k} \Rightarrow \hat{v}_k$  fresh batches estimate variance without sacrificing adaptivity or acceleration.
- > The framework is estimator-agnostic: any good filtration-adapted variance estimator can be used.

## Assumptions I: unbiased estimate

---

**A1. Conditional unbiasedness.** Given  $x_k$  for the main update samples  $\xi_{k+1}$ ,

$$\mathbb{E}_{\xi_{k+1}}[G(x_k, \xi_{k+1}) \mid \mathcal{F}_k] = \nabla f(x_k). \quad \text{classic}$$

## Assumptions I: unbiased estimate

**A1. Conditional unbiasedness.** Given  $x_k$  for the main update samples  $\xi_{k+1}$ ,

$$\mathbb{E}_{\xi_{k+1}}[G(x_k, \xi_{k+1}) \mid \mathcal{F}_k] = \nabla f(x_k). \quad \text{classic}$$

For the stepsize selection samples  $\bar{\xi}_k$  and  $\hat{\xi}_k$ , it holds that

$$\begin{aligned} \mathbb{E}_{\bar{\xi}_k} \left[ G(x_k, \bar{\xi}_k) \mid \mathcal{F}_{k-\frac{2}{3}} \right] &= \nabla f(x_k), & \mathbb{E}_{\hat{\xi}_k} \left[ G(x_k, \hat{\xi}_k) \mid \mathcal{F}_{k-\frac{1}{3}} \right] &= \nabla f(x_k), \\ \mathbb{E}_{\hat{\xi}_k} \left[ F(x_k, \hat{\xi}_k) \mid \mathcal{F}_{k-\frac{1}{3}} \right] &= f(x_k), & \mathbb{E}_{\hat{\xi}_k} \left[ F(x_{k-1}, \hat{\xi}_k) \mid \mathcal{F}_{k-\frac{1}{3}} \right] &= f(x_{k-1}). \end{aligned}$$

## Assumptions I: unbiased estimate

**A1. Conditional unbiasedness.** Given  $x_k$  for the main update samples  $\xi_{k+1}$ ,

$$\mathbb{E}_{\xi_{k+1}}[G(x_k, \xi_{k+1}) \mid \mathcal{F}_k] = \nabla f(x_k). \quad \text{classic}$$

For the stepsize selection samples  $\bar{\xi}_k$  and  $\hat{\xi}_k$ , it holds that

$$\begin{aligned} \mathbb{E}_{\bar{\xi}_k} \left[ G(x_k, \bar{\xi}_k) \mid \mathcal{F}_{k-\frac{2}{3}} \right] &= \nabla f(x_k), & \mathbb{E}_{\hat{\xi}_k} \left[ G(x_k, \hat{\xi}_k) \mid \mathcal{F}_{k-\frac{1}{3}} \right] &= \nabla f(x_k), \\ \mathbb{E}_{\hat{\xi}_k} \left[ F(x_k, \hat{\xi}_k) \mid \mathcal{F}_{k-\frac{1}{3}} \right] &= f(x_k), & \mathbb{E}_{\hat{\xi}_k} \left[ F(x_{k-1}, \hat{\xi}_k) \mid \mathcal{F}_{k-\frac{1}{3}} \right] &= f(x_{k-1}). \end{aligned}$$

Classical unbiasedness extends naturally to the new filtrations induced by the parameter-free setting.

## Assumptions II: Bounded local conditional variance

---

**A2. Bounded local conditional variance.** Given  $x_k$ ,  $\exists \sigma_k \geq 0$  s.t. for a fresh sample  $\xi_{k+1}$ ,

$$\mathbb{E}_{\xi_{k+1}} [\|G(x_k, \xi_{k+1}) - \nabla f(x_k)\|^2 \mid \mathcal{F}_k] \leq \sigma_k^2. \quad \text{classic}$$

## Assumptions II: Bounded local conditional variance

**A2. Bounded local conditional variance.** Given  $x_k$ ,  $\exists \sigma_k \geq 0$  s.t. for a fresh sample  $\xi_{k+1}$ ,

$$\mathbb{E}_{\xi_{k+1}} [\|G(x_k, \xi_{k+1}) - \nabla f(x_k)\|^2 \mid \mathcal{F}_k] \leq \sigma_k^2. \quad \text{classic}$$

$\exists \delta_k \geq 0$  such that, for fresh batch-size and stepsize selection samples  $\xi_k^b$  and  $\bar{\xi}_k$ ,

$$\mathbb{E}_{\xi_k^b} [\|G(x_k, \xi_k^b) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \delta_k^2,$$

$$\mathbb{E}_{\bar{\xi}_k} [\|G(x_k, \bar{\xi}_k) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \delta_k^2.$$

## Assumptions II: Bounded local conditional variance

**A2. Bounded local conditional variance.** Given  $x_k$ ,  $\exists \sigma_k \geq 0$  s.t. for a fresh sample  $\xi_{k+1}$ ,

$$\mathbb{E}_{\xi_{k+1}} [\|G(x_k, \xi_{k+1}) - \nabla f(x_k)\|^2 \mid \mathcal{F}_k] \leq \sigma_k^2. \quad \text{classic}$$

$\exists \delta_k \geq 0$  such that, for fresh batch-size and stepsize selection samples  $\xi_k^b$  and  $\bar{\xi}_k$ ,

$$\mathbb{E}_{\xi_k^b} [\|G(x_k, \xi_k^b) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \delta_k^2,$$

$$\mathbb{E}_{\bar{\xi}_k} [\|G(x_k, \bar{\xi}_k) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \delta_k^2.$$

$\exists v_k \geq 0$  such that for a fresh stepsize selection sample  $\hat{\xi}_k$ ,

$$\mathbb{E}_{\hat{\xi}_k} [|\tilde{L}_k(\hat{\xi}_k) - L_k|^2 \mid \mathcal{F}_{k-\frac{2}{3}}] \leq v_k^2,$$

where  $L_k$  is defined is the local smoothness estimator.

## Assumptions II: Bounded local conditional variance

**A2. Bounded local conditional variance.** Given  $x_k$ ,  $\exists \sigma_k \geq 0$  s.t. for a fresh sample  $\xi_{k+1}$ ,

$$\mathbb{E}_{\xi_{k+1}} [\|G(x_k, \xi_{k+1}) - \nabla f(x_k)\|^2 \mid \mathcal{F}_k] \leq \sigma_k^2. \quad \text{classic}$$

$\exists \delta_k \geq 0$  such that, for fresh batch-size and stepsize selection samples  $\xi_k^b$  and  $\bar{\xi}_k$ ,

$$\mathbb{E}_{\xi_k^b} [\|G(x_k, \xi_k^b) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \delta_k^2,$$

$$\mathbb{E}_{\bar{\xi}_k} [\|G(x_k, \bar{\xi}_k) - \nabla f(x_k)\|^2 \mid \mathcal{G}_k] \leq \delta_k^2.$$

$\exists v_k \geq 0$  such that for a fresh stepsize selection sample  $\hat{\xi}_k$ ,

$$\mathbb{E}_{\hat{\xi}_k} \left[ |\tilde{L}_k(\hat{\xi}_k) - L_k|^2 \mid \mathcal{F}_{k-\frac{2}{3}} \right] \leq v_k^2,$$

where  $L_k$  is defined is the local smoothness estimator.

Classical bounded variance extends naturally to the new filtrations induced by the parameter-free setting;  $v_k$  is due to the bias of stepsize estimator.

## Assumptions III: Finite-sample cocoercivity–smoothness

---

**A3.** Given a query pair  $(x_{k-1}, x_k)$ ,  $\exists \mathcal{L}(\bar{\xi}_k, \hat{\xi}_k)$ ,  $\mathcal{L}(\hat{\xi}_k) > 0$  s.t.

$$\frac{\|\Delta G(x_k, \bar{\xi}_k)\|^2}{2\mathcal{L}(\bar{\xi}_k, \hat{\xi}_k)} \leq T(x_k, \hat{\xi}_k) \leq \frac{\mathcal{L}(\hat{\xi}_k)}{2} \|x_k - x_{k-1}\|^2, \quad \text{a.s.}$$

## Assumptions III: Finite-sample cocoercivity–smoothness

**A3.** Given a query pair  $(x_{k-1}, x_k)$ ,  $\exists \mathcal{L}(\bar{\xi}_k, \hat{\xi}_k)$ ,  $\mathcal{L}(\hat{\xi}_k) > 0$  s.t.

$$\frac{\|\Delta G(x_k, \bar{\xi}_k)\|^2}{2\mathcal{L}(\bar{\xi}_k, \hat{\xi}_k)} \leq T(x_k, \hat{\xi}_k) \leq \frac{\mathcal{L}(\hat{\xi}_k)}{2} \|x_k - x_{k-1}\|^2, \quad \text{a.s.}$$

A finite-sample analogue of cocoercivity + smoothness. Does **not** require each  $F(\cdot, \hat{\xi})$  to be convex or smooth — only locally well-behaved along  $(x_{k-1}, x_k)$ .

## Adaptivity to the Lipschitz constant

### Theorem 1 (Ji & Lan, 2026)

Under A1–A3, with the previously defined adaptive stepsize and minibatches rules, to reach  $\varepsilon$  solution such that  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$ , stochastic AC-FGM requires

$$\mathcal{O} \left( \sqrt{\frac{\mathcal{L}D_0^2}{\varepsilon} \cdot \max \left\{ \frac{v_{\max}}{v_0}, 1 \right\}} \right) \text{ iterations.}$$

## Adaptivity to the Lipschitz constant

### Theorem 1 (Ji & Lan, 2026)

Under A1–A3, with the previously defined adaptive stepsize and minibatches rules, to reach  $\varepsilon$  solution such that  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$ , stochastic AC-FGM requires

$$\mathcal{O} \left( \sqrt{\frac{\mathcal{L}D_0^2}{\varepsilon} \cdot \max \left\{ \frac{v_{\max}}{v_0}, 1 \right\}} \right) \text{ iterations.}$$

- > Matches the  $L$ -known AC-SA rate [Lan \(2012\)](#) — optimal in  $\varepsilon$  [Nemirovski \(1985\)](#)

## Adaptivity to the Lipschitz constant

### Theorem 1 (Ji & Lan, 2026)

Under A1–A3, with the previously defined adaptive stepsize and minibatches rules, to reach  $\varepsilon$  solution such that  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$ , stochastic AC-FGM requires

$$\mathcal{O} \left( \sqrt{\frac{\mathcal{L}D_0^2}{\varepsilon} \cdot \max \left\{ \frac{v_{\max}}{v_0}, 1 \right\}} \right) \text{ iterations.}$$

- > Matches the  $L$ -known AC-SA rate [Lan \(2012\)](#) — optimal in  $\varepsilon$  [Nemirovski \(1985\)](#)
- > Pessimistic  $\mathcal{L}, v_{\max}$ , due to the randomness of the stepsize when seeking in expectation guarantee; see also [Lan, Li, & Xu 2024](#);

## Adaptivity to the Lipschitz constant

### Theorem 1 (Ji & Lan, 2026)

Under A1–A3, with the previously defined adaptive stepsize and minibatches rules, to reach  $\varepsilon$  solution such that  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$ , stochastic AC-FGM requires

$$\mathcal{O} \left( \sqrt{\frac{\mathcal{L}D_0^2}{\varepsilon} \cdot \max \left\{ \frac{v_{\max}}{v_0}, 1 \right\}} \right) \text{ iterations.}$$

- > Matches the  $L$ -known AC-SA rate [Lan \(2012\)](#) — optimal in  $\varepsilon$  [Nemirovski \(1985\)](#)
- > Pessimistic  $\mathcal{L}, v_{\max}$ , due to the randomness of the stepsize when seeking in expectation guarantee; see also [Lan, Li, & Xu 2024](#);
- > Under light-tail assumptions,  $\mathcal{L}, v_{\max} \rightarrow$  local quantities  $\hat{\mathcal{L}}_N$  and  $v_N^{\max}$ .

## Sample complexity

To reach  $\varepsilon$  solution s.t.  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$ , stochastic AC-FGM requires

$$\mathcal{O} \left( \sqrt{\frac{\mathcal{L}D_0^2}{\varepsilon} \cdot \frac{v_{\max}}{v_0}} + \frac{r_N \mathcal{L}^2 D_0^2}{\varepsilon^2} \cdot \frac{v_{\max}^2}{v_0^2} \right) \text{ samples}$$

where  $r_N$  is the *average variance-to-smoothness ratio* along the trajectory

$$r_N^2 := \frac{1}{N} \sum_{k=1}^N \frac{v_{k-1}^{\max} D_0^2 + \sigma_{k-1}^2}{L_k^2} \implies \frac{\sigma^2}{L^2}$$

## Sample complexity

To reach  $\varepsilon$  solution s.t.  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$ , stochastic AC-FGM requires

$$\mathcal{O} \left( \sqrt{\frac{\mathcal{L}D_0^2}{\varepsilon} \cdot \frac{v_{\max}}{v_0}} + \frac{r_N \mathcal{L}^2 D_0^2}{\varepsilon^2} \cdot \frac{v_{\max}^2}{v_0^2} \right) \text{ samples}$$

where  $r_N$  is the *average variance-to-smoothness ratio* along the trajectory

$$r_N^2 := \frac{1}{N} \sum_{k=1}^N \frac{v_{k-1}^{\max} D_0^2 + \sigma_{k-1}^2}{L_k^2} \implies \frac{\sigma^2}{L^2}$$

To reach  $\varepsilon$  solution s.t.  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$ , AC-SA [Lan 2012](#) requires

$$\mathcal{O} \left( \sqrt{\frac{LD_0^2}{\varepsilon}} + \frac{\sigma^2 D_0^2}{\varepsilon^2} \right) \text{ samples}$$

## Iteration limit $N$ free: anchored regularization

---

The proximal subproblem includes an extra anchor term:

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle G_k, z \rangle + h(z) + \frac{1}{2\eta_k} \|y_{k-1} - z\|^2 + \frac{\gamma_k}{2\eta_k} \|y_0 - z\|^2 \right\}$$

## Iteration limit $N$ free: anchored regularization

---

The proximal subproblem includes an extra anchor term:

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle G_k, z \rangle + h(z) + \frac{1}{2\eta_k} \|y_{k-1} - z\|^2 + \frac{\gamma_k}{2\eta_k} \|y_0 - z\|^2 \right\}$$

>  $\gamma_k$  chosen small (decaying) so convergence rate is unchanged

## Iteration limit $N$ free: anchored regularization

The proximal subproblem includes an extra anchor term:

$$z_k = \operatorname{argmin}_{z \in X} \left\{ \langle G_k, z \rangle + h(z) + \frac{1}{2\eta_k} \|y_{k-1} - z\|^2 + \frac{\gamma_k}{2\eta_k} \|y_0 - z\|^2 \right\}$$

- >  $\gamma_k$  chosen small (decaying) so convergence rate is unchanged
- > Enables horizon-free analysis — no prior knowledge of total iterations  $N$

## Adaptive to the iteration limit

### Theorem 2 (Ji & Lan, 2026)

Under A1–A3, with anchored regularization  $\gamma_k = \frac{1}{k}$ ,  $\tau_k = \frac{k+2-\beta}{2}$ , and suppose  $\eta_1 > 0$  and

$$\eta_2 = \min \left\{ \frac{1}{16L_1}, \frac{2(1-\beta)}{3-\beta} \eta_1 \right\}, \quad \eta_k = \min \left\{ \frac{k-1}{16L_{k-1}}, \frac{(k-1)(k+2-\beta)}{k^2} \eta_{k-1} \right\}, \quad \forall k \geq 3.$$

Furthermore, suppose the mini batch size satisfies

$$m_k = \mathcal{O} \left( \left\lceil \frac{(k+1)k^2}{L_{k-1}^2} \max \left\{ v_{k-1}^{\max}, \frac{\sigma_{k-1}^2}{D_0^2} \right\} \right\rceil \right), \quad n_k = \mathcal{O} \left( \left\lceil \frac{(k+1)k^2}{L_{k-1}^2} \max \left\{ v_{k-1}^{\max}, \frac{\delta_k^2}{D_0^2} \right\} \right\rceil \right).$$

To reach  $\varepsilon$  solution such that  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$  on  $A$ , stochastic AC-FGM requires

$$\mathcal{O} \left( \sqrt{\frac{\mathcal{L}D_0^2}{\varepsilon} \cdot \max \left\{ \frac{v_{\max}}{v_0}, 1 \right\}} \right) \text{ iterations.}$$

## Adaptivity to the local variances

### Theorem 3 (Ji & Lan, 2026)

Under A1-A3, the same parameter choices as Theorem 2, suppose

$$A := \{\forall k, \widehat{v}_{k-1} \geq v_{k-1}, \widehat{\sigma}_{k-1}^2 \geq \sigma_{k-1}^2, \widehat{\delta}_k^2 \geq \delta_k^2\}, \quad \mathbb{P}(A) \geq 1 - \delta.$$

Furthermore, suppose the mini batch size satisfies

$$m_k = \mathcal{O} \left( \left\lceil \frac{(k+1)k^2}{L_{k-1}^2} \max \left\{ \widehat{v}_{k-1}^{\max}, \frac{\widehat{\sigma}_{k-1}^2}{D_0^2} \right\} \right\rceil \right), \quad n_k = \mathcal{O} \left( \left\lceil \frac{(k+1)k^2}{L_{k-1}^2} \max \left\{ \widehat{v}_{k-1}^{\max}, \frac{\widehat{\delta}_k^2}{D_0^2} \right\} \right\rceil \right).$$

To reach  $\varepsilon$  solution such that  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$  on  $A$ , stochastic AC-FGM requires

$$\mathcal{O} \left( \sqrt{\frac{\mathcal{L}D_0^2}{\varepsilon} \cdot \max \left\{ \frac{v_{\max}}{v_0}, 1 \right\}} \right) \text{ iterations.}$$

## High-Probability Guarantee

### Theorem 4 (Ji & Lan, 2026)

Under additional light tail assumption, w.p. at least  $1 - (N+1)e^{-\Lambda^2/3} - 4(N+1)e^{-\Lambda}$ :

$$\Phi(x_N) - \Phi(x^*) \leq \varepsilon \quad \text{in}$$

$$\mathcal{O} \left( \sqrt{\frac{\widehat{L}_N D_0^2}{\varepsilon} \cdot \max \left\{ \frac{v_{N+1}^{\max}}{v_0}, 1 \right\}} \right) \quad \text{iterations.}$$

## High-Probability Guarantee

### Theorem 4 (Ji & Lan, 2026)

Under additional light tail assumption, w.p. at least  $1 - (N+1)e^{-\Lambda^2/3} - 4(N+1)e^{-\Lambda}$ :

$$\Phi(x_N) - \Phi(x^*) \leq \varepsilon \quad \text{in}$$

$$\mathcal{O} \left( \sqrt{\frac{\widehat{L}_N D_0^2}{\varepsilon} \cdot \max \left\{ \frac{v_{N+1}^{\max}}{v_0}, 1 \right\}} \right) \quad \text{iterations.}$$

- > **Trajectory-dependent:**  $\widehat{L}_N, v_{N+1}^{\max}$  depend only on iterates *actually visited*
- > **Prior work** (Kreidler et al. 2024): bound depends on  $\sigma^* = \max_{x \in \mathbb{B}(x^*, 2D_0)} \sigma_x$  — supremum over entire ball
- > **Ours:** no global worst-case penalty; pays only for noise and curvature encountered

## High-Probability sample complexity

Stochastic AC-FGM requires

$$\mathcal{O} \left( \sqrt{\frac{L_N D_0^2}{\varepsilon} \cdot \frac{v_{N+1}^{\max}}{v_0}} + \frac{r_N \mathcal{L}^2 D_0^2}{\varepsilon^2} \cdot \frac{(v_{N+1}^{\max})^2}{v_0^2} \right) \text{ samples}$$

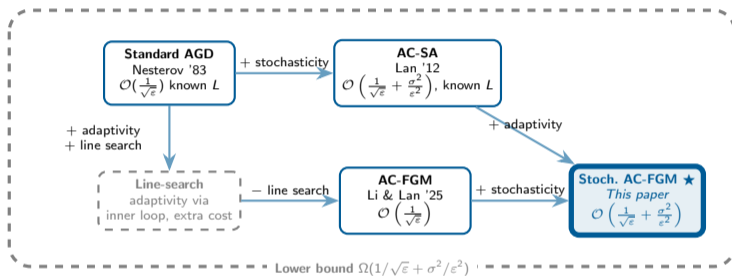
where  $r_N$  is the *average variance-to-smoothness ratio* along the trajectory

$$r_N^2 := \frac{1}{N} \sum_{k=1}^N \frac{v_{k-1}^{\max} D_0^2 + \sigma_{k-1}^2}{\bar{L}_k^2} \implies \frac{\sigma^2}{L^2}$$

To reach  $\varepsilon$  solution s.t.  $\mathbb{E}[\Phi(x_N) - \Phi(x^*)] \leq \varepsilon$ , AC-SA [Lan 2012](#) requires

$$\mathcal{O} \left( \sqrt{\frac{LD_0^2}{\varepsilon}} + \frac{\sigma^2 D_0^2}{\varepsilon^2} \right) \text{ samples}$$

## Take-Home Message



Stochastic AC-FGM achieves:

- > Adaptive to the Lipschitz constant  $L$ , the iteration horizon  $N$ , and the local variances
- > Optimal convergence rate without boundedness assumption

# Thank You

## Stochastic Auto-conditioned Fast Gradient Methods with Op

Yao Ji & Guanghui (George) Lan

**Preprint:** Stay tuned for arXiv soon!

**Contact:** [yaoji@gatech.edu](mailto:yaoji@gatech.edu) [george.lan@isye.gatech.edu](mailto:george.lan@isye.gatech.edu)