

High-dimensional DFO: Stochastic Subspace Descent and improvements

Stephen Becker, University of Colorado Boulder

Collaborators: Alireza Doostan, Nuojin Chen, Killian Wood, Drona Khurana (CU)
David Kozak, Luis Tenorio (Colorado School of Mines)

Slides will be posted at <https://stephenbeckr.github.io/papers/#talks>



ICCOPT, July 2025, USC (Los Angeles)

Outline

Outline and learning objectives of this talk

- Introduce a class of 0th order optimization methods, “SSD”
 - survey our prior analysis results
- Argue that **stepsize** selection is a key issue
 - suggest two ideas to help
 - show some ML **examples** where these methods make sense
- Discuss ongoing directions

Motivation: 0th order methods / derivative-free optimization

$\min_{\boldsymbol{x} \in \mathbb{R}^d} f(\boldsymbol{x})$ where we do not have access to $\nabla f(\boldsymbol{x})$

Traditional applications: PDE constrained optimization, when the **adjoint state method** or **automatic differentiation** is inapplicable

- e.g., multiphysics codes with complicated adjoints (and hard to code in HPC); memory issues in AD for time-dependent problems, etc.

Motivation: 0th order methods / derivative-free optimization

$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$ where we do not have access to $\nabla f(\mathbf{x})$

Traditional applications: PDE constrained optimization, when the **adjoint state method** or **automatic differentiation** is inapplicable

- e.g., multiphysics codes with complicated adjoints (and hard to code in HPC); memory issues in AD for time-dependent problems, etc.

Machine learning applications

- hyper parameter tuning
- black-box attacks, e.g., adversarial attacks on a model (attacker doesn't have access to source code)

Motivation: 0th order methods / derivative-free optimization

$$\min_{\boldsymbol{x} \in \mathbb{R}^d} f(\boldsymbol{x}) \text{ where we do not have access to } \nabla f(\boldsymbol{x})$$

Traditional applications: PDE constrained optimization, when the **adjoint state method** or **automatic differentiation** is inapplicable

- e.g., multiphysics codes with complicated adjoints (and hard to code in HPC); memory issues in AD for time-dependent problems, etc.

Machine learning applications

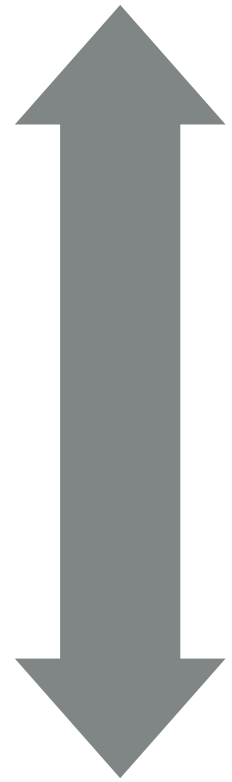
- hyper parameter tuning
- black-box attacks, e.g., adversarial attacks on a model (attacker doesn't have access to source code)

Recent theme of my work: exploit **multi-fidelity** models

- in applied math, for PDE simulations, we often have several physics based models (with different approximations), different discretizations, different numerical precisions, reduced-order models, etc
- **machine learning** also has examples of multi-fidelity models

Context

High dimensions,
low accuracy



Low dimensions,
high accuracy

A guiding principle:

the time spent in the **optimization method** (creating and solving surrogate models)
should equal

the time spent in the **function evaluation** (e.g., solving the PDE, training neural net...)

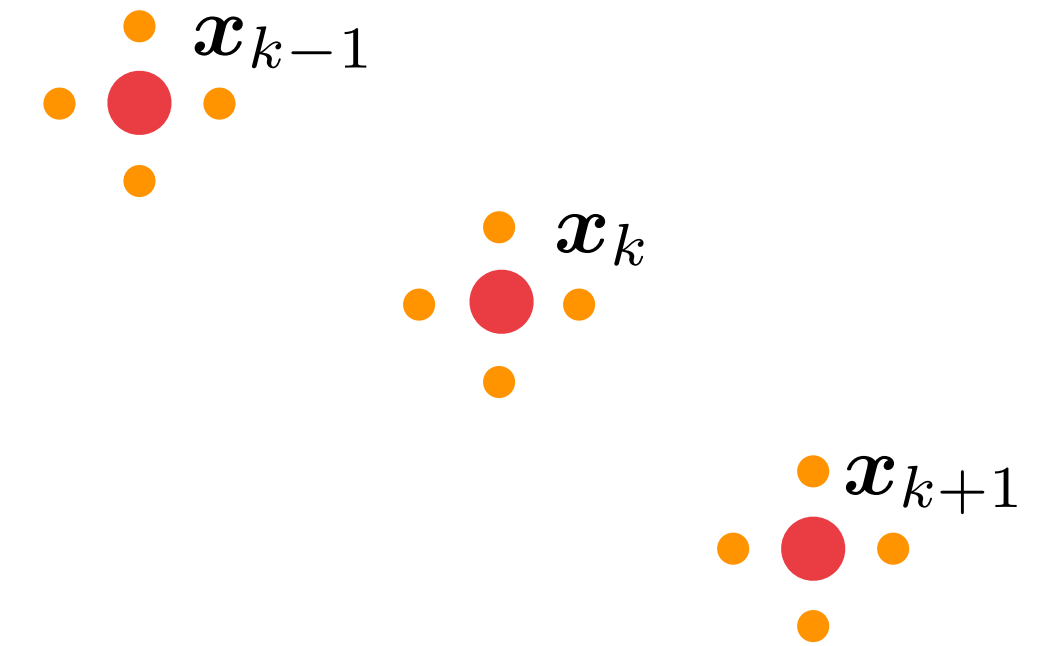
... hence, the best method to use depends a lot on the problem.

Local methods

“zeroth order”: approximate $\nabla f(x)$ (e.g., w/ finite differences)

stochastic zeroth order: approximate g such that $\mathbb{E}[g] = \nabla f(x)$
today's talk

polling methods: Nelder-Mead, coordinate descent, etc.

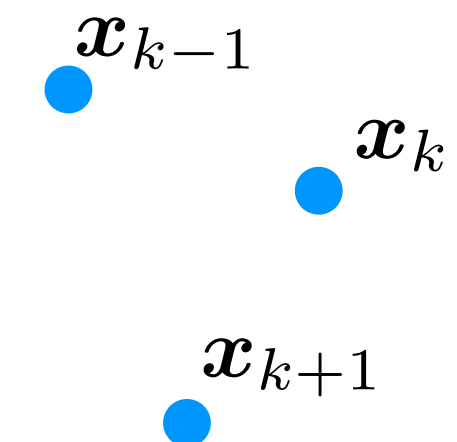


Global / model-based methods

polynomial model, minimize with a trust-region (“DFO-TR”)

see our new paper “A Unified Framework for Entropy Search and Expected Improvement in Bayesian Optimization”, N. Cheng et al. ICML ‘25

Gaussian process model, use acquisition function to tradeoff exploration and exploitation (“Bayesian Optimization”)



Misc. / heuristics

genetic algorithms, particle swarm, CMA-ES, simulated annealing, etc.

A stochastic zeroth-order method: SSD

Assume we can compute/approximate $\mathbf{q}\mathbf{q}^\top \nabla f(\mathbf{x}) = \left(\lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{q}) - f(\mathbf{x})}{h} \right) \mathbf{q}$ e.g., finite differences, or *forward-mode* AutoDiff

directional derivative, “two-point” estimator in ZO literature

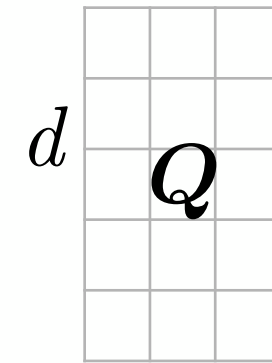
another perspective: $\mathbf{q}^\top \nabla f(\mathbf{x}) = \left. \frac{d}{dt} \varphi(t) \right|_{t=0}$ where $\varphi(t) = f(\mathbf{x} + t\mathbf{q})$

A stochastic zeroth-order method: SSD

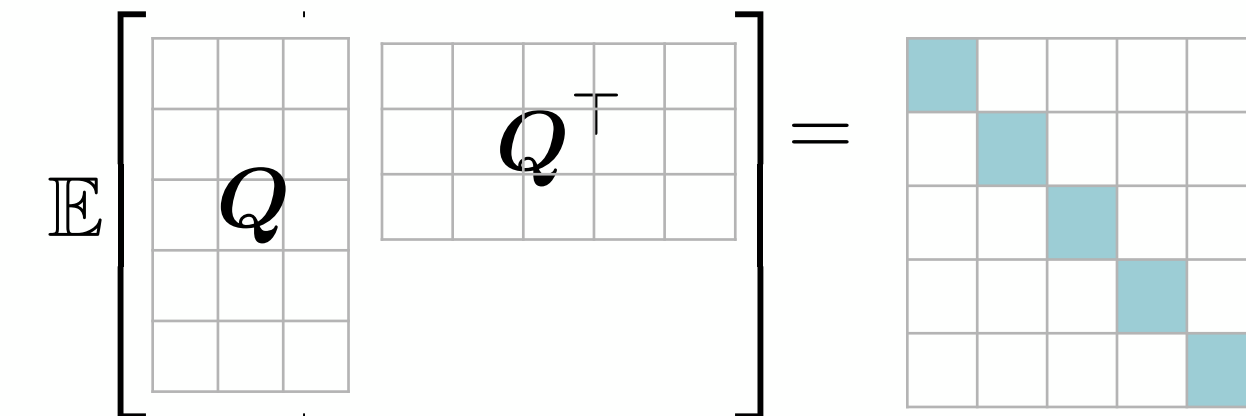
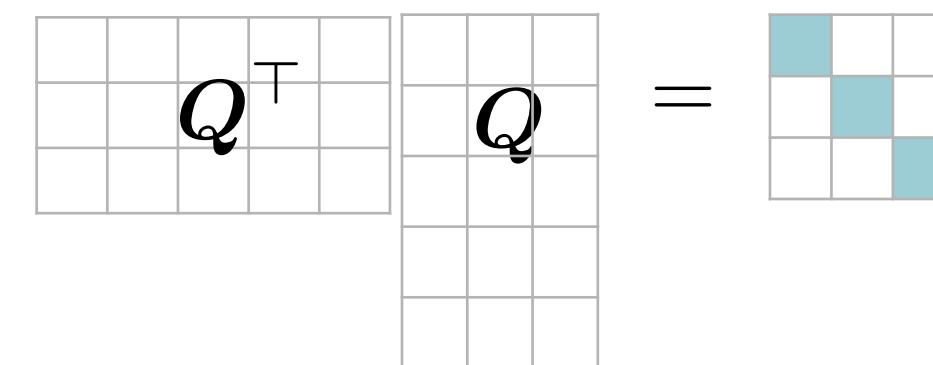
Assume we can compute/approximate $\mathbf{q}\mathbf{q}^\top \nabla f(\mathbf{x}) = \left(\lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{q}) - f(\mathbf{x})}{h} \right) \mathbf{q}$

Even better, average a few copies to reduce variance

$$\mathbf{Q}\mathbf{Q}^\top \nabla f(\mathbf{x}) = \sum_{i=1}^{\ell} \mathbf{q}_i \mathbf{q}_i^\top \nabla f(\mathbf{x}) \quad \mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_\ell] \quad \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_{\ell \times \ell}, \quad \mathbb{E} \left(\frac{d}{\ell} \mathbf{Q}\mathbf{Q}^\top \right) = \mathbf{I}_{d \times d}$$



columns **not** independent!



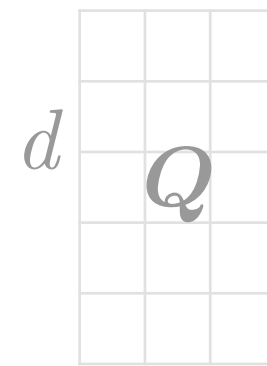
A stochastic zeroth-order method: SSD

Assume we can compute/approximate $\mathbf{q}\mathbf{q}^\top \nabla f(\mathbf{x}) = \left(\lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{q}) - f(\mathbf{x})}{h} \right) \mathbf{q}$

Even better, average a few copies to reduce variance

$$\mathbf{Q}\mathbf{Q}^\top \nabla f(\mathbf{x}) = \sum_{i=1}^{\ell} \mathbf{q}_i \mathbf{q}_i^\top \nabla f(\mathbf{x})$$

$$\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_\ell]$$



$$\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_{\ell \times \ell}, \quad \mathbb{E} \left(\frac{d}{\ell} \mathbf{Q}\mathbf{Q}^\top \right) = \mathbf{I}_{d \times d}$$

Algorithm: stochastic subspace descent (SSD)

Repeat:

Draw a \mathbf{Q} $\mathcal{O}(\ell)$ oracle calls
 $\mathbf{x} \leftarrow \mathbf{x} - \eta \frac{d}{\ell} \mathbf{Q}\mathbf{Q}^\top \nabla f(\mathbf{x})$

stepsize

It's a type of “stochastic gradient method” (direction is unbiased)
 but has much stronger guarantees due to its structure

e.g., the direction is **descent direction** with prob. 1 if \mathbf{Q} is continuous

Important: draw new (independent) \mathbf{Q} every iteration

A stochastic zeroth-order method: SSD

Assume we can compute/approximate $\mathbf{q}\mathbf{q}^\top \nabla f(\mathbf{x}) = \left(\lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{q}) - f(\mathbf{x})}{h} \right) \mathbf{q}$ e.g., finite differences, or *forward-mode AutoDiff*

directional derivative, “two-point” estimator in ZO literature

Even better, average a few copies to reduce variance

$$\mathbf{Q}\mathbf{Q}^\top \nabla f(\mathbf{x}) = \sum_{i=1}^{\ell} \mathbf{q}_i \mathbf{q}_i^\top \nabla f(\mathbf{x}) \quad \mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_\ell] \quad \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_{\ell \times \ell}, \quad \mathbb{E} \left(\frac{d}{\ell} \mathbf{Q}\mathbf{Q}^\top \right) = \mathbf{I}_{d \times d}$$

Equivalent formulation... and why we call it “SSD”

Columns of \mathbf{Q} form a basis for \mathcal{V}

Algorithm: stochastic subspace descent (SSD)

Repeat:

Draw a \mathbf{Q} $\mathcal{O}(\ell)$ oracle calls
 $\mathbf{x} \leftarrow \mathbf{x} - \eta \frac{d}{\ell} \mathbf{Q}\mathbf{Q}^\top \nabla f(\mathbf{x})$

stepsize

Algorithm 1 Stochastic Subspace Descent (SSD)

Require: η ▷ Stepsize

Require: $\mathbf{x}_0 \in \mathbb{R}^d$ ▷ Initial point

```

1: for  $k = 0, 1, 2, \dots$  do
2:   Choose subspace  $\mathcal{V}_k$  of dimension  $\ell \leq d$ 
3:    $\mathbf{g}_k \leftarrow \text{proj}_{\mathcal{V}_k} (\nabla f(\mathbf{x}_k))$  ▷ Project onto subspace
4:    $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \eta \mathbf{g}_k$ 
5: end for
    
```

Stochastic 0th-order methods

Variants have been investigated for a long time...

Under various names, e.g., “random gradient”, “random pursuit”, “directional search”, “random search”

- J. Matyas, *Random Optimization*, Automation & Remote Control, **1965**
- M. Gaviano, *Some general results on convergence of random search algorithms in minimization problems*, Towards Global Optimisation, **1975**.
- F.J. Solis and R. J-B. Wets, *Minimization by random search techniques*, Math. of Operations Research, **1981** (no rate)
- B. Polyak **1987**
- Yu. Ermoliev and R.J.-B. Wets, *Numerical techniques for stochastic optimization*, chapter 6, Springer-Verlag, **1988**.

For recent results (2015-2025) on zeroth-order ML, see:

- “Zeroth-order Machine Learning” AAAI tutorial 2024
 - Wotao Yin, Sijia Liu, Pin-Yu Chen
 - <https://sites.google.com/view/zo-tutorial-aaai-2024/>
- Chen et al., *DeepZero: Scaling up Zeroth-Order Optimization for Deep Model Training*, ICLR '23
- Yu et al., *SubZero: Random Subspace Zeroth-Order Optimization for Memory-Efficient LLM Fine-Tuning*, '24
- Liu et al., *Sparse MeZO: Less Parameters for Better Performance in Zeroth-Order LLM Fine-Tuning*, '24
- Zhang et al., *Revisiting Zeroth-Order Optimization for Memory-Efficient LLM Fine-Tuning: A Benchmark*, ICML '24

Much recent work on variants, 2011—2020

- D. Leventhal and A.S. Lewis, *Randomized Hessian estimation and directional search*, Optimization (2011)
- Z. Zhang, *Scalable Derivative-Free Optimization Algorithms with Low-Dimensional Subspace Techniques* (thesis in Chinese 2011, preprint with proof sketches in English 2025, arXiv 2501.04536), “NEWUOAs”
- S. U. Stich, C. Muller, and B. Gartner, *Optimization of convex functions with random pursuit*, SIAM J. Opt. (2013)
- Yu. Nesterov, *Random gradient-free minimization of convex functions*, '11 / Yu. Nesterov and V. Spokoiny, FoCM 2017
- P. Dvurechensky, A. Gasnikov, and A. Tiurin, *Randomized similar triangles method: A unifying framework for accelerated randomized optimization methods (coordinate descent, directional search, derivative-free method)*, arXiv:1707.08486
- P. Dvurechensky, A. Gasnikov, and E. Gorbunov, *An accelerated directional derivative method for smooth stochastic convex optimization*, arXiv:1804.02394
- S. Ghadimi and G. Lan, *Stochastic first- and zeroth-order methods for nonconvex stochastic programming*, SIAM J. Opt. (2013)
- R. Chen and S. Wild, *Randomized derivative-free optimization of noisy convex functions*, arXiv:1507.03332 (2015).
- K. Choromanski, M. Rowland, V. Sindhwani, R. E. Turner, and A. Weller, *Structured evolution with compact architectures for scalable policy optimization*, ICML, 2018.
- T. Salimans, J. Ho, X. Chen, S. Sidor, and I. Sutskever, *Evolution strategies as a scalable alternative to reinforcement learning*, arXiv:1703.03864 (2017).
- J. Duchi, M. Jordan, M. Wainwright, A. Wibisono, *Optimal Rates for Zero-Order Convex Optimization: The Power of Two Function Evaluations*, IEEE Trans Info Theory (2015)
- A. S. Berahas, L. Cao, K. Choromanski, K. Scheinberg, *A Theoretical and Empirical Comparison of Gradient Approximations in Derivative-Free Optimization*, 2019 (published in FoCM '22)
- F. Hanzely, K. Mishchenko, P. Richtarik, *SEGA: Variance Reduction via Gradient Sketching*, NeurIPS 2018
- *cousin of “direct search” methods*, cf. S. Gratton, C. W. Royer, L. N. Vicente, Z. Zhang, *Direct Search Based on Probabilistic Descent*, SIAM J. Opt. (2015)

Since 2020...

- D. Kozak, S. B., A. Doostan, and L. Tenorio. *A stochastic subspace approach to gradient-free optimization in high dimensions*. COAP, 2021.
- D. Kozak, C. Molinari, L. Rosasco, L. Tenorio, S. Villa. *Zeroth-order optimization with orthogonal random directions*. Math. Program., 2023
- M. Rando, C. Molinari, L. Rosasco, S. Villa. *An optimal structured zeroth-order algorithm for nonsmooth optimization*. NeurIPS 2024.
- M. Rando, C. Molinaro, S. Villa, L. Rosasco. *Stochastic zeroth order descent with structured directions*. COAP, 2024.
- R. Nozawa, P.-L. Poirion, A. Takeda. *Zeroth-Order Random Subspace Algorithm for Non-smooth Convex Optimization*, JOTA, 2025.
- H.Q. Cai, D. McKenzie, W. Yin, Z. Zhang. *Zeroth-Order Regularized Optimization (ZORO): Approximately Sparse Gradients and Adaptive Sampling*, SIOPT, 2022

Stochastic 0th-order methods

Variants have been investigated for a long time...

Under various names, e.g., “random gradient”, “random pursuit”, “directional search”, “random search”

- J. Matyas, *Random Optimization*, Automation & Remote Control, **1965**
- M. Gaviano, *Some general results on convergence of random search algorithms in minimization problems*, Towards Global Optimisation, **1975**.
- F.J. Solis and R. J-B. Wets, *Minimization by random search techniques*, Math. of Operations Research, **1981** (no rate)
- B. Polyak **1987**
- Yu. Ermoliev and R.J.-B. Wets, *Numerical techniques for stochastic optimization*, chapter 6, Springer-Verlag, **1988**.

Most literature focuses on $\ell = 1$

How is \mathbf{q} typically chosen?

$$\|\mathbf{q}\|_2 = 1$$
$$\mathbb{E}[\mathbf{q}\mathbf{q}^\top] = \frac{1}{d}\mathbf{I}_{d \times d}$$

- spherical (or Gaussian, scaled appropriately)

or

- canonical basis vector, $\mathbf{q} \sim \text{uniform}[\mathbf{e}_1, \dots, \mathbf{e}_d]$

(hence SSD reduces to randomized coordinate descent)

Much recent work on variants, 2011—2020

- D. Leventhal and A.S. Lewis, *Randomized Hessian estimation and directional search*, Optimization (2011)
- Z. Zhang, *Scalable Derivative-Free Optimization Algorithms with Low-Dimensional Subspace Techniques* (thesis in Chinese 2011, preprint with proof sketches in English 2025, arXiv 2501.04536), “NEWUOAs”
- S. U. Stich, C. Muller, and B. Gartner, *Optimization of convex functions with random pursuit*, SIAM J. Opt. (2013)
- Yu. Nesterov, *Random gradient-free minimization of convex functions*, '11 / Yu. Nesterov and V. Spokoiny, FoCM 2017
- P. Dvurechensky, A. Gasnikov, and A. Tiurin, *Randomized similar triangles method: A unifying framework for accelerated randomized optimization methods (coordinate descent, directional search, derivative-free method)*, arXiv:1707.08486
- P. Dvurechensky, A. Gasnikov, and E. Gorbunov, *An accelerated directional derivative method for smooth stochastic convex optimization*, arXiv:1804.02394
- S. Ghadimi and G. Lan, *Stochastic first- and zeroth-order methods for nonconvex stochastic programming*, SIAM J. Opt. (2013)
- R. Chen and S. Wild, *Randomized derivative-free optimization of noisy convex functions*, arXiv:1507.03332 (2015).
- K. Choromanski, M. Rowland, V. Sindhwani, R. E. Turner, and A. Weller, *Structured evolution with compact architectures for scalable policy optimization*, ICML, 2018.
- T. Salimans, J. Ho, X. Chen, S. Sidor, and I. Sutskever, *Evolution strategies as a scalable alternative to reinforcement learning*, arXiv:1703.03864 (2017).
- J. Duchi, M. Jordan, M. Wainwright, A. Wibisono, *Optimal Rates for Zero-Order Convex Optimization: The Power of Two Function Evaluations*, IEEE Trans Info Theory (2015)
- A. S. Berahas, L. Cao, K. Choromanski, K. Scheinberg, *A Theoretical and Empirical Comparison of Gradient Approximations in Derivative-Free Optimization*, 2019 (published in FoCM '22)
- F. Hanzely, K. Mishchenko, P. Richtarik, *SEGA: Variance Reduction via Gradient Sketching*, NeurIPS 2018
- *cousin of “direct search” methods*, cf. S. Gratton, C. W. Royer, L. N. Vicente, Z. Zhang, *Direct Search Based on Probabilistic Descent*, SIAM J. Opt. (2015)

Since 2020...

- D. Kozak, S. B., A. Doostan, and L. Tenorio. *A stochastic subspace approach to gradient-free optimization in high dimensions*. COAP, 2021.
- D. Kozak, C. Molinari, L. Rosasco, L. Tenorio, S. Villa. *Zeroth-order optimization with orthogonal random directions*. Math. Program., 2023
- M. Rando, C. Molinari, L. Rosasco, S. Villa. *An optimal structured zeroth-order algorithm for nonsmooth optimization*. NeurIPS 2024.
- M. Rando, C. Molinaro, S. Villa, L. Rosasco. *Stochastic zeroth order descent with structured directions*. COAP, 2024.
- R. Nozawa, P.-L. Poirion, A. Takeda. *Zeroth-Order Random Subspace Algorithm for Non-smooth Convex Optimization*, JOTA, 2025.
- H.Q. Cai, D. McKenzie, W. Yin, Z. Zhang. *Zeroth-Order Regularized Optimization (ZORO): Approximately Sparse Gradients and Adaptive Sampling*, SIOPT, 2022

Stochastic 0th-order methods

Variants have been investigated for a long time...

Under various names, e.g. “random gradient”, “random

pursuit”, “directional

- J. Matyas, *Random Optimization*
- M. Gaviano, *Some general minimization problems*
- F.J. Solis and R. J-B. *Operations Research*, 1980
- B. Polyak **1987**
- Yu. Ermoliev and R.J. *chapter 6*, Springer-Verlag

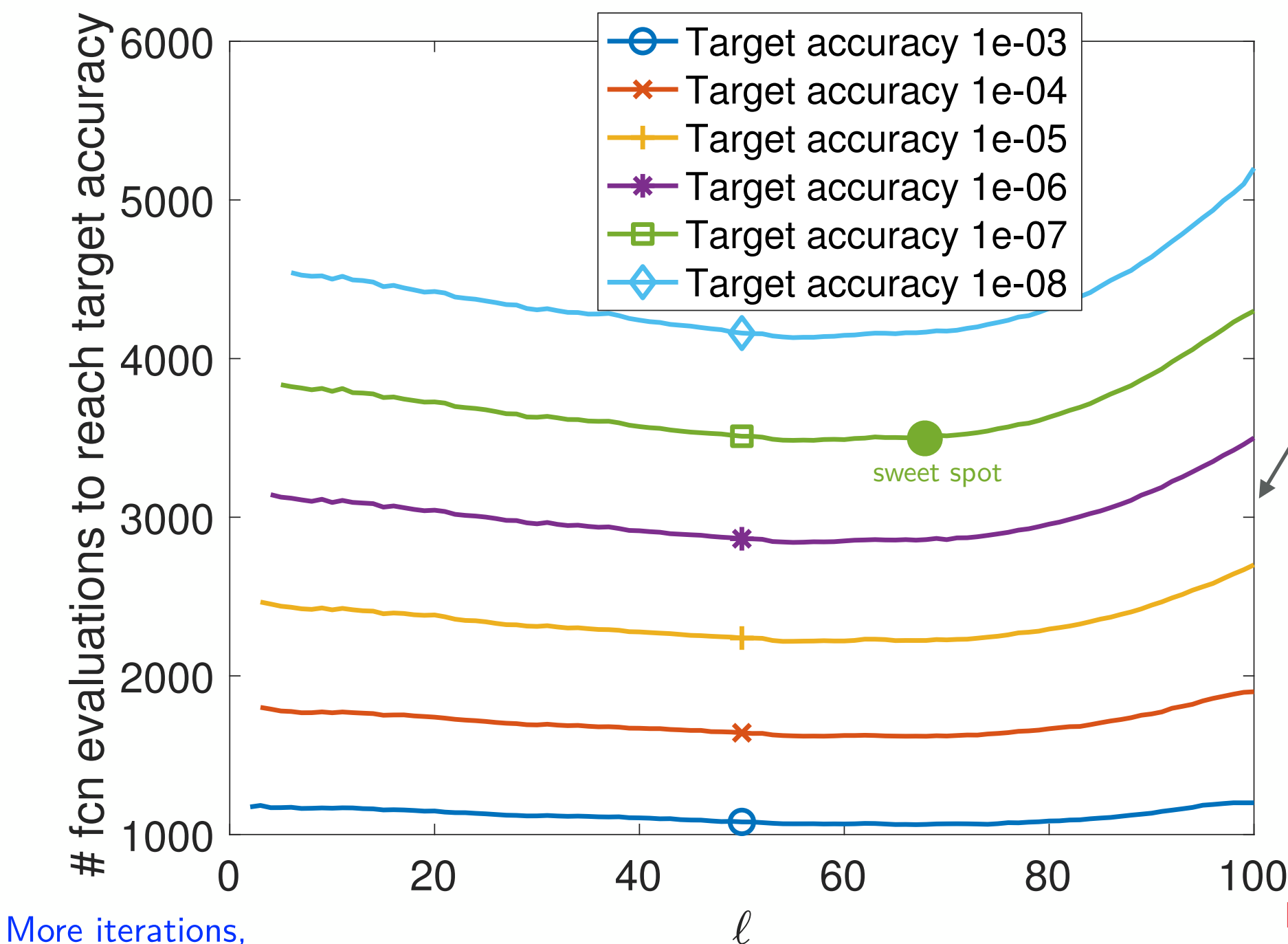
Most literature focuses on

How is q typically

- spherical (or Gaussian)
- or

- canonical basis vector, $q \sim \text{uniform}[e_1, \dots, e_d]$
- (hence SSD reduces to randomized coordinate descent)

Most work has focused on a single directional derivative, $\ell = 1$



100 dimensional quadratic test problem, using exact linesearch, averaged over 200 experiments

$\ell = d$ is gradient descent, and non-stochastic

Much recent work on variants, 2011—2020

- D. Leventhal and A.S. Lewis, *Randomized Hessian estimation and directional search*, Optimization (2011)
- Z. Zhang, *Scalable Derivative-Free Optimization Algorithms with Low-Dimensional Subspace Techniques* (thesis in Chinese 2011, preprint with proof sketches in English 2025, arXiv 2501.04536), “NEWUOAs”
- *Stochastic zeroth order descent with random pursuit*, SIAM J. Opt. (2013)
- Yu. Nesterov and V. Spokoiny, *FoCM* 2017
- *Method: A unifying framework for accelerated derivative-free method*, arXiv:1707.08486
- *Derivative method for smooth stochastic convex*
- *Stochastic programming*, SIAM J. Opt. (2013)
- *Functions*, arXiv:1507.03332 (2015).
- *Directed evolution with compact architectures for*
- *scalable alternative to reinforcement learning*,
- *Convex Optimization: The Power of Two*
- *Cal Comparison of Gradient Approximations in*
- *Sketching*, NeurIPS 2018
- Zhang, *Direct Search Based on Probabilistic*
- *gradient-free optimization in high dimensions.*
- *with orthogonal random directions*. Math.
- *er algorithm for nonsmooth optimization*. NeurIPS
- M. Rando, C. Molinaro, S. Villa, L. Rosasco. *Stochastic zeroth order descent with structured directions*. COAP, 2024.
- R. Nozawa, P.-L. Poirion, A. Takeda. *Zeroth-Order Random Subspace Algorithm for Non-smooth Convex Optimization*, JOTA, 2025.
- H.Q. Cai, D. McKenzie, W. Yin, Z. Zhang. *Zeroth-Order Regularized Optimization (ZORO): Approximately Sparse Gradients and Adaptive Sampling*, SIOPT, 2022

Our first analysis

Theorem (Kozak, Becker, Tenorio, Doostan '20)

Assume: minimizer attained, gradient Lipschitz, stepsize η_k chosen appropriately.

1. If f is **convex**,

$$\mathbb{E} f(\mathbf{x}_k) - f^* \leq 2 \frac{d}{\ell} \frac{L}{k} R^2 = \mathcal{O}(k^{-1})$$

2. If f is **not necessarily convex** but satisfies the **Polyak-Lojasiewicz** inequality,

$$\mathbb{E} f(\mathbf{x}_k) - f^* \leq \rho^k (f(\mathbf{x}_0) - f^*) = \mathcal{O}(\rho^k) \quad \text{and} \quad f(\mathbf{x}_k) \xrightarrow{\text{a.s.}} f^*$$

3. If f is **strongly convex**, statements of 2 above hold, and also

$$\mathbf{x}_k \xrightarrow{\text{a.s.}} \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x})$$

4. If f is **not convex** (nor PL),

$$\min_{k' \in \{0, \dots, k\}} \mathbb{E} \|\nabla f(\mathbf{x}_{k'})\|^2 \leq \frac{d}{\ell} \frac{2L(f(\mathbf{x}_0) - f^*)}{k+1}$$

Generic SSD

$$\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_{\ell \times \ell}, \quad \mathbb{E} \left(\frac{d}{\ell} \mathbf{Q} \mathbf{Q}^\top \right) = \mathbf{I}_{d \times d}$$



David Kozak

$$f^* \stackrel{\text{def}}{=} \min_{\mathbf{x}} f(\mathbf{x})$$

$$\rho = 1 - \frac{\mu}{L} \frac{\ell}{d}$$

d = ambient dimension

ℓ = # directional derivs

$\frac{d}{\ell} = 1$ is gradient descent

μ is PL constant

L is gradient Lipschitz constant

$\eta = \frac{\ell}{d} \frac{1}{L}$ is stepsize

Our first analysis

Our analysis is comparable to analysis of similar algorithms

Assume f obtains its minimum and ∇f is L -Lipschitz continuous.

Theorem 1 (Kozak, Becker, Tenorio, Doostan '19, Thm. 2.4). *The SSD algorithm with stepsize $\eta = \frac{1}{L} \frac{\ell}{d}$ gives*

$$\mathbb{E} f(x_k) - f^* \leq \boxed{2 \frac{d}{\ell} \frac{L}{k} R^2}$$

where

$$R = \sup_{x|f(x) \leq f(x_0)} \inf_{x^* \in \arg\min f} \|x - x^*\|$$

(e.g., f is coercive $\implies R < \infty$).

SSD

$$1 \leq \ell \leq d$$

Theorem 2 (Nesterov, Spokoiny '17, Thm. 8). *Take stepsize $\eta = \frac{1}{4(d+4)L}$, then the random gradient method with a Gaussian direction converges as*

$$\frac{1}{k} \sum_{i=0}^{k-1} \mathbb{E} f(x_i) - f^* \leq \boxed{\frac{4(d+4)L}{k} \|x_0 - x^*\|^2}$$

where x^* is any optimal solution.

Gaussian

$$\ell = 1$$

convex, not necessarily strongly convex, scenario

Our first analysis

Our analysis is comparable to analysis of similar algorithms

Assume f obtains its minimum, ∇f is L -Lipschitz continuous, and f is μ PL or strongly convex.

Theorem 3 (Kozak, Becker, Tenorio, Doostan '19, Cor. 2.3). *The SSD algorithm with stepsize $\eta = \frac{1}{L} \frac{\ell}{d}$ gives*

$$\mathbb{E} f(x_k) - f^* \leq \rho^k (f(x_0) - f^*) \quad \text{with} \quad \rho = \left[1 - \frac{\mu \ell}{L d} \right]. \quad 1 \leq \ell \leq d$$

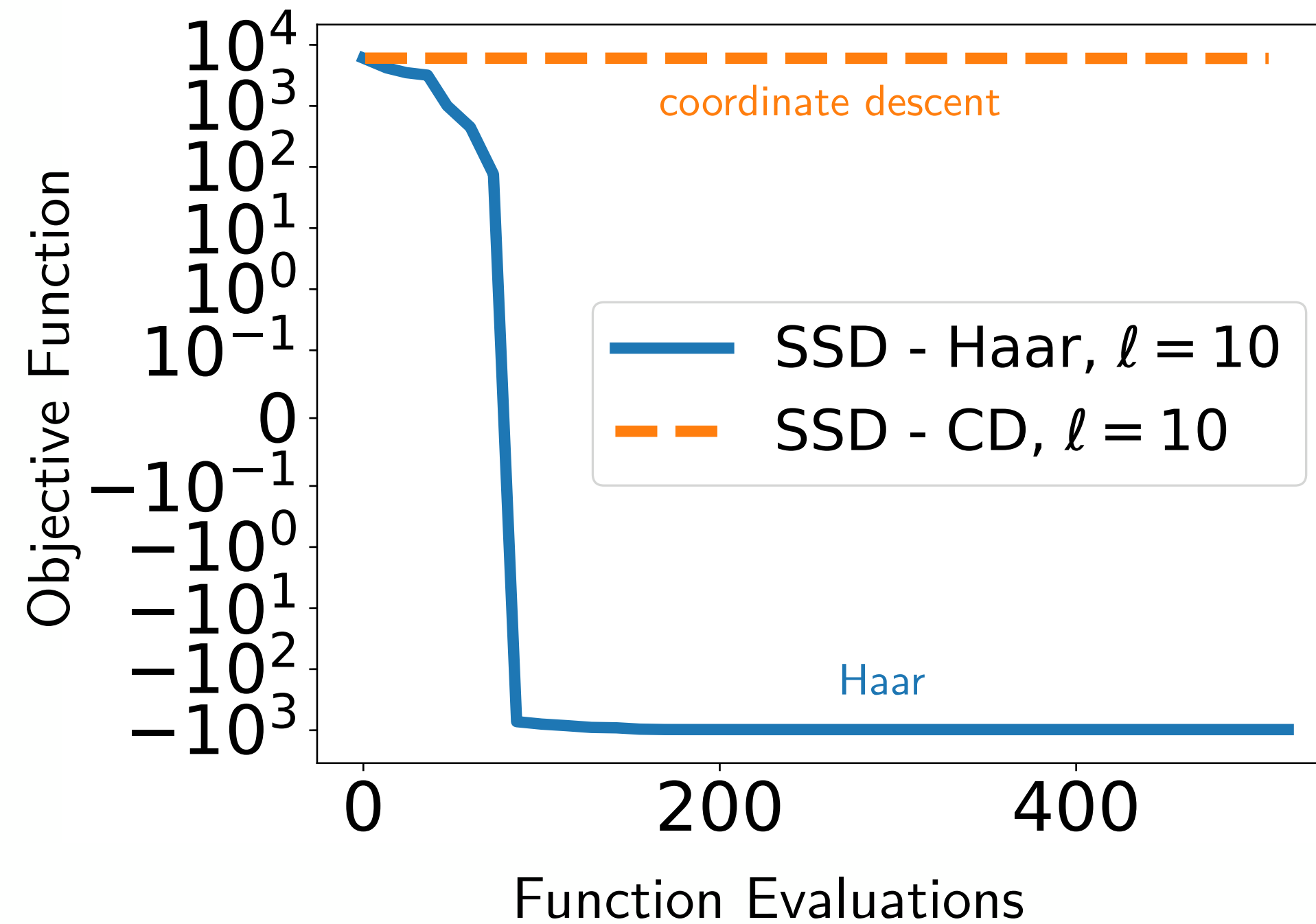
Theorem 4 (Nesterov, Spokoiny '17, Thm. 8). *Take stepsize $\eta = \frac{1}{4(d+4)L}$, then the random gradient method with a Gaussian direction converges as*

$$\mathbb{E} f(x_k) - f^* \leq \frac{L}{2} \rho^k \|x_0 - x^*\|^2 \quad \text{with} \quad \rho = \left[1 - \frac{\mu}{L} \frac{1}{8(d+4)} \right]$$

where x^* is any optimal solution.

strongly convex or PL scenario

... but that theory doesn't capture the full story



Observation: sometimes SSD (with Haar) drastically outperforms randomized coordinate descent (CD)

Both of them are valid instantiations of SSD

So, the details of Q matter!

Algorithm: “Haar” SSD

Draw the random matrix Q as follows:

$\{\mathbf{q}_1, \dots, \mathbf{q}_\ell\}$ a basis for $\text{span}\{\tilde{\mathbf{q}}_1, \dots, \tilde{\mathbf{q}}_\ell\}$ $\tilde{\mathbf{q}}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$
(or any uniformly random subspace)

Equivalently, draw from the Haar distribution

Satisfies

$$Q^\top Q = \mathbf{I}_{\ell \times \ell}, \quad \mathbb{E} \left(\frac{d}{\ell} Q Q^\top \right) = \mathbf{I}_{d \times d}$$

but goes even further.

... but that theory doesn't capture the full story

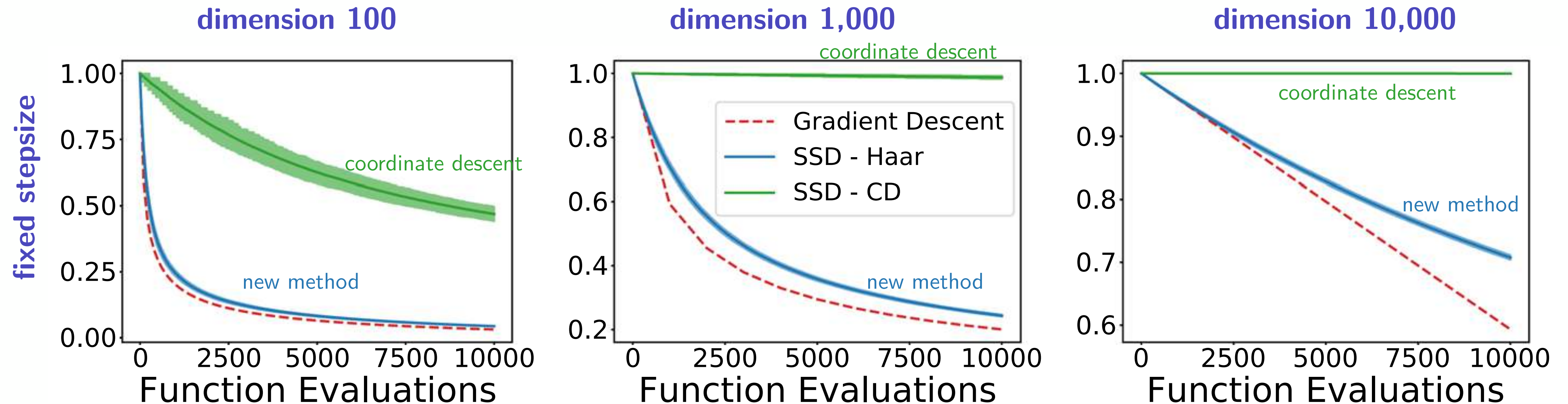
SSD drastically outperforms
randomized coordinate descent (CD)

We can force it to happen by making a
problem with low “intrinsic” dimension, e.g.,
Nesterov’s “worst function in the world”

$$f_{\lambda,r}(\mathbf{x}) = \lambda((x_1^2 + \sum_{i=1}^{r-1} (x_i - x_{i+1})^2 + x_r^2)/2 - x_1)/4,$$

This has **intrinsic dimension** of r

$r = 20, \ell = 3$



Improved theory (specialized to SSD-Haar)

Recall...
 d = ambient dimension
 ℓ = # directional derivs

Tighter analysis using concentration-of-measure:

Lemma 2 (Johnson-Lindenstrauss style embedding, from Kozak, Becker, Tenorio '19, Lemma 1).

$\forall \epsilon \in (0, 1)$, if $\ell \gtrsim \epsilon^{-2}$, $Q \sim \text{Haar}(d \times \ell)$, then $\forall 0 \neq g \in \mathbb{R}^d$,

$$1 - \epsilon \leq \frac{d}{\ell} \frac{\|Q^T g\|^2}{\|g\|^2} \leq 1 + \epsilon \quad w/ \text{ prob. } \delta \geq 0.8$$

Note: *coordinate descent* style projections do **not** have similar nice embedding properties

Improved theory (specialized to SSD-Haar)

Recall...
 d = ambient dimension
 ℓ = # directional derivs

Tighter analysis using concentration-of-measure:

Lemma 2 (Johnson-Lindenstrauss style embedding, from Kozak, Becker, Tenorio '19, Lemma 1).

$\forall \epsilon \in (0, 1)$, if $\ell \gtrsim \epsilon^{-2}$, $Q \sim \text{Haar}(d \times \ell)$, then $\forall 0 \neq g \in \mathbb{R}^d$,

$$1 - \epsilon \leq \frac{d}{\ell} \frac{\|Q^T g\|^2}{\|g\|^2} \leq 1 + \epsilon \quad \text{w/ prob. } \delta \geq 0.8$$

... and putting it all together

Theorem 3 (Kozak, Becker, Tenorio '19, Thm. 1). *If f is strongly convex and ∇f is Lipschitz continuous, then for an appropriate stepsize η_k , the sequence (x_k) generated by SSD (with $Q \sim \text{Haar}$), for $k > 100$, satisfies*

$$f(x_k) - f^* \leq (1 + (1 - \epsilon)\rho)^{k/2} (f(x_0) - f^*) \quad \text{with probability } \geq 0.998,$$

where $\rho < 1$ depends on ℓ , d and the Lipschitz and strong convexity parameters.

error in JL embedding

due to possibility of failure of JL

Extension: Variance Reduction

Algorithm SVRG-style Variance Reduced SSD method, “VRSSD”

```

1: for  $k = 1, 2, \dots$  do                                 $\triangleright k$  is the “epoch”
2:    $\bar{z} \leftarrow \nabla f(x_k)$                                  $\triangleright$  Expensive, but not done often
3:    $w_0 \leftarrow x_k$ 
4:   for  $t = 1, 2, \dots, T$  do                                 $\triangleright$  Typically  $T = \mathcal{O}(d)$ 
5:     Draw  $Q \sim \text{Haar}(d \times \ell)$ 
6:      $w_{t+1} \leftarrow w_t - \eta \left( \underbrace{\frac{d}{\ell} Q Q^T \nabla f(w_t)}_{\text{regular SSD term}} - \alpha_k \underbrace{\left( \frac{d}{\ell} Q Q^T - I \right) \bar{z}}_{\text{orthogonal projection}} \right)$   $\triangleright \alpha_k$  to be estimated
7:    $x_{k+1} \leftarrow w_T$ 

```

only use control variate in orthogonal subspace
(since we know gradient in main subspace)

Theorem 4 (Kozak, Becker, Tenorio, Doostan 2019; Thm. 2.7). *If f is strongly convex and ∇f is Lipschitz continuous, then for an appropriate stepsize η_k , the sequence (x_k) generated by VRSSD converges almost surely to the (unique) minimizer of f and at a linear rate (the rate depends on η_k and α_k).*



We do not require the ERM structure!

from the literature:

Algorithm SAGA (Defazio, Bach, Lacoste-Julien '14) for solving the ERM model

```

1:  $\forall i = 1, \dots, N, x^{(i)} \stackrel{\text{def}}{=} x_0$ ; store  $\{\nabla f_i(x^{(i)})\}_{i=1}^N$  in table
2: for  $k = 1, 2, \dots$  do
3:   Draw  $j \sim \text{Uniform}([1, \dots, N])$ 
4:    $\bar{z} \leftarrow \frac{1}{N} \sum_{i=1}^N \nabla f_i(x^{(i)})$   $\triangleright$  From table
5:    $x_{k+1} \leftarrow x_k - \eta \left( \nabla f_j(x_k) - \nabla f_j(x^{(j)}) + \bar{z} \right)$ 
6:   Re-define  $x^{(j)} \leftarrow x_k$  and update table with  $\nabla f_j(x^{(j)})$ 

```

our variant:

Algorithm SAGA-style Variance Reduced SSD method

```

1: Pre-compute  $\bar{z} \leftarrow \nabla f(x_0)$ 
2: for  $k = 1, 2, \dots$  do
3:   Draw  $Q \sim \text{Haar}(p \times r)$ 
4:    $x_{k+1} \leftarrow x_k - \eta \left( \frac{d}{\ell} Q Q^T \nabla f(x_k) - \frac{d}{\ell} Q Q^T \bar{z} + \bar{z} \right)$ 
5:    $\bar{z} \leftarrow \bar{z} + Q Q^T (\nabla f(x_k) - \bar{z})$   $\triangleright$  Update of  $\bar{z}$  is low-memory, unlike original SAGA

```

key: update control variate in the subspace

$$f_c(x) \approx f(x)$$

generic (non-algorithmic) control variates

control variate, coarse approximation, cheap to evaluate

Algorithm Proposed coarse-model variance reduced SSD/Random-Gradient

```

1: for  $k = 1, 2, \dots$  do
2:    $\bar{z} \leftarrow \nabla f_c(x_k)$   $\triangleright$  Full coarse-grid gradient
3:   Draw  $Q \sim \text{Haar}(d \times \ell)$ 
4:    $x_{k+1} \leftarrow x_k - \eta_k \left( \frac{d}{\ell} Q Q^T \nabla f(x_k) + \alpha_k \left( \frac{d}{\ell} Q Q^T \bar{z} - \bar{z} \right) \right)$ 

```

Key idea: easy to do **orthogonal projection**

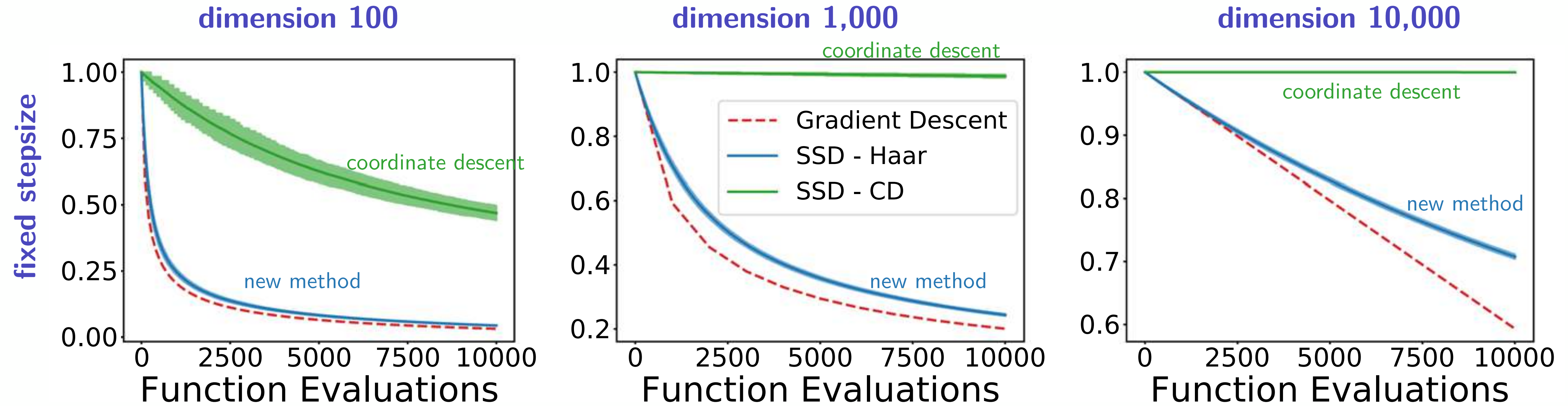
Extension: stepsize selection

$$r = 20, \ell = 3$$

Nesterov's "worst function in the world"

$$f_{\lambda,r}(\mathbf{x}) = \lambda((x_1^2 + \sum_{i=1}^{r-1} (x_i - x_{i+1})^2 + x_r^2)/2 - x_1)/4,$$

This has **intrinsic dimension** of r



Recall previous example

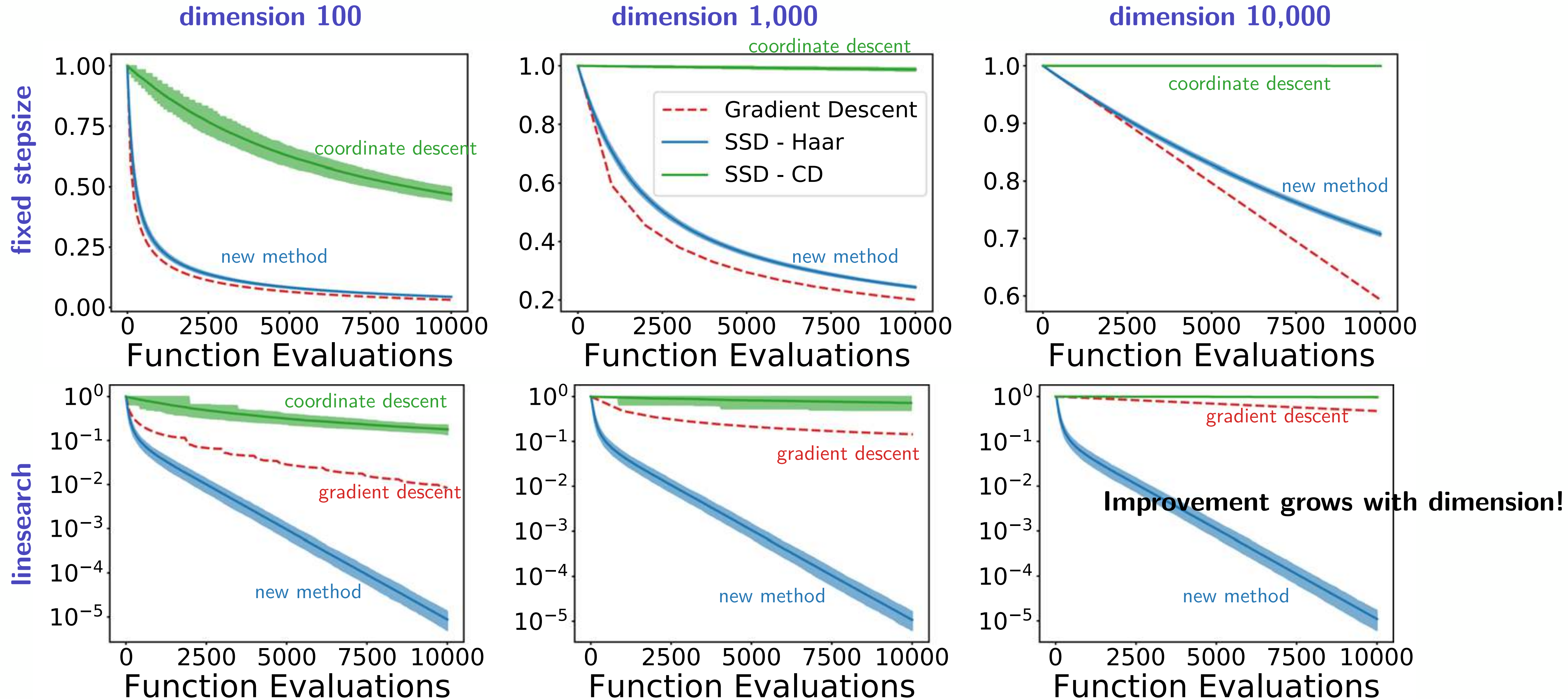
Extension: stepsize selection

$$r = 20, \ell = 3$$

Nesterov's "worst function in the world"

$$f_{\lambda,r}(\mathbf{x}) = \lambda((x_1^2 + \sum_{i=1}^{r-1} (x_i - x_{i+1})^2 + x_r^2)/2 - x_1)/4,$$

This has **intrinsic dimension** of r



Extension: stepsize selection

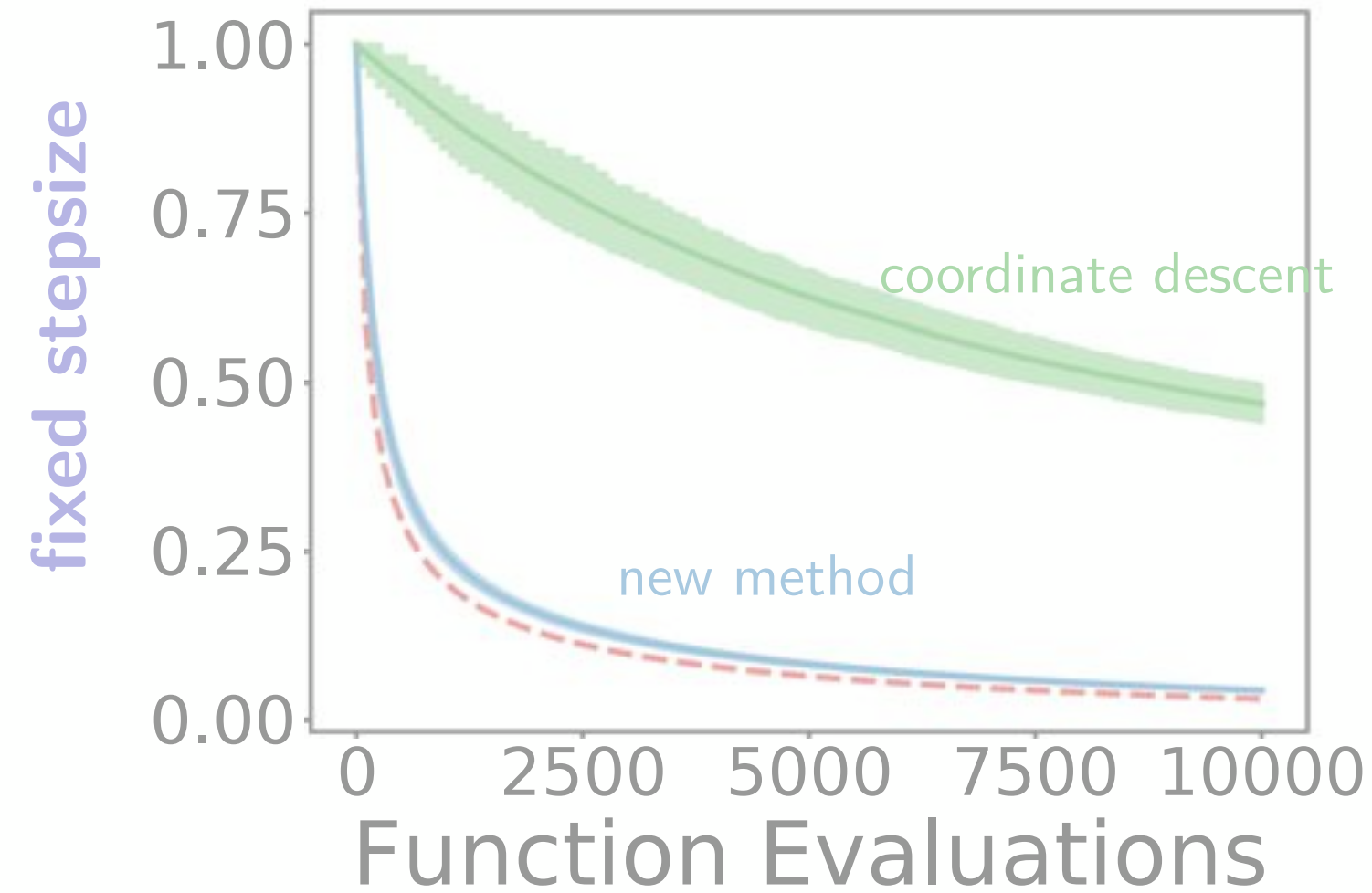
$$r = 20, \ell = 3$$

Nesterov's "worst function in the world"

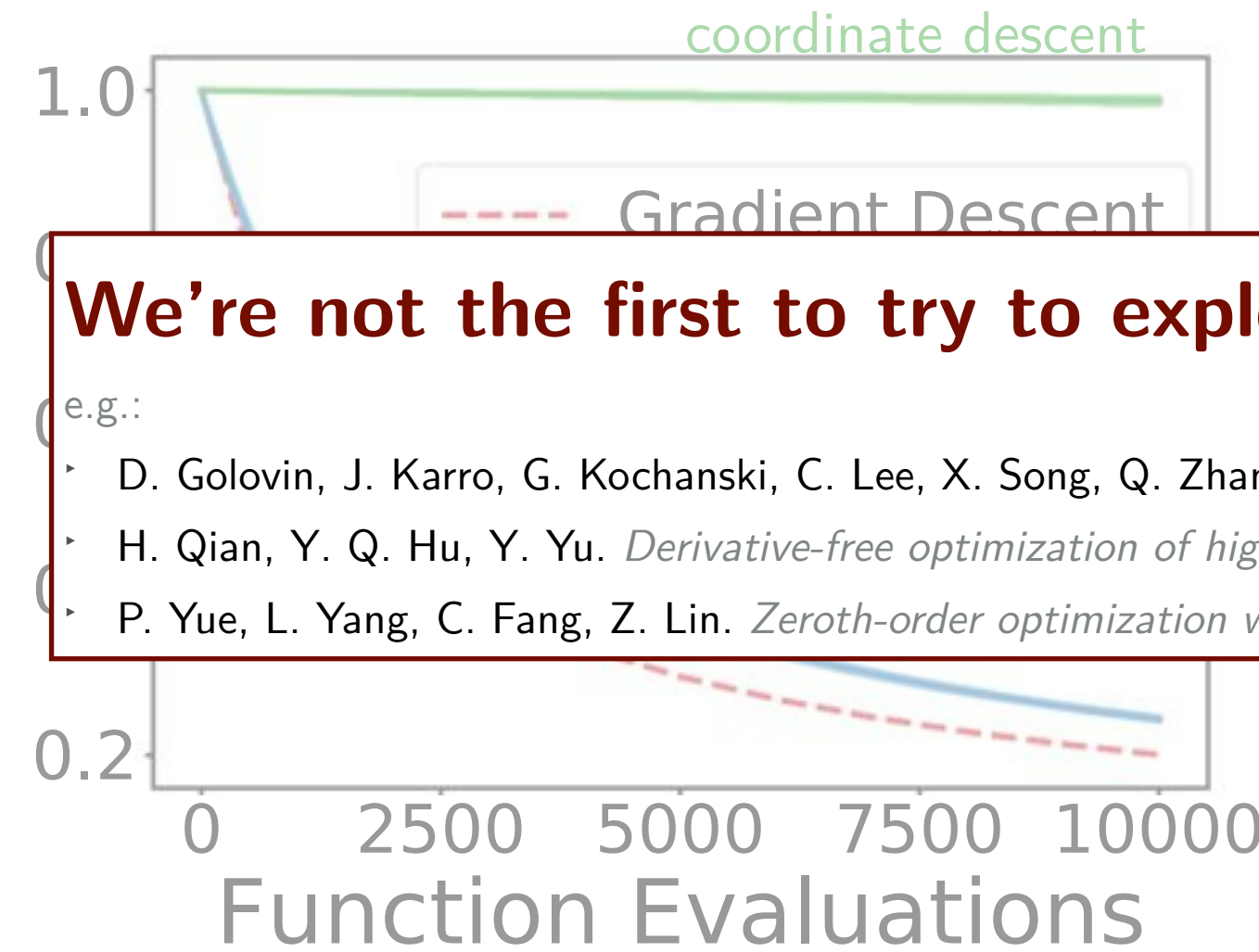
$$f_{\lambda,r}(\mathbf{x}) = \lambda((x_1^2 + \sum_{i=1}^{r-1} (x_i - x_{i+1})^2 + x_r^2)/2 - x_1)/4,$$

This has **intrinsic dimension** of r

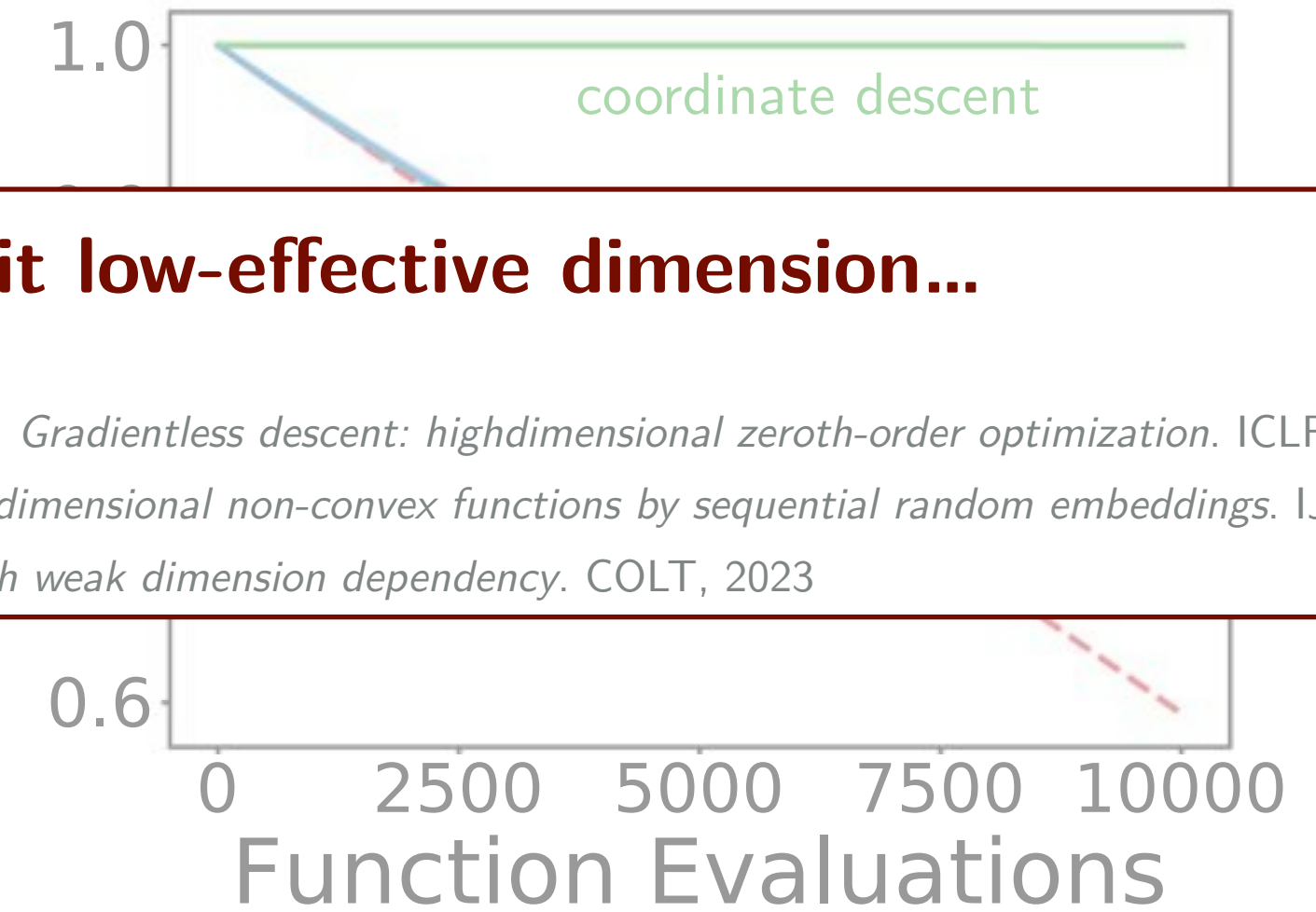
dimension 100



dimension 1,000



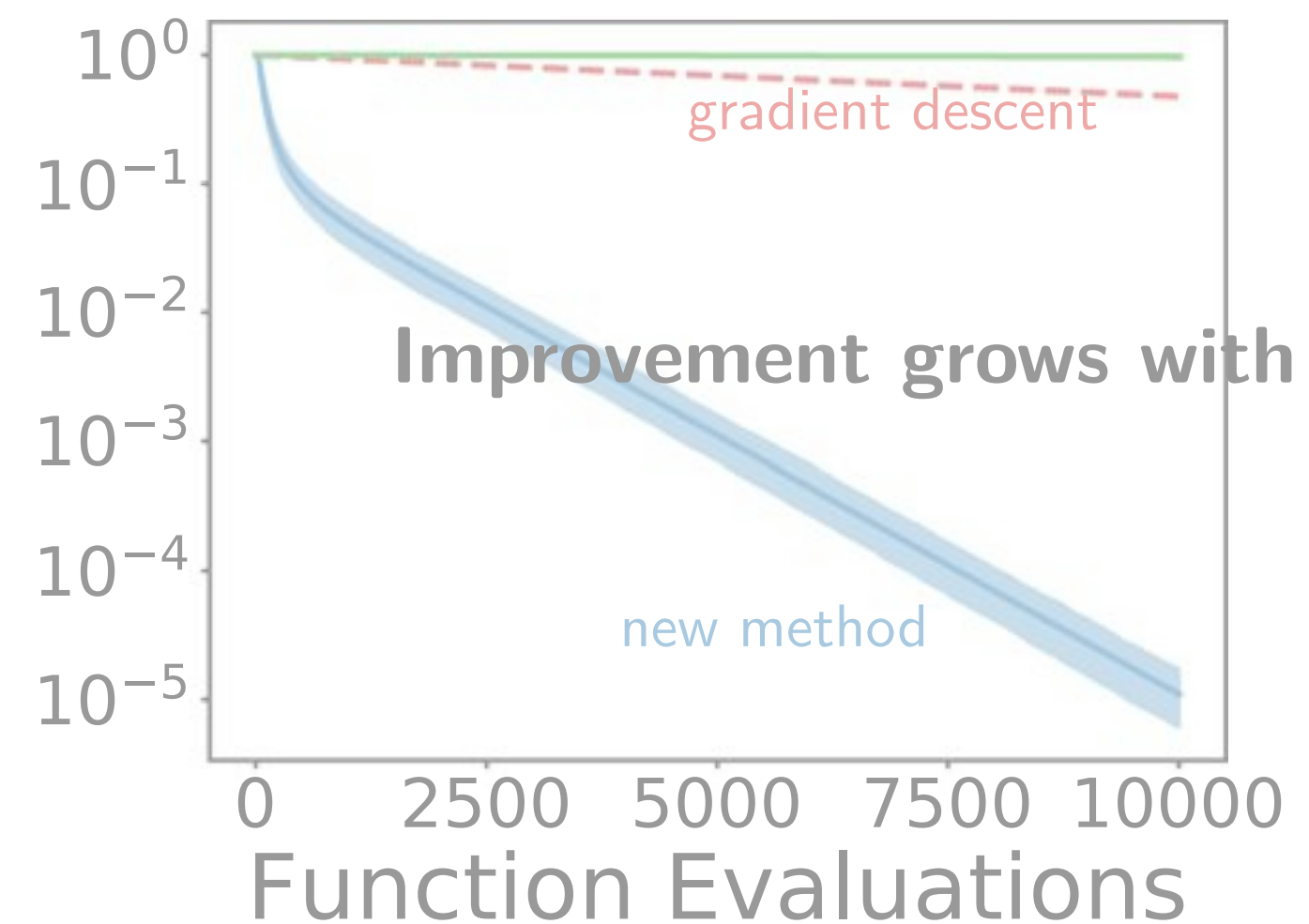
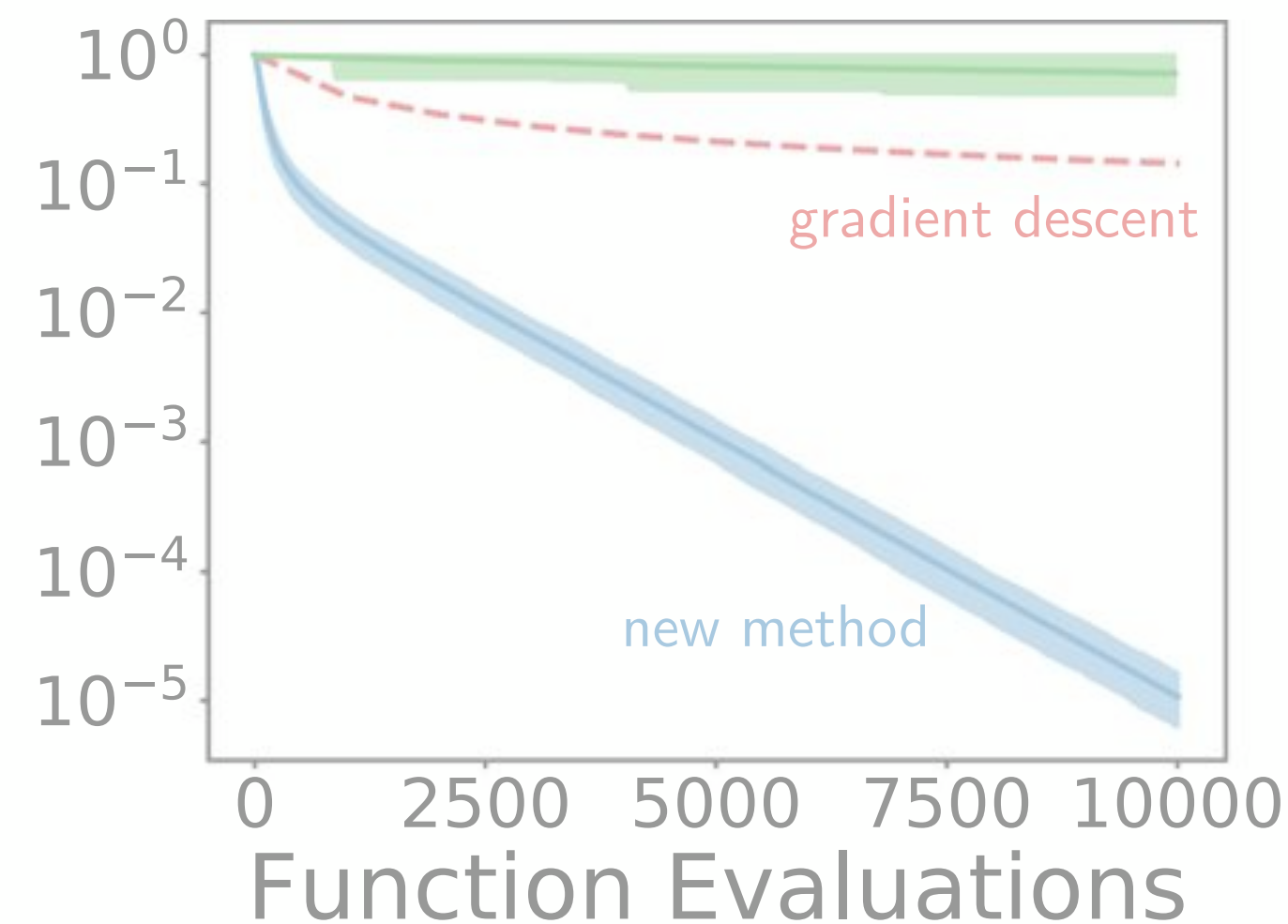
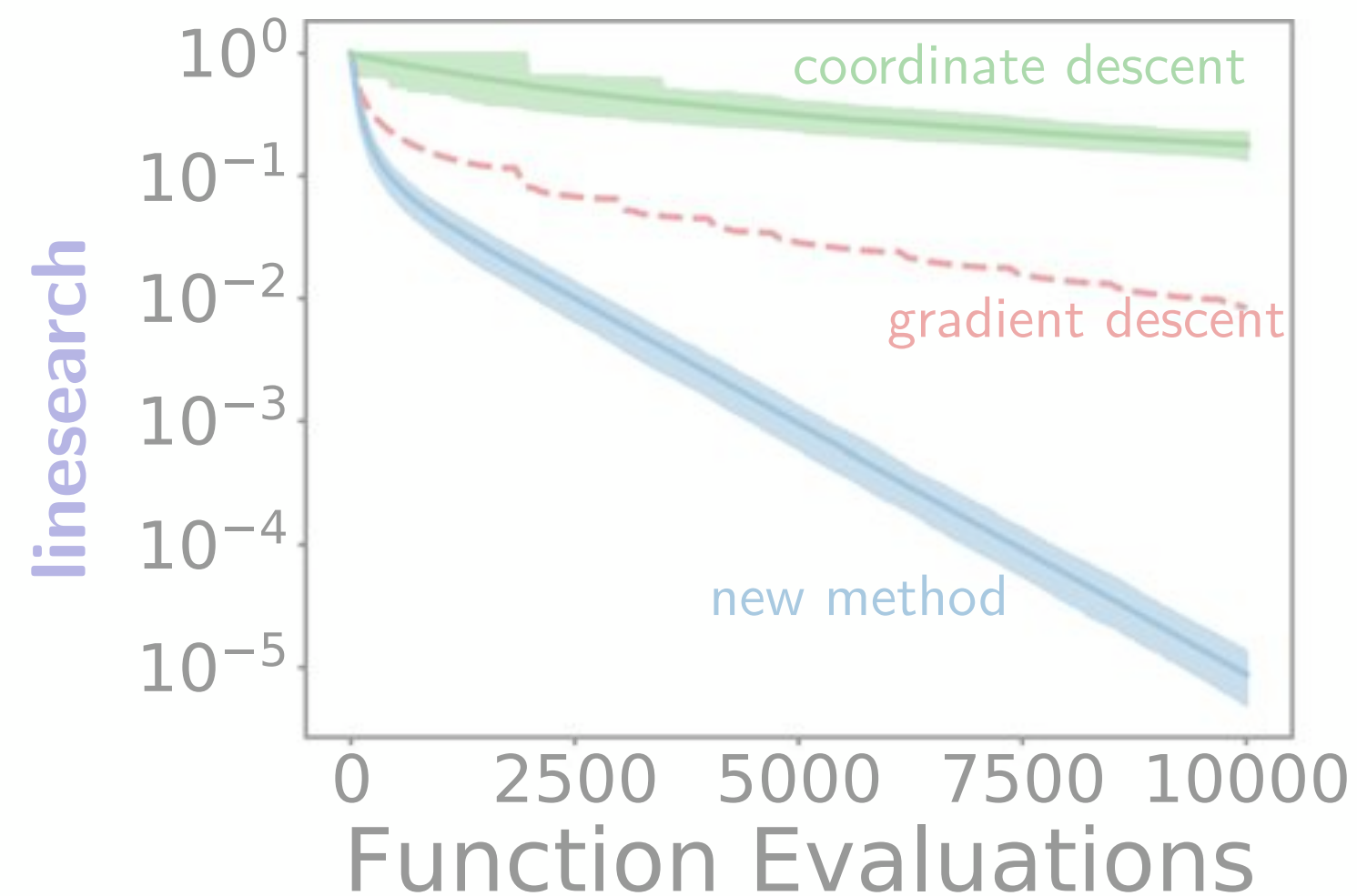
dimension 10,000



We're not the first to try to exploit low-effective dimension...

e.g.:

- D. Golovin, J. Karro, G. Kochanski, C. Lee, X. Song, Q. Zhang. *Gradientless descent: highdimensional zeroth-order optimization*. ICLR, 2020
- H. Qian, Y. Q. Hu, Y. Yu. *Derivative-free optimization of high-dimensional non-convex functions by sequential random embeddings*. IJCAI, 2016
- P. Yue, L. Yang, C. Fang, Z. Lin. *Zeroth-order optimization with weak dimension dependency*. COLT, 2023



Extension: stepsize selection

$$r = 20, \ell = 3$$

Nesterov's "worst function in the world"

$$\frac{r-1}{r}$$

$$(x_{i+1})^2 + x_r^2)/2 - x_1)/4,$$

Working hypothesis: SSD-Haar nicely exploits low-dimensional structure...

... if we have an aggressive stepsize

So the main question is, **how to choose the stepsize?**

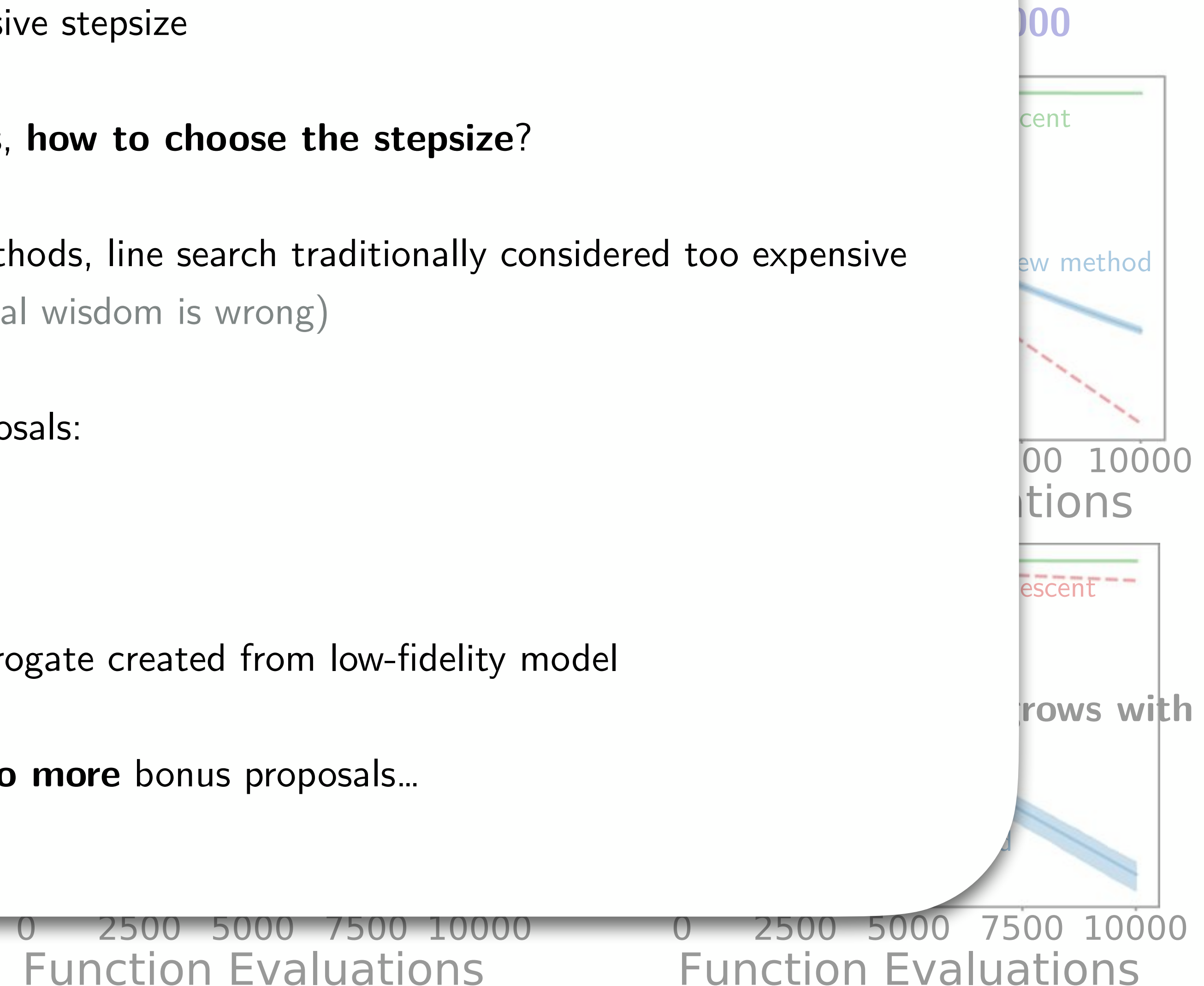
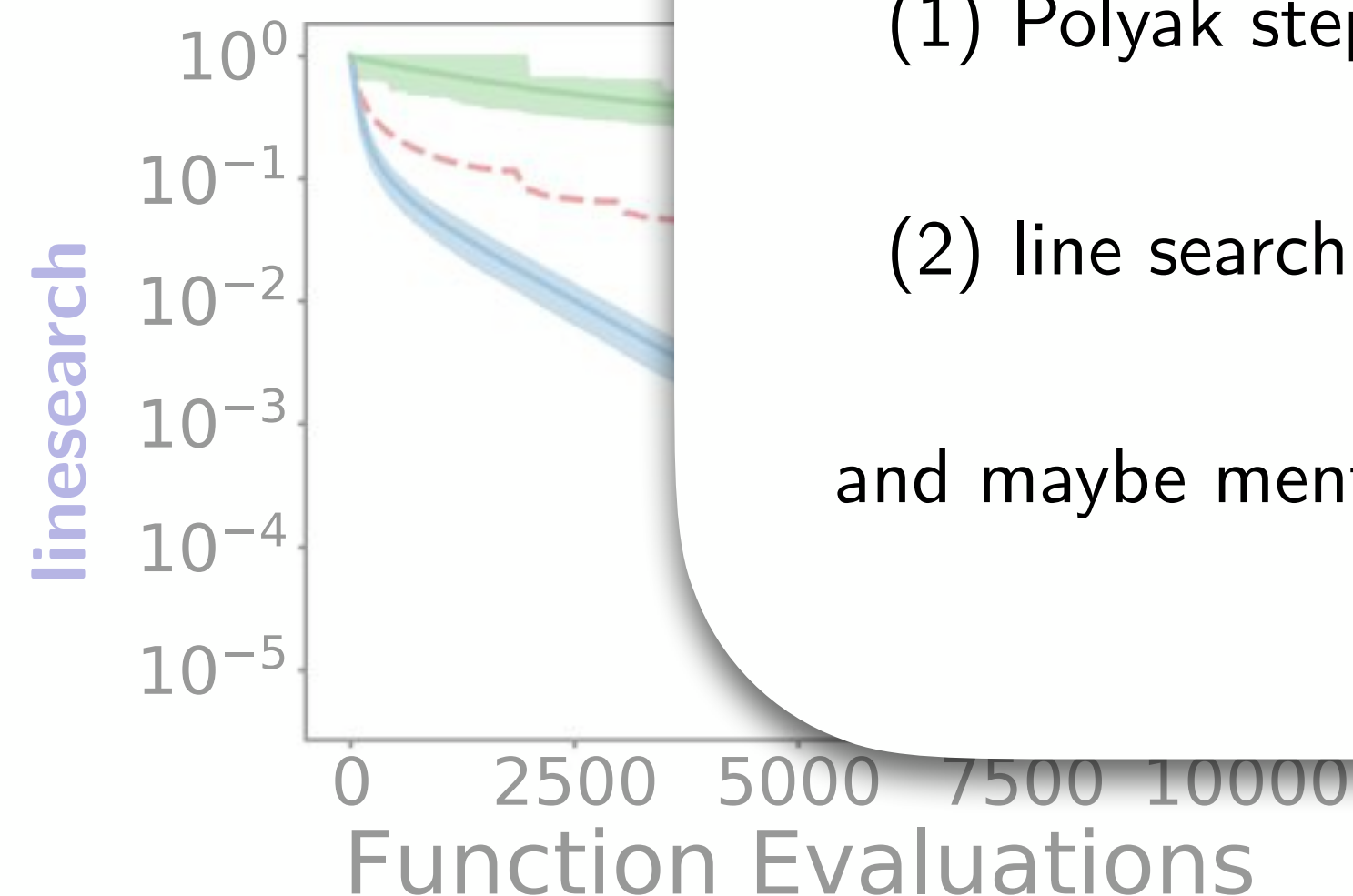
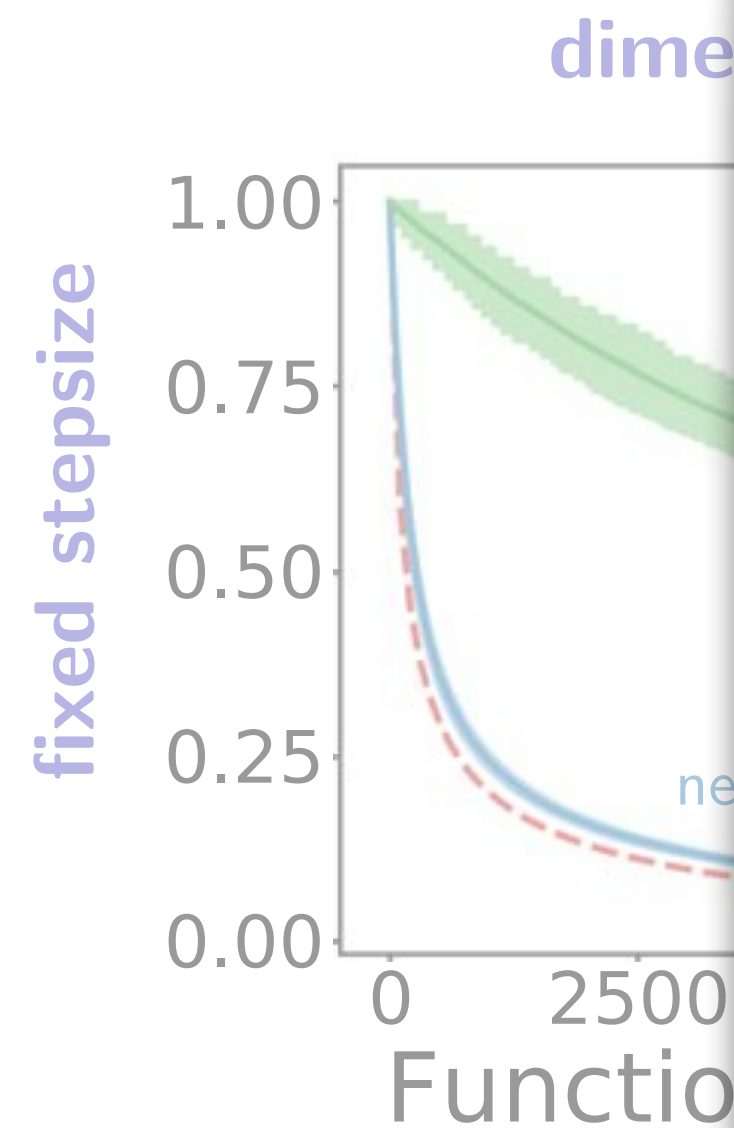
For DFO/0th order methods, line search traditionally considered too expensive
(perhaps conventional wisdom is wrong)

We'll provide **two** proposals:

(1) Polyak stepsize

(2) line search on surrogate created from low-fidelity model

and maybe mention **two more** bonus proposals...



Stepsize selection: Polyak stepsize

Joint project with Killian Wood,
Drona Khurana



Polyak stepsize for gradient descent (1983)

Recently revisited a lot in literature

$$\eta_k^{\text{Polyak}} = \frac{f(\mathbf{x}_k) - f^*}{\|\nabla f(\mathbf{x}_k)\|^2} \quad f^* = \min_{\mathbf{x}} f(\mathbf{x})$$
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k^{\text{Polyak}} \nabla f(\mathbf{x}_k)$$

well, to be precise, actually

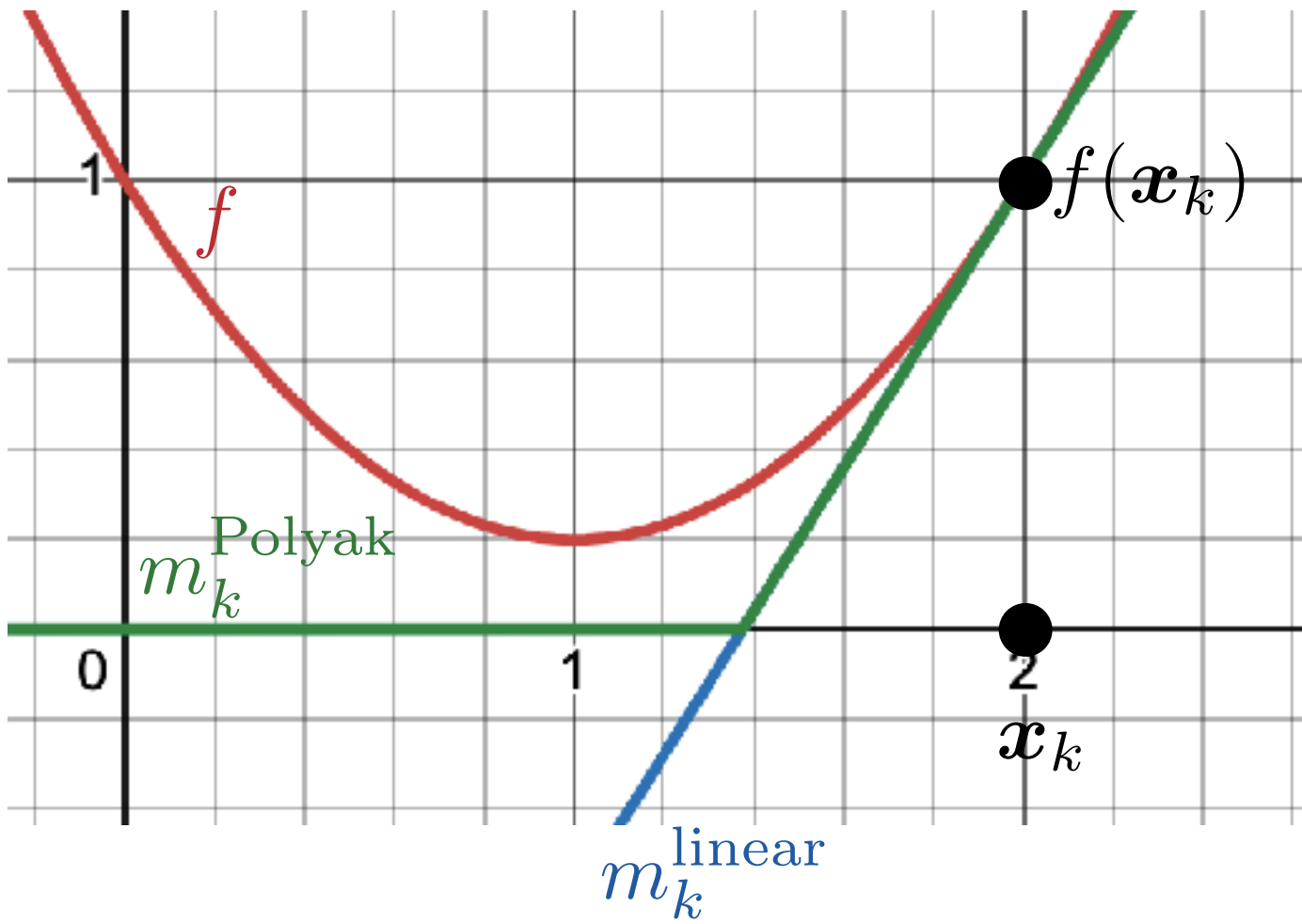
$$\eta_k^{\text{Polyak}} = \min \left\{ \frac{f(\mathbf{x}_k) - f^*}{\|\nabla f(\mathbf{x}_k)\|^2}, \eta_{\max} \right\}$$

cf. John Duchi’s semi-plenary July 22

Derivation from a *model-based view point*: $\mathbf{x}_{k+1} = \underset{\mathbf{y}}{\operatorname{argmin}} \ m_k(\mathbf{y}) + \frac{1}{2\eta_{\max}} \|\mathbf{y} - \mathbf{x}_k\|^2$

Asi & Duchi '19; Loyizou et al. '21; Berrada et al. '19; Davis & Drusvyatskiy '19; Schaipp et al. '23

Model	Algorithm
$m_k(\mathbf{y}) = f(\mathbf{y})$	Proximal Point
$m_k(\mathbf{y}) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{y} - \mathbf{x}_k \rangle$	Gradient descent, fixed stepsize
$m_k(\mathbf{y}) = \max \{ f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{y} - \mathbf{x}_k \rangle, f^* \}$	Gradient descent, Polyak stepsize



Stepsize selection: Polyak stepsize

Polyak stepsize for gradient descent (1983)

Recently revisited a lot in literature

$$\eta_k^{\text{Polyak}} = \frac{f(\mathbf{x}_k) - f^*}{\|\nabla f(\mathbf{x}_k)\|^2} \quad f^* = \min_{\mathbf{x}} f(\mathbf{x})$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k^{\text{Polyak}} \nabla f(\mathbf{x}_k)$$

Our extension to SSD case:

$$\eta_k^{\text{Polyak-SSD}} = \frac{f(\mathbf{x}_k) - f_k^*}{\|\mathbf{Q}^\top \nabla f(\mathbf{x}_k)\|^2} \quad f_k^* = \min_{\mathbf{x} \in \{\mathbf{x}_k\} + \text{col } \mathbf{Q}} f(\mathbf{x})$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k^{\text{Polyak-SSD}} \mathbf{Q} \mathbf{Q}^\top \nabla f(\mathbf{x}_k)$$

Derivation from a *model-based view point*: $\mathbf{x}_{k+1} = \underset{\mathbf{y}}{\text{argmin}} \ m_k(\mathbf{y}) + \frac{1}{2\eta_{\max}} \|\mathbf{y} - \mathbf{x}_k\|^2$

keep the same model... $m_k(\mathbf{y}) = \max \{f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{y} - \mathbf{x}_k \rangle, f_k^*\}$

... but change the context: $\mathbf{x}_{k+1} = \underset{\mathbf{y} \in \{\mathbf{x}_k\} + \text{col } \mathbf{Q}}{\text{argmin}} \ m_k(\mathbf{y}) + \frac{1}{2\eta} \|\mathbf{y} - \mathbf{x}_k\|^2$

Stepsize selection: Polyak stepsize

Polyak stepsize for gradient descent (1983)

Recently revisited a lot in literature

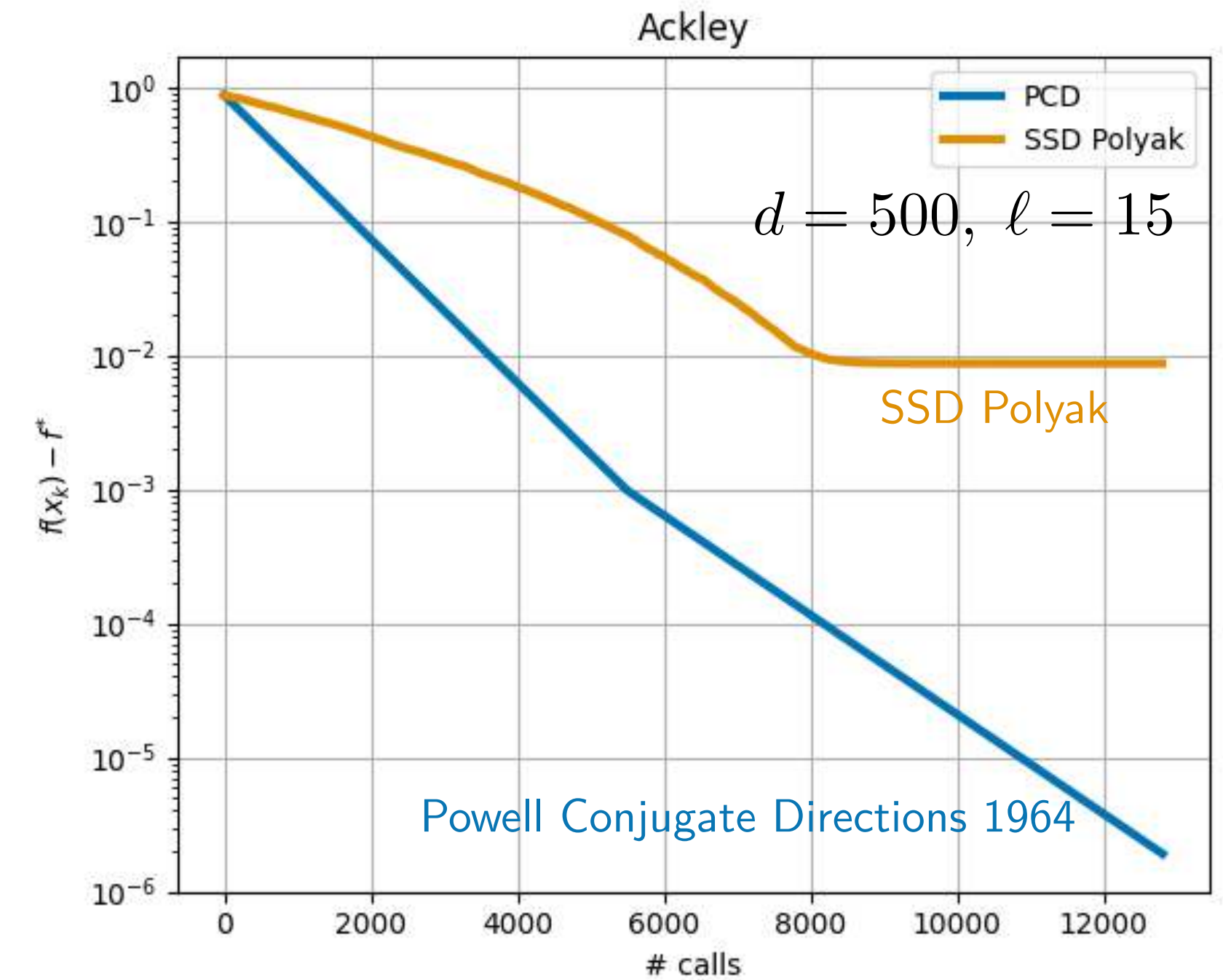
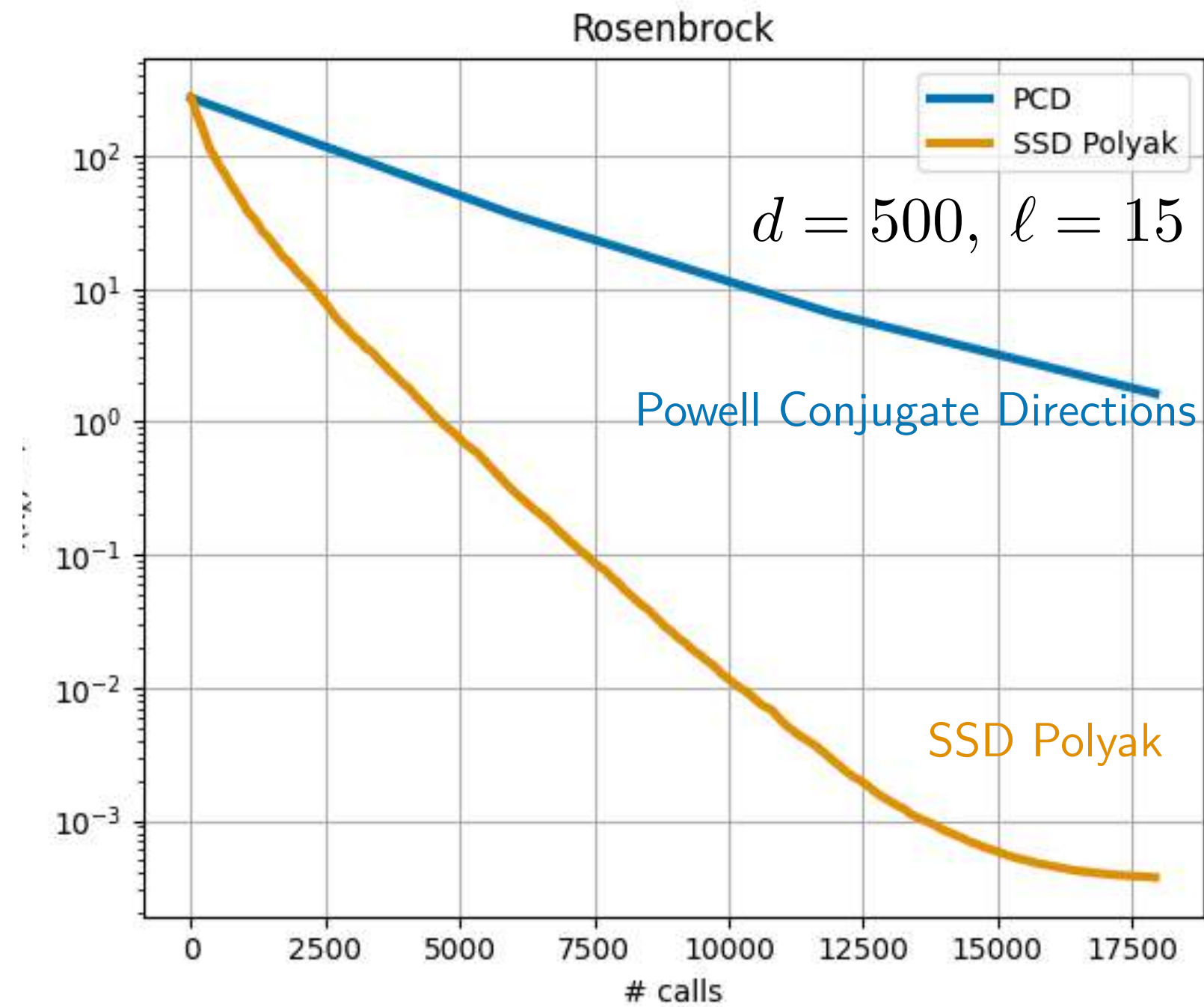
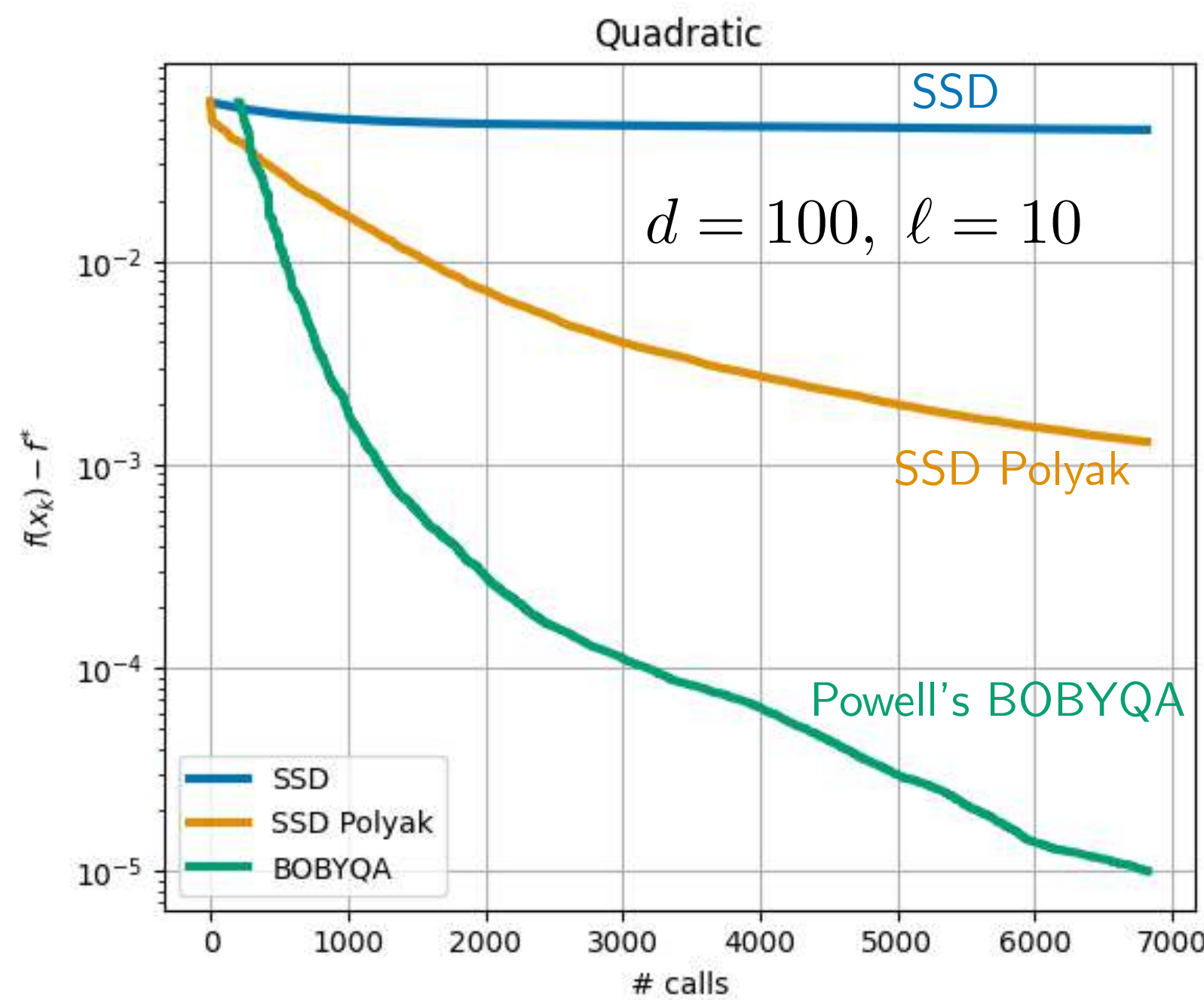
$$\eta_k^{\text{Polyak}} = \frac{f(\mathbf{x}_k) - f^*}{\|\nabla f(\mathbf{x}_k)\|^2} \quad f^* = \min_{\mathbf{x}} f(\mathbf{x})$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k^{\text{Polyak}} \nabla f(\mathbf{x}_k)$$

Our extension to SSD case:

$$\eta_k^{\text{Polyak-SSD}} = \frac{f(\mathbf{x}_k) - f_k^*}{\|\mathbf{Q}^\top \nabla f(\mathbf{x}_k)\|^2} \quad f_k^* = \min_{\mathbf{x} \in \{\mathbf{x}_k\} + \text{col} \mathbf{Q}} f(\mathbf{x})$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k^{\text{Polyak-SSD}} \mathbf{Q} \mathbf{Q}^\top \nabla f(\mathbf{x}_k)$$



Works fine in practice, our analysis is ongoing

Stepsize selection: bifidelity surrogate

Joint project with Nuojin (Noki) Cheng
(Google)



Classic exact linesearch

$$\mathbf{g}_k = \mathbf{Q}\mathbf{Q}^\top \nabla f(\mathbf{x}_k)$$

$$\eta^* = \operatorname{argmin} \varphi(\eta) \quad \varphi(\eta) \stackrel{\text{def}}{=} f(\mathbf{x}_k - \eta \mathbf{g}_k)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta^* \mathbf{g}_k$$

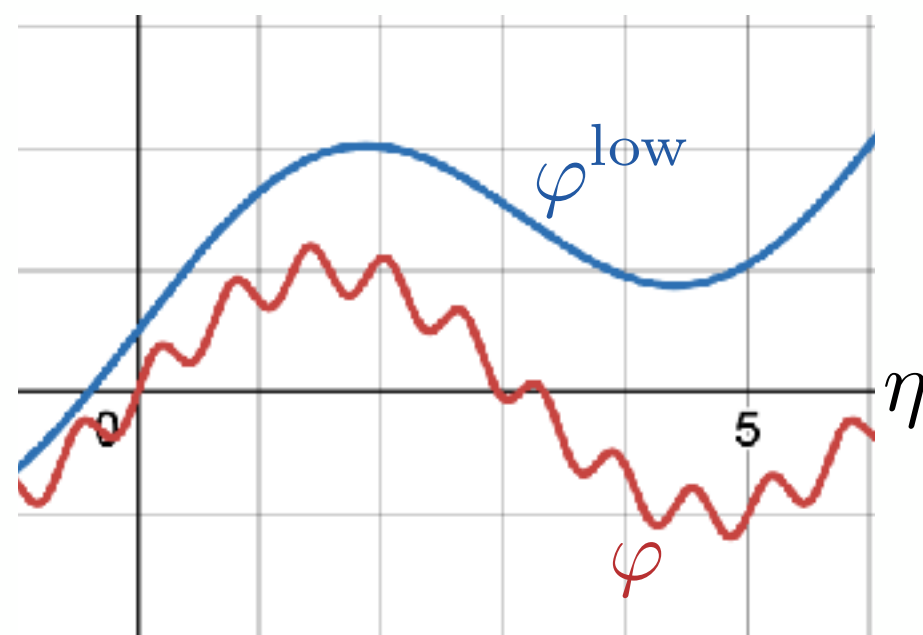
ideally would do line search
on this but too expensive



Premise: suppose we have a cheap, inaccurate approximation f^{low}

Expensive $\varphi(\eta) \stackrel{\text{def}}{=} f(\mathbf{x}_k - \eta \mathbf{g}_k)$

Inaccurate $\varphi^{\text{low}}(\eta) \stackrel{\text{def}}{=} f^{\text{low}}(\mathbf{x}_k - \eta \mathbf{g}_k)$



Stepsize selection: bifidelity surrogate

Joint project with Nuojin (Noki) Cheng
(Google)



Classic exact linesearch

$$\mathbf{g}_k = \mathbf{Q}\mathbf{Q}^\top \nabla f(\mathbf{x}_k)$$

$$\eta^* = \operatorname{argmin} \varphi(\eta) \quad \varphi(\eta) \stackrel{\text{def}}{=} f(\mathbf{x}_k - \eta \mathbf{g}_k)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta^* \mathbf{g}_k$$

ideally would do line search
on this but too expensive



Premise: suppose we have a cheap, inaccurate approximation f^{low}

Expensive

$$\varphi(\eta) \stackrel{\text{def}}{=} f(\mathbf{x}_k - \eta \mathbf{g}_k)$$

Inaccurate

$$\varphi^{\text{low}}(\eta) \stackrel{\text{def}}{=} f^{\text{low}}(\mathbf{x}_k - \eta \mathbf{g}_k)$$

data (function evaluations)

$$\{\varphi(0), \varphi(\eta_{\max})\}$$

$$\{\varphi^{\text{low}}(\eta_i)\}_{i=1}^{20}$$



co-kriging (1D)



surrogate model

$$\psi(\eta)$$

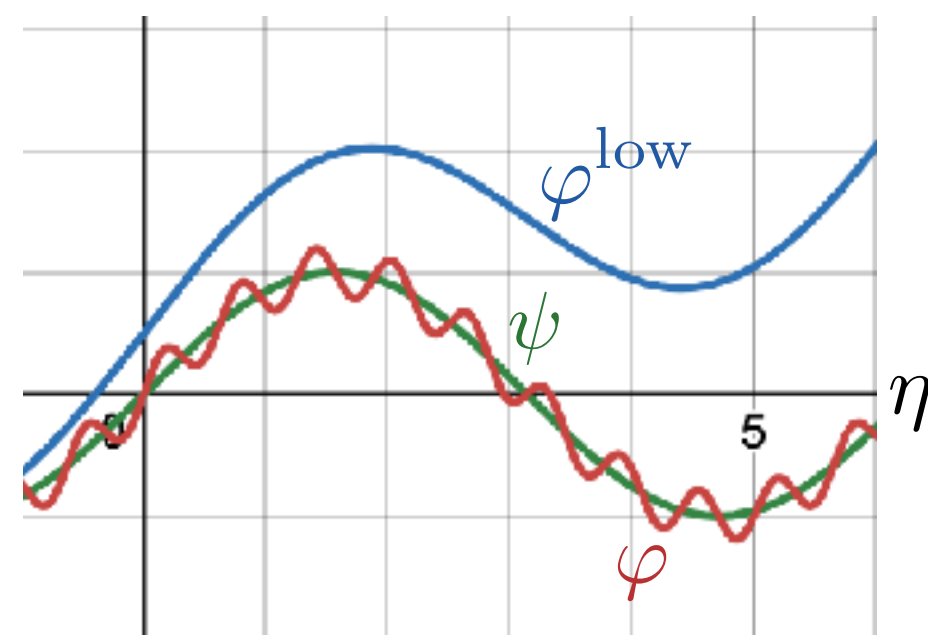


$$\eta^* = \operatorname{argmin} \psi(\eta)$$

traditional line search
on surrogate model

(computationally “free”)

e.g., calibrate low-fidelity model



Convergence analysis in our preprint “Stochastic Subspace Descent Accelerated via Bi-fidelity Line Search”
arxiv.org/abs/2505.00162, Nuojin Chen, Alireza Doostan, Stephen Becker

ML bifidelity example 1

Context:

- black-box model
- high-dimensional, low accuracy

Premise: suppose we have a cheap, inaccurate approximation f^{low}

Example: black-box adversarial attack Carlini & Wagner '17, black-box extension Chen et al. '17

For a given sample, find a small perturbation such that the machine learning algorithm misclassifies it

$$\min_{\epsilon} -f_{\text{cross-entropy}}(g(x^{\dagger} + \epsilon), y^{\dagger}) + \tau \|\epsilon\|^2$$

particular training example
(features and true label)

encourages small perturbation

model output: vector with probability of different classes

switching to ML notation!

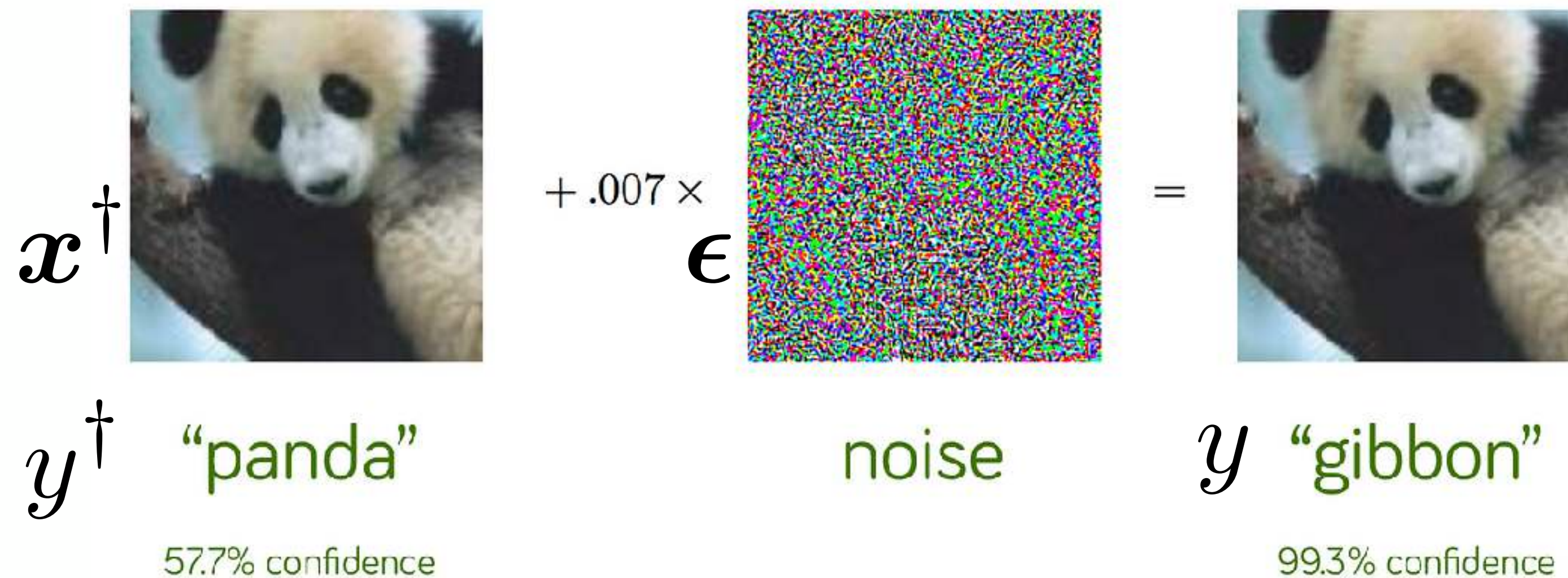


Image source: Explaining and Harnessing Adversarial Examples, Goodfellow et al, ICLR 2015

ML bifidelity example 1

Premise: suppose we have a cheap, inaccurate approximation f^{low}

- Context:
- black-box model
 - high-dimensional, low accuracy

Example: black-box adversarial attack Carlini & Wagner '17, black-box extension Chen et al. '17
For a given sample, find a small perturbation such that the machine learning algorithm misclassifies it

particular training example
(features and true label)

$$\min_{\epsilon} -f_{\text{cross-entropy}}(g(x^\dagger + \epsilon), y^\dagger) + \tau \|\epsilon\|^2$$

encourages small perturbation

model output: vector with probability of different classes

switching to ML notation!

MNIST is 28 x 28 images so $d = 784$



Train two models on MNIST data: (60k training, 10k test)

f is output of **large model**, trained conventionally

- convolution (32 filters) -> convolution (64 filters) ->
- max-pooling/flatten, fully connected (1024 neurons)
- > 10 class output. ReLU activation, 5x5 kernels

f^{low} is output of **small model**

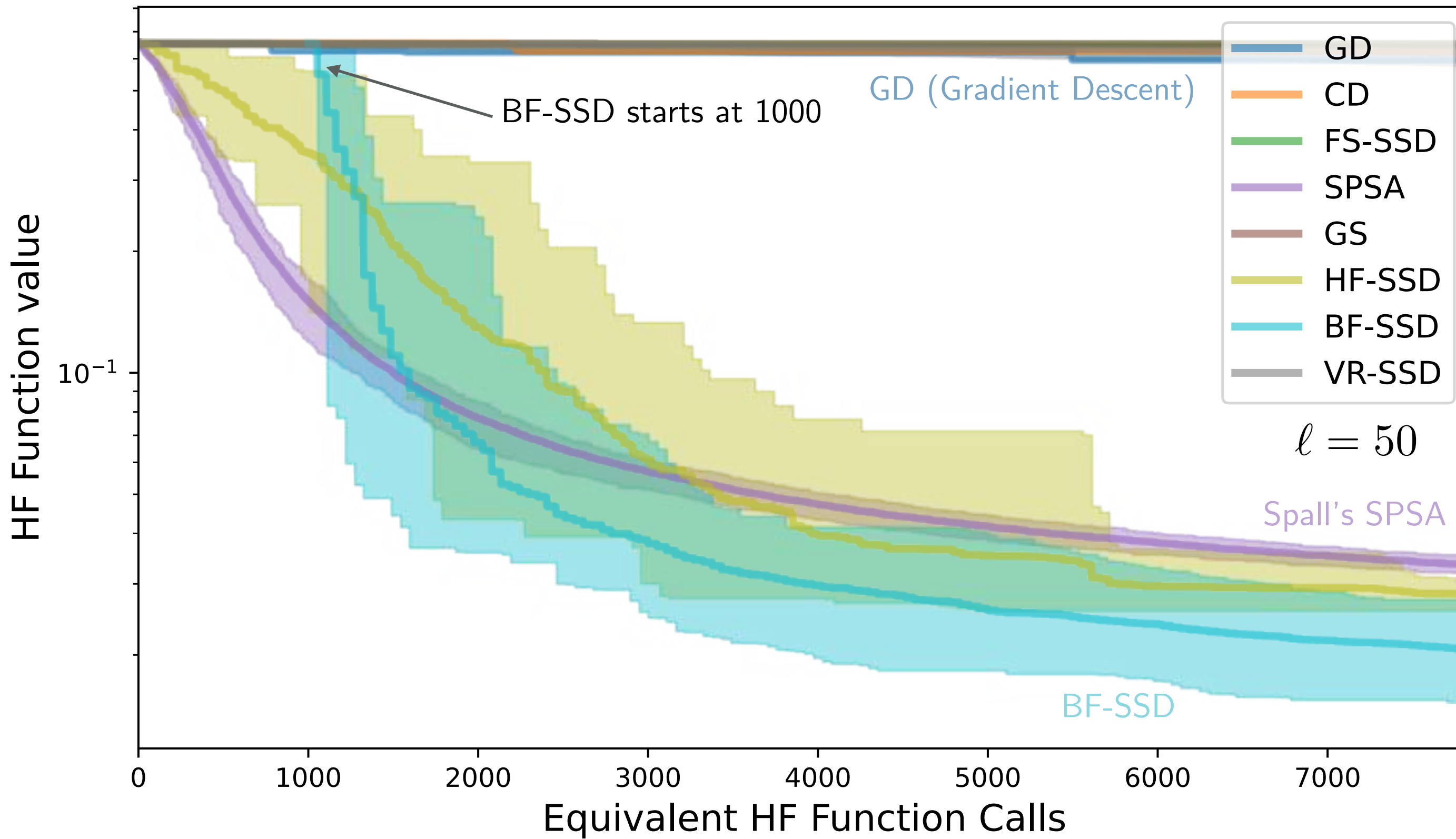
- trained not on MNIST but on output of large model
- (**knowledge distillation**), 1000 samples
- convolution (2 filters) -> max-pooling/flatten, fully connected (16 neurons)
- > 10 class output. ReLU activation, 2x3 kernels

	119x larger	
	Large model	Small model
# parameters	3,274,634	27,562
Test Accuracy	99.02%	82.21%

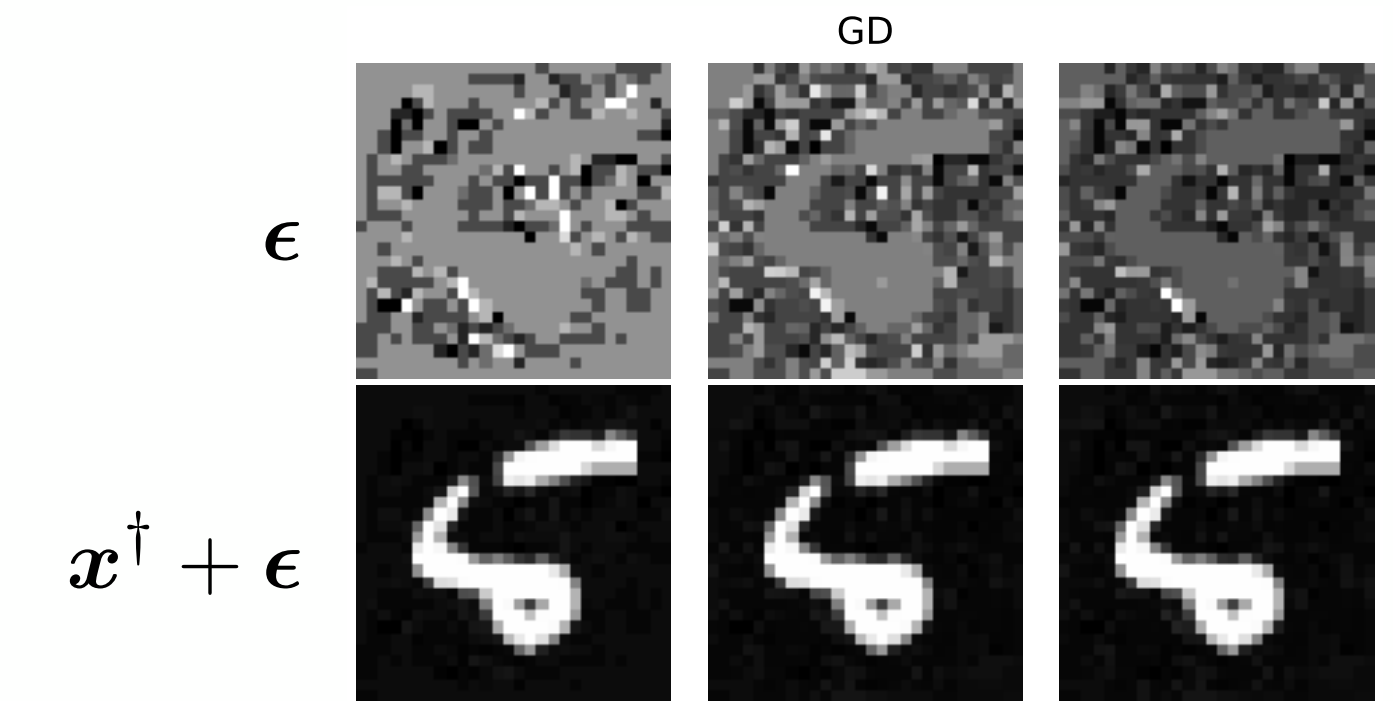
In some scenarios, small model
is not just cheap but “free”

ML bifidelity example 1: Results

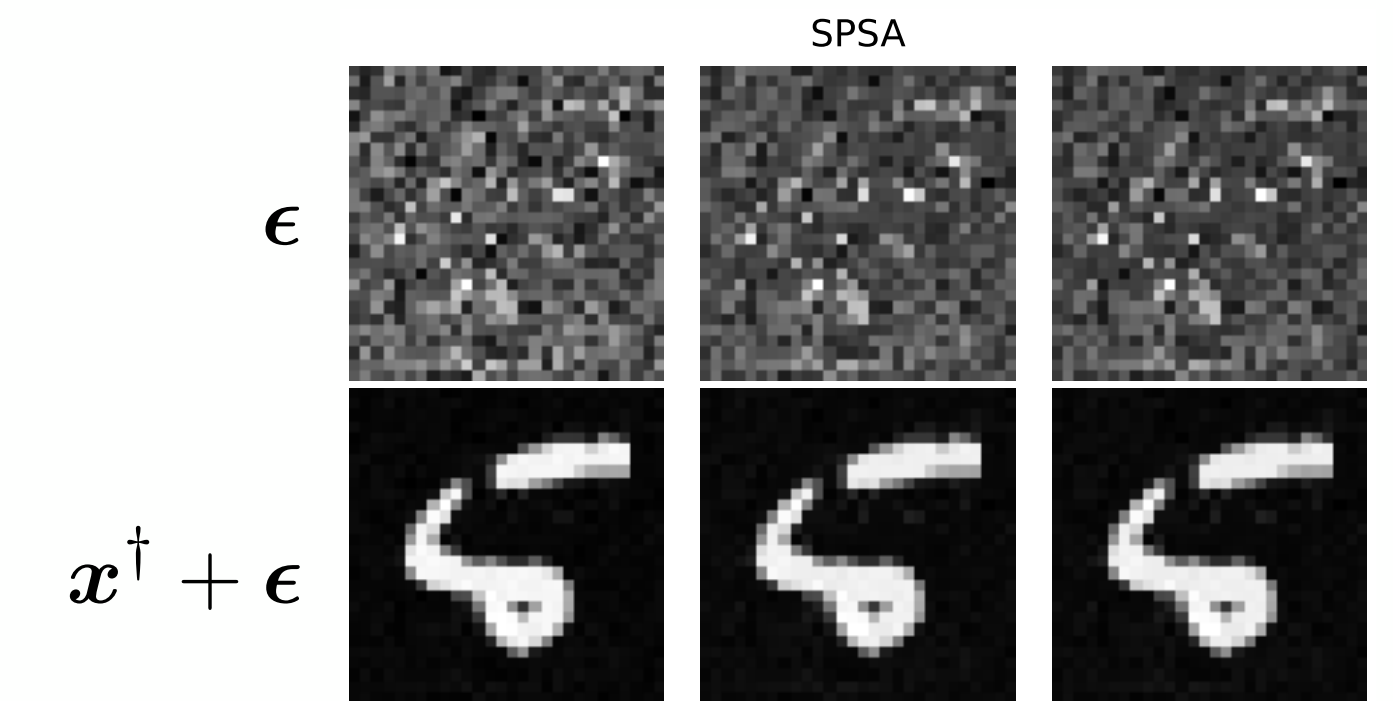
Test Case 1



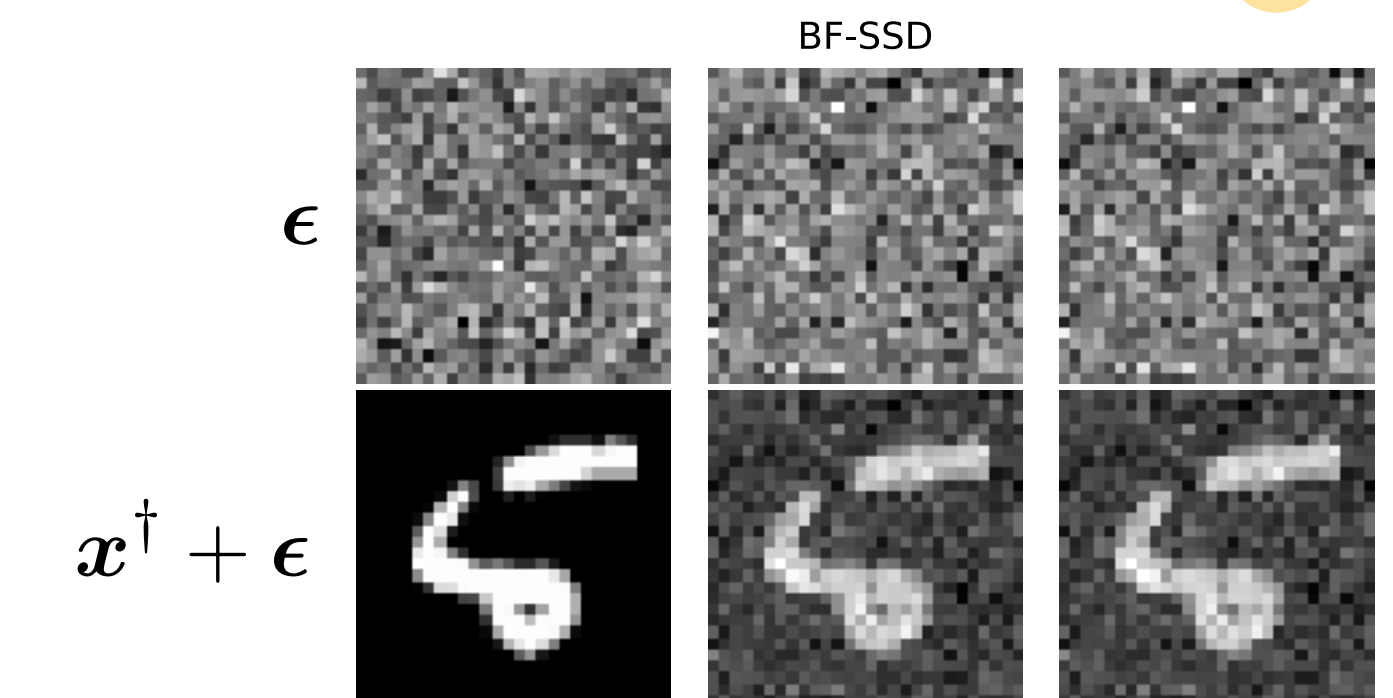
HF fun calls: 2000 5000 7000



(a) GD (Label=5, Predict=5)



(b) SPSA (Label=5, Predict=6)



(d) BF-SSD (Label=5, Predict=6)

ML bifidelity example 2

Premise: suppose we have a cheap, inaccurate approximation f^{low}

Context:

- black-box model
- high-dimensional, low accuracy

Example: soft prompting black-box LLM

We want to fine-tune a LLM like BERT or GPT

Instead of modifying network, lightweight alternative is to learn embeddings that are prepended to input sequence

ML bifidelity example 2

Context:

- black-box model
- high-dimensional, low accuracy

Premise: suppose we have a cheap, inaccurate approximation f^{low}

Example: soft prompting black-box LLM

We want to fine-tune a LLM like BERT or GPT

Instead of modifying network, lightweight alternative is to learn embeddings that are prepended to input sequence

Task: binary sentiment analysis (*classify a movie review as positive or negative*)

$$\text{Pretrained: } \left\{ \begin{array}{ll} f_{\text{token}} : \text{str} \rightarrow \mathbb{R}^{L_t \times d} & \text{tokenizer converts strings of any length to an embedding} \\ f_c : \mathbb{R}^{L_t \times d} \rightarrow [0, 1] & \text{classifier (we use small DistilBERT, small version of BERT)} \\ (z, y) \in \text{str} \times \{0, 1\} & \text{data from aclIMDB database} \end{array} \right. \quad d = 784$$

transformer

ML bifidelity example 2

Context:

- black-box model
- high-dimensional, low accuracy

Premise: suppose we have a cheap, inaccurate approximation f^{low}

Example: soft prompting black-box LLM

We want to fine-tune a LLM like BERT or GPT

Instead of modifying network, lightweight alternative is to learn embeddings that are **prepended to input sequence**

Task: binary sentiment analysis (*classify a movie review as positive or negative*)

Pretrained: $\begin{cases} f_{\text{token}} : \text{str} \rightarrow \mathbb{R}^{L_t \times d} & \text{tokenizer converts strings of any length to an embedding} \\ f_c : \mathbb{R}^{L_t \times d} \rightarrow [0, 1] & \text{classifier (we use small DistilBERT, small version of BERT)} \\ (z, y) \in \text{str} \times \{0, 1\} & \text{data from ac1IMDB database} \end{cases}$

$d = 784$

$$\mathbf{x}^* \in \underset{\mathbf{x} \in \mathbb{R}^d}{\operatorname{argmin}} \mathbb{E}_{(z, y)} [\text{CE}(f_c(\text{cat}[\mathbf{x}, f_{\text{token}}(z)]), y)]$$

cross-entropy loss

risk, replaced by empirical risk for training

$$f(\mathbf{x}) = \frac{1}{10} \sum_{i=1}^{10} \text{CE}(f_c(\text{cat}[\mathbf{x}, f_{\text{token}}(z_i)]), y_i)$$

ML bifidelity example 2

- Context:
- black-box model
 - high-dimensional, low accuracy

Premise: suppose we have a cheap, inaccurate approximation f^{low}

Example: soft prompting black-box LLM

We want to fine-tune a LLM like BERT or GPT

Instead of modifying network, lightweight alternative is to learn embeddings that are prepended to input sequence

Task: binary sentiment analysis (*classify a movie review as positive or negative*)

Pretrained: $\begin{cases} f_{\text{token}} : \text{str} \rightarrow \mathbb{R}^{L_t \times d} & \text{tokenizer converts strings of any length to an embedding} \\ f_c : \mathbb{R}^{L_t \times d} \rightarrow [0, 1] & \text{classifier (we use small DistilBERT, small version of BERT)} \\ (z, y) \in \text{str} \times \{0, 1\} & \text{data from a1IMDB database} \end{cases}$

transformer

$d = 784$

f uses a sample size of 10 **High-Fidelity**
 f^{low} uses a sample size of 2 **Low-Fidelity**

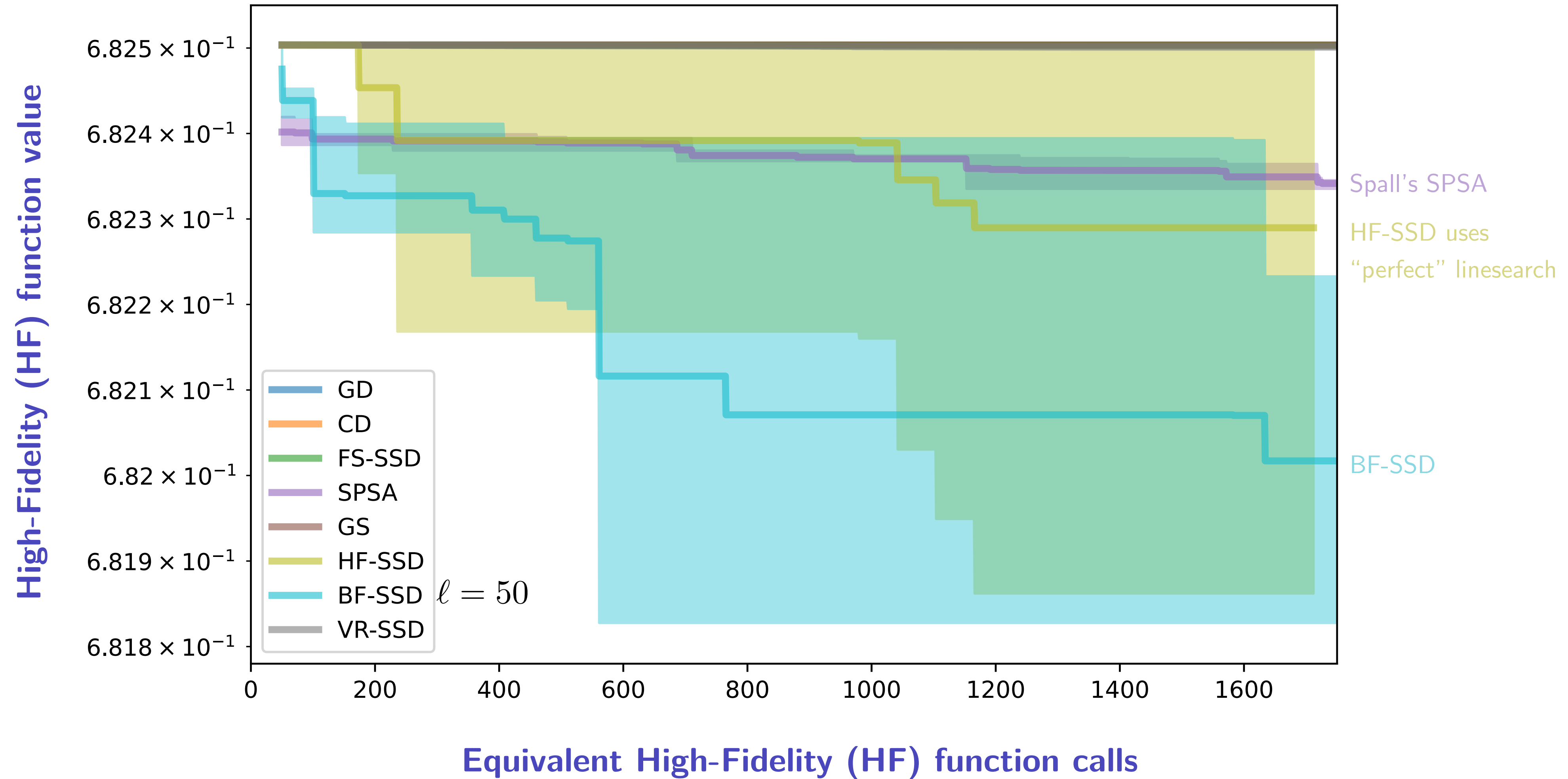
cross-entropy loss

$$\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} \mathbb{E}_{(z, y)} [\text{CE}(f_c(\text{cat}[\mathbf{x}, f_{\text{token}}(z)]), y)]$$

risk, replaced by empirical risk for training

$$f(\mathbf{x}) = \frac{1}{10} \sum_{i=1}^{10} \text{CE}(f_c(\text{cat}[\mathbf{x}, f_{\text{token}}(z_i)]), y_i)$$

ML bifidelity example 2: Results



Bonus ideas

3. Just do a line search!

e.g.: C. Paquette, K. Scheinberg. *A Stochastic Line Search Method with Expected Complexity Analysis*. SIOPT, 2020.

Probably works great, but we haven't analyzed in our context

Bonus ideas

3. Just do a line search!

e.g.: C. Paquette, K. Scheinberg. *A Stochastic Line Search Method with Expected Complexity Analysis*. SIOPT, 2020.

Probably works great, but we haven't analyzed in our context

4. Barzilai-Borwein stepsize

has been used before in SGD methods, e.g., C. Tan, S. Ma, Y.-H. Dai, Y. Qian. *Barzilai-Borwein Step Size for Stochastic Gradient Descent*. NIPS, 2016.

Popular scalar approximation to attempt to solve secant equation $\mathbf{B}\mathbf{s} = \mathbf{y}$ where

$$\mathbf{s} = \mathbf{x}_k - \mathbf{x}_{k-1}$$
$$\mathbf{y} = \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1})$$

$$\eta_{\text{BB1}} = \frac{\|\mathbf{s}\|^2}{\mathbf{s}^\top \mathbf{y}}$$

$$\eta_{\text{BB2}} = \frac{\mathbf{s}^\top \mathbf{y}}{\|\mathbf{y}\|^2}$$

In a DFO setting, the BB1 stepsize is reasonable: we have \mathbf{s} explicitly, and $\mathbf{s}^\top \mathbf{y}$ can be computed with two directional derivatives

Downside: even in deterministic full-gradient settings, convergence isn't guaranteed

What's next? Biased sampling

Suppose we have a crude estimate, \mathbf{g} , of the gradient
e.g., from past iterations, or low-fidelity model, etc.

try quasi-Newton ideas?

How to exploit this?

we already saw a few ways — here is one more:

$$\mathbf{v} = \frac{1}{\ell} \mathbf{\Sigma}^{-1} \mathbf{Q} \mathbf{Q}^{\top} \nabla f(\mathbf{x}_k) \quad \text{and} \quad \mathbf{x}_{k+1} = \mathbf{x}_k - \eta \mathbf{v}$$

where each column \mathbf{q} of \mathbf{Q} is iid $\mathbf{q} \sim \mathcal{N}(0, \mathbf{\Sigma})$

$$\text{and } \mathbf{\Sigma} = \sigma^2 \mathbf{I} + \mathbf{g} \mathbf{g}^{\top} / \|\mathbf{g}\|^2$$

adjust based on our *confidence* in gradient estimate

Then still unbiased, $\mathbb{E}[\mathbf{v}] = \nabla f(\mathbf{x}_k)$

but hopefully much lower variance.

inspired by Boullé and Townsend '22

What's next? Biased sampling

Suppose we have a crude estimate, \mathbf{g} , of the gradient
e.g., from past iterations, or low-fidelity model, etc.

try quasi-Newton ideas?

How to exploit this?

we already saw a few ways — here is one more:

$$\mathbf{v} = \frac{1}{\ell} \Sigma^{-1} \mathbf{Q} \mathbf{Q}^T \nabla f(\mathbf{x}_k) \quad \text{and} \quad \mathbf{x}_{k+1} = \mathbf{x}_k - \eta \mathbf{v}$$

where each column \mathbf{q} of \mathbf{Q} is iid $\mathbf{q} \sim \mathcal{N}(0, \Sigma)$

$$\text{and } \Sigma = \sigma^2 \mathbf{I} + \mathbf{g} \mathbf{g}^T / \|\mathbf{g}\|^2$$

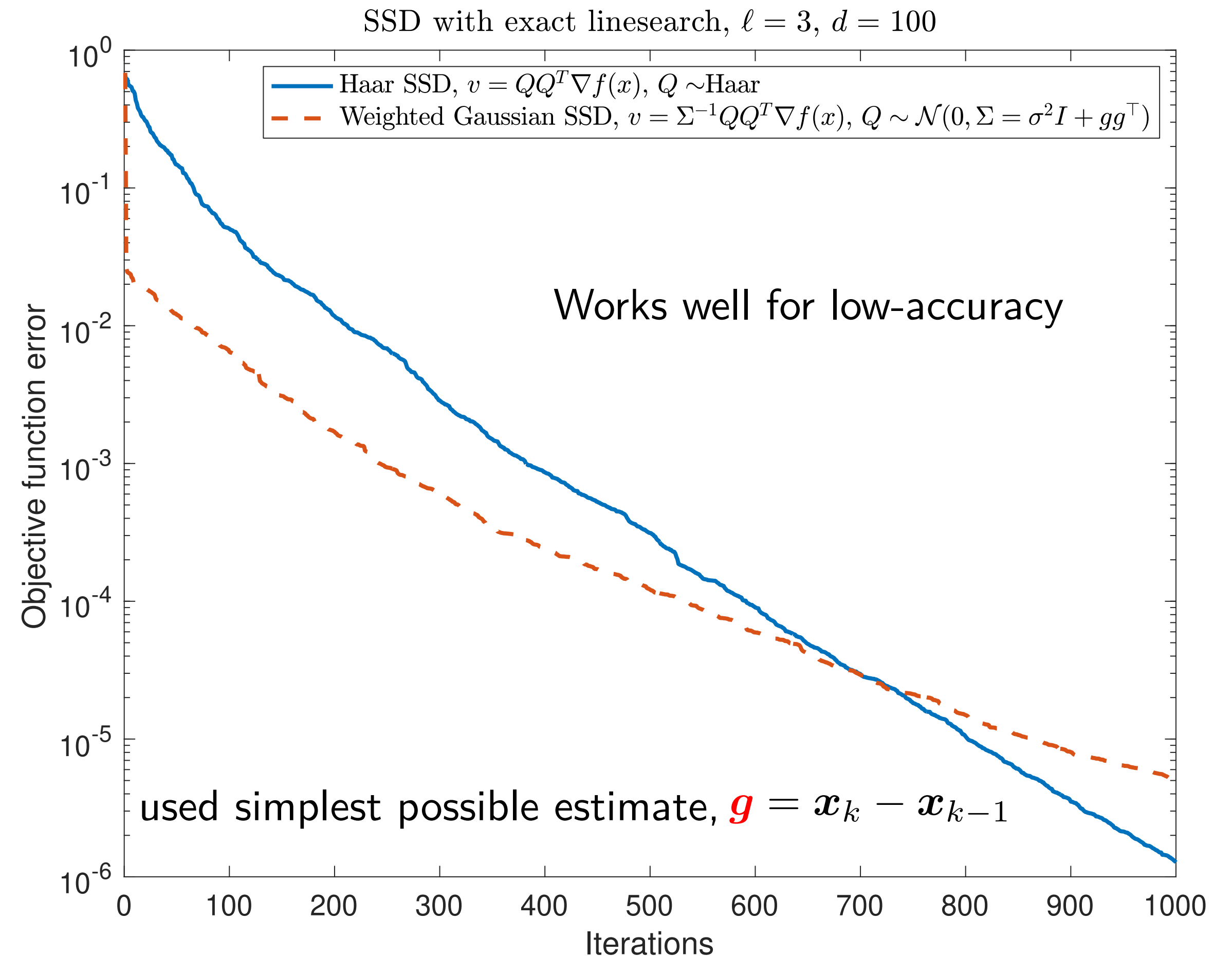
adjust based on our *confidence* in gradient estimate

Then still unbiased, $\mathbb{E}[\mathbf{v}] = \nabla f(\mathbf{x}_k)$

but hopefully much lower variance.

inspired by Boullé and Townsend '22

Preliminary example on a quadratic function



works poorly when ℓ is large,
since then Gaussian is inferior to Haar

What's next? Understand dependence on Hessian eigenvalues

Restrict to convex quadratic functions for simplicity, $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top H\mathbf{x} + \mathbf{c}^\top \mathbf{x}$

The **min** and **max** eigenvalue of H control worst-case behavior (if using fixed stepsize)

... but the **interior** eigenvalues play a role with *average-case analysis* and/or *exact line search*

What's next? Understand dependence on Hessian eigenvalues

Restrict to convex quadratic functions for simplicity, $f(x) = \frac{1}{2}x^\top Hx + c^\top x$

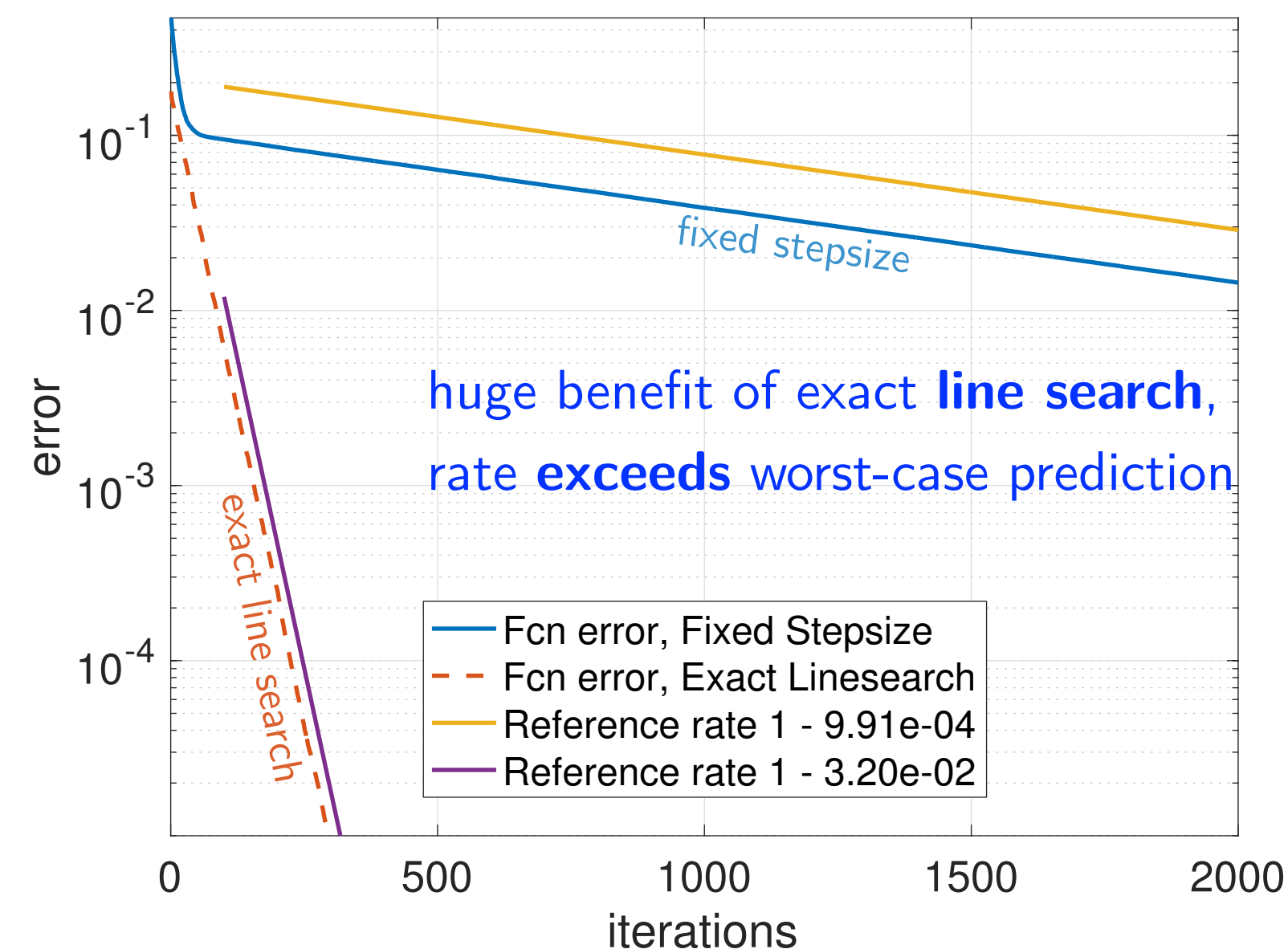
The **min** and **max** eigenvalue of H control worst-case behavior (if using fixed stepsize)

$d = 100, \ell = 5$

... but the **interior** eigenvalues play a role with *average-case analysis* and/or *exact line search*

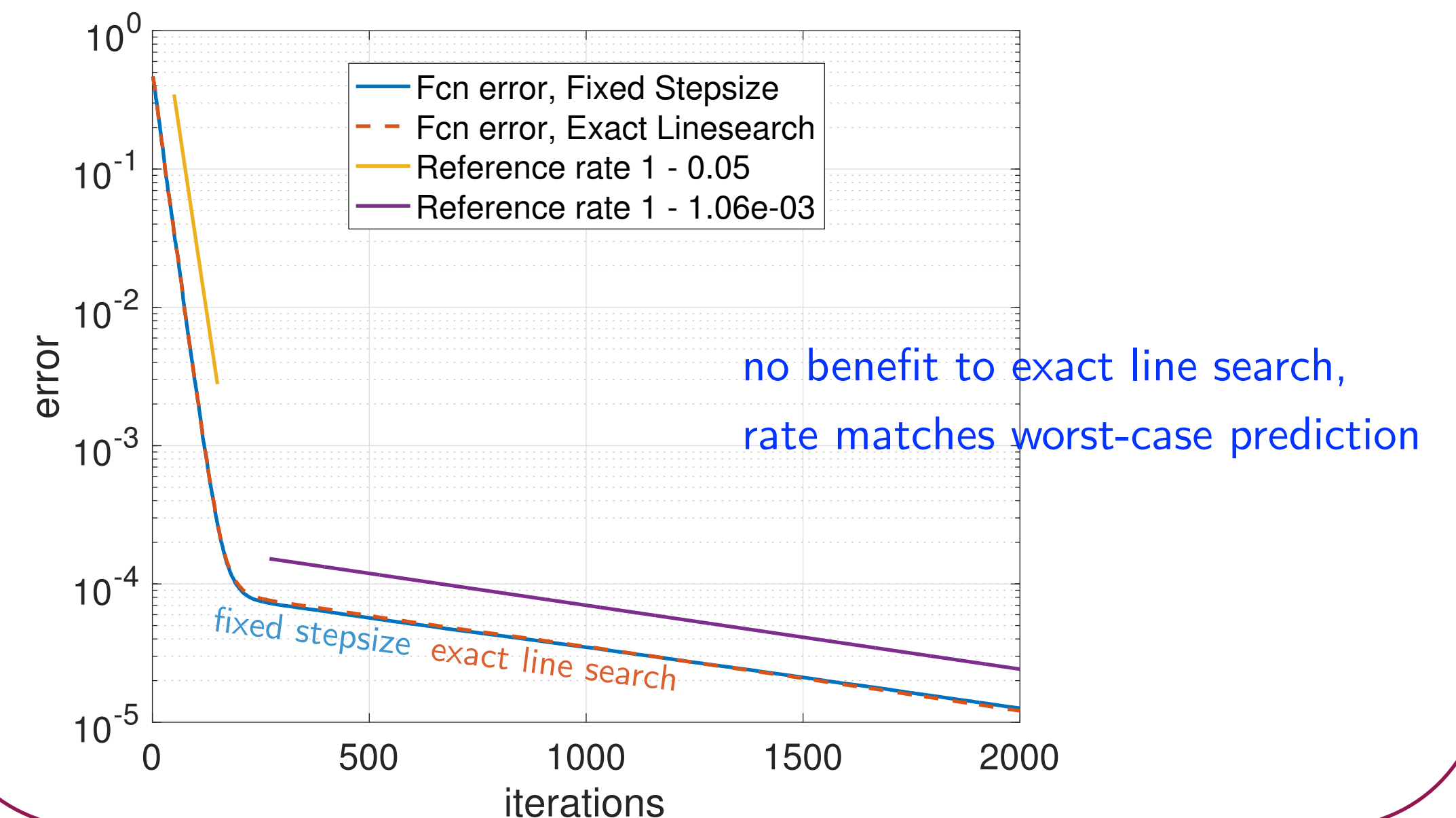
Case A: 1 big eigenvalue, the rest small

(Small “numerical rank”) $\kappa = \lambda_{\max}/\lambda_{\min} = 100$



Case B: 1 small eigenvalue, the rest large

(Full “numerical rank”) $\kappa = \lambda_{\max}/\lambda_{\min} = 100$



Conclusion

- 0th order methods have their roles
- stepsize selection (and/or line search) is important
- multi fidelity is useful

Thanks for your attention

Slides will be posted at <https://stephenbeckr.github.io/papers/#talks>

