

# Approximate Proximal-Gradient Methods

Anis Hamadouche, Yun Wu, Andrew M. Wallace,  
and Joaõ F. C. Mota

The Institute of Sensors, Signals & Systems (ISSS)

*Heriot-Watt University, Edinburgh, UK*



# Composite Optimization Problems

Many problems in science, engineering, and *defense* can be written as

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad g(x) + h(x)$$

$\left| \begin{array}{l} \text{non-differentiable, can } \textit{encode constraints}: h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\} \\ \text{differentiable, real-valued: } g : \mathbb{R}^n \rightarrow \mathbb{R} \end{array} \right.$



**Example:** State-space model of a drone

$$\left| x[k] = (\text{position time } k, \text{ velocity time } k, \% \text{ mission completed}, \dots) \right.$$

$$x[k+1] \simeq Ax[k] + Bu[k]$$

$\left| \text{input at time } k \right.$

**Goal:** given  $x[0]$ , drive state to  $x_f$  in  $T$  time steps, while minimizing energy

$$\underset{\substack{x[1], \dots, x[T] \\ u[0], \dots, u[T-1]}}{\text{minimize}} \quad (x[T] - x_f)^2 + \sum_{k=0}^{T-1} u^2[k] = g(x)$$

subject to  $x[k+1] = Ax[k] + Bu[k], \quad k = 0, \dots, T-1$  ——— can be encoded in  $h(x)$

# Approximate Proximal Methods

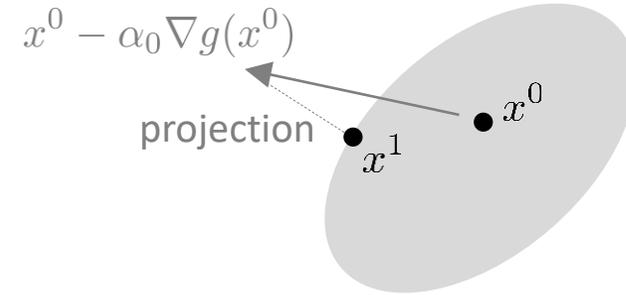
$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

$$x^{k+1} = \text{prox}_{\alpha_k h} \left( \underbrace{x^k - \alpha_k \nabla g(x^k)}_{\text{gradient descent along } g \text{ (differentiable)}} \right)$$

*gradient descent* along  $g$  (differentiable)

$$\text{prox}_f(x) := \arg \min_y f(x) + \frac{1}{2} \|y - x\|_2^2$$

generalization of a *projection*



*Errors* may be introduced to *save power*

*hardware, software, linear algebra, or algorithmic approximations*

*how to model ?*



$$x^{k+1} = \text{prox}_{\alpha_k h} \left( x^k - \alpha_k \nabla g(x^k) + \epsilon_1^k \right) + \epsilon_2^k$$

*Existing convergence proofs hold for this type of errors ?*

*Tradeoffs between power savings / accuracy / execution time ?*

# Related Work

*Approximate PGD with decreasing errors*

*Schmidt et al. 2010*

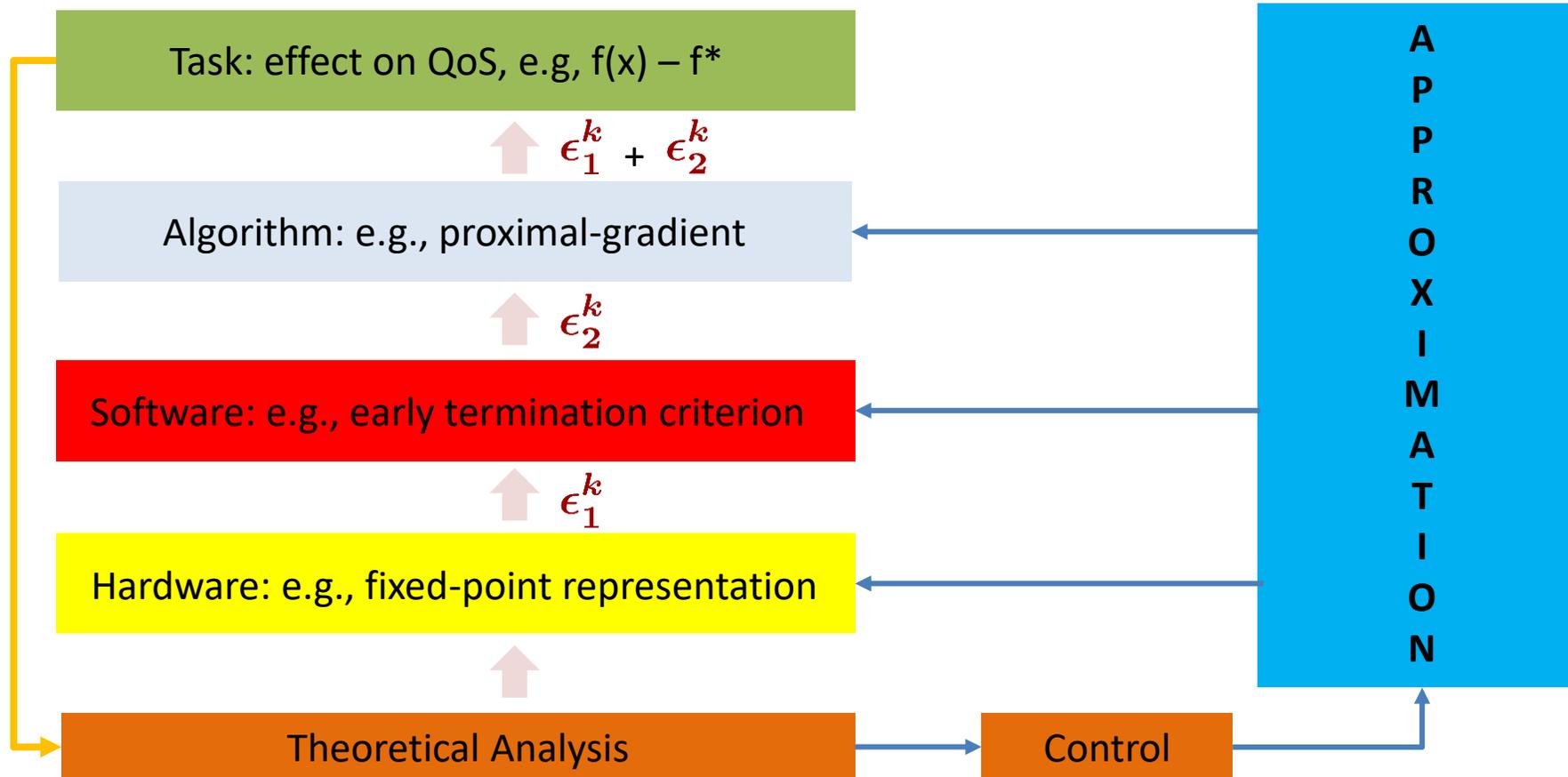
Criterion	Bound	Constants
$f\left(\frac{1}{k}\sum_{i=1}^k x_i\right) - f(x^*)$	$\frac{L}{2k} \left(\ x_0 - x^*\  + 2A_k + \sqrt{2B_k}\right)^2$	$A_k = \sum_{i=1}^k \left(\frac{\ e_i\ }{L} + \sqrt{\frac{2\varepsilon_i}{L}}\right), \quad B_k = \sum_{i=1}^k \frac{\varepsilon_i}{L}$

*Approximate Accelerated PGD with square summable (weighted) errors*

*Schmidt et al. 2010, Aujol et al. 2015*

Criterion	Bound	Constants
$f(x_k) - f(x^*)$	$\frac{2L}{(k+1)^2} \left(\ x_0 - x^*\  + 2\tilde{A}_k + \sqrt{2\tilde{B}_k}\right)^2$	$\tilde{A}_k = \sum_{i=1}^k i \left(\frac{\ e_i\ }{L} + \sqrt{\frac{2\varepsilon_i}{L}}\right), \quad \tilde{B}_k = \sum_{i=1}^k \frac{i^2 \varepsilon_i}{L}$
$t_N^2 w_N + \sum_{n=2}^N \rho_n w_{n-1} + \frac{1}{2\gamma} \ u_N - x^*\ ^2$	$\frac{1}{2\gamma} \left(\ u_0 - x^*\  + 2A_{i,N} + \sqrt{2B_N}\right)^2$	$A_{1,n} = \sum_{k=1}^n t_k \left(\gamma \ e_k\  + \sqrt{2\gamma \varepsilon_k}\right)$
$u_n = x_{n-1} + t_n(x_n - x_{n-1})$		$A_{2,n} = \gamma \sum_{k=1}^n t_k \ e_k\  \quad B_n = \gamma \sum_{k=1}^n t_k^2 \varepsilon_k$
$w_n := F(x_n) - F(x^*)$		

# Analysis of Error Propagation



# Error Models

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

Error Type	Probabilistic Model	
<b>Gradient computation (linear)</b> $\nabla^{\epsilon_1^k} g(x^k) := \nabla g(x^k) + \epsilon_1^k$	<ol style="list-style-type: none"> <li>Centered &amp; CMI</li> <li>Bounded</li> <li>CMI of the iterates</li> </ol>	$\mathbb{E}[\epsilon_{1\Omega}^k \mid \epsilon_{1\Omega}^1, \dots, \epsilon_{1\Omega}^{k-1}] = \mathbb{E}[\epsilon_{1\Omega}^k] = 0,$ $\mathbb{P}( \epsilon_{1\Omega j}^k  \leq \delta) = 1, \quad \text{for all } j = 1, \dots, n,$ $\mathbb{E}[\epsilon_{1\Omega}^k \top x_{\Omega}^k \mid \epsilon_{1\Omega}^1, \dots, \epsilon_{1\Omega}^{k-1}, x_{1\Omega}^1, \dots, x_{1\Omega}^{k-1}] = \mathbb{E}[\epsilon_{1\Omega}^k \top x_{\Omega}^k] = 0,$
<b>Proximal computation (nonlinear)</b> $r^i = x^i - \bar{x}^i$ $h(x^{k+1}) + \frac{1}{2s} \ x^{k+1} - x^k + s \nabla^{\epsilon_1^k} g(x^k)\ _2^2 \leq$ $\epsilon_2^k + h(\bar{x}^{k+1}) + \frac{1}{2s} \ \bar{x}^{k+1} - x^k + s \nabla^{\epsilon_1^k} g(x^k)\ _2^2$	<ol style="list-style-type: none"> <li>Centered and CMI</li> <li>CMI of the iterates</li> <li>Bounded</li> </ol>	$\mathbb{E}[r_{\Omega}^k \mid r_{\Omega}^1, \dots, r_{\Omega}^{k-1}] = \mathbb{E}[r_{\Omega}^k] = 0,$ $\mathbb{E}[r_{\Omega}^k \top x_{\Omega}^k \mid r_{\Omega}^1, \dots, r_{\Omega}^{k-1}, x_{1\Omega}^1, \dots, x_{1\Omega}^{k-1}] = \mathbb{E}[r_{\Omega}^k \top x_{\Omega}^k] = 0,$ $\mathbb{P}( \epsilon_{2\Omega}^k  \leq \epsilon_0) = 1$

Independence  $\longrightarrow$  Conditional Mean Independence (CMI)  $\longrightarrow$  Uncorrelatedness

**Probabilistic analysis** is a hybrid of worst-case and average-case analyses that inherits advantages of both. It measures the expected performance of algorithms under slight random perturbations of worst-case inputs

# Convergence Results

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

Scheme	Analysis	Bound
Proximal-Gradient Method	Deterministic	$f\left(\frac{1}{k+1} \sum_{i=0}^k x^{i+1}\right) - f(x^*) \leq \underbrace{\frac{1}{k+1} \left[ \sum_{i=0}^k \epsilon_2^i + \sum_{i=1}^k \left( \ \epsilon_1^i\ _2 + \sqrt{\frac{2\epsilon_2^i}{s}} \right) \ x^* - x^0\ _2 \right]}_{\text{Error term}} + \underbrace{\frac{1}{2s} \ x^* - x^0\ _2^2}_{\text{Error-free}}$
	Probabilistic $2 - 4 \exp(-\frac{\gamma^2}{2})$	$f\left(\frac{1}{k} \sum_{i=1}^k x_{\Omega}^i\right) - f(x^*) \leq \underbrace{\frac{1}{k} \sum_{i=1}^k \epsilon_{2\Omega}^i + \frac{\gamma}{\sqrt{k}} \left( \sqrt{n} \delta  + \sqrt{\frac{2\epsilon_0}{s}} \right) \ x^* - x^0\ _2}_{\text{Error term}} + \underbrace{\frac{1}{2sk} \ x^* - x^0\ _2^2}_{\text{Error-free}}$

# Convergence Results

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

Assumption on error for convergence	Scheme	Analysis	Error Bounds	Rate
$\epsilon_1 \propto O(1/k^{1+\lambda})$ $\epsilon_2 \propto O(1/k^{2+\lambda})$	Proximal-Gradient Method	Deterministic	$\frac{1}{k+1} \left[ \sum_{i=0}^k \epsilon_2^i + \sum_{i=1}^k \left( \ \epsilon_1^i\ _2 + \sqrt{\frac{2\epsilon_2^i}{s}} \right) \ x^* - x^0\ _2 \right]$	$O(1/k)$
$\epsilon_1 \propto O(1/k^{0.5+\lambda})$ $\epsilon_2 \propto O(1/k^{1+\lambda})$		Probabilistic $2 - 4 \exp(-\frac{\gamma^2}{2})$	$\frac{1}{k} \sum_{i=1}^k \epsilon_{2\Omega}^i + \frac{\gamma}{\sqrt{k}} \left( \sqrt{n} \delta  + \sqrt{\frac{2\epsilon_0}{s}} \right) \ x^* - x^0\ _2$	$O(1/k)$

Schmidt et al. 2010

Assumption on error	Scheme	Analysis	Error Bounds	Rate
$\epsilon_1 \propto O(1/k^{1+\lambda})$ $\epsilon_2 \propto O(1/k^{2+\lambda})$	Proximal-Gradient Method	Deterministic Probabilistic	$\frac{1}{2sk} \left[ \left\  x^* - x^0 \right\ _2 + 2A_k + \sqrt{2B_k} \right]^2$ $A_k = \sum_{i=1}^k \left( \frac{\ \epsilon_1^i\ _2}{L} + \sqrt{\frac{2\epsilon_2^i}{L}} \right), \quad B_k = \sum_{i=1}^k \frac{\epsilon_2^i}{L}$	$O(1/k)$

# Convergence Results

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

Assumption on error for convergence	Scheme	Analysis	Error Bounds	Rate
$\epsilon_1 \propto O(1/k)$ $\epsilon_2 \propto O(1/k^2)$	Proximal-Gradient Method	Deterministic	$\frac{1}{k+1} \left[ \sum_{i=0}^k \epsilon_2^i + \sum_{i=1}^k \left( \ \epsilon_1^i\ _2 + \sqrt{\frac{2\epsilon_2^i}{s}} \right) \ x^* - x^0\ _2 \right]$	$O(\log k/k)$
$\epsilon_1 \propto O(1/k^{0.5})$ $\epsilon_2 \propto O(1/k)$		Probabilistic $2 - 4 \exp(-\frac{\gamma^2}{2})$	$\frac{1}{k} \sum_{i=1}^k \epsilon_{2\Omega}^i + \frac{\gamma}{\sqrt{k}} \left( \sqrt{n} \delta  + \sqrt{\frac{2\epsilon_0}{s}} \right) \ x^* - x^0\ _2$	$O(\log k/k)$

Schmidt et al. 2010

Assumption on error	Scheme	Analysis	Error Bounds	Rate
$\epsilon_1 \propto O(1/k)$ $\epsilon_2 \propto O(1/k^2)$	Proximal-Gradient Method	Deterministic Probabilistic	$\frac{1}{2sk} \left[ \left\  x^* - x^0 \right\ _2 + 2A_k + \sqrt{2B_k} \right]^2$ $A_k = \sum_{i=1}^k \left( \frac{\ \epsilon_1^i\ _2}{L} + \sqrt{\frac{2\epsilon_2^i}{L}} \right), \quad B_k = \sum_{i=1}^k \frac{\epsilon_2^i}{L}$	$O(\log^2 k/k)$

# Convergence Results

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

Assumption on error for convergence	Scheme	Analysis	Error Bounds	Rate
$\epsilon_1 \leq \delta$ $\epsilon_2 \leq \epsilon_0$	Proximal-Gradient Method	Deterministic	$\frac{1}{k+1} \left[ \sum_{i=0}^k \epsilon_2^i + \sum_{i=1}^k \left( \ \epsilon_1^i\ _2 + \sqrt{\frac{2\epsilon_2^i}{s}} \right) \ x^* - x^0\ _2 \right]$	$O(1)$
<b>Stationary</b> $\epsilon_1 \leq \delta$ $\epsilon_2 \leq \epsilon_0$		Probabilistic $2 - 4 \exp(-\frac{\gamma^2}{2})$	$E(\epsilon_{2\Omega}) + \frac{\gamma}{\sqrt{k}} \left( \frac{\epsilon_0}{2} + \sqrt{n} \delta  \ x^* - x^0\ _2 \right)$ <p><i>E(ϵ<sub>2Ω</sub>)-optimality if stationary</i></p>	$O(1/\sqrt{k})$

Schmidt et al. 2010

Assumption on error	Scheme	Analysis	Error Bounds	Rate
$\epsilon_1 \leq \delta$ $\epsilon_2 \leq \epsilon_0$	Proximal-Gradient Method	Deterministic	$\frac{1}{2sk} \left[ \left\  x^* - x^0 \right\ _2 + 2A_k + \sqrt{2B_k} \right]^2$ $A_k = \sum_{i=1}^k \left( \frac{\ \epsilon_1^i\ _2}{L} + \sqrt{\frac{2\epsilon_2^i}{L}} \right), \quad B_k = \sum_{i=1}^k \frac{\epsilon_2^i}{L}$	$O(k)$

# Experimental Setup:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

$$x^{k+1} = \text{prox}_{\alpha_k h} \left( x^k - \alpha_k \nabla g(x^k) \right)$$

- LASSO with 600 random examples and 100 features using fixed-point representation (round-off error) and finite solver precision (CVX).

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|Ax - y\|_2^2 + \lambda \|x\|_1,$$

- Quantization according to Q-format:

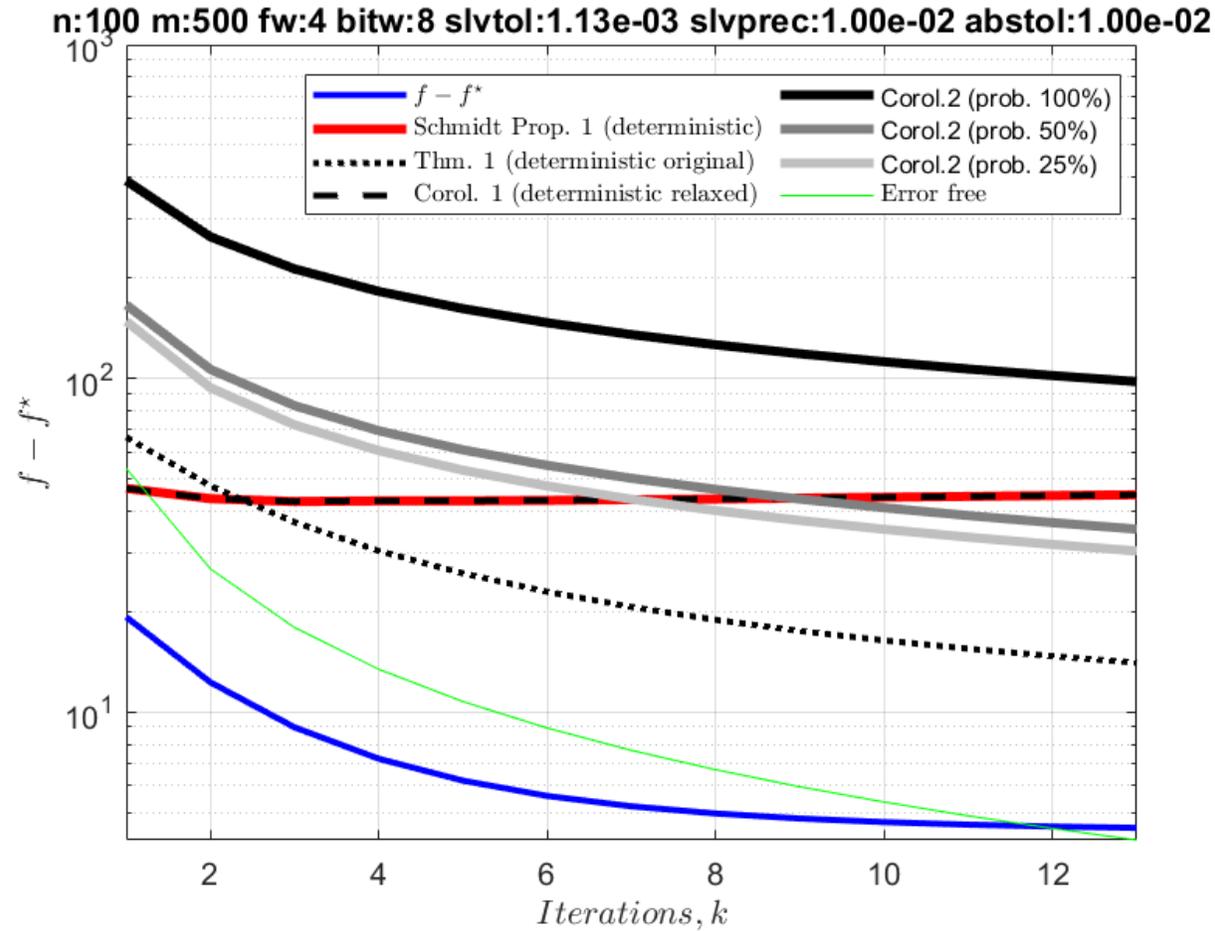
$$sI.F(x) = \sum_{i=0}^{W-2} b_i 2^{i-F} - b_{W-1} 2^{I-1}.$$

precision	bit	frac.	ABSTOL
2.22e-16	8	4	2.22e-16
	16	6	0.001
		8	2.22e-16
0.001	8	4	0.001
	16	8	0.01
		8	2.22e-16
0.01	8	4	0.001
	16	6	0.01
		8	0.01

# Experimental Results:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

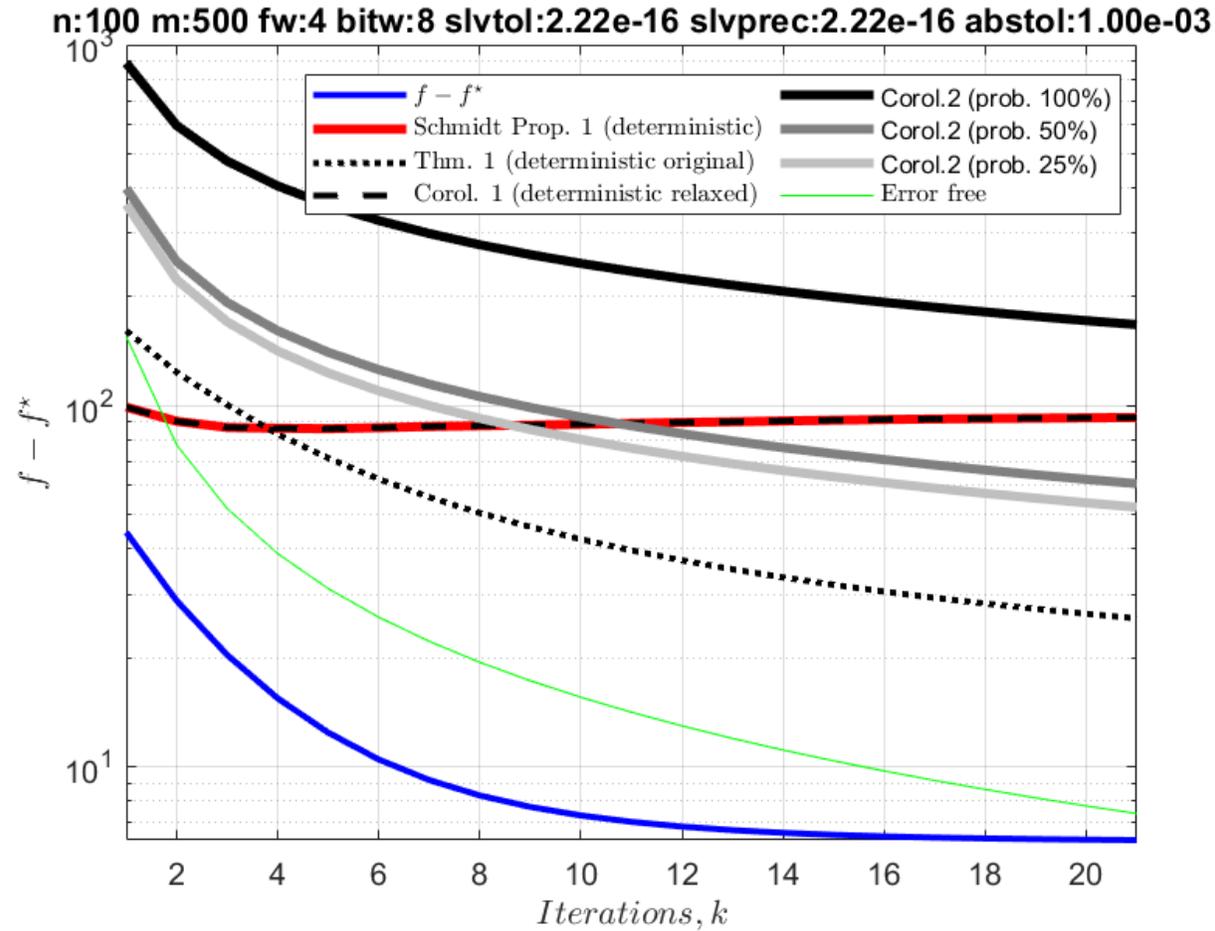
$$x^{k+1} = \text{prox}_{\alpha_k h} \left( x^k - \alpha_k \nabla g(x^k) \right)$$



# Experimental Results:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

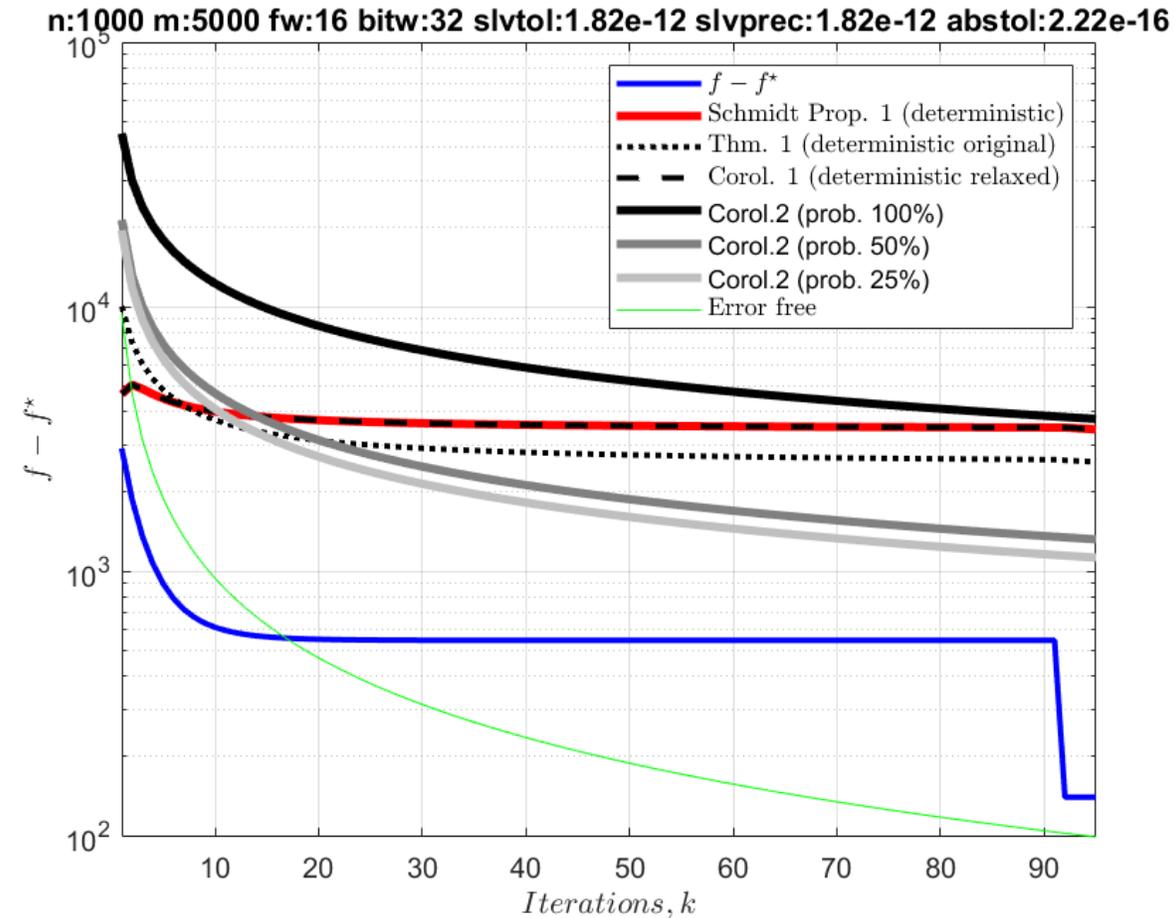
$$x^{k+1} = \text{prox}_{\alpha_k h} \left( x^k - \alpha_k \nabla g(x^k) \right)$$



# Experimental Results:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) := g(x) + h(x)$$

$$x^{k+1} = \text{prox}_{\alpha_k h} \left( x^k - \alpha_k \nabla g(x^k) \right)$$



# Conclusion:

- We obtained new tighter deterministic bounds and we demonstrated their validity on a practical optimization example (LASSO) solved on a reduced-precision machine combined with reduced-precision solver.

# Conclusion:

- We obtained new tighter deterministic bounds and we demonstrated their validity on a practical optimization example (LASSO) solved on a reduced-precision machine combined with reduced-precision solver.
- We also derived probabilistic upper bounds.

# Conclusion:

- We obtained new tighter deterministic bounds and we demonstrated their validity on a practical optimization example (LASSO) solved on a reduced-precision machine combined with reduced-precision solver.
- We also derived probabilistic upper bounds.
- Worst-case running time can be much worse than the observed running time in practice.

# Conclusion:

- We obtained new tighter deterministic bounds and we demonstrated their validity on a practical optimization example (LASSO) solved on a reduced-precision machine combined with reduced-precision solver.
- We also derived probabilistic upper bounds.
- Worst-case running time can be much worse than the observed running time in practice.
- Probabilistic bounds are more practical.

# Conclusion:

- We obtained new tighter deterministic bounds and we demonstrated their validity on a practical optimization example (LASSO) solved on a reduced-precision machine combined with reduced-precision solver.
- We also derived probabilistic upper bounds.
- Worst-case running time can be much worse than the observed running time in practice.
- Probabilistic bounds are more practical.
- More relaxations on the assumptions are needed in order to incorporate more general perturbations into the analysis.

