# Stochastic Optimisation
# for Large-Scale Inverse Problems

Matthias J. Ehrhardt

Department of Mathematical Sciences, University of Bath, UK

6 September, 2024

# Main Aim and Outline

$$x^\sharp \in \arg \min_x \left\{ \sum_{i=1}^n f_i(A_i x) + g(x) \right\}$$

► proper, convex and lower semi-continuous
► *n* large and/or $A_i x$ expensive

**Outline:**
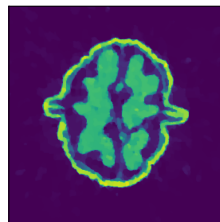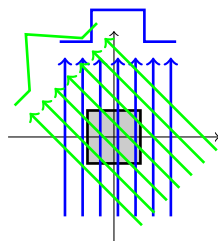
1) **Why?** Inverse Problems and Optimization
2) **How?** Randomized Algorithms for Convex Optimization
3) **So what?** Applications: PET, CT, ...

# CT Reconstruction with TV

**Total variation (TV)**

Rudin, Osher, Fatemi '92

$$\mathcal{R}(x) = \|Dx\|_1$$



$$\min_x \left\{ \sum_{j=1}^{s} \|K_j x - b_j\|^2 + \lambda\|Dx\|_1 + \imath_+(x) \right\}$$

$$\min_x \left\{ \sum_{i=1}^{n} f_i(A_i x) + g(x) \right\}$$
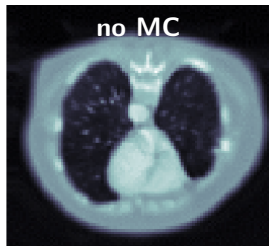
$n = s$
$f_i(y) = \|y - b_i\|^2 \quad i \in [n]$
$A_i = K_i \quad i \in [n]$
$g(x) = \lambda\|Dx\|_1 + \imath_+(x)$

# Motion corrected CT reconstruction

$$\min_x \left\{ \sum_{i=1}^{s} \| K M_i x - b_i \|^2 + \mathcal{R}(x) \right\}$$
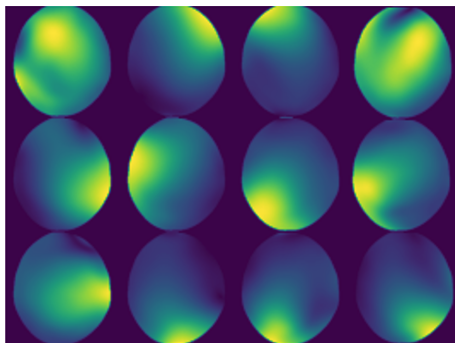
- ▶ $M_i$ motion transformation
- ▶ here $s = 10$ motion gates; computations are a bottleneck
- ▶ No motion correction: $M_i = I$



e.g. Delplancke, Thielemans, Ehrhardt '21

# Parallel MRI

$$\min_x \left\{ \sum_{i=1}^{s} \| SF C_i x - b_i \|^2 + \mathcal{R}(x) \right\}$$

▶ $C_i$ sensitivity map for $i$th MR coil, $s = 12$



Pruessmann et al. '99

# Stochastic Optimisation Algorithms

# Building blocks for Convex Optimisation

Template:

$$\min_x \{f(Ax) + g(x) = F(x) + g(x)\}$$

- ▶ Ingredient 1 (gradient descent)

$$x^+ = x - \tau \nabla F(x)$$

- ▶ Ingredient 2 (proximal point algorithm)

$$x^+ = \text{prox}_{\tau g}(x) = \arg \min_z \left\{ \frac{1}{2} \|z - x\|^2 + \tau g(z) \right\}$$

- ▶ Ingredient 3 (conjugation)
  if $f$ is prox-friendly, but $f \circ A$ is not: split $f$ and $A$
  $f(Ax) = f^{**}(Ax) = \sup_y \{\langle Ax, y \rangle - f^*(x)\}$

  Dual: $\min_y \{f^*(y) + g^*(-A^*y)\}$
  Primal-Dual: $\min_x \max_y \{\langle Ax, y \rangle - f^*(y) + g(x)\}$

# Building Algorithms

Template: $\min_x \{f(Ax) + g(x) = F(x) + g(x)\}$

**New algorithms** are designed by mix-and-match:

**Proximal Gradient Descent**: Combettes and Wajs '05

$$x^+ = \text{prox}_{\tau g}(x - \tau \nabla F(x))$$

**Primal-Dual Hybrid Gradient** Chambolle and Pock '11

$$x^+ = \text{prox}_{\tau g}(x - \tau A^* y)$$
$$\overline{x} = x + \theta(x^+ - x)$$
$$y^+ = \text{prox}_{\sigma f^*}(y + \sigma A\overline{x})$$

**GD**

$$x^+ = x - \tau \nabla F(x)$$

# Revisiting Gradient Descent: SGD and its variants ($g = 0$)

**GD**

$$x^+ = x - \tau \sum_{i=1}^{n} \nabla F_i(x)$$

# Revisiting Gradient Descent: SGD and its variants ($g = 0$)

**GD**

$$x^+ = x - \tau \sum_{i=1}^{n} \nabla F_i(x)$$

**SGD** and variants

Uniformly at random select $j$

$$x^+ = x - \tau \tilde{\nabla}^j F(x)$$

▶ SGD: randomly choose $j$,

$$\tilde{\nabla}^j F(x) = n \nabla F_j(x)$$

nonconvergence for fixed $\tau$, "slow" convergence for carefully decreasing $\tau$ Robbins and Monro '51

# Revisiting Gradient Descent: SGD and its variants ($g = 0$)

**GD**

$$x^+ = x - \tau \sum_{i=1}^n \nabla F_i(x)$$

**SGD** and variants

Uniformly at random select $j$

$$x^+ = x - \tau \tilde{\nabla}^j F(x)$$

▶ SGD: randomly choose $j$,

$$\tilde{\nabla}^j F(x) = n \nabla F_j(x)$$

nonconvergence for fixed $\tau$, "slow" convergence for carefully decreasing $\tau$ Robbins and Monro '51

▶ SAGA/SVRG: randomly choose $j$,

$$\tilde{\nabla}^j F(x) = n(\nabla F_j(x) - G_j) + G$$

$G$ historic gradient, $G_j$ historic stochastic gradient Defazio et al. '14, Johnsen and Zhang '13, SAGA converges for $\tau \leq 1/(3nL_{\max})$

# Revisiting Gradient Descent: SGD and its variants ($g = 0$)

**GD**

$$x^+ = x - \tau \sum_{i=1}^n \nabla F_i(x)$$

**SGD** and variants

Uniformly at random select $j$

$$x^+ = x - \tau \tilde{\nabla}^j F(x)$$

▶ SGD: randomly choose $j$,

$$\tilde{\nabla}^j F(x) = n \nabla F_j(x)$$

nonconvergence for fixed $\tau$, "slow" convergence for carefully decreasing $\tau$ Robbins and Monro '51

▶ SAGA/SVRG: randomly choose $j$,

$$\tilde{\nabla}^j F(x) = n(\nabla F_j(x) - G_j) + G$$

$G$ historic gradient, $G_j$ historic stochastic gradient Defazio et al. '14, Johnsen and Zhang '13, SAGA converges for $\tau \leq 1/(3nL_{\max})$

▶ Similar algorithms for proximal point Bianchi '16, Traore et al. '23

## Revisiting PDHG

**PDHG**:

$$x^+ = \text{prox}_{\tau g}(x - \tau A^* y)$$
$$\overline{x} = x^+ + \theta(x^+ - x)$$
$$y^+ = \text{prox}_{\sigma f^*}(y + \sigma A\overline{x})$$

# Revisiting PDHG

**PDHG**:

$$x^+ = \text{prox}_{\tau g}(x - \tau A^* y)$$
$$\overline{x} = x^+ + \theta(x^+ - x)$$
$$y^+ = \text{prox}_{\sigma f^*}(y + \sigma A \overline{x})$$

**PDHG (dual extrapolation)**:

$$y^+ = \text{prox}_{\sigma f^*}(y + \sigma A x)$$
$$\overline{y} = y^+ + \theta(y^+ - y)$$
$$x^+ = \text{prox}_{\tau g}(x - \tau A^* \overline{y})$$

# Revisiting PDHG

**PDHG**:

$$x^+ = \text{prox}_{\tau g}(x - \tau A^* y)$$
$$\overline{x} = x^+ + \theta(x^+ - x)$$
$$y^+ = \text{prox}_{\sigma f^*}(y + \sigma A\overline{x})$$

**PDHG (dual extrapolation)**:

$$y^+ = \text{prox}_{\sigma f^*}(y + \sigma Ax)$$
$$\overline{y} = y^+ + \theta(y^+ - y)$$
$$x^+ = \text{prox}_{\tau g}(x - \tau A^*\overline{y})$$

**PDHG (dual extrapolation with $f = \sum_i f_i$)**:

$$y_i^+ = \text{prox}_{\sigma f_i^*}(y_i + \sigma A_i x), \quad i = 1, \ldots, n$$
$$\overline{y}_i = y_i^+ + \theta(y_i^+ - y_i), \quad i = 1, \ldots, n$$
$$x^+ = \text{prox}_{\tau g}(x - \tau \sum_{i=1}^n A_i^* \overline{y}_i)$$

# From PDHG to SPDHG

**PDHG (dual extrapolation with $f = \sum_i f_i$):**

$$y_i^+ = \text{prox}_{\sigma f_i^*}(y_i + \sigma A_i x), \quad i = 1, \dots, n$$

$$\bar{y}_i = y_i^+ + \theta(y_i^+ - y_i), \quad i = 1, \dots, n$$

$$x^+ = \text{prox}_{\tau g}\left(x - \tau \sum_{i=1}^n A_i^* \bar{y}_i\right)$$

# From PDHG to SPDHG

**PDHG (dual extrapolation with $f = \sum_i f_i$):**

$$y_i^+ = \text{prox}_{\sigma f_i^*}(y_i + \sigma A_i x), \quad i = 1, \ldots, n$$

$$\overline{y}_i = y_i^+ + \theta(y_i^+ - y_i), \quad i = 1, \ldots, n$$

$$x^+ = \text{prox}_{\tau g}\left(x - \tau \sum_{i=1}^n A_i^* \overline{y}_i\right)$$

**Stochastic PDHG (SPDHG):** Chambolle, Ehrhardt, Richtárik, Schönlieb '18

Uniform at randomly select $j$

$$y_i^+ = \text{prox}_{\sigma f_i^*}(y_i + \sigma A_i x), \quad i = j$$

$$\overline{y}_i = y_i^+ + \theta n(y_i^+ - y_i), \quad i = j; \quad \overline{y}_i = y_i \text{ else}$$

$$x^+ = \text{prox}_{\tau g}\left(x - \tau \sum_{i=1}^n A_i^* \overline{y}_i\right)$$

▶ convergence for $\sigma \tau < 1/(n \max_i \|A_i\|^2)$, $\theta = 1$

Chambolle, Ehrhardt, Richtárik, Schönlieb '18, Gutiérrez, Delplancke, Ehrhardt '21, Alacaoglu, Fercoq, Cevher '22

# SPDHG as SAGA

**SPDHG**:

Uniform at randomly select $j$

$y_i^+ = \text{prox}_{\sigma f_i^*}(y_i + \sigma A_i x), \quad i = j$

$\overline{y}_i = y_i^+ + \theta n(y_i^+ - y_i), \quad i = j; \quad \overline{y}_i = y_i$ else

$x^+ = \text{prox}_{\tau g}\left(x - \tau \sum_{i=1}^n A_i^* \overline{y}_i\right)$

# SPDHG as SAGA

**SPDHG**: Chambolle, Ehrhardt, Richtárik, Schönlieb '18

Uniform at randomly select $j$

$y_i^+ = \text{prox}_{\sigma f_i^*}(y_i + \sigma A_i x), \quad i = j$

$\bar{y}_i = y_i^+ + \theta n(y_i^+ - y_i), \quad i = j; \quad \bar{y}_i = y_i \text{ else}$

$x^+ = \text{prox}_{\tau g}(x - \tau \sum_{i=1}^n A_i^* \bar{y}_i)$

**SPDHG as SAGA (new)**:

Uniform at randomly select $j$

$y_i^+ = \text{prox}_{\sigma f_i^*}(y_i + \sigma A_i x), \quad i = j$

$\tilde{\nabla}^j = (1 + \theta n)A_j^*(y_j^+ - y_j) + \sum_{i=1}^n A_i^* y_i$

$x^+ = \text{prox}_{\tau g}(x - \tau \tilde{\nabla}^j)$

▶ essentially SAGA version of SPDHG

▶ for $\sigma = 1$, step size bound $\tau < 1/(n \max_i \|A_i\|^2)$ $3\times$ larger

# Numerical Results

# Subsets / minibatching

Forward Operator: $K : X \to \mathbb{R}^s$

$$\min_x \left\{ \sum_{j=1}^{s} \|K_j x - b_j\|^2 + \lambda \|Dx\|_1 + \imath_+(x) \right\}$$



- Choose subsets $S_i$
- $A_i = (K_j)_{j \in S_i} : X \to \mathbb{R}^{|S_i|}$
- $f_i(y) = \sum_{j \in S_i} \|K_j x - b_j\|^2$
- $n$ depends on the size of the subsets $S_i$
- $g(x) = \lambda \|Dx\|_1 + \imath_+(x)$

$$\min_x \left\{ \sum_{i=1}^{n} f_i(A_i x) + g(x) \right\}$$

# PET: Sanity Check, Convergence to Saddle Point (TV)



**saddle point (5000 iter PDHG)**

**SPDHG (20 epochs, 252 subsets)**

Ehrhardt, Markiewicz, Schönlieb '19

# PET: Faster than PDHG, TV, 20 epochs

**PDHG**

**SPDHG (252 subsets)**



Ehrhardt, Markiewicz, Schönlieb '19

# PET: Faster than PDHG, TV, 5 epochs



**PDHG**

**SPDHG (252 subsets)**

# PET: Faster than PDHG, TV, 1 epochs

**PDHG**



**SPDHG (252 subsets)**



Ehrhardt, Markiewicz, Schönlieb '19

# PET: More subsets are faster

$n = 1, 21, 100, 252$



Ehrhardt, Markiewicz, Schönlieb '19

# Step-size condition of SPDHG

$$\sigma\tau < 1/(n \max_i \|A_i\|^2)$$

▶ Is a large-product $\sigma\tau$ good? Empirically yes

# Step-size condition of SPDHG

$$\sigma\tau < 1/(n \max_i \|A_i\|^2)$$

▶ Is a large-product $\sigma\tau$ good? Empirically yes
▶ Is upper bound tight? No, e.g. for PDHG $\sigma\tau\|A\|^2 < 4/3$ is possible Ma et al. '23 (and in fact optimal). Empirically observed for SPDHG, e.g. Schramm and Holler '22

# Step-size condition of SPDHG

$$\sigma\tau < 1/(n \max_i \|A_i\|^2)$$

- ▶ Is a large-product $\sigma\tau$ good? Empirically yes
- ▶ Is upper bound tight? No, e.g. for PDHG $\sigma\tau\|A\|^2 < 4/3$ is possible Ma et al. '23 (and in fact optimal). Empirically observed for SPDHG, e.g. Schramm and Holler '22
- ▶ Is the ratio $\sigma/\tau$ important? Yes Delplancke et al. '20



(a) synthetic data                    (b) real data

# Step-size condition of SPDHG

$$\sigma\tau < 1/(n \max_i \|A_i\|^2)$$

- ▶ Is a large-product $\sigma\tau$ good? Empirically yes
- ▶ Is upper bound tight? No, e.g. for PDHG $\sigma\tau\|A\|^2 < 4/3$ is possible Ma et al. '23 (and in fact optimal). Empirically observed for SPDHG, e.g. Schramm and Holler '22
- ▶ Is the ratio $\sigma/\tau$ important? Yes Delplancke et al. '20



(a) synthetic data      (b) real data

- ▶ How to choose the ratio $\sigma/\tau$? Open question

# Adaptive step-sizes

- Idea: let $\sigma$ and $\tau$ vary with iterations
- PDHG: a bit of theory + emprical results <span style="font-size:smaller">Goldstein et al. '15</span>
- SPDHG: empirical results for MPI <span style="font-size:smaller">Zdun and Brandt '21</span>

# Adaptive step-sizes

- Idea: let $\sigma$ and $\tau$ vary with iterations
- PDHG: a bit of theory + emprical results Goldstein et al. '15
- SPDHG: empirical results for MPI Zdun and Brandt '21
- SPDHG: theory + numerics for CT Chambolle, Ehrhardt et al. '24
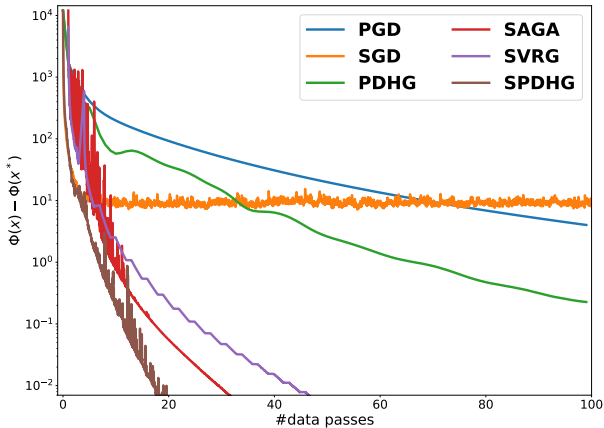
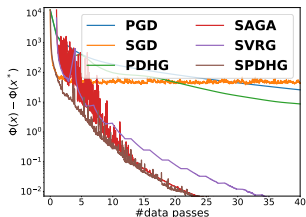# CT: 3 epochs <span>Ehrhardt, Kereta, Liang, Tang '24</span>
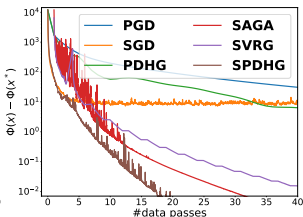
# CT: Quantitative Comparison
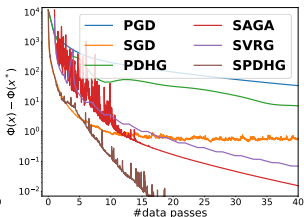


Ehrhardt, Kereta, Liang, Tang '24
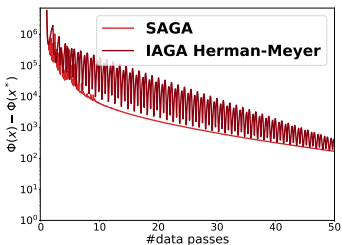
# CT: Quantitative Comparison, Noise



**high noise**     **medium noise (shown)**     **low noise**

▶ Speed seems to depend on noise in the data

▶ Gradient based methods more effected

Ehrhardt, Kereta, Liang, Tang '24

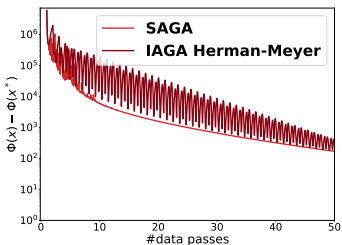# CT: Random v Deterministic
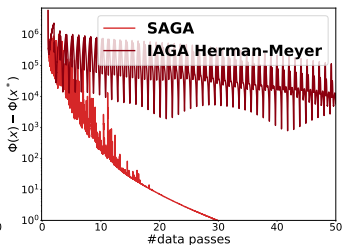


**30 subsets**

▶ similar convergence for 30 subsets (similar to literature)

Herman and Meyer '93, Ehrhardt, Kereta, Liang, Tang '24

# CT: Random v Deterministic



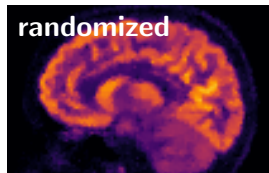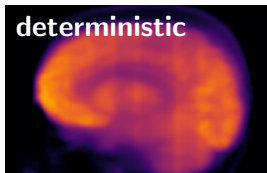**30 subsets**                    **240 subsets**

- ▶ similar convergence for 30 subsets (similar to literature)
- ▶ big difference for 240 subsets

Herman and Meyer '93, Ehrhardt, Kereta, Liang, Tang '24

# Conclusions and Outlook

**Conclusions:**

- ▶ **Zoo** of stochastic algorithms exists (gets larger and larger)
- ▶ **Randomness** seems important in general and not just mathematical convenience
- ▶ **Speeds up** reconstruction of inverse problems; e.g. PET, listmode PET (randomize over events), CT, parallel MRI, motion-corrected CT, magnetic particle imaging



deterministic



randomized

**Future directions:**

- ▶ Tighter analysis
- ▶ Inverse problems specific analysis
- ▶ Learned algorithms