# A Randomized Algorithm for Convex Optimization and Medical Imaging Applications

Matthias J. Ehrhardt

Department of Mathematical Sciences, University of Bath, UK

March 1, 2022

The Leverhulme Trust

UK RI Engineering and Physical Sciences Research Council

UNIVERSITY OF BATH

# Main Aim and Outline

**Main aim:**

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

▶ proper, convex and lower semi-continuous
▶ non-smooth
▶ $n$ is large and/or $\mathbf{B}_i x$ expensive

# Main Aim and Outline

**Main aim:**

$$x^{\sharp} \in \arg\min_{x} \left\{ \sum_{i=1}^{n} f_i(\mathbf{B}_i x) + g(x) \right\}$$

- ▶ proper, convex and lower semi-continuous
- ▶ non-smooth
- ▶ $n$ is large and/or $\mathbf{B}_i x$ expensive

**Outline:**

1) From Inverse Problems to Optimization (**Why?**)
2) Randomized Algorithm for Convex Optimization (**How?**)
3) Application: Medical Imaging (PET, CT, MRI)

From Inverse Problems to Optimization

# What is an inverse problem? Inverse to what?

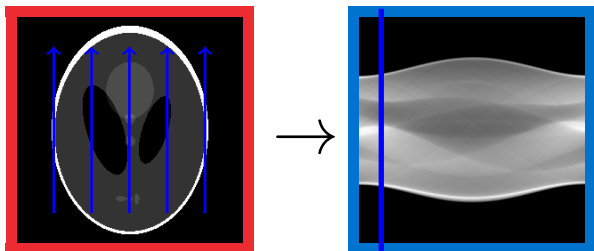**Forward problem:** given $u$, compute $\mathbf{A}u = v$. Evaluate $\mathbf{A}$

▶ $\mathbf{A} : U \to V$ (non-)linear operator between spaces $U$ and $V$

# What is an inverse problem? Inverse to what?

**Forward problem:** given $u$, compute $\mathbf{A}u = v$. Evaluate $\mathbf{A}$

▶ $\mathbf{A} : U \to V$ (non-)linear operator between spaces $U$ and $V$
▶ Example: Radon / X-ray transform (used in CT, PET, ...)

$$\mathbf{A}u(L) = \int_L u(r)\,dr$$

# What is an inverse problem? Inverse to what?

**Forward problem:** given $u$, compute $\mathbf{A}u = v$. Evaluate $\mathbf{A}$

- $\mathbf{A} : U \rightarrow V$ (non-)linear operator between spaces $U$ and $V$
- Example: Radon / X-ray transform (used in CT, PET, ...)
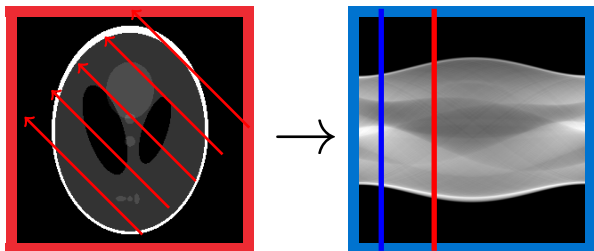
$$\mathbf{A}u(L) = \int_L u(r)\,dr$$

# What is an inverse problem? Inverse to what?

**Forward problem:** given $u$, compute $\mathbf{A}u = v$. Evaluate $\mathbf{A}$

- ▶ $\mathbf{A} : U \to V$ (non-)linear operator between spaces $U$ and $V$
- ▶ Example: Radon / X-ray transform (used in CT, PET, ...)
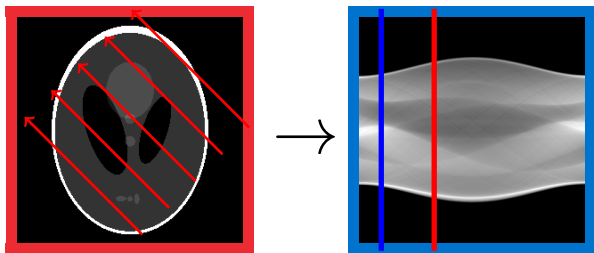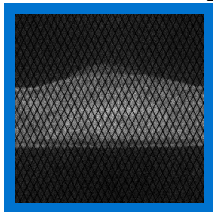
$$\mathbf{A}u(L) = \int_L u(r)\, dr$$



**Inverse problem:** given $v$, solve $\mathbf{A}u = v$. "Invert" $\mathbf{A}$

# What is the problem with inverse problems?

▶ PET example: $\mathbf{A}u(L) = \int_L u(r)\,dr$
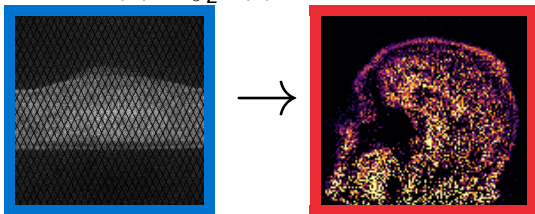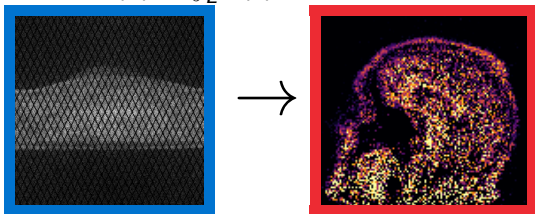
 $\longrightarrow$

# What is the problem with inverse problems?

▶ PET example: $\mathbf{A}u(L) = \int_L u(r)dr$

# What is the problem with inverse problems?

▶ PET example: $\mathbf{A}u(L) = \int_L u(r)dr$



**Definition (Hadamard, 1902):** We call an inverse problem $\mathbf{A}u = v$ **well-posed** if

    (1) a solution $u^*$ **exists**

    (2) the solution $u^*$ is **unique**

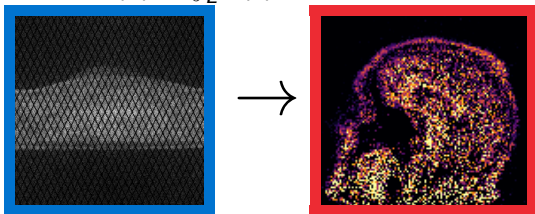    (3) $u^*$ depends **continuously** on data $v$.

Otherwise, it is called **ill-posed**.

Jacques Hadamard

# What is the problem with inverse problems?

▶ PET example: $\mathbf{A}u(L) = \int_L u(r)dr$



**Definition (Hadamard, 1902):** We call an inverse problem $\mathbf{A}u = v$ **well-posed** if

    (1) a solution $u^*$ **exists**

    (2) the solution $u^*$ is **unique**

    (3) $u^*$ depends **continuously** on data $v$.

Otherwise, it is called **ill-posed**.
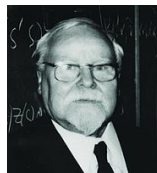
Jacques Hadamard

Most interesting problems are ill-posed, in particular (3) is violated.

# A way to solve inverse problems

**Tikhonov regularization (1943)**

Approximate a solution $u^*$ of $\mathbf{A}u = v$ via

$$u_\lambda = \arg\min_u \left\{ \|\mathbf{A}u - v\|^2 + \lambda\|u\|^2 \right\}$$
$$= (\mathbf{A}^*\mathbf{A} + \lambda I)^{-1}\mathbf{A}^*v$$



Andrey Tikhonov

# A way to solve inverse problems

**Variational regularization**

Approximate a solution $u^*$ of $\mathbf{A}u = v$ via

$$u_\lambda = \arg\min_u \left\{ D(\mathbf{A}u, v) + \lambda R(u) \right\}$$

▶ **data fit** $D$: quantify fit of prediction $\mathbf{A}u$ to data $v$. Usually a "divergence", i.e. $D(x, y) \geq 0$ and $D(x, y) = 0$ iff $x = y$

$$D(x, y) = \|x - y\|_2^2, \|x - y\|_1, \int x - y + y \log(y/x), \ldots$$

▶ **regularizer** $R$: penalize unwanted features, ensures stability

$$R(x) = \|x\|_2^2, \|x\|_1, \mathsf{TV}(x) = \|\nabla x\|_1, \mathsf{TGV}, \ldots$$

# PET Modelling

$$b_i \sim \text{Poisson}(a_i^T u + r_i)$$

- data $b_i \in \mathbb{N}$
- forward model $a_i^T u \approx \gamma_i \int_{L_i} u$ (x-ray transform)
- multiplicative correction $\gamma_i > 0$ (attenuation, normalisation)
- background $r_i > 0$ (scatter, randoms)
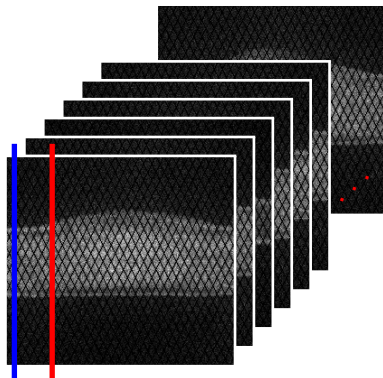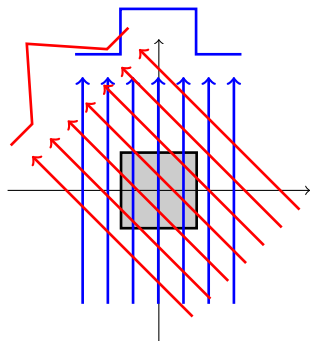
# PET Modelling

$$b_i \sim \text{Poisson}(a_i^T u + r_i)$$

- data $b_i \in \mathbb{N}$
- forward model $a_i^T u \approx \gamma_i \int_{L_i} u$ (x-ray transform)
- multiplicative correction $\gamma_i > 0$ (attenuation, normalisation)
- background $r_i > 0$ (scatter, randoms)
- number of data / rays: 2D $N = 86k$, 3D $N = 355M$

# PET Reconstruction[1]

$$u_\lambda \in \arg\min_u \left\{ \sum_{j=1}^m \mathcal{D}_j(\mathbf{A}_j u + r_j) + \lambda \mathcal{R}(u) + \imath_+(u) \right\}$$

▶ Partition data in "subsets" $\mathbb{S}_1, \ldots, \mathbb{S}_m$

$$\mathcal{D}_j(y) := \sum_{i \in \mathbb{S}_j} \mathsf{KL}(y_i; b_i)$$

▶ Kullback–Leibler divergence

$$\mathsf{KL}(y; b) = \begin{cases} y - b + b \log\left(\frac{b}{y}\right) & \text{if } y > 0 \\ \infty & \text{else} \end{cases}$$

▶ Regularizer $\mathcal{R}$, see next page
▶ Constraint

$$\imath_+(u) = \begin{cases} 0, & \text{if } u_i \geq 0 \text{ for all } i \\ \infty, & \text{else} \end{cases}$$

---

[1] Brune '10, Brune et al. '10, Setzer et al. '10, Müller et al. '11, Anthoine et al. '12, Knoll et al. '16, Ehrhardt et al. '16, Hohage and Werner '16, Schramm et al. '17, Rasch et al. '17, Ehrhardt et al. '17, Mehranian et al. '17 and many, many more

# PET Reconstruction with TV



**Total variation (TV)**

Rudin, Osher, Fatemi 1992

$$\mathcal{R}(u) = \|\nabla u\|_1$$

$$\min_u \left\{ \sum_{j=1}^m \mathcal{D}_j(\mathbf{A}_j u) + \lambda \|\nabla u\|_1 + \imath_+(u) \right\}$$

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

| | |
|---|---|
| $n = m + 1$ | $g(x) = \imath_+(x)$ |
| $\mathbf{B}_i = \mathbf{A}_i$ | $f_i = \mathcal{D}_i \quad i \in [m]$ |
| $\mathbf{B}_n = \nabla$ | $f_n = \lambda \|\cdot\|_1$ |

# PET Reconstruction with TGV

**Total generalized variation (TGV)**

Bredies, Kunisch, Pock 2010

$$\mathcal{R}(u) = \min_v \|\nabla u - v\|_1 + \beta \|\mathbf{D}v\|_1$$



$$\min_{u,v} \left\{ \sum_{j=1}^{m} \mathcal{D}_j(\mathbf{A}_j u) + \lambda \|\nabla u - v\|_1 + \lambda\beta \|\mathbf{D}v\|_1 + \imath_+(u) \right\}$$

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

| | |
|---|---|
| $n = m + 2$ | |
| $x = (u; v)$ | $g(x) = \imath_+(u)$ |
| $\mathbf{B}_i = (\mathbf{A}_i, 0)$ | $f_i = \mathcal{D}_i \quad i \in [m]$ |
| $\mathbf{B}_{n-1} = (\nabla, -\mathbf{I})$ | $f_{n-1} = \lambda \|\cdot\|_1$ |
| $\mathbf{B}_n = (0, \mathbf{D})$ | $f_n = \lambda\beta \|\cdot\|_1$ |

# Observations

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

▶ **Proper:** Extended valued $f : X \mapsto \mathbb{R} \cup \{\infty\}$ and $f \not\equiv \infty$

▶ **Convex:** e.g. $C$ convex $\Rightarrow \imath_C$ convex

▶ **Lower semi-continuous (lsc):** $x_k \to x$ then

$$f(x) \leq \liminf_{k \to \infty} f(x_k)$$

   ▶ continuous $\Rightarrow$ lsc
   ▶ $C$ closed $\Rightarrow \imath_C$ lsc

▶ $f(z) = \sum_i f_i(z_i)$ is "**separable**". Not separable in $x$.

# Observations

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

- **Proper:** Extended valued $f : X \mapsto \mathbb{R} \cup \{\infty\}$ and $f \not\equiv \infty$
- **Convex:** e.g. $C$ convex $\Rightarrow \imath_C$ convex
- **Lower semi-continuous (lsc):** $x_k \to x$ then

$$f(x) \le \liminf_{k \to \infty} f(x_k)$$

  - continuous $\Rightarrow$ lsc
  - $C$ closed $\Rightarrow \imath_C$ lsc
- $f(z) = \sum_i f_i(z_i)$ is "**separable**". Not separable in $x$.

Problem 1: The functions $f_i, g$ are non-smooth but "simple"
Problem 2: $n$ is large and/or $\mathbf{B}_i x$ expensive

Optimization

# Subgradient

If $f$ is convex and smooth, then for all $x, y \in X$ we have

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$$

# Subgradient

From now on: $X = \mathbb{R}^d$

If $f$ is convex and smooth, then for all $x, y \in X$ we have

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$$

Extend definition to non-differentiable functions:

**Definition:** $p \in X$ is called a **subgradient** of $f : X \mapsto \mathbb{R} \cup \{\infty\}$ at $x \in X$ if for all $y \in X$

$$f(y) \geq f(x) + \langle p, y - x \rangle$$

holds. The set of all subgradients at $x \in X$ is called the **subdifferential** and denoted by $\partial f(x)$.

Example: $f(x) = |x|$

$$\partial f(x) = \begin{cases} \{1\} & \text{if } x > 0 \\ [-1, 1] & \text{if } x = 0 \\ \{-1\} & \text{if } x < 0 \end{cases}$$

# Proximal Operators: A **gradient descent** point of view

**(Sub-)Gradient descent:** $p \in \partial f(x)$   $(= \{\nabla f(x)\}$ if $f$ is diff.$)$

$$x^+ = x - p$$

# Proximal Operators: A **gradient descent** point of view

**(Sub-)Gradient descent:** $p \in \partial f(x)$ $\quad (= \{\nabla f(x)\}$ if $f$ is diff.$)$

$$x^+ = x - p$$

**Implicit (Sub-)Gradient descent:** $p^+ \in \partial f(x^+)$

$$x^+ = x - p^+ \in x - \partial f(x^+)$$

# Proximal Operators: A **gradient descent** point of view

**(Sub-)Gradient descent:** $p \in \partial f(x)$   $(= \{\nabla f(x)\}$ if $f$ is diff.$)$

$$x^+ = x - p$$

**Implicit (Sub-)Gradient descent:** $p^+ \in \partial f(x^+)$

$$x^+ = x - p^+ \in x - \partial f(x^+)$$
$$\Leftrightarrow \quad x \in (I + \partial f)x^+$$

# Proximal Operators: A **gradient descent** point of view

**(Sub-)Gradient descent:** $p \in \partial f(x)$ $\quad (= \{\nabla f(x)\}$ if $f$ is diff.)

$$x^+ = x - p$$

**Implicit (Sub-)Gradient descent:** $p^+ \in \partial f(x^+)$

$$x^+ = x - p^+ \in x - \partial f(x^+)$$
$$\Leftrightarrow \quad x \in (I + \partial f)x^+$$
$$\Leftrightarrow \quad x^+ = (I + \partial f)^{-1}x \quad =: \mathsf{prox}_f(x)$$

**Definition:** The **proximal operator** of $f$ is defined as
$$\mathsf{prox}_f(x) := (I + \partial f)^{-1}(x).$$

$\mathsf{prox}_f$ has *many* names:
*prox / proximal / proximity / resolvent operator*

# Proximal Operators: A **minimization** point of view

**Definition:** The **proximal operator** of $f$ is defined as
$$\text{prox}_f(x) := \arg \min_z \left\{ \frac{1}{2} \|z - x\|^2 + f(z) \right\}$$

# Proximal Operators: A **minimization** point of view

**Definition:** The **proximal operator** of $f$ is defined as
$$\text{prox}_f(x) := \arg\min_z \left\{ \frac{1}{2}\|z - x\|^2 + f(z) \right\}$$

**Proposition:** $(I + \partial f)^{-1}(x) = \arg\min_z \left\{ \frac{1}{2}\|z - x\|^2 + f(z) \right\}$

# Proximal Operators: A **minimization** point of view

**Definition:** The **proximal operator** of $f$ is defined as
$$\text{prox}_f(x) := \arg\min_z \left\{ \frac{1}{2}\|z - x\|^2 + f(z) \right\}$$

**Proposition:** $(I + \partial f)^{-1}(x) = \arg\min_z \left\{ \frac{1}{2}\|z - x\|^2 + f(z) \right\}$

"Proof":
$$x^+ = \arg\min_z \left\{ \frac{1}{2}\|z - x\|^2 + f(z) \right\}$$
$$\Leftrightarrow \quad 0 \in \partial \left\{ \frac{1}{2}\|x^+ - x\|^2 + f(x^+) \right\}$$
$$\Leftrightarrow \quad 0 \in x^+ - x + \partial f(x^+)$$
$$\Leftrightarrow \quad x \in (I + \partial f)x^+$$
$$\Leftrightarrow \quad x^+ = (I + \partial f)^{-1}x$$

# Proximal operator: properties and examples

$$\mathsf{prox}_f(x) = \arg\min_z \left\{ \frac{1}{2}\|z - x\|^2 + f(z) \right\}$$

**Many rules:** e.g.

**Proposition:** Let $f$ be separable, i.e. $f(x) = \sum_i f_i(x_i)$. Then
$$\mathsf{prox}_f(x)_i = \mathsf{prox}_{f_i}(x_i).$$

Examples:
- $f(x) = \frac{1}{2}\|x\|_2^2$:   $\mathsf{prox}_f(x) = \frac{1}{2}x$
- $f(x) = \|x\|_1$:
$$\mathsf{prox}_f(x)_i = \begin{cases} x_i - 1 & \text{if } x_i > 1 \\ 0 & |x_i| \leq 1 \\ x_i + 1 & \text{if } x_i < -1 \end{cases}$$
- $f = \imath_C$:   $\mathsf{prox}_f(x) = \mathsf{proj}_C(x)$
- $f = \imath_{\geq 0}$:   $\mathsf{prox}_f(x)_i = \max(x_i, 0)$

# Proximal operator: properties and examples

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

**Many rules:** e.g.

**Proposition:** Let $f$ be separable, i.e. $f(x) = \sum_i f_i(x_i)$. Then
$$\text{prox}_f(x)_i = \text{prox}_{f_i}(x_i).$$

Examples:

- $f(x) = \frac{1}{2}\|x\|_2^2$:  $\text{prox}_f(x) = \frac{1}{2}x$
- $f(x) = \|x\|_1$:
$$\text{prox}_f(x)_i = \begin{cases} x_i - 1 & \text{if } x_i > 1 \\ 0 & |x_i| \leq 1 \\ x_i + 1 & \text{if } x_i < -1 \end{cases}$$

- $f = \iota_C$:  $\text{prox}_f(x) = \text{proj}_C(x)$
- $f = \iota_{\geq 0}$:  $\text{prox}_f(x)_i = \max(x_i, 0)$

**Problem:** What is the proximal operator of $f(x) = \|\mathbf{C}x\|_1$?

# The way out: Saddle Point Problems

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

- $f(y) := \sum_i f_i(y_i)$, $\mathbf{B} = [\mathbf{B}_1; \ldots; \mathbf{B}_n]$

$$x^\sharp \in \arg\min_x \left\{ f(\mathbf{B}x) + g(x) \right\}$$

# The way out: Saddle Point Problems

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

▶ $f(y) := \sum_i f_i(y_i)$, $\mathbf{B} = [\mathbf{B}_1; \ldots; \mathbf{B}_n]$

$$x^\sharp \in \arg\min_x \left\{ f(\mathbf{B}x) + g(x) \right\}$$

**Definition:** The **convex conjugate** of $f$ is given by
$$f^*(y) := \sup_z \langle z, y \rangle - f(z).$$

**Theorem:** Let $f$ be proper, convex and lsc, then
$$f(z) = (f^*)^*(z) = \sup_y \langle z, y \rangle - f^*(y).$$

# The way out: Saddle Point Problems

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

► $f(y) := \sum_i f_i(y_i)$, $\mathbf{B} = [\mathbf{B}_1; \ldots; \mathbf{B}_n]$

$$x^\sharp \in \arg\min_x \left\{ f(\mathbf{B}x) + g(x) \right\}$$

**Definition:** The **convex conjugate** of $f$ is given by
$$f^*(y) := \sup_z \langle z, y \rangle - f(z).$$

**Theorem:** Let $f$ be proper, convex and lsc, then
$$f(z) = (f^*)^*(z) = \sup_y \langle z, y \rangle - f^*(y).$$

$$(x^\sharp, y^\sharp) \in \arg\min_x \sup_y \left\{ \langle \mathbf{B}x, y \rangle - f^*(y) + g(x) \right\}$$

# Primal-Dual Hybrid Gradient (PDHG) Algorithm[1]

Given $x^0, y^0, \overline{y}^0 = y^0$

(1) $x^{k+1} = \mathrm{prox}_{\tau g}(x^k - \tau \mathbf{B}^* \overline{y}^k)$

(2) $y^{k+1} = \mathrm{prox}_{\sigma f^*}(y^k + \sigma \mathbf{B} x^{k+1})$

(3) $\overline{y}^{k+1} = y^{k+1} + \theta(y^{k+1} - y^k)$

- ► evaluation of $\mathbf{B}$ and $\mathbf{B}^*$
- ► proximal operator
- ► convergence: $\theta = 1, \sigma\tau\|\mathbf{B}\|^2 < 1$

[1]Pock, Cremers, Bischof, Chambolle '09, Chambolle and Pock '11

# Primal-Dual Hybrid Gradient (PDHG) Algorithm[1]

Given $x^0, y^0, \overline{y}^0 = y^0$

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

    (2) $y_i^{k+1} = \text{prox}_{\sigma f_i^*}(y_i^k + \sigma \mathbf{B}_i x^{k+1}) \quad i = 1, \dots, n$

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

- $f(y) := \sum_i f_i(y_i), \; [\text{prox}_{f^*}(y)]_i = \text{prox}_{f_i^*}(y_i)$
- $\mathbf{B} = [\mathbf{B}_1; \dots; \mathbf{B}_n]^T, \; \mathbf{B}^* y = \sum_{i=1}^n \mathbf{B}_i^* y_i$

---

[1] Pock, Cremers, Bischof, Chambolle '09, Chambolle and Pock '11

# Primal-Dual Hybrid Gradient (PDHG) Algorithm[1]

Given $x^0, y^0, \overline{y}^0 = y^0$

$\quad$ (1) $x^{k+1} = \text{prox}_{\tau g}(x^k - \tau \sum_{i=1}^n \mathbf{B}_i^* \overline{y}_i^k)$

$\quad$ (2) $y_i^{k+1} = \text{prox}_{\sigma f_i^*}(y_i^k + \sigma \mathbf{B}_i x^{k+1})$ $\quad i = 1, \ldots, n$

$\quad$ (3) $\overline{y}_i^{k+1} = y_i^{k+1} + \theta(y_i^{k+1} - y_i^k)$ $\quad i = 1, \ldots, n$

- $f(y) := \sum_i f_i(y_i)$, $[\text{prox}_{f^*}(y)]_i = \text{prox}_{f_i^*}(y_i)$
- $\mathbf{B} = [\mathbf{B}_1; \ldots; \mathbf{B}_n]^T$, $\mathbf{B}^* y = \sum_{i=1}^n \mathbf{B}_i^* y_i$

---

[1]Pock, Cremers, Bischof, Chambolle '09, Chambolle and Pock '11

# Stochastic PDHG Algorithm[1]

Given $x^0, y^0, \overline{y}^0 = y^0$

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

Select $\mathbb{S}^{k+1} \subset \{1, \ldots, n\}$ randomly.

(2) $y_i^{k+1} = \begin{cases} \text{prox}_{\sigma_i f_i^*}(y_i^k + \sigma_i \mathbf{B}_i x^{k+1}) & i \in \mathbb{S}^{k+1} \\ y_i^k & \text{else} \end{cases}$

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

▶ probabilities $p_i := \mathbb{P}(i \in \mathbb{S}^{k+1}) > 0$ (**proper** sampling)

▶ $\sum_{i=1}^n \mathbf{B}_i^* \overline{y}_i^k$ can be computed using only $\mathbf{B}_i^*$ for $i \in \mathbb{S}^k$

▶ evaluation of $\mathbf{B}_i$ and $\mathbf{B}_i^*$ only for $i \in \mathbb{S}^{k+1}$.

---

[1] Chambolle, Ehrhardt, Richtárik, Schönlieb '18

# Convergence Guarantees

# Step Size Condition with ESO[1]

Tall matrix $\mathbf{C} = [\mathbf{C}_1; \ldots; \mathbf{C}_n]$, $\mathbf{C}^* h = \sum_{i=1}^{n} \mathbf{C}_i^* h_i$

**Definition (Expected Separable Overapproximation, ESO):**
Random subset $\mathbb{S} \subset \{1, \ldots, n\}$. The ESO parameters $v_i$ fulfill the ESO inequality if for all $h$
$$\mathbb{E}_{\mathbb{S}} \left\| \sum_{i \in \mathbb{S}} \mathbf{C}_i^* h_i \right\|^2 \leq \sum_{i=1}^{n} p_i v_i \|h_i\|^2 .$$

---

[1] Qu, Richtárik, Zhang '14

# Step Size Condition with ESO[1]

Tall matrix $\mathbf{C} = [\mathbf{C}_1; \ldots; \mathbf{C}_n]$, $\mathbf{C}^* h = \sum_{i=1}^{n} \mathbf{C}_i^* h_i$

**Definition (Expected Separable Overapproximation, ESO):**
Random subset $\mathbb{S} \subset \{1, \ldots, n\}$. The ESO parameters $v_i$ fulfill the ESO inequality if for all $h$

$$\mathbb{E}_{\mathbb{S}} \left\| \sum_{i \in \mathbb{S}} \mathbf{C}_i^* h_i \right\|^2 \leq \sum_{i=1}^{n} p_i v_i \|h_i\|^2 .$$

**Example (Full Sampling):** $\mathbb{S} = \{1, \ldots, n\}, p_i = 1, v_i = \|\mathbf{C}\|^2$

$$LHS = \|\mathbf{C}^* h\|^2$$

---

[1] Qu, Richtárik, Zhang '14

# Step Size Condition with ESO[1]

Tall matrix $\mathbf{C} = [\mathbf{C}_1; \ldots; \mathbf{C}_n]$, $\mathbf{C}^* h = \sum_{i=1}^{n} \mathbf{C}_i^* h_i$

**Definition (Expected Separable Overapproximation, ESO):**
Random subset $\mathbb{S} \subset \{1, \ldots, n\}$. The ESO parameters $v_i$ fulfill the ESO inequality if for all $h$
$$\mathbb{E}_{\mathbb{S}} \left\| \sum_{i \in \mathbb{S}} \mathbf{C}_i^* h_i \right\|^2 \leq \sum_{i=1}^{n} p_i v_i \|h_i\|^2 .$$

**Example (Full Sampling):** $\mathbb{S} = \{1, \ldots, n\}$, $p_i = 1$, $v_i = \|\mathbf{C}\|^2$
$$LHS = \|\mathbf{C}^* h\|^2 \leq \|\mathbf{C}^*\|^2 \|h\|^2$$

---

[1] Qu, Richtárik, Zhang '14

# Step Size Condition with ESO[1]

Tall matrix $\mathbf{C} = [\mathbf{C}_1; \ldots; \mathbf{C}_n]$, $\mathbf{C}^* h = \sum_{i=1}^{n} \mathbf{C}_i^* h_i$

**Definition (Expected Separable Overapproximation, ESO):**
Random subset $\mathbb{S} \subset \{1, \ldots, n\}$. The ESO parameters $v_i$ fulfill the ESO inequality if for all $h$

$$\mathbb{E}_{\mathbb{S}} \left\| \sum_{i \in \mathbb{S}} \mathbf{C}_i^* h_i \right\|^2 \leq \sum_{i=1}^{n} p_i v_i \|h_i\|^2.$$

**Example (Full Sampling):** $\mathbb{S} = \{1, \ldots, n\}$, $p_i = 1$, $v_i = \|\mathbf{C}\|^2$

$$LHS = \|\mathbf{C}^* h\|^2 \leq \|\mathbf{C}^*\|^2 \|h\|^2 = \sum_{i=1}^{n} \|\mathbf{C}^*\|^2 \|h_i\|^2$$

---

[1] Qu, Richtárik, Zhang '14

# Step Size Condition with ESO[1]

Tall matrix $\mathbf{C} = [\mathbf{C}_1; \ldots; \mathbf{C}_n]$, $\mathbf{C}^* h = \sum_{i=1}^{n} \mathbf{C}_i^* h_i$

**Definition (Expected Separable Overapproximation, ESO):**
Random subset $\mathbb{S} \subset \{1, \ldots, n\}$. The ESO parameters $v_i$ fulfill the ESO inequality if for all $h$

$$\mathbb{E}_{\mathbb{S}} \left\| \sum_{i \in \mathbb{S}} \mathbf{C}_i^* h_i \right\|^2 \leq \sum_{i=1}^{n} p_i v_i \|h_i\|^2.$$

**Example (Full Sampling):** $\mathbb{S} = \{1, \ldots, n\}, p_i = 1, v_i = \|\mathbf{C}\|^2$

$$LHS = \|\mathbf{C}^* h\|^2 \leq \|\mathbf{C}^*\|^2 \|h\|^2 = \sum_{i=1}^{n} \|\mathbf{C}^*\|^2 \|h_i\|^2$$

**Example (Serial Sampling):** $\mathbb{S} = \{i\}, v_i = \|\mathbf{C}_i\|^2$
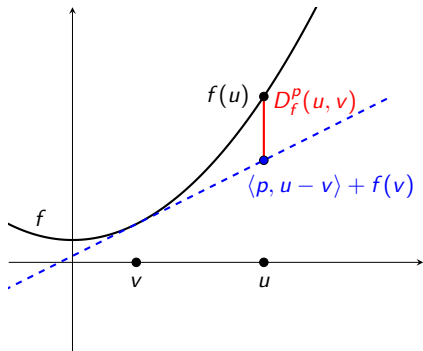
$$LHS = \sum_{i=1}^{n} p_i \|\mathbf{C}_i^* h_i\|^2$$

[1]Qu, Richtárik, Zhang '14

# Step Size Condition with ESO[1]

Tall matrix $\mathbf{C} = [\mathbf{C}_1; \ldots; \mathbf{C}_n]$, $\mathbf{C}^* h = \sum_{i=1}^n \mathbf{C}_i^* h_i$

**Definition (Expected Separable Overapproximation, ESO):**
Random subset $\mathbb{S} \subset \{1, \ldots, n\}$. The ESO parameters $v_i$ fulfill the ESO inequality if for all $h$

$$\mathbb{E}_{\mathbb{S}} \left\| \sum_{i \in \mathbb{S}} \mathbf{C}_i^* h_i \right\|^2 \leq \sum_{i=1}^n p_i v_i \|h_i\|^2 .$$

**Example (Full Sampling):** $\mathbb{S} = \{1, \ldots, n\}, p_i = 1,\ v_i = \|\mathbf{C}\|^2$

$$LHS = \|\mathbf{C}^* h\|^2 \leq \|\mathbf{C}^*\|^2 \|h\|^2 = \sum_{i=1}^n \|\mathbf{C}^*\|^2 \|h_i\|^2$$

**Example (Serial Sampling):** $\mathbb{S} = \{i\},\ v_i = \|\mathbf{C}_i\|^2$

$$LHS = \sum_{i=1}^n p_i \|\mathbf{C}_i^* h_i\|^2 \leq \sum_{i=1}^n p_i \|\mathbf{C}_i^*\|^2 \|h_i\|^2$$

[1] Qu, Richtárik, Zhang '14

# Bregman Distance

**Definition:** The **Bregman distance** of $f$ is defined as
$$D_f^p(u, v) = f(u) - f(v) - \langle p, u - v \rangle, \qquad p \in \partial f(v).$$

# Convergence of SPDHG

Let $\theta = 1$ and choose $\sigma_i, \tau$ such that there exist ESO parameters $v_i$ of $\mathbf{C} = [\mathbf{C}_1; \ldots, \mathbf{C}_n]$ with $\mathbf{C}_i = \sqrt{\sigma_i \tau} \mathbf{B}_i$ which satisfy $v_i < p_i$.

**Theorem:** Chambolle, Ehrhardt, Richtárik, Schönlieb '18
Let $(x^\sharp, y^\sharp)$ be a saddle point. Then

- Almost surely: $D_g^{\mathbf{B}^* y^\sharp}(x^k, x^\sharp) + D_{f^*}^{-\mathbf{B} x^\sharp}(y^k, y^\sharp) \to 0$

- Rate for ergodic sequence $(\hat{x}^K, \hat{y}^K) = \frac{1}{K} \sum_{k=1}^{K} (x^k, y^k)$

$$\mathbb{E}\left\{ D_g^{\mathbf{B}^* y^\sharp}(\hat{x}^K x^\sharp) + D_{f^*}^{-\mathbf{B} x^\sharp}(\hat{y}^K, y^\sharp) \right\} \leq \frac{C}{K}$$

# Convergence of SPDHG

Let $\theta = 1$ and choose $\sigma_i, \tau$ such that there exist ESO parameters $v_i$ of $\mathbf{C} = [\mathbf{C}_1; \ldots, \mathbf{C}_n]$ with $\mathbf{C}_i = \sqrt{\sigma_i \tau} \mathbf{B}_i$ which satisfy $v_i < p_i$.

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

Let $(x^\sharp, y^\sharp)$ be a saddle point. Then

▶ Almost surely: $D_g^{\mathbf{B}^* y^\sharp}(x^k, x^\sharp) + D_{f^*}^{-\mathbf{B}x^\sharp}(y^k, y^\sharp) \to 0$

▶ Rate for ergodic sequence $(\hat{x}^K, \hat{y}^K) = \frac{1}{K}\sum_{k=1}^{K}(x^k, y^k)$

$$\mathbb{E}\left\{ D_g^{\mathbf{B}^* y^\sharp}(\hat{x}^K x^\sharp) + D_{f^*}^{-\mathbf{B}x^\sharp}(\hat{y}^K, y^\sharp) \right\} \leq \frac{C}{K}$$

**Theorem:** Gutiérrez, Delplancke, Ehrhardt '21

Let $\mathbb{S}$ be serial sampling. Then there exists a saddle point $(x^\sharp, y^\sharp)$ such that almost surely

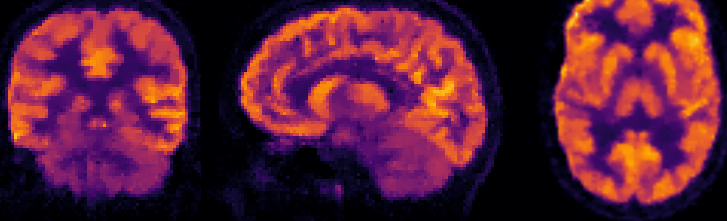$$(x^k, y^k) \to (x^\sharp, y^\sharp).$$

# Convergence of SPDHG

Let $\theta = 1$ and choose $\sigma_i, \tau$ such that there exist ESO parameters $v_i$ of $\mathbf{C} = [\mathbf{C}_1; \ldots, \mathbf{C}_n]$ with $\mathbf{C}_i = \sqrt{\sigma_i \tau} \mathbf{B}_i$ which satisfy $v_i < p_i$.

**Theorem:** Chambolle, Ehrhardt, Richtárik, Schönlieb '18
Let $(x^\sharp, y^\sharp)$ be a saddle point. Then

► Almost surely: $D_g^{\mathbf{B}^* y^\sharp}(x^k, x^\sharp) + D_{f^*}^{-\mathbf{B} x^\sharp}(y^k, y^\sharp) \to 0$

► Rate for ergodic sequence $(\hat{x}^K, \hat{y}^K) = \frac{1}{K} \sum_{k=1}^{K} (x^k, y^k)$

$$\mathbb{E} \left\{ D_g^{\mathbf{B}^* y^\sharp}(\hat{x}^K x^\sharp) + D_{f^*}^{-\mathbf{B} x^\sharp}(\hat{y}^K, y^\sharp) \right\} \leq \frac{C}{K}$$

**Theorem:** Gutiérrez, Delplancke, Ehrhardt '21
Let $\mathbb{S}$ be serial sampling. Then there exists a saddle point $(x^\sharp, y^\sharp)$ such that almost surely

$$(x^k, y^k) \to (x^\sharp, y^\sharp).$$

Remark: See also Alacaoglu, Fercoq, Cevher '20 for a similar result with a different proof.

Applications

# Sanity Check: Convergence to Saddle Point (TV)



**saddle point (5000 iter PDHG)**

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

# More subsets are faster

$m = 1, 21, 100, 252$



Ehrhardt, Markiewicz, Schönlieb '19

# "Balanced sampling" is faster

uniform sampling: $p_i = 1/n$

balanced sampling: $p_i = \begin{cases} \frac{1}{2m} & i < n \\ \frac{1}{2} & i = n \end{cases}$



Ehrhardt, Markiewicz, Schönlieb '19

# Faster than PDHG, TV



**PDHG (20 epochs)**

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

Ehrhardt, Markiewicz, Schönlieb '19

# Faster than PDHG, TV

**PDHG (5 epochs)**



**SPDHG (252 subsets, balanced, 5 epochs)**



Ehrhardt, Markiewicz, Schönlieb '19

# Faster than PDHG, TV

**PDHG (1 epoch)**



**SPDHG (252 subsets, balanced, 1 epoch)**



Ehrhardt, Markiewicz, Schönlieb '19

# Total Generalized Variation

**saddle point (PDHG, 5000 iterations)**



**SPDHG (252 subsets, balanced, 10 epochs)**



Ehrhardt, Markiewicz, Schönlieb '19

# Total Generalized Variation

**PDHG (10 epochs)**



**SPDHG (252 subsets, balanced, 10 epochs)**

Ehrhardt, Markiewicz, Schönlieb '19

# Motion corrected CT reconstruction

$$\min_x \sum_{i=1}^{n} \|AM_i x - b_i\|^2 + R(x)$$

▶ Here $n = 10$ motion gates
▶ No motion correction: $M_i = I$



Delplancke, Thielemans, Ehrhardt '21

# Motion corrected CT reconstruction

$$\min_x \sum_{i=1}^{n} \|A M_i x - b_i\|^2 + R(x)$$

- Here $n = 10$ motion gates
- No motion correction: $M_i = I$



PDHG converged | SPDHG 30 epochs | PDHG 30 epochs

Delplancke, Thielemans, Ehrhardt '21

# Parallel MRI

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

▶ Here $n = 8$ coils



PDHG converged  SPDHG 100 epochs  PDHG 100 epochs

Gutiérrez, Delplancke, Ehrhardt '21

# Conclusions and Outlook

- **Randomized** optimisation for cost functionals with "separable structure"
- **Generalisation** of PDHG and its convergence results
- **Much faster** PET reconstruction: making advanced models feasible for clinical data

**Not shown today:**

- Convergence theorems: 1) $\mathcal{O}(1/k^2)$ acceleration, 2) linear convergence
  Chambolle, Ehrhardt, Richtárik, Schönlieb '18

**Future work:**

- adaptive and optimal sampling
- adaptive step-sizes