

---

# Sketchy Decisions

---



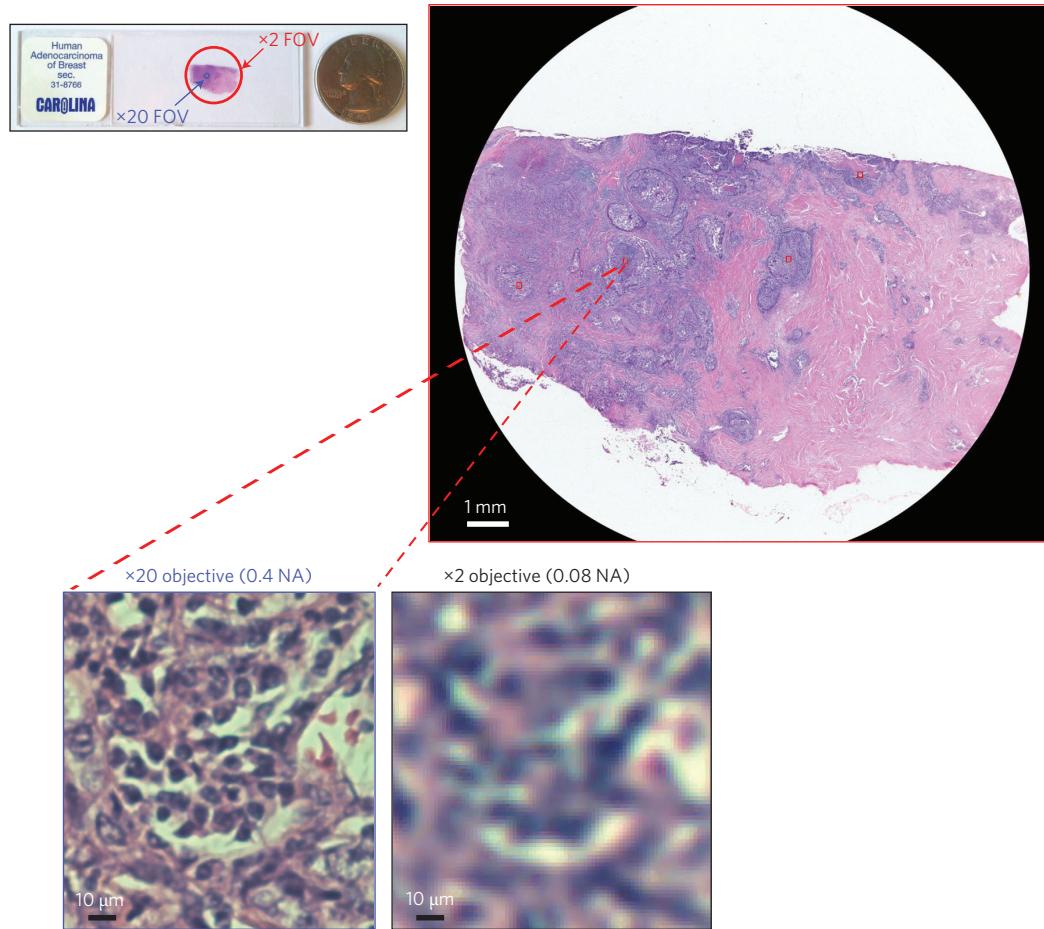
Joel A. Tropp

Computing + Mathematical Sciences  
California Institute of Technology  
jtropp@cms.caltech.edu

Collaborators: Volkan Cevher (EPFL), Roarke Horstmeyer (Duke),  
Quoc Tran-Dinh (UNC), Madeleine Udell (Cornell), Alp Yurtsever (EPFL)

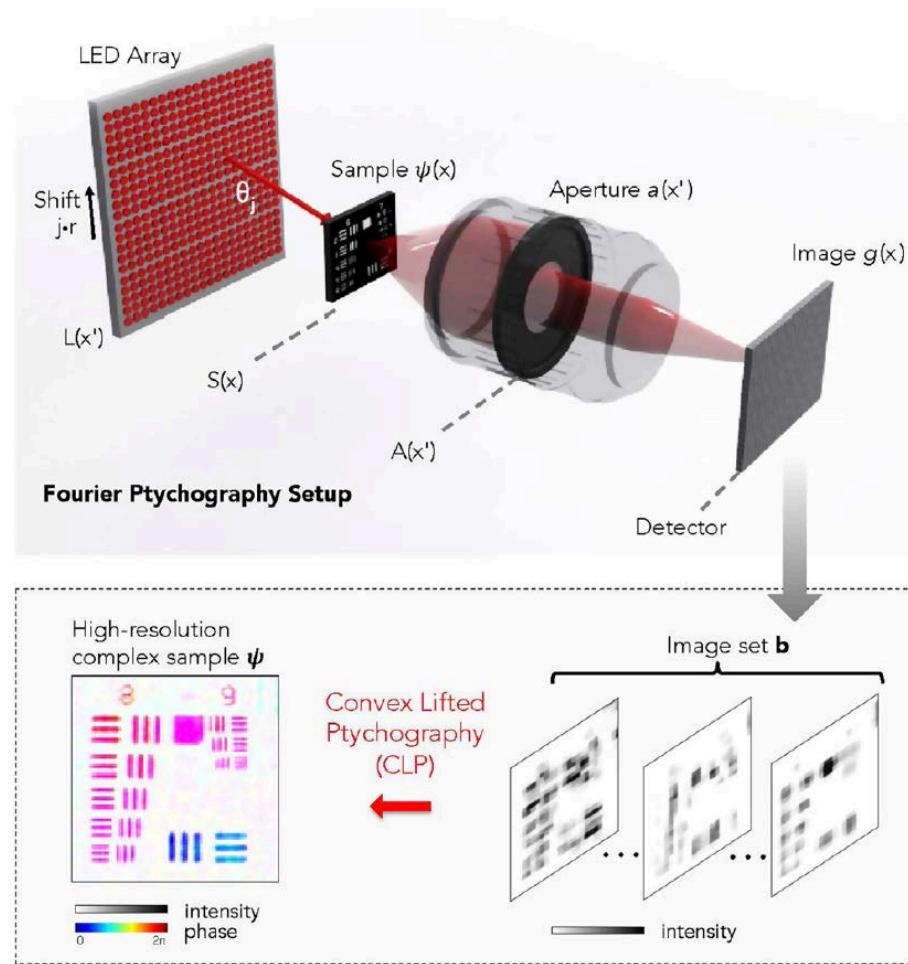
# Fourier Ptychography

# Microscopy: Field of View / Resolution



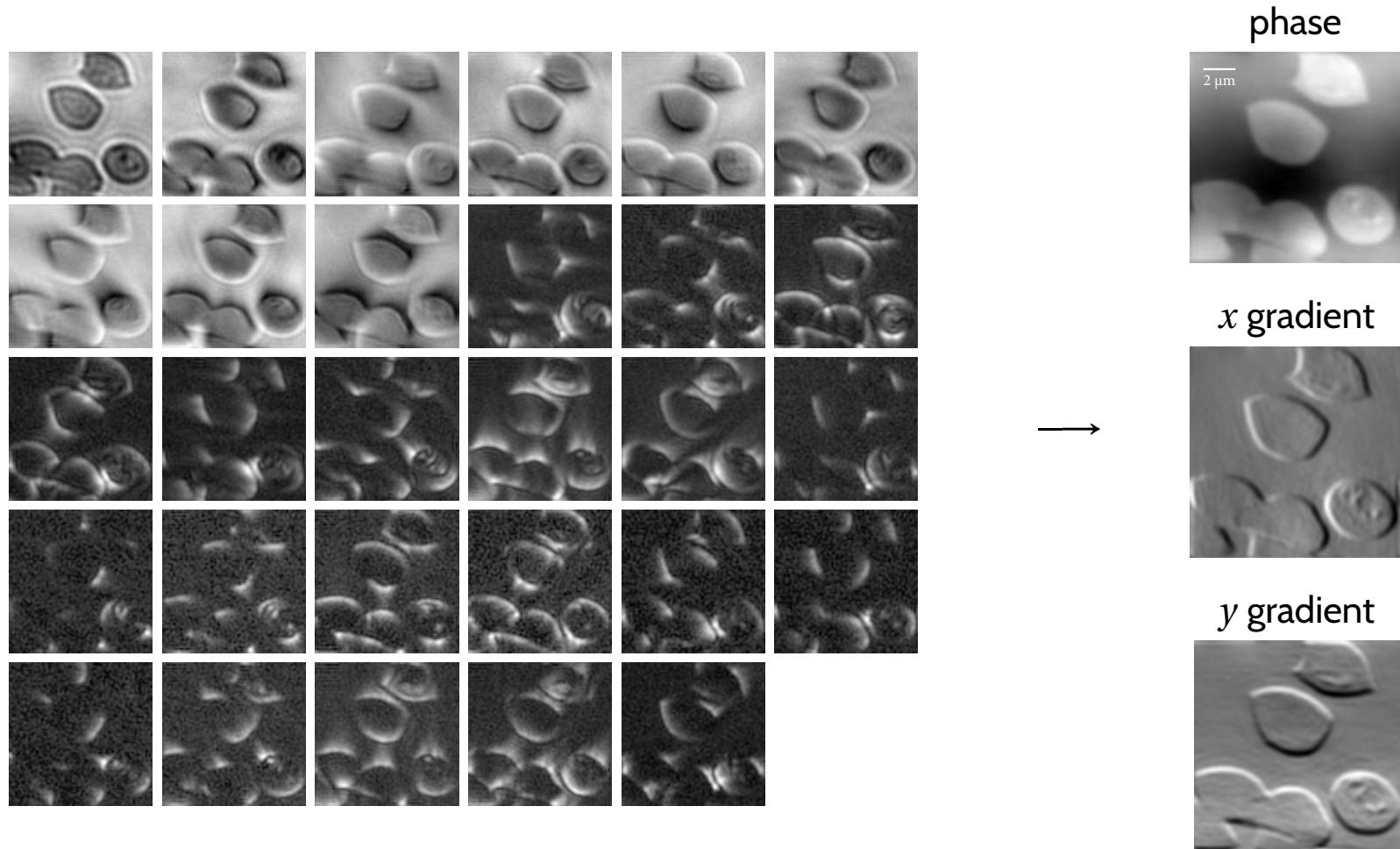
Source: Adapted from Zhang et al. 2013.

# Fourier Ptychography: Field of View + Resolution



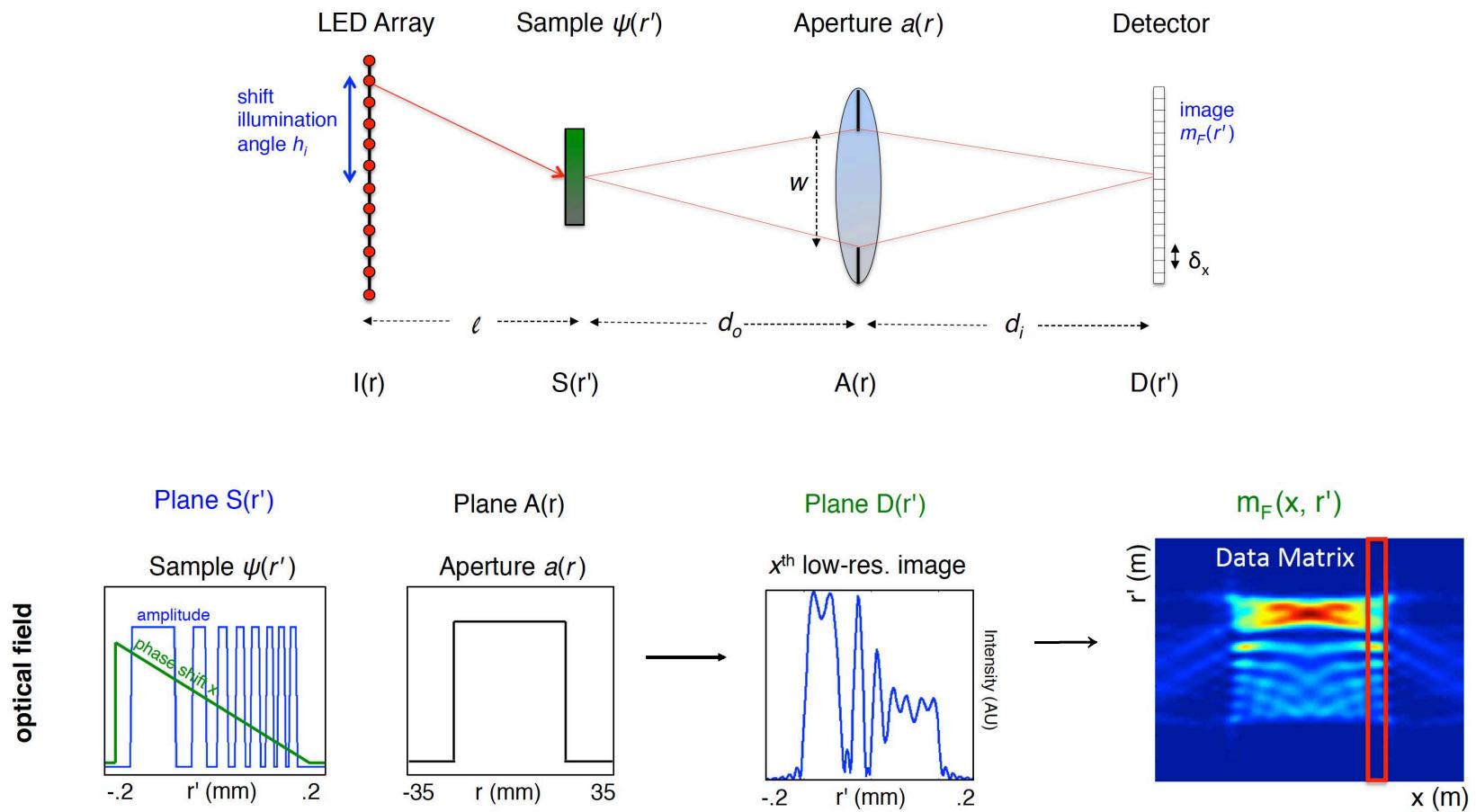
Sources: Zhang et al. 2013; Horstmeyer & Yang 2014; Ou et al. 2014; Horstmeyer et al. 2015.

# Fourier Ptychography: Malaria Example



Source: Yurtsever et al. 2017.

# Fourier Ptychography: Schematic



Source: Adapted from Horstmeyer & Yang 2014.

---

# Fourier Ptychography: Reconstruction

---

- Acquire a family of noisy measurements:

$$b_i = |\langle \mathbf{a}_i, \boldsymbol{\psi} \rangle|^2 + \xi_i \quad \text{for } i = 1, \dots, d$$

- $\mathbf{a}_i \in \mathbb{C}^n$  are known measurement vectors that model FP system
- $\boldsymbol{\psi} \in \mathbb{C}^n$  is the unknown sample transmission function
- $\xi_i \in \mathbb{R}$  is unknown noise
- Reconstruction via optimization:

$$\underset{\mathbf{x} \in \mathbb{C}^n}{\text{minimize}} \quad \sum_{i=1}^d \text{loss}\left(|\langle \mathbf{a}_i, \mathbf{x} \rangle|^2; b_i\right)$$

- Assume  $\text{loss}(\cdot; b)$  is a convex function
- **Malaria example:**  $n = 25,600$  and  $d = 185,600$

Sources: Zhang et al. 2013; Horstmeyer & Yang 2014; Horstmeyer et al. 2015.

---

# Fourier Ptychography: Convex Reconstruction

---

- ✿ **Observe:**  $|\langle \mathbf{a}, \mathbf{x} \rangle|^2 = \mathbf{a}^*(\mathbf{x}\mathbf{x}^*)\mathbf{a} = \mathbf{a}^*\mathbf{X}\mathbf{a}$  where  $\mathbf{X}$  is rank-one, psd

- ✿ Lift to matrix optimization problem:

$$\underset{\mathbf{X} \in \mathbb{C}^{n \times n}}{\text{minimize}} \quad \sum_{i=1}^d \text{loss}(\mathbf{a}_i^* \mathbf{X} \mathbf{a}_i; b_i) \quad \text{subject to} \quad \text{rank}(\mathbf{X}) = 1; \quad \mathbf{X} \text{ psd}$$

- ✿ Replace rank constraint with trace constraint to obtain convex problem:

$$\underset{\mathbf{X} \in \mathbb{C}^{n \times n}}{\text{minimize}} \quad \sum_{i=1}^d \text{loss}(\mathbf{a}_i^* \mathbf{X} \mathbf{a}_i; b_i) \quad \text{subject to} \quad \text{trace}(\mathbf{X}) = \alpha; \quad \mathbf{X} \text{ psd}$$

- ✿ Return maximum eigenvector  $\mathbf{x}_*$  of a solution  $\mathbf{X}_*$

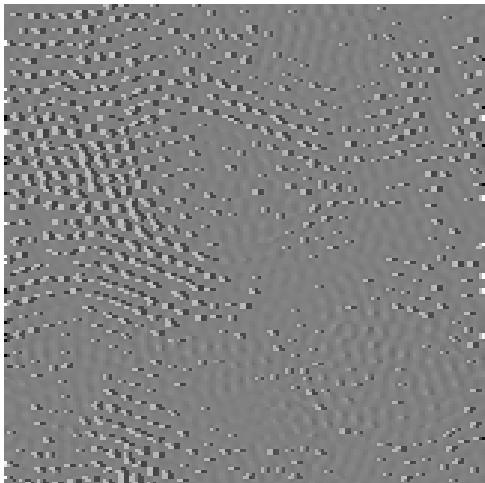
- ✿ **Malaria example:** Matrix  $\mathbf{X}$  has  $n^2 = 6.55 \cdot 10^8$  real dof

Sources: AIM Frames Workshop 2008; Edidin et al. 2009; Chai et al. 2011; Candès et al. 2013; Horstmeyer et al. 2015.

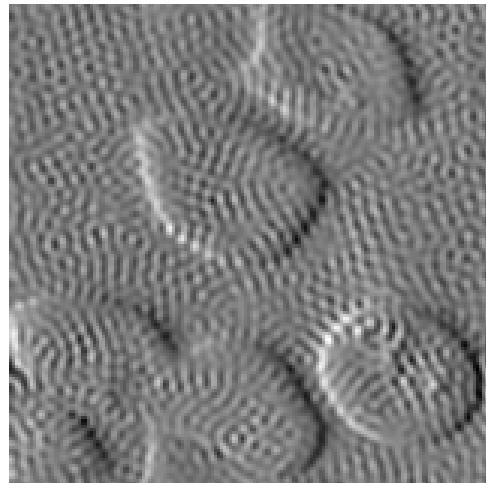
---

# Convexity: Why Bother?

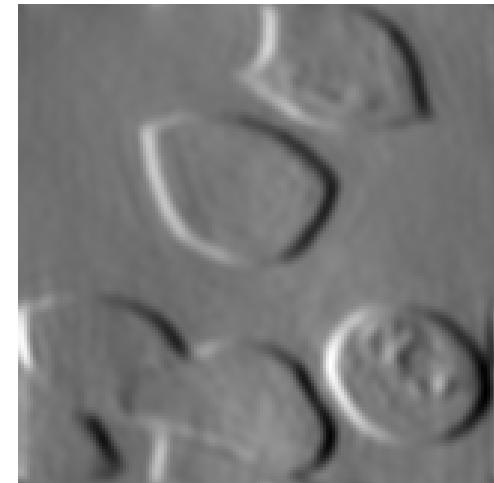
---



Wirtinger Flow  
(not convex)



Burer-Monteiro  
(sort of convex)



???  
(convex)

images of  $x$  phase gradient

**Challenge:** How to solve the convex ptychography problem at scale?

Sources: Burer & Monteiro 2003; Candès et al. 2014; Horstmeyer et al. 2015; Yurtsever et al. 2017.

# Optimization with Optimal Storage

---

# Optimization with Optimal Storage

---

Can we develop algorithms  
that reliably solve an optimization problem  
using **storage** that does not exceed  
the size of the problem data  
or the size of the solution?

---

# Convex Low-Rank Matrix Optimization

---

$$\underset{X \in \mathbb{H}_n}{\text{minimize}} \quad f(\mathcal{A}X) \quad \text{subject to} \quad \text{trace}(X) = \alpha; \quad X \text{ psd}$$

## Details:

- $\mathcal{A} : \mathbb{H}_n \rightarrow \mathbb{R}^d$  is a real-linear map on  $n \times n$  Hermitian matrices
- $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex and differentiable
- In many applications,
  - $\mathcal{A}$  extracts  $d$  linear measurements of  $n \times n$  matrix
  - $f = \text{loss}(\cdot; \mathbf{b})$  for data  $\mathbf{b} \in \mathbb{R}^d$
  - $d \ll n^2$
  - $\alpha$  modulates rank of solution
- Models problems in signal processing, statistics, and machine learning (e.g., convex ptychography)

---

# Optimal Storage

---

What kind of storage bounds can we hope for?

- Assume black-box implementation of operations with linear map:

$$\begin{array}{ll} \mathbf{u} \mapsto \mathcal{A}(\mathbf{u}\mathbf{u}^*) & (\mathbf{u}, \mathbf{z}) \mapsto (\mathcal{A}^*\mathbf{z})\mathbf{u} \\ \mathbb{C}^n \rightarrow \mathbb{R}^d & \mathbb{C}^n \times \mathbb{R}^d \rightarrow \mathbb{C}^n \end{array}$$

- Need  $\Theta(n + d)$  storage for output of black-box operations
- Need  $\Theta(r n)$  storage for rank- $r$  approximate solution of model problem

**Definition.** An algorithm for the model problem has **optimal storage** if its working storage is  $\Theta(d + r n)$  rather than  $\Theta(n^2)$ .

Source: Yurtsever et al. 2017; Cevher et al. 2017.

---

# So Many Algorithms...

---

- ❖ 1990s: **Interior-point methods**
  - ❖ Storage cost  $\Theta(n^4)$  for Hessian
- ❖ 2000s: **Convex first-order methods**
  - ❖ (Accelerated) proximal gradient, spectral bundle methods, and others
  - ❖ Store matrix variable  $\Theta(n^2)$
- ❖ 2008–Present: **Storage-efficient convex first-order methods**
  - ❖ Conditional gradient method (CGM) and extensions
  - ❖ Store matrix in low-rank form  $\mathcal{O}(tn)$ ; no storage guarantees
- ❖ 2009–Present: **Nonconvex heuristics**
  - ❖ Burer–Monteiro factorization idea + various nonlinear programming methods
  - ❖ Store low-rank matrix factors  $\Theta(r n)$
  - ❖ For guaranteed solution, need unrealistic + unverifiable statistical assumptions

Sources: Interior-point: Nemirovski & Nesterov 1994; ... First-order: Rockafellar 1976; Helmberg & Rendl 1997; Auslender & Teboulle 2006; ... CGM: Frank & Wolfe 1956; Levitin & Poljak 1967; Jaggi 2013; ... Heuristics: Burer & Monteiro 2003; Keshavan et al. 2009; Jain et al. 2012; Candès et al. 2014; Bhojanapalli et al. 2015; Boumal et al. 2016; ....

---

# The Challenge

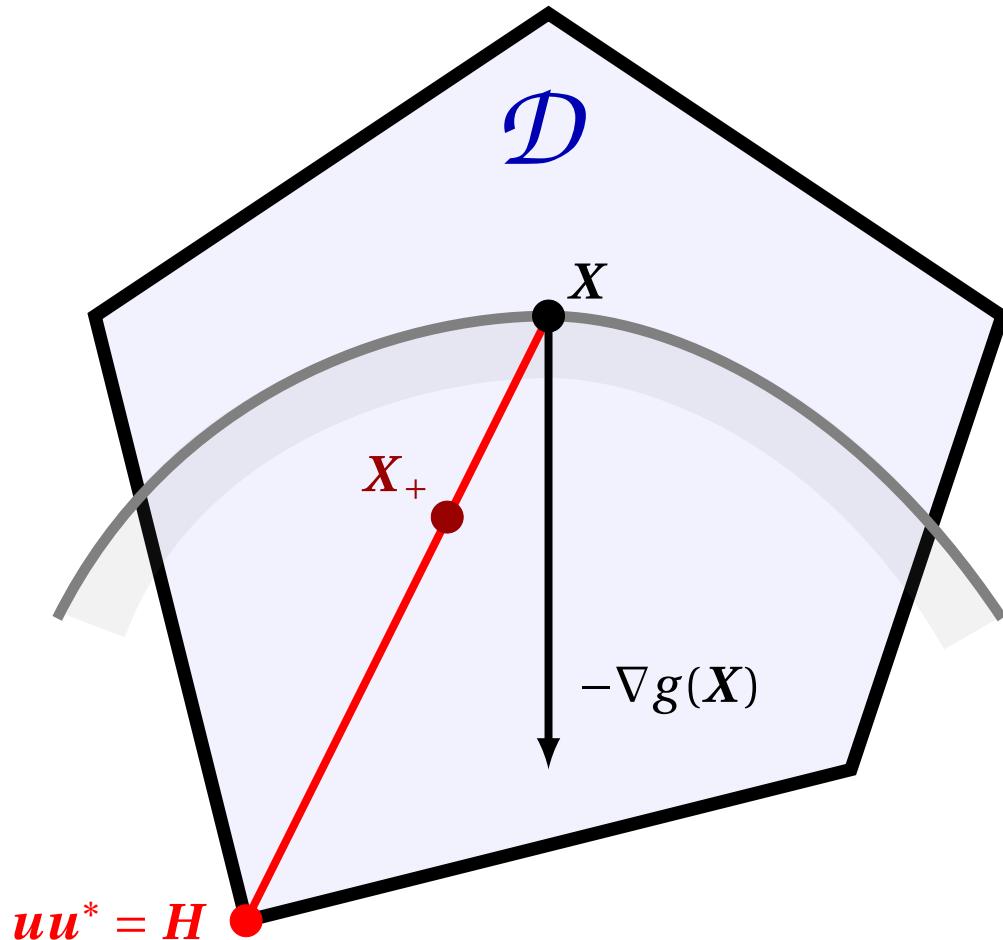
---

- Some algorithms provably solve the model problem...
- Some algorithms have optimal storage guarantees...

Is there an **algorithm**  
that provably computes  
a **low-rank approximation**  
to a solution of the model problem  
+ has **optimal storage** guarantees?

# Conditional Gradient Method

# Geometry of CGM



$$H = \arg \max_{Y \in \mathcal{D}} \langle Y, -\nabla g(X) \rangle$$

$$X_+ = (1 - \eta)X + \eta H$$

$$\{Y : g(Y) \leq g(X)\}$$

$$\min_{X \in \mathcal{D}} g(X)$$

$$\mathcal{D} = \{Y \text{ psd} : \text{trace}(Y) = 1\}$$

---

# CGM for the Model Problem

---

**Input:** Problem data; suboptimality  $\varepsilon$

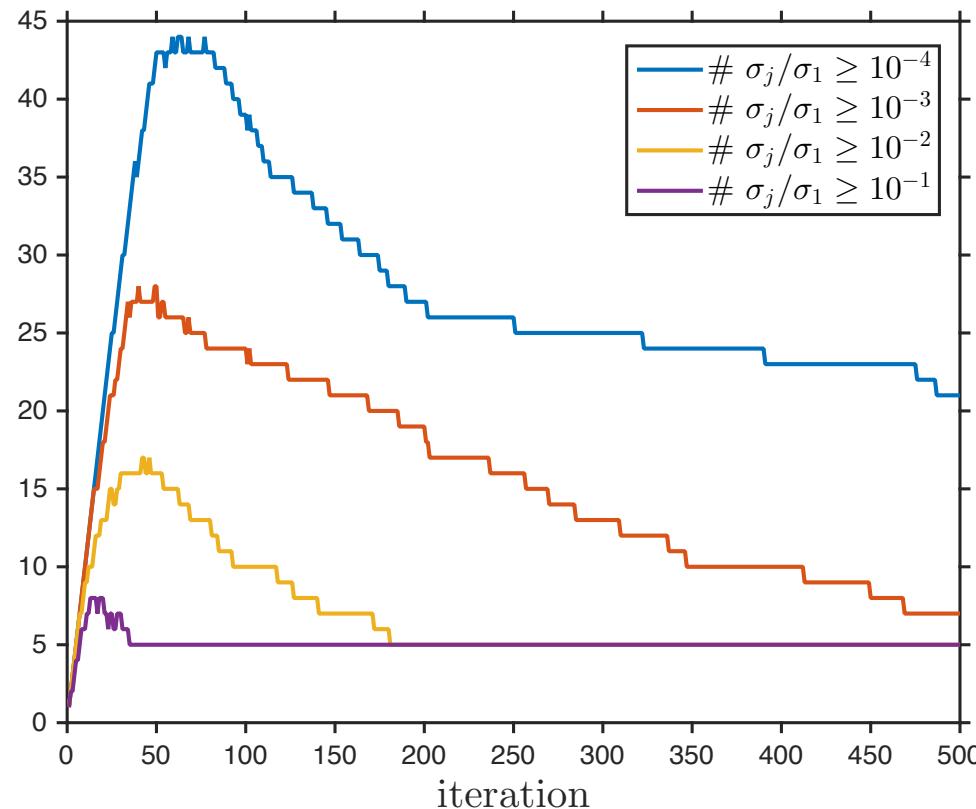
**Output:** Approximate solution  $X_{\text{cgm}}$

```
1  function CGM
2       $X \leftarrow \mathbf{0}$                                  $\triangleright$  Initialize variable
3      for  $t \leftarrow 0, 1, 2, 3, \dots$  do
4           $u \leftarrow \text{MinEigVec}(\mathcal{A}^*(\nabla f(\mathcal{A}X)))$      $\triangleright$  Lanczos!
5           $H \leftarrow -\alpha uu^*$                                  $\triangleright$  Form update direction
6          if  $\langle X - H, \mathcal{A}^*(\nabla f(\mathcal{A}X)) \rangle \leq \varepsilon$ 
7              then break for                                 $\triangleright$  Stop when  $\varepsilon$ -suboptimal
8           $\eta \leftarrow 2/(t + 2)$                              $\triangleright$  Update learning rate
9           $X \leftarrow (1 - \eta)X + \eta H$                      $\triangleright$  Update variable
10     return  $X$ 
```

**Comment:** In notation of last slide,  $g = f \circ \mathcal{A}$ . The gradient  $\nabla g = \mathcal{A}^* \circ \nabla f \circ \mathcal{A}$ .

**Sources:** Frank & Wolfe 1956; Levitin & Poljak 1967; Hazan 2008; Clarkson 2010; Jaggi 2013.

# Evolution of $\varepsilon$ -Rank of CGM Iterates



Comments: Malaria data, quadratic loss,  $\alpha = 1,400$ .

# SketchyCGM

---

# Crisis / Opportunity

---

## Crisis:

- ❖ CGM needs many iterations to converge to a near-low-rank solution
- ❖ The  $\varepsilon$ -rank of the CGM iterates can increase without bound
- ❖ CGM requires high + unpredictable storage
- ❖ Typically involves dynamic memory allocation

## Opportunity:

- ❖ **Modify CGM to work with optimal storage!**
- ❖ Drive the CGM iteration with small “dual” variable  $z = \mathcal{A}^T X$
- ❖ Maintain small randomized sketch of primal matrix variable  $X$
- ❖ After iteration terminates, reconstruct matrix variable  $X$  from sketch

Source: Yurtsever et al. 2017.

---

# CGM, Redux

---

**Input:** Problem data; suboptimality  $\varepsilon$

**Output:** Approximate solution  $X_{\text{cgm}}$

```
1  function CGM
2       $X \leftarrow \mathbf{0}_{n \times n}$ 
3      for  $t \leftarrow 0, 1, 2, 3, \dots$  do
4           $u \leftarrow \text{MinEigVec}(\mathcal{A}^*(\nabla f(\mathcal{A}X)))$ 
5           $H \leftarrow -\alpha uu^*$ 
6          if  $\langle \mathcal{A}(X - H), \nabla f(\mathcal{A}X) \rangle \leq \varepsilon$ 
7              then break for
8           $\eta \leftarrow 2/(t + 2)$ 
9           $X \leftarrow (1 - \eta)X + \eta H$ 
10     return  $X$ 
```

**Idea:** Apply  $\mathcal{A}$  to all expressions involving  $X$  and  $H$ ...

---

# A Dual Formulation of CGM

---

**Input:** Problem data; suboptimality  $\varepsilon$

**Output:** Approximate dual solution  $\mathcal{A}\mathbf{X}_{\text{cgm}}$

```
1  function DUALCGM
2       $\mathcal{A}\mathbf{X} \leftarrow \mathcal{A}\mathbf{0}_{n \times n}$ 
3      for  $t \leftarrow 0, 1, 2, 3, \dots$  do
4           $\mathbf{u} \leftarrow \text{MinEigVec}(\mathcal{A}^*(\nabla f(\mathcal{A}\mathbf{X})))$ 
5           $\mathcal{A}\mathbf{H} \leftarrow \mathcal{A}(-\alpha \mathbf{u} \mathbf{u}^*)$ 
6          if  $\langle \mathcal{A}(\mathbf{X} - \mathbf{H}), \nabla f(\mathcal{A}\mathbf{X}) \rangle \leq \varepsilon$ 
7              then break for
8           $\eta \leftarrow 2/(t + 2)$ 
9           $\mathcal{A}\mathbf{X} \leftarrow (1 - \eta)\mathcal{A}\mathbf{X} + \eta\mathcal{A}\mathbf{H}$ 
10     return  $\mathcal{A}\mathbf{X}$ 
```

**Idea:** Change variables  $\mathbf{z} = \mathcal{A}\mathbf{X}$  and  $\mathbf{h} = \mathcal{A}\mathbf{H}$ ...

---

# A Dual Formulation of CGM

---

**Input:** Problem data; suboptimality  $\varepsilon$

**Output:** Approximate dual solution  $\mathbf{z}_{\text{cgm}}$

```
1  function DUALCGM
2       $\mathbf{z} \leftarrow \mathbf{0}_d$ 
3      for  $t \leftarrow 0, 1, 2, 3, \dots$  do
4           $\mathbf{u} \leftarrow \text{MinEigVec}(\mathcal{A}^*(\nabla f(\mathbf{z})))$ 
5           $\mathbf{h} \leftarrow \mathcal{A}(-\alpha \mathbf{u} \mathbf{u}^*)$ 
6          if  $\langle \mathbf{z} - \mathbf{h}, \nabla f(\mathbf{z}) \rangle \leq \varepsilon$ 
7              then break for
8           $\eta \leftarrow 2/(t + 2)$ 
9           $\mathbf{z} \leftarrow (1 - \eta)\mathbf{z} + \eta\mathbf{h}$ 
10     return  $\mathbf{z}$ 
```

---

# A Dual Formulation of CGM

---

**Input:** Problem data; suboptimality  $\varepsilon$

**Output:** Approximate dual solution  $z_{\text{cgm}}$

```
1  function DUALCGM
2       $z \leftarrow \mathbf{0}_d$                                  $\triangleright$  Initialize dual variable
3      for  $t \leftarrow 0, 1, 2, 3, \dots$  do
4           $\mathbf{u} \leftarrow \text{MinEigVec}(\mathcal{A}^*(\nabla f(z)))$          $\triangleright$  Lanczos!
5           $\mathbf{h} \leftarrow \mathcal{A}(-\alpha \mathbf{u} \mathbf{u}^*)$                            $\triangleright$  Form dual update direction
6          if  $\langle z - \mathbf{h}, \nabla f(z) \rangle \leq \varepsilon$ 
7              then break for
8           $\eta \leftarrow 2/(t + 2)$ 
9           $z \leftarrow (1 - \eta)z + \eta \mathbf{h}$                                  $\triangleright$  Update dual variable
10     return  $z$ 
```

**Benefit:** Only uses storage  $\Theta(n + d)$ !

**Problem:** Where do we get  $X_{\text{cgm}}$ ?

---

# Sketching the Decision Variable

---

- Idea: Maintain small sketch of primal variable  $\mathbf{X}$ !
- Fix target rank  $r$  of solution
- Draw Gaussian dimension reduction map

$$\boldsymbol{\Omega} \in \mathbb{C}^{n \times k} \quad \text{where } k = 2r$$

- Sketch takes the form

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\Omega} \in \mathbb{C}^{n \times k}$$

- Can perform linear update  $\mathbf{X} \leftarrow (1 - \eta)\mathbf{X} + \eta\mathbf{H}$  by operating on sketch
- Can compute provably good rank- $r$  approximation  $\hat{\mathbf{X}}$  from sketch
- Only needs additional storage  $\Theta(rn)$ !

Sources: Woolfe et al. 2008; Clarkson & Woodruff 2009; Halko et al. 2009; Gittens 2011, 2013; Woodruff 2014; Cohen et al. 2015; Boutsidis et al. 2015; Tropp et al. 2016, 2017; ....

---

# SketchyCGM for the Model Problem

---

**Input:** Problem data; suboptimality  $\varepsilon$ ; target rank  $r$

**Output:** Rank- $r$  approximate solution  $\hat{X} = V\Lambda V^*$  in factored form

```
1  function SKETCHYCGM
2      SKETCH.INIT( $n, r$ )                                ▷ Initialize sketch to zero
3       $z \leftarrow \mathbf{0}$ 
4      for  $t \leftarrow 0, 1, 2, 3, \dots$  do
5           $\mathbf{u} \leftarrow \text{MinEigVec}(\mathcal{A}^*(\nabla f(z)))$ 
6           $\mathbf{h} \leftarrow \mathcal{A}(-\alpha \mathbf{u} \mathbf{u}^*)$ 
7          if  $\langle z - \mathbf{h}, \nabla f(z) \rangle \leq \varepsilon$  then break for
8           $\eta \leftarrow 2/(t + 2)$ 
9           $z \leftarrow (1 - \eta)z + \eta \mathbf{h}$ 
10         SKETCH.CGMUPDATE( $-\sqrt{\alpha} \mathbf{u}, \eta$ )           ▷ Update sketch of  $X$ 
11          $(V, \Lambda) \leftarrow \text{SKETCH.RECONSTRUCT}()$        ▷ Approx. eigendecomps of  $X$ 
12         return  $(V, \Lambda)$ 
```

Source: Yurtsever et al. 2017.

---

# Methods for SKETCH Object

---

```
1 function SKETCH.INIT( $n, r$ )                                ▷ Rank- $r$  approx of  $n \times n$  psd matrix
2    $k \leftarrow 2r$ 
3    $\Omega \leftarrow \text{randn}(\mathbb{C}, n, k)$ 
4    $Y \leftarrow \text{zeros}(n, k)$ 

5 function SKETCH.CGMUPDATE( $s, \theta$ )
6    $Y \leftarrow (1 - \theta)Y + \theta s(s^* \Omega)$                       ▷ Average  $ss^*$  into sketch

7 function SKETCH.RECONSTRUCT()
8    $C \leftarrow \text{chol}(\Omega^* Y)$                                      ▷ Cholesky decomposition
9    $Z \leftarrow Y/C$                                                  ▷ Solve least-squares problems
10   $(U, \Sigma, \sim) \leftarrow \text{svds}(Z, r)$                          ▷ Compute  $r$ -truncated SVD
11  return ( $U, \Sigma^2$ )                                         ▷ Return eigenvalue factorization
```

Sources: Yurtsever et al. 2017; Tropp et al. 2017.

---

## Less Filling / Great Taste

---

**Theorem 1** (YUTC 2016). *SKETCHYCGM has the following properties:*

- *SKETCHYCGM has optimal storage guarantee  $\Theta(d + r n)$*
- *SKETCHYCGM produces an  $\varepsilon$ -suboptimal objective value after  $O(\varepsilon^{-1})$  iterations*
- *Suppose CGM produces iterates  $X_t$  that converge to  $X_{\text{cgm}}$ . Then SKETCHYCGM produces rank- $r$  iterates  $\hat{X}_t$  that satisfy*

$$\limsup_{t \rightarrow \infty} \mathbb{E} \left\| \hat{X}_t - X_{\text{cgm}} \right\|_{S_1} \leq \text{const} \cdot \left\| X_{\text{cgm}} - [X_{\text{cgm}}]_r \right\|_{S_1}$$

*In particular, if  $\text{rank}(X_{\text{cgm}}) \leq r$ , then  $\mathbb{E} \left\| \hat{X}_t - X_{\text{cgm}} \right\|_{S_1} \rightarrow 0$*

Source: “Everything you always wanted in an algorithm. And less.”

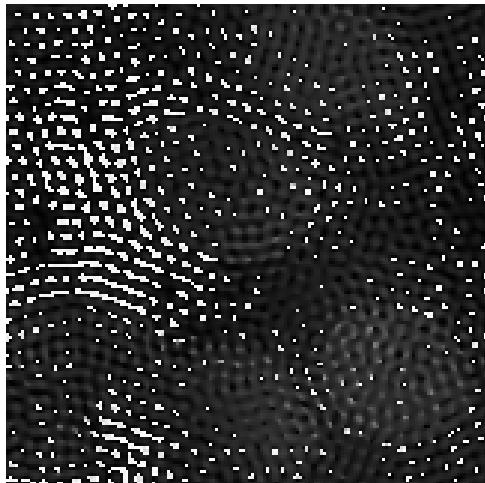
<https://www.youtube.com/watch?v=0agZEMEpiVI>.

# Performance of SketchyCGM

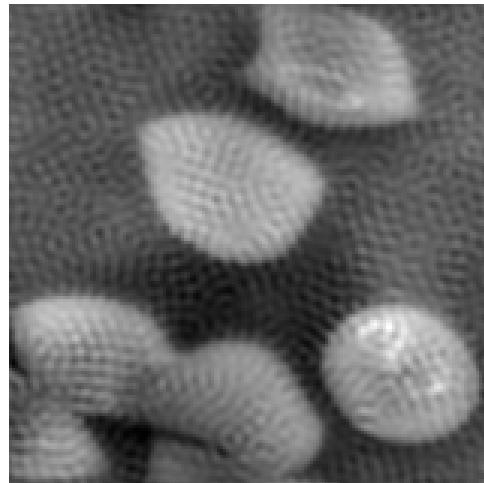
---

# Fourier Ptychography, Redux

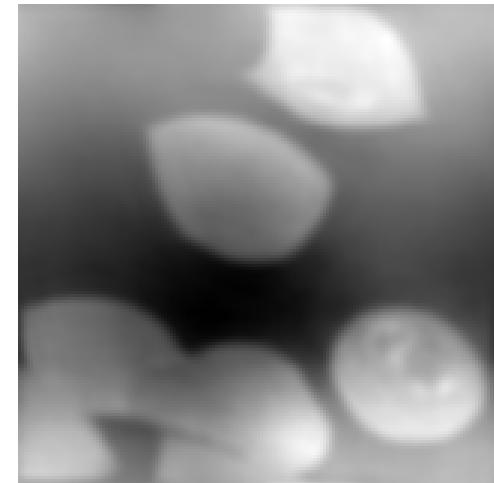
---



Wirtinger Flow



Burer-Monteiro



SKETCHYCGM

29 illuminations;  $80 \times 80$  pixels each;  $d = 1.86 \cdot 10^5$  measurements

image size  $n = 160 \times 160$  pixels; matrix size  $n^2 = 6.55 \cdot 10^8$

SKETCHYCGM storage (rank  $r = 1$ ):  $6.53 \cdot 10^5$   
quadratic loss

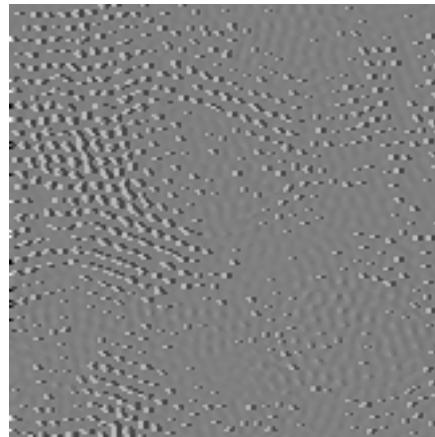
Sources: Burer & Monteiro 2003; Candès et al. 2014; Horstmeyer et al. 2015; Yurtsever et al. 2017.

---

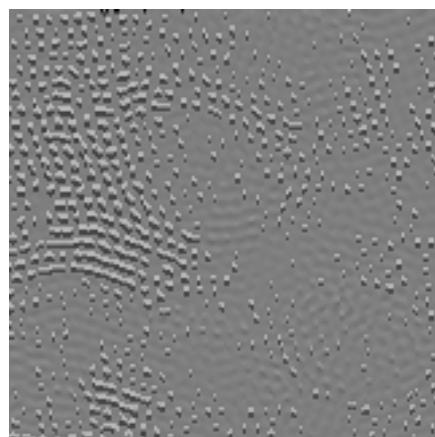
# Fourier Ptychography: Malaria Phase Gradients

---

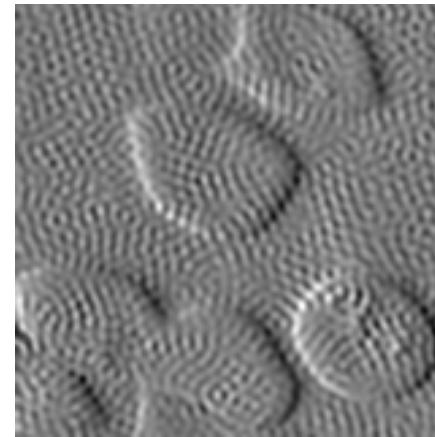
$\Delta_x$



$\Delta_y$



Wirtinger Flow



Burer-Monteiro



SKETCHYCGM

---

# Convex Low-Rank Matrix Completion

---

- Suppose  $\mathbf{X}_\natural \in \mathbb{R}^{m \times n}$  is a rank- $r$  matrix
- Observe a subset of entries + noise:

$$b_{ij} = (\mathbf{X}_\natural)_{ij} + \xi_{ij} \quad \text{for } (i, j) \in E$$

- Matrix completion via convex programming:

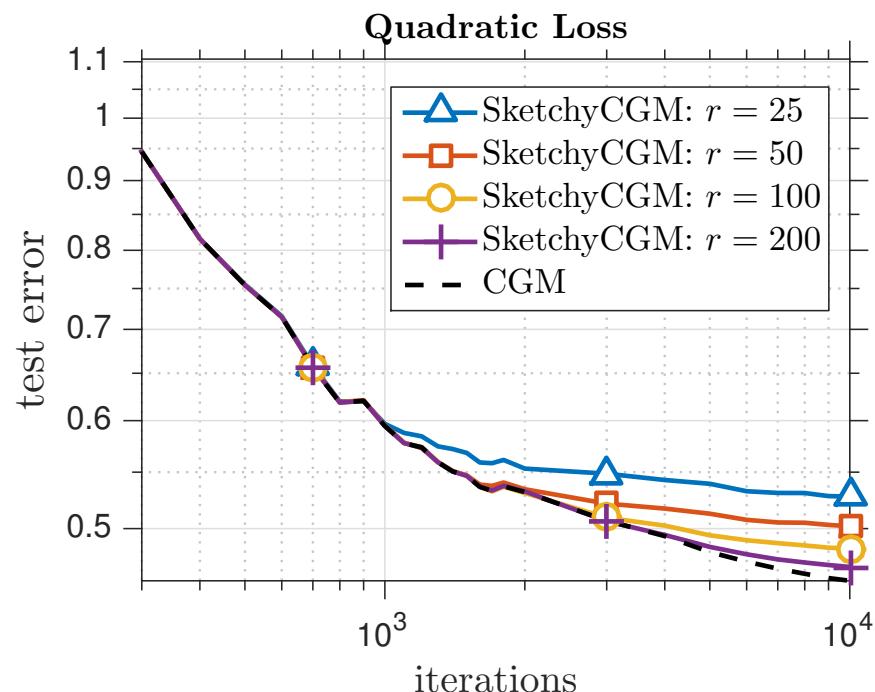
$$\underset{\mathbf{X} \in \mathbb{R}^{m \times n}}{\text{minimize}} \quad \sum_{(i, j) \in E} \text{loss}(x_{ij}; b_{ij}) \quad \text{subject to} \quad \|\mathbf{X}\|_{S_1} \leq \alpha$$

- Matrix  $\mathbf{X}_\natural$  has about  $r(m + n)$  degrees of freedom
- Convex method often effective when  $\#E = \Theta(r(m + n))$
- But decision variable  $\mathbf{X}$  has  $mn$  degrees of freedom!
- SKETCHYCGM works with storage  $\Theta(\#E + r(m + n))$

Sources: Srebro et al. 2004; Candès & Recht 2010; Yurtsever et al. 2017.

# MovieLens 10M

•  $m = 71,567$  users,  $n = 10,681$  movies,  $d = 10^7$  ratings, dim.  $mn = 7.64 \cdot 10^8$



Approximate storage costs

Rank ( $r$ )	SKETCHYCGM
25	$3.28 \cdot 10^7$
50	$4.51 \cdot 10^7$
100	$6.98 \cdot 10^7$
200	$1.19 \cdot 10^8$

Source: Harper & Konstan 2015; Yurtsever et al. 2017.

# Denouement

---

## Beyond...

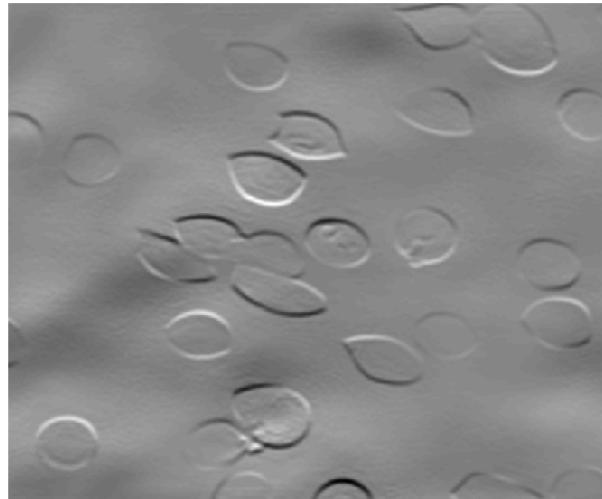
---

- ✿ Other low-rank matrix optimization problems!
- ✿ More effective algorithms!
- ✿ Many applications!
- ✿ Beyond low-rank matrix optimization!

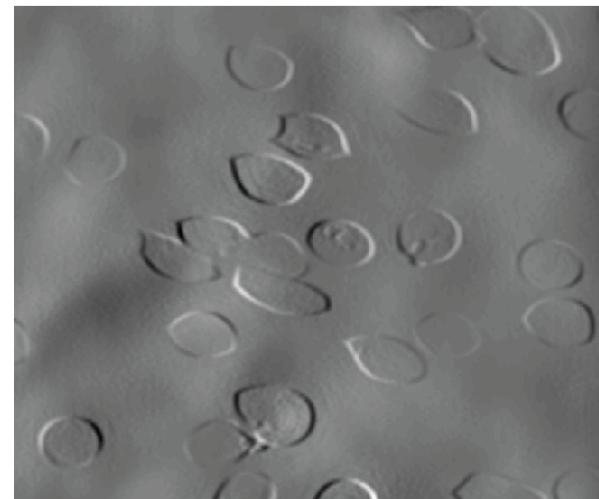
---

# A Large-Scale Problem

---



$x$  phase gradient



$y$  phase gradient

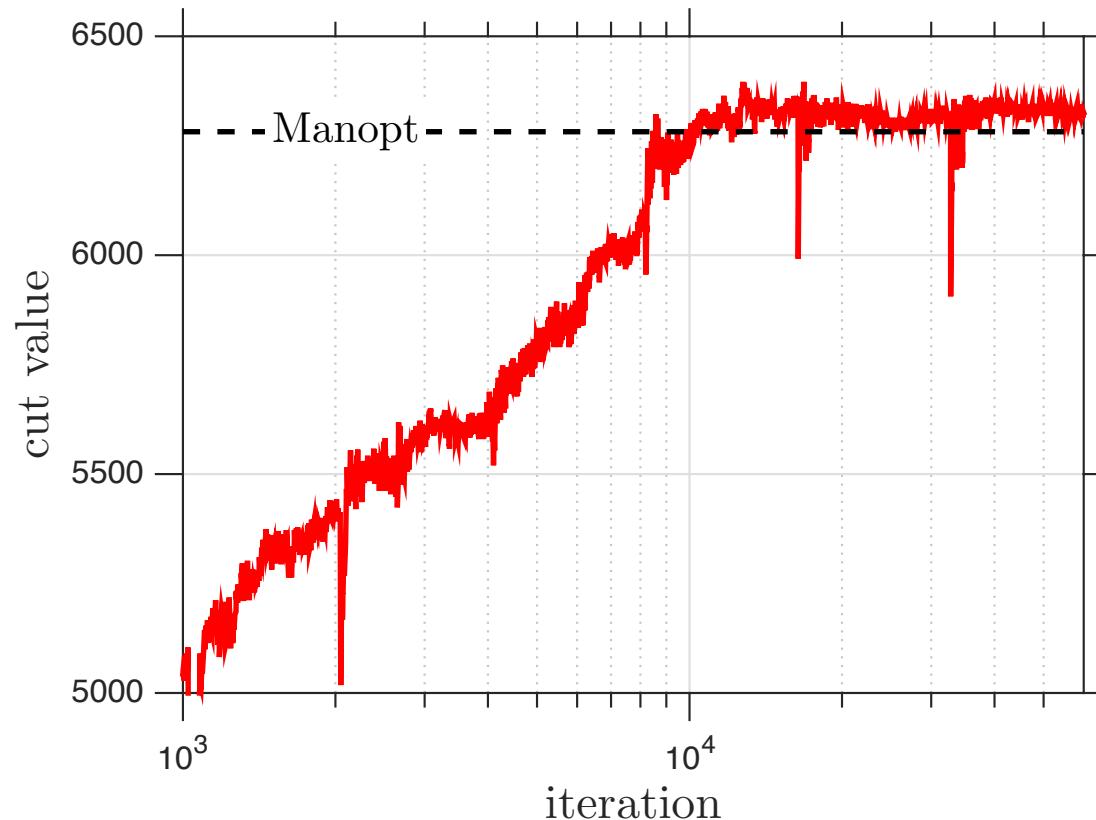
29 illuminations;  $250 \times 250$  pixels each;  $d = 1.81 \cdot 10^6$  measurements  
image size  $n = 501 \times 501$  pixels; matrix size  $n^2 = 6.25 \cdot 10^{10}$   
solution via **SKETCHYUPD**

Sources: Cevher et al. 2017.

---

## The MAXCUT SDP

---



MAXCUT SDP; sparse graph with 10,000 nodes (G67); via SKETCHYUPD

Sources: Goemans & Williamson 1995; Boumal 2015; Cevher et al. 2017.

---

## To learn more...

---

**E-mail:** [jtropp@cms.caltech.edu](mailto:jtropp@cms.caltech.edu)

**Web:** <http://users.cms.caltech.edu/~jtropp>

**Papers:**

- Halko, Martinsson, & Tropp, “[Finding structure with randomness: Probabilistic algorithms for computing approximate matrix decompositions](#),” *SIAM Review*, 2011
- Horstmeyer et al. “[Solving ptychography with a convex relaxation](#),” *New J. Physics*, 2015
- Tropp, Yurtsever, Udell, & Cevher, “[Randomized single-view algorithms for low-rank matrix approximation](#),” Caltech ACM TR 2017-01, arXiv cs.NA 1609.00048, [updates coming soon!](#)
- Yurtsever, Udell, Tropp, & Cevher, “[Sketchy decisions: Convex low-rank matrix optimization with optimal storage](#),” AISTATS 2017, arXiv ma.OC arXiv:1702.06838
- **More to come!**