# Sketchy Decisions

❦

## Joel A. Tropp

Steele Family Professor
of Applied & Computational Mathematics

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)

# 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:

- $\mathcal{A} : \mathbb{H}_n \to \mathbb{R}^d$ is a real-linear map on $n \times n$ Hermitian matrices
- $f : \mathbb{R}^d \to \mathbb{R}$ is convex and differentiable

- In many applications,
  - $\mathcal{A}$ extracts $d$ linear measurements of $n \times n$ matrix
  - $f = \text{loss}(\,\cdot\,; b)$ for data $b \in \mathbb{R}^d$
  - $d \ll n^2$
  - $\alpha$ modulates rank of solution

- Models problems in signal processing, statistics, and machine learning

# Optimal Storage

## What kind of storage bounds can we hope for?

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

$$\boldsymbol{u} \mapsto \mathcal{A}(\boldsymbol{uu}^*) \qquad\qquad (\boldsymbol{u}, \boldsymbol{z}) \mapsto (\mathcal{A}^*\boldsymbol{z})\boldsymbol{u}$$

$$\mathbb{C}^n \to \mathbb{R}^d \qquad\qquad \mathbb{C}^n \times \mathbb{R}^d \to \mathbb{C}^n$$

- Need $\Theta(n+d)$ storage for output of black-box operations

- Need $\Theta(rn)$ storage for rank-$r$ approximate solution of model problem

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

# 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 $O(tn)$; no storage guarantees

- 2009–Present: **Nonconvex heuristics**
  - Burer–Monteiro factorization idea + various nonlinear programming methods
  - Store low-rank matrix factors $\Theta(rn)$
  - 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; ....
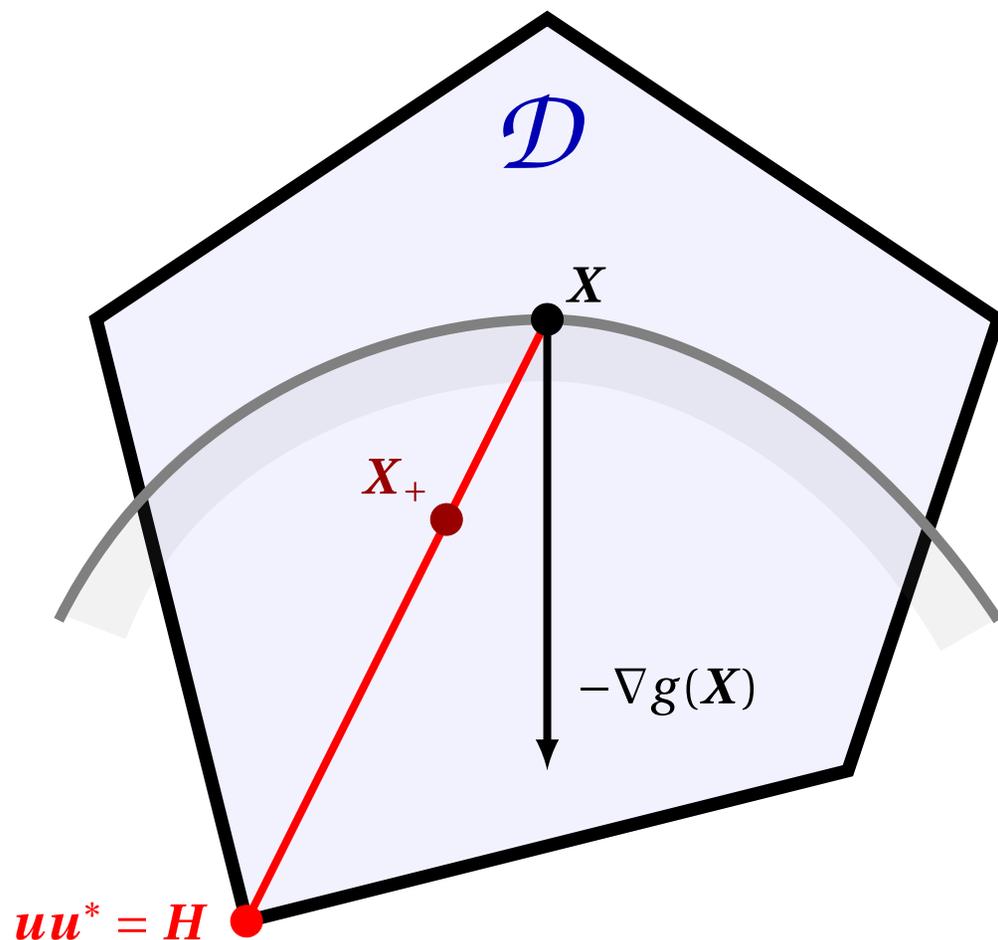
# The Challenge

❧ 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?

# SketchyCGM

# 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 $X_{\mathrm{cgm}}$

1  **function** CGM
2      $X \leftarrow \mathbf{0}$                                                            ▷ Initialize variable
3      **for** $t \leftarrow 0, 1, 2, 3, \ldots$ **do**
4          $u \leftarrow \mathtt{MinEigVec}(\mathcal{A}^*(\nabla f(\mathcal{A}X)))$                     ▷ Lanczos!
5          $H \leftarrow -\alpha u u^*$                                          ▷ Form update direction
6          **if** $\langle X - H, \, \mathcal{A}^*(\nabla f(\mathcal{A}X)) \rangle \leq \varepsilon$
7                  **then break for**                              ▷ Stop when $\varepsilon$-suboptimal
8          $\eta \leftarrow 2/(t+2)$                                          ▷ Update step size
9          $X \leftarrow (1-\eta)X + \eta H$                                 ▷ 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.

# 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 $z = \mathcal{A}X$
- Maintain small randomized sketch of primal matrix variable $X$
- After iteration terminates, reconstruct matrix variable $X$ from sketch

**Source:** Yurtsever et al. 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$(n, r)$               $\triangleright$ Initialize SKETCH to zero

3      $z \leftarrow 0$

4      **for** $t \leftarrow 0, 1, 2, 3, \dots$ **do**

5        $u \leftarrow \text{MinEigVec}(\mathcal{A}^*(\nabla f(z)))$          $\triangleright$ Lanczos!

6        $h \leftarrow \mathcal{A}(-\alpha u u^*)$

7        **if** $\langle z - h, \nabla f(z) \rangle \leq \varepsilon$ **then break for**

8        $\eta \leftarrow 2/(t+2)$

9        $z \leftarrow (1-\eta)z + \eta h$

10        SKETCH.CGMUPDATE$(-\sqrt{\alpha}\, u, \eta)$      $\triangleright$ Update sketch of $X$

11      $(V, \Lambda) \leftarrow$ SKETCH.RECONSTRUCT$()$      $\triangleright$ Approx. eigendecomp 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      $k \leftarrow 2r$                                                        ▷ Increase $k$ for better quality

3      $\mathbf{\Omega} \leftarrow \mathsf{randn}(\mathbb{C}, n, k)$

4      $Y \leftarrow \mathsf{zeros}(n, k)$

5  **function** SKETCH.CGMUPDATE($s$, $\theta$)

6      $Y \leftarrow (1 - \theta)\,Y + \theta\,s(s^*\mathbf{\Omega})$            ▷ Average $ss^*$ into sketch

7  **function** SKETCH.RECONSTRUCT( )

8      $C \leftarrow \mathsf{chol}(\mathbf{\Omega}^* Y)$                         ▷ Cholesky decomposition

9      $Z \leftarrow Y / C$                                                     ▷ Solve least-squares problems

10     $(U, \mathbf{\Sigma}, \sim) \leftarrow \mathsf{svds}(Z, r)$              ▷ Compute $r$-truncated SVD

11     **return** $(U, \mathbf{\Sigma}^2)$                                      ▷ Return eigenvalue factorization

**Sources:** Williams & Seeger 2001; Drineas & Mahoney 2005; Gittens 2011, 2013; Pourkamali-Anaraki & Becker 2016; Wang, Gittens, & Mahoney 2017; Tropp et al. 2017.

# Less Filling / Great Taste

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

- *SKETCHYCGM has optimal storage guarantee $\Theta(d + rn)$*

- *SKETCHYCGM produces an $\varepsilon$-suboptimal objective value after $O(\varepsilon^{-1})$ iterations*

- *Suppose CGM produces iterates $X_t$ that converge to a rank-$r$ matrix $X_{\mathrm{cgm}}$. Then SKETCHYCGM produces rank-$r$ iterates $\hat{X}_t$ that satisfy*
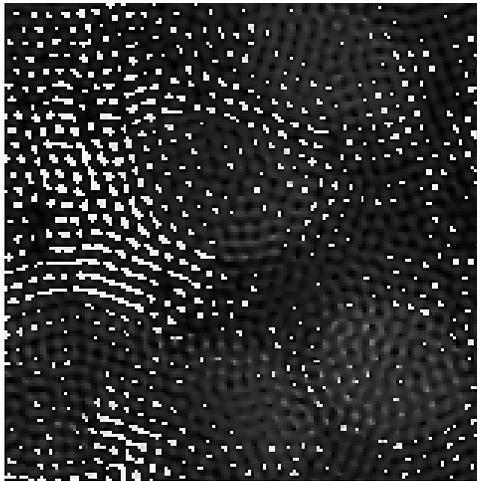
$$\mathbb{E} \left\| \hat{X}_t - X_{\mathrm{cgm}} \right\|_{S_1} \to 0.$$

**Source:** "Everything you always wanted in an algorithm. And less."
https://www.youtube.com/watch?v=0agZEMEpiVI.

# Performance of SketchyCGM

# Fourier Ptychography



Wirtinger Flow   Burer–Monteiro   SKETCHYCGM

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

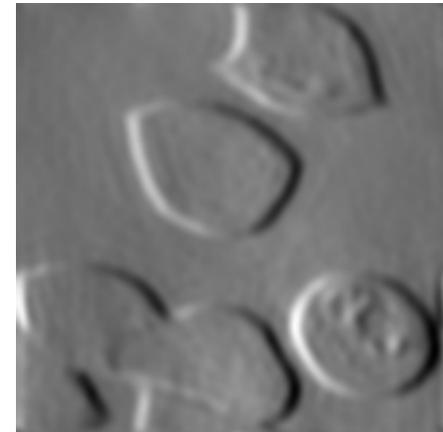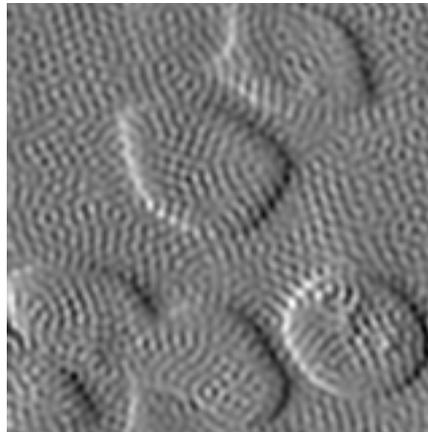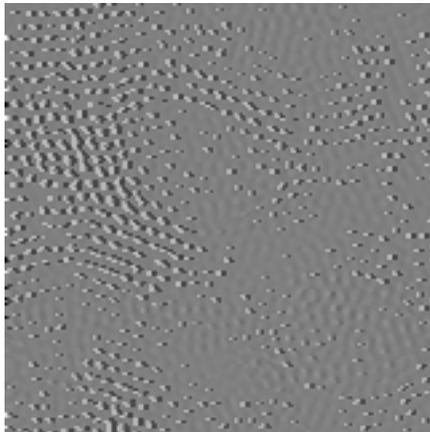image size $n = 160 \times 160$ pixels; matrix size $n^2 = 6.55 \cdot 10^8$

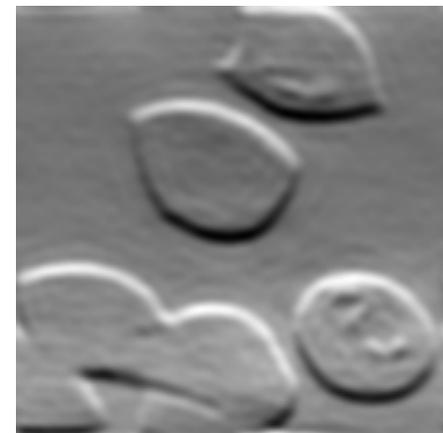SKETCHYCGM storage (rank $r = 1$): $6.53 \cdot 10^5$
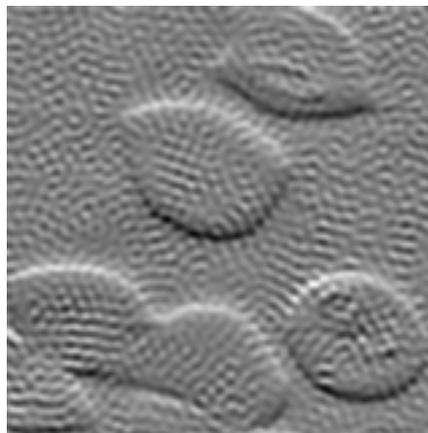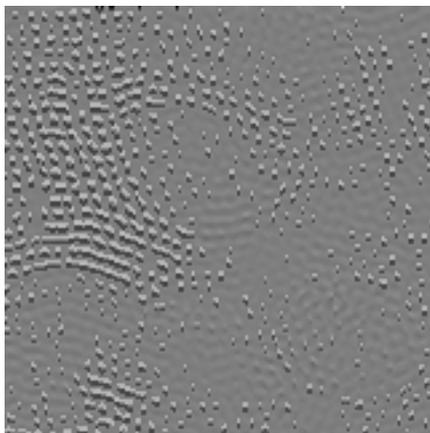
quadratic loss

**Sources:** Burer & Monteiro 2003; Balan et al. 2008; Chai et al. 2011; Zheng, Horstmeyer, & Yang 2013; Horstmeyer & Yang 2014; Candès et al. 2014; Horstmeyer et al. 2015; Yurtsever et al. 2017.
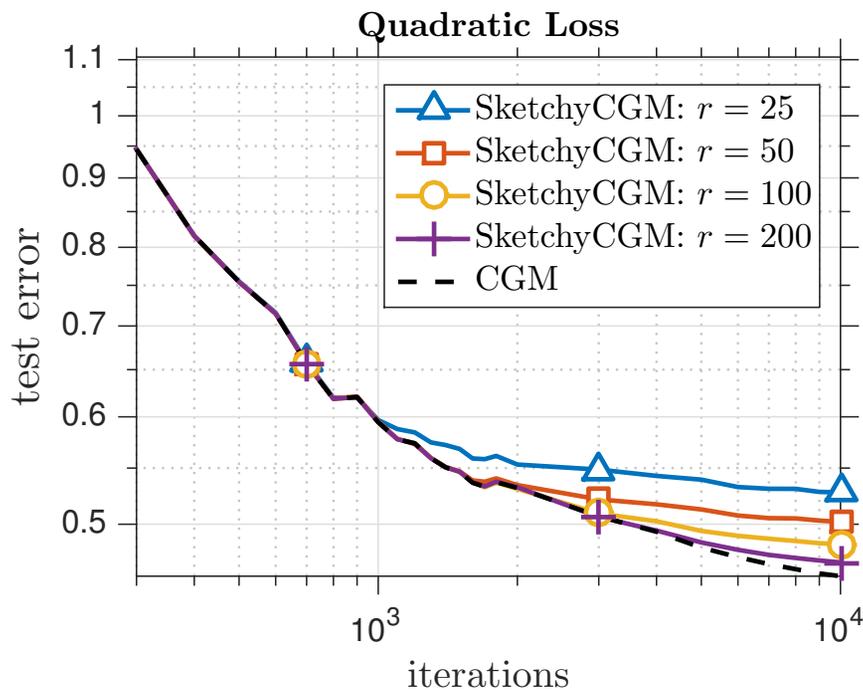
# Fourier Ptychography: *Malaria Phase Gradients*



$\Delta_x$

$\Delta_y$

Wirtinger Flow     Burer–Monteiro     SKETCHYCGM

# 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.

# To learn more...

**E-mail:** `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, "Fixed-rank approximation of a positive-semidefinite matrix from streaming data," NIPS, 2017
- Tropp, Yurtsever, Udell, & Cevher, "Practical sketching algorithms for low-rank matrix approximation," *SIMAX*, 2017
- Tropp, Yurtsever, Udell, & Cevher, "More practical sketching algorithms for low-rank matrix approximation," soon!
- Yurtsever, Udell, Tropp, & Cevher, "Sketchy decisions: Convex low-rank matrix optimization with optimal storage," AISTATS, 2017
- Cevher, Tropp, & Yurtsever, "Storage-optimal algorithms for semidefinite programming," eventually!