

Pursuit of Low-dimensional Structures in High-dimensional (Visual) Data

Yi Ma

**School of Information Science & Technology
ShanghaiTech University, China**



CONTEXT – *Data increasingly massive, high-dimensional...*



Images

↓ ➤ **1M pixels**

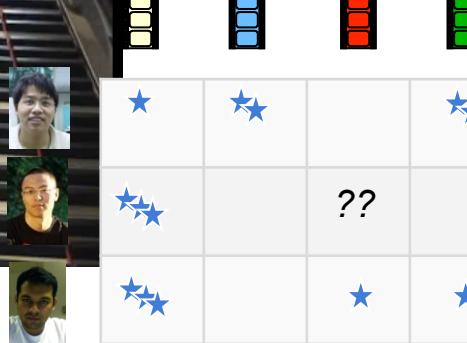
Compression
De-noising
Super-resolution
Recognition...



Videos

↓ ➤ **1B voxels**

Streaming
Tracking
Stabilization...



User data

↓ ➤ **1B users**

Clustering
Classification
Collaborative filtering...

U.S. COMMERCE'S ORTNER SAYS YEN UNDervalued

Commerce Dept. undersecretary of economic affairs Robert Ortner said that he believed the dollar at current levels was fairly priced against most European currencies.

In a wide ranging address sponsored by the Export-Import Bank, Ortner, the bank's senior economist also said he believed that the yen was undervalued and could go up by 10 or 15 pct.

"I do not regard the dollar as undervalued at this point against the yen," he said.

On the other hand, Ortner said that he thought that "the yen is still a little bit undervalued," and "could go up another 10 or 15 pct."

In addition, Ortner, who said he was speaking personally, said he thought that the dollar against most European currencies was "fairly priced."

Ortner said his analysis of the various exchange rate values was based on such economic particulars as wage rate differentials.

Ortner said there had been little impact on U.S. trade deficit by the decline of the dollar because at the time of the Plaza Accord, the dollar was extremely overvalued and that the first 15 pct decline had little impact.

He said there were indications now that the trade deficit was beginning to level off.

Turning to Brazil and Mexico, Ortner made it clear that it would be almost impossible for those countries to earn enough foreign exchange to pay the service on their debts. He said the best way to deal with this was to use the policies outlined in Treasury Secretary James Baker's debt initiative.

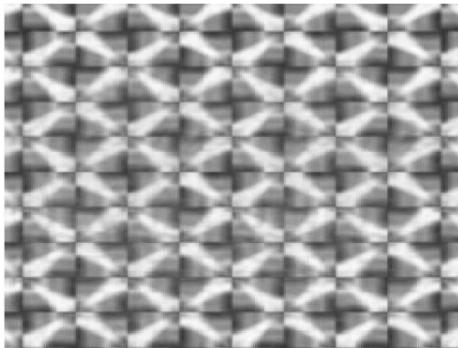
Web data

↓ ➤ **100B webpages**

Indexing
Ranking
Search...

How to extract *low-dim structures* from such *high-dim data?*

CONTEXT – *Low dimensional structures in visual data*



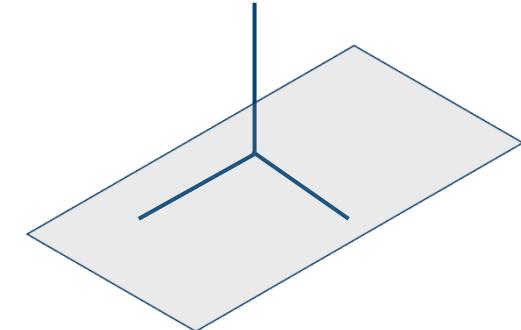
(71)) which turns out in the end to be mathematically equivalent to maximum entropy. The problem is interesting also in that we can see a continuous gradation from decision problems so simple that common sense tells us the answer instantly, with no need for any mathematical theory, through problems more and more involved so that common sense becomes more and more difficultly in making a decision, until finally we reach a point where nobody has yet claimed to be able to see the right decision intuitively, and we require the mathematics to tell us what to do.

Finally, the widget problem turns out to be very close to an important real problem faced by oil prospectors. The details of the real problem are shrouded in proprietary caution; but without giving away any secrets to report that, a few years ago, the writer spent a week at research laboratories of one of our large oil companies, lecturing for over 20 hours on the widget problem. We went through every part of the calculation in excruciating detail – in a room full of engineers armed with calculators, checking up on every stage of the numerical work.

Here is the problem: Mr A is in charge of a widget factory, which proudly advertises that it can make delivery in 24 hours on any size order. This, of course, is not really true, and Mr A's job is to protect, as best he can, the advertising manager's reputation for veracity. This means that each morning he must decide whether the day's run of 200 widgets will be painted red or green. (For complex technological reasons, not relevant to the present problem, only one color can be produced per day.) We follow his problem of decision through several



Visual data exhibit ***low-dimensional structures*** due to rich ***local*** regularities, ***global*** symmetries, ***repetitive*** patterns, or ***redundant*** sampling.



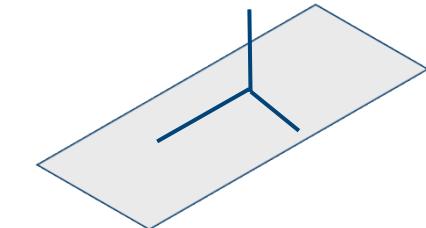
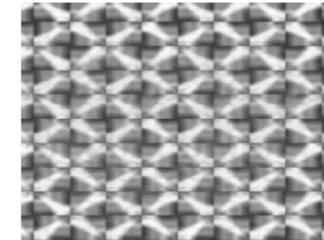
CONTEXT – PCA: *Fitting Data with a Low-dim. Subspace*

If we view the data (image) as a matrix

$$\textcolor{red}{A} = [\mathbf{a}_1 \mid \cdots \mid \mathbf{a}_n] \in \mathbb{R}^{m \times n}$$

then

$$r \doteq \text{rank}(\textcolor{red}{A}) \ll m.$$

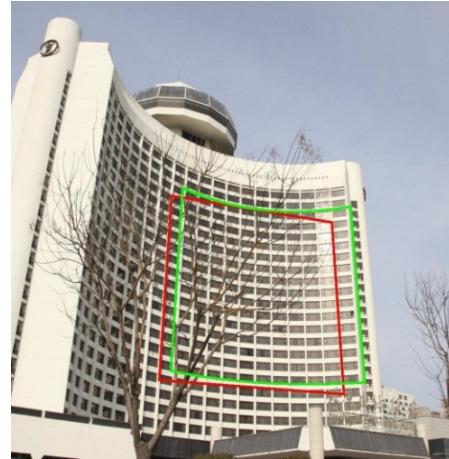


Principal Component Analysis (PCA) via singular value decomposition (SVD):

- Optimal estimate of $\textcolor{red}{A}$ under iid Gaussian noise $D = \textcolor{red}{A} + \textcolor{red}{Z}$
- Efficient and scalable computation
- Fundamental statistical tool, with huge impact in image processing, vision, web search, bioinformatics...

But... **PCA breaks down under even a single corrupted observation.**

CONTEXT – *But life is not so easy...*



*Real application data often contain **missing observations**, **corruptions**, or subject to unknown **deformation** or **misalignment**.*

Classical methods (e.g., PCA, least square regression) break down...

Everything old ...

A long and rich history of robust estimation with error correction and missing data imputation:



R. J. Boscovich. *De calculo probabilitatum que respondent diversis valoribus summe errorum post plures observationes ...*, before 1756



A. Legendre. *Nouvelles methodes pour la determination des orbites des cometes*, 1806



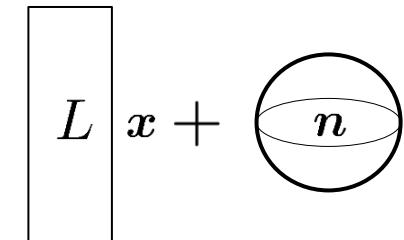
C. Gauss. *Theory of motion of heavenly bodies*, 1809



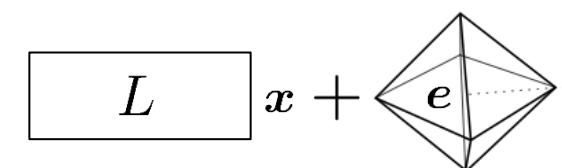
A. Beurling. *Sur les intégrales de Fourier absolument convergentes et leur application à une transformation fonctionnelle*, 1938

⋮

B. Logan. *Properties of High-Pass Signals*, 1965



over-determined
+ dense, Gaussian



underdetermined
+ sparse, Laplacian

CONTEXT – *Recent related progress*

Sparse recovery: Given $y = L\mathbf{x}_0$, $L \in \mathbb{R}^{m \times n}$, $m \ll n$, recover \mathbf{x}_0 .

$$y \in \mathbb{R}^m \quad \left[\begin{array}{c} \text{blue} \\ \text{yellow} \\ \text{green} \\ \text{red} \\ \text{purple} \end{array} \right] = L \in \mathbb{R}^{m \times n} \quad \left[\begin{array}{c} \text{red} \\ \text{blue} \\ \text{yellow} \\ \text{green} \\ \text{purple} \end{array} \right] \quad \mathbf{x} \in \mathbb{R}^n$$

Impossible in general ($m \ll n$)

Well-posed if \mathbf{x}_0 is structured (sparse), but still **NP-hard**

Tractable via convex optimization: $\min \|\mathbf{x}\|_1$ s.t. $y = L\mathbf{x}$

... if L is “nice” (random, incoherent, RIP)

Hugely active area: Donoho+Huo '01, Elad+Bruckstein '03, Candès+Tao '04, '05, Tropp '04, '06, Donoho '04, Fuchs '05, Zhao+Yu '06, Meinshausen+Buhlmann '06, Wainwright '09, Donoho+Tanner '09 ... and many others

CONTEXT – *Recent related progress*

Robust recovery: Given $y = L\mathbf{x}_0 + \mathbf{e}_0$, $L \in \mathbb{R}^{m \times n}$, $m \ll n$, recover \mathbf{x}_0 and \mathbf{e}_0 .

$$y = L\mathbf{x} + \mathbf{e}$$

Impossible in general ($m \ll n + m$)

Well-posed if \mathbf{x}_0 is *sparse*, errors \mathbf{e}_0 not too dense, but still **NP-hard**

Tractable: via convex optimization: $\min \|\mathbf{x}\|_1 + \|\mathbf{e}\|_1$ s.t. $y = L\mathbf{x} + \mathbf{e}$

... if L is “nice” (*cross and bouquet*)

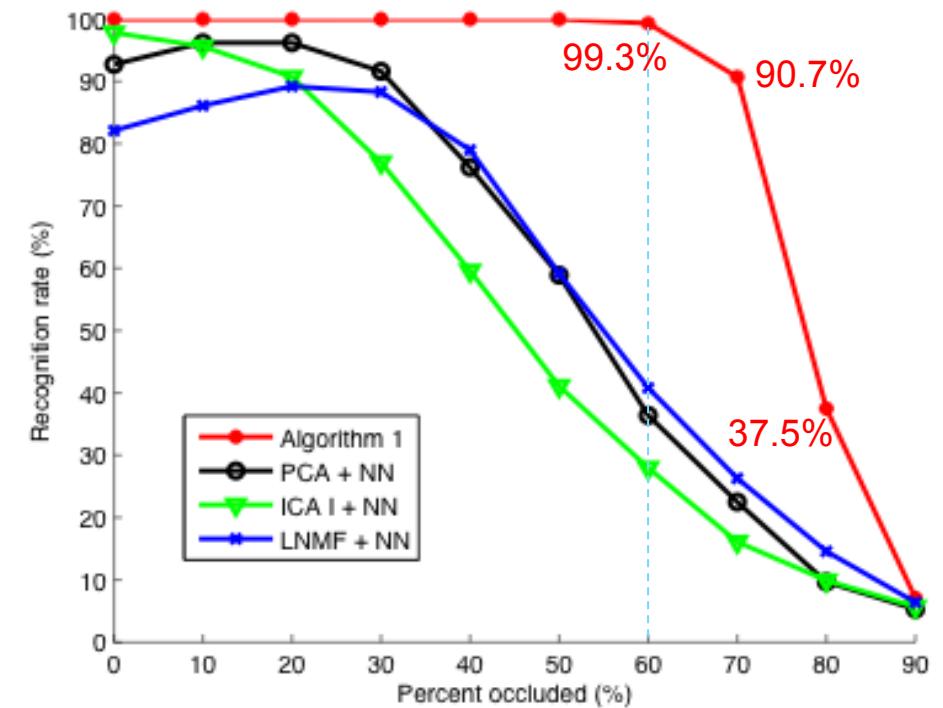
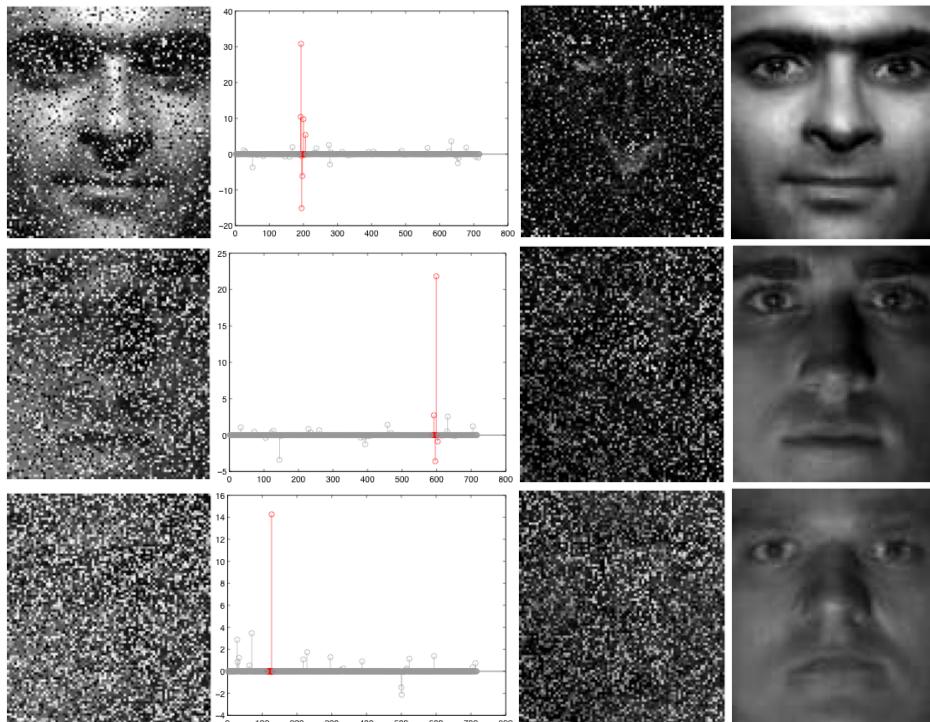
Hugely active area: Candès+Tao '05, Wright+Ma '10, Nguyen+Tran '11, Li '11,
also Zhang, Yang, Huang '11, etc...

EXPERIMENTS – *Varying Level of Random Corruption*

Extended Yale B Database
(38 subjects)

Training: subsets 1 and 2 (717 images)
Testing: subset 3 (453 images)

y \hat{x}_1 \hat{e}_1 $\hat{y}_0 = A\hat{x}_1$



CONTEXT – *Recent related progress*

Low-rank recovery: Given $y = \mathcal{L}[A_0]$, $\mathcal{L} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$, recover A_0 .

$$y \in \mathbb{R}^p \quad \left\| \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right\| = \mathcal{L} \left\langle \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right\rangle , \quad A \in \mathbb{R}^{m \times n} \quad \left\langle \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right\rangle i = 1 \dots p$$

Impossible in general ($p \ll mn$)

Well-posed if A_0 is structured (*low-rank*), but still **NP-hard**

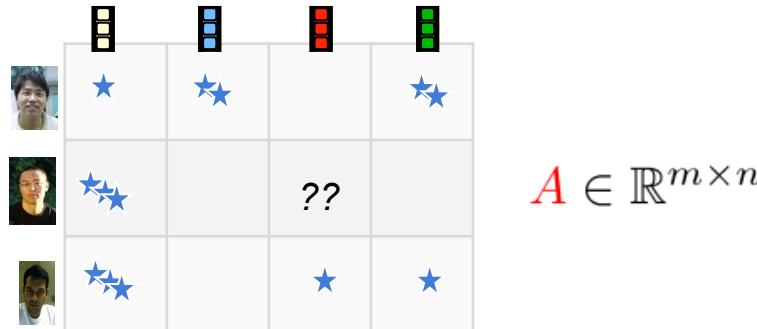
Tractable via convex optimization: $\min \|A\|_*$ s.t. $y = \mathcal{L}(A)$

... if \mathcal{L} is “nice” (*random, rank-RIP*)

Hugely active area: Recht+Fazel+Parillo ‘07, Candès+Plan ‘10, Mohan+Fazel ‘10, Recht+Xu+Hassibi ‘11, Chandrasekaran+Recht+Parillo+Willsky ‘11, Negahban+Wainwright ‘11 ...

CONTEXT – *Recent related progress*

Matrix completion: Given $y = \mathcal{P}_\Omega(\mathbf{A}_0)$, $\Omega \subset [m] \times [n]$, recover \mathbf{A}_0 .



Impossible in general ($|\Omega| \ll mn$)

Well-posed if \mathbf{A}_0 is structured (*low-rank*), but still **NP-hard**

Tractable via convex optimization: $\min \|A\|_*$ s.t. $y = \mathcal{P}_Q(A)$

... if Ω is “nice” (*random subset*) ...

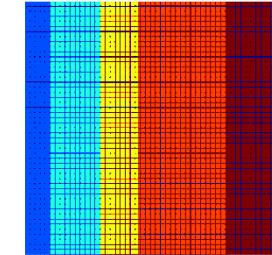
... and \mathbf{A}_0 interacts “nicely” with \mathcal{P}_Ω (\mathbf{A}_0 *incoherent – not “spiky”*).

Hugely active area: Candès+Recht ‘08, Keshavan+Oh+Montanari ‘09, Candès+Tao ‘09, Gross ‘10, Recht ‘10, Negahban+Wainwright ‘10

CONTEXT – *Low-dimensional Models*

The data should be **low-dimensional (low-rank)**:

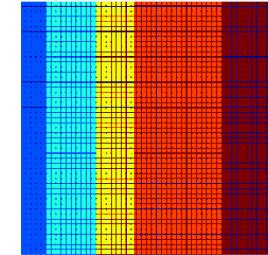
$$\textcolor{red}{A} = [\mathbf{a}_1 \mid \cdots \mid \mathbf{a}_n] \in \mathbb{R}^{m \times n}, \quad \text{rank}(\textcolor{red}{A}) \ll m.$$



CONTEXT – *Low-dimensional Models*

*The data should be **low-dimensional**:*

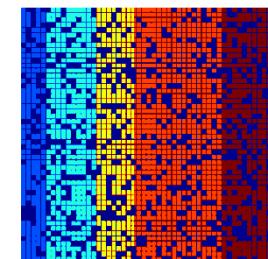
$$A = [\mathbf{a}_1 \mid \cdots \mid \mathbf{a}_n] \in \mathbb{R}^{m \times n}, \quad \text{rank}(A) \ll m.$$



*... but some of the observations are **grossly corrupted**:*

$$\mathbf{A} + \mathbf{E}, \quad |E_{ij}|$$

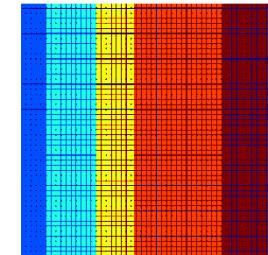
E_{ij} arbitrarily large, but most are zero.



CONTEXT – Low-dimensional Models

*The data should be **low-dimensional**:*

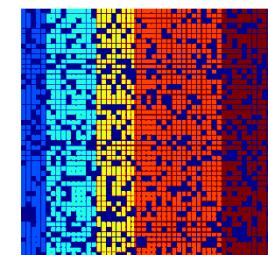
$$A = [\mathbf{a}_1 \mid \cdots \mid \mathbf{a}_n] \in \mathbb{R}^{m \times n}, \quad \text{rank}(A) \ll m.$$



*... but some of the observations are **grossly corrupted**:*

$$A + E, \quad |E_{ij}|$$

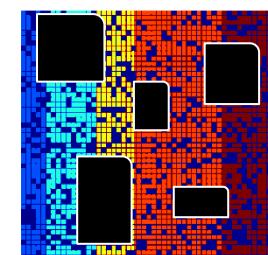
E_{ij} arbitrarily large, but most are zero.



*... and some of them can be **missing** too:*

$$D = \mathcal{P}_{\Omega}[A + E],$$

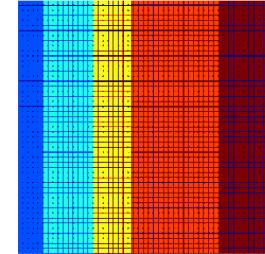
$\Omega \subset [m] \times [n]$ the set of observed entries.



CONTEXT – *Low-dimensional Models*

*The data should be **low-dimensional**:*

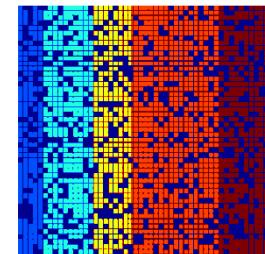
$$A = [\mathbf{a}_1 \mid \cdots \mid \mathbf{a}_n] \in \mathbb{R}^{m \times n}, \quad \text{rank}(A) \ll m.$$



*... but some of the observations are **grossly corrupted**:*

$$A + E, \quad |E_{ij}|$$

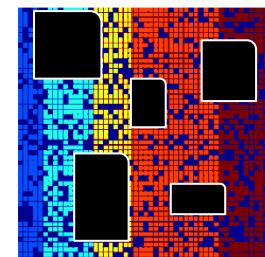
E_{ij} arbitrarily large, but most are zero.



*... and some of them can be **missing** too:*

$$D = \mathcal{P}_{\Omega}[A + E],$$

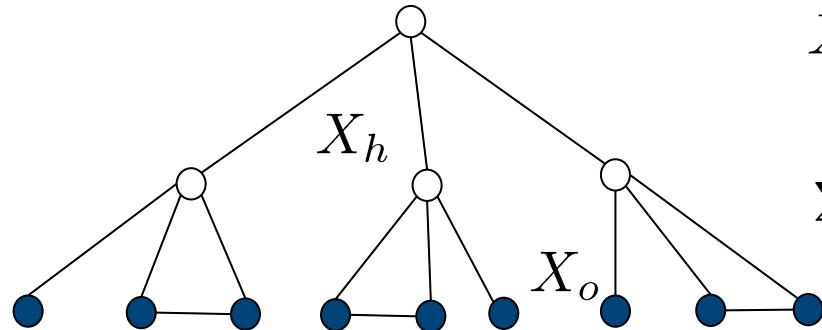
$\Omega \subset [m] \times [n]$ the set of observed entries.



... special cases of a more general problem:

$$D = \mathcal{L}_1(\mathbf{A}) + \mathcal{L}_2(\mathbf{E}) + \mathbf{Z} \quad \mathbf{A}, \mathbf{E} \text{ either sparse or low-rank}$$

CONTEXT: Learning Graphical Models



$$X = (X_o, X_h) \sim \mathcal{N}(0, \Sigma)$$

$$\Sigma = \begin{bmatrix} \Sigma_o & \Sigma_{oh} \\ \Sigma_{ho} & \Sigma_h \end{bmatrix} \Rightarrow \Sigma^{-1} = \begin{bmatrix} J_o & J_{oh} \\ J_{ho} & J_h \end{bmatrix}$$

$$X_i, X_j \text{ cond. indep. given other variables} \Leftrightarrow (\Sigma^{-1})_{ij} = 0$$

Separation Principle:

$$\begin{array}{rcl} \Sigma_o^{-1} & = & J_o - J_{oh} J_h^{-1} J_{ho} \\ \text{observed} & = & \text{sparse} + \text{low-rank} \end{array}$$

- sparse pattern \rightarrow conditional (in)dependence
- rank of second component \rightarrow number of hidden variables

THIS TALK

Given observations $D = \mathcal{P}_Q[\mathbf{A} + \mathbf{E} + \mathbf{Z}]$, with

\mathbf{A} low-rank,

\mathbf{E} sparse,

\mathbf{Z} small, dense noise,

recover a good estimate of \mathbf{A} and \mathbf{E} .

Theory and Algorithm

- Provably Correct and Tractable Solution
- Provably Optimal and Efficient Algorithms

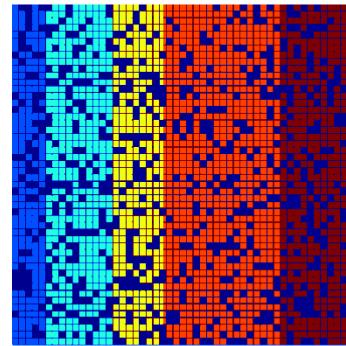
Potential Applications

- Visual Data (Restoration, Reconstruction, Recognition)
- Other Data

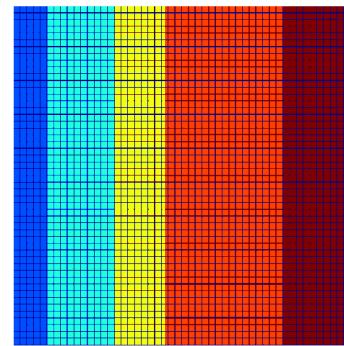
Extensions and Conclusions

ROBUST PCA – *Problem Formulation*

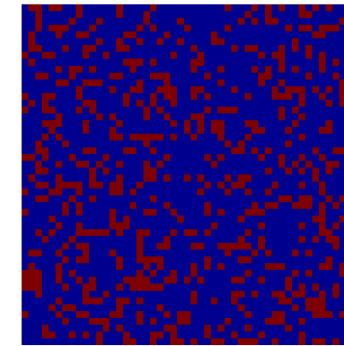
D - observation



A_0 – low-rank



E_0 – sparse



$$= +$$

Problem: Given $D = \underbrace{A_0}_{\text{Low-rank component}} + \underbrace{E_0}_{\text{Sparse component (gross errors)}}$, recover A_0 and E_0 .

Numerous approaches in the literature:

- Multivariate trimming [Gnanadesikan and Kettering '72]
- Power Factorization [Wieber '70s]
- Random sampling [Fischler and Bolles '81]
- Alternating minimization [Shum & Ikeuchi '96, Ke and Kanade '03]
- Influence functions [de la Torre and Black '03]

Key question: ***can guarantee correctness with an efficient algorithm?***

ROBUST PCA – Convex Surrogates for Sparsity and Rank

Seek the lowest-rank A that agrees with the data up to some sparse error E :

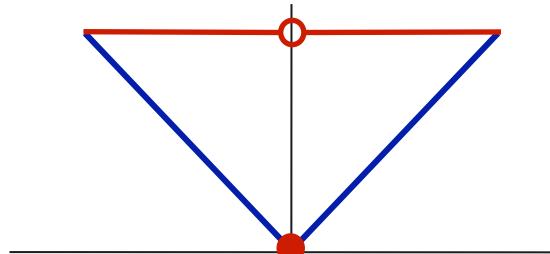
$$\min \text{rank}(A) + \gamma \|E\|_0 \quad \text{subj } A + E = D.$$

But INTRACTABLE! Relax with convex surrogates:

$$\|E\|_0 = \#\{E_{ij} \neq 0\} \quad \rightarrow \quad \|E\|_1 = \sum_{ij} |E_{ij}|. \quad L_1 \text{ norm}$$

$$\text{rank}(A) = \#\{\sigma_i(A) \neq 0\} \quad \rightarrow \quad \|A\|_* = \sum_i \sigma_i(A). \quad \text{Nuclear norm}$$

Convex envelope over $B_{2,2} \times B_{1,\infty}$



ROBUST PCA – *By Convex Optimization*

Seek the lowest-rank A that agrees with the data up to some sparse error E :

$$\min \text{rank}(A) + \gamma \|E\|_0 \quad \text{subj } A + E = D.$$

But INTRACTABLE! Relax with convex surrogates:

$$\|E\|_0 = \#\{E_{ij} \neq 0\} \quad \rightarrow \quad \|E\|_1 = \sum_{ij} |E_{ij}|. \quad L_1 \text{ norm}$$

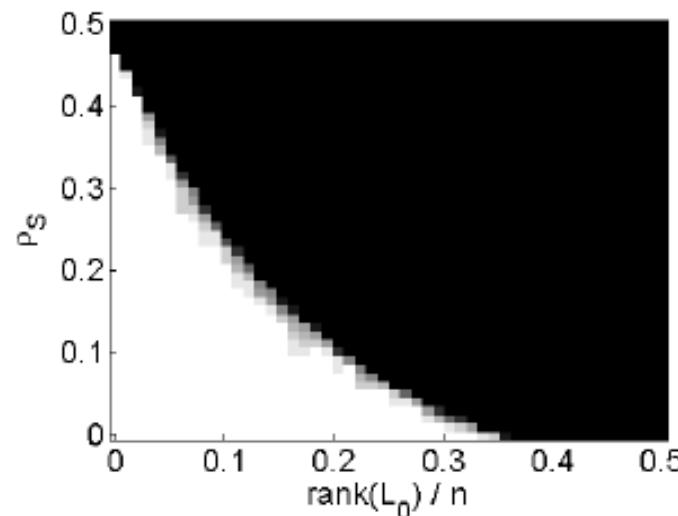
$$\text{rank}(A) = \#\{\sigma_i(A) \neq 0\} \quad \rightarrow \quad \|A\|_* = \sum_i \sigma_i(A). \quad \text{Nuclear norm}$$

$$\min \|A\|_* + \lambda \|E\|_1 \quad \text{subj } A + E = D.$$

Semidefinite program, solvable in polynomial time

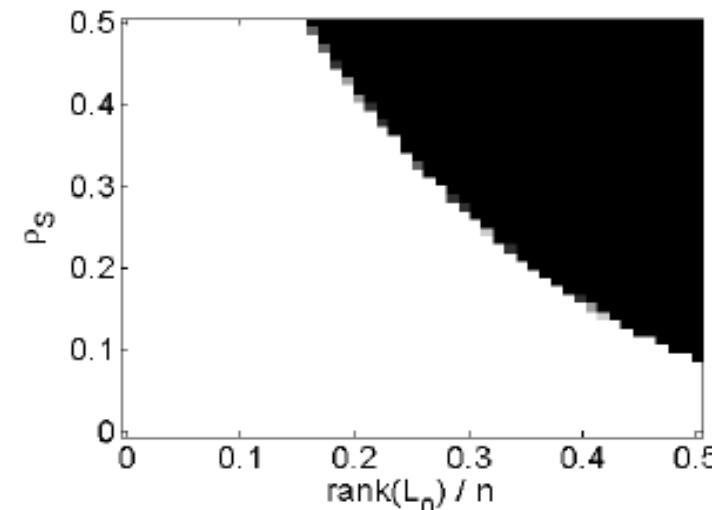
ROBUST PCA – *When the Convex Program Works?*

$$D = \textcolor{red}{A} + \textcolor{green}{E}$$



Robust PCA, Random Signs

$$D = \mathcal{P}_\Omega[\textcolor{red}{A}]$$



Matrix Completion

White regions are instances with perfect recovery.

Correct recovery when $\textcolor{red}{A}$ is indeed **low-rank** and $\textcolor{green}{E}$ is indeed **sparse**?

MAIN THEORY – *Exact Solution by Convex Optimization*

Theorem 1 (Principal Component Pursuit). If $A_0 \in \mathbb{R}^{m \times n}$, $m \geq n$ has rank

Non-adaptive weight factor

and E_0 has Bernoulli support with error probability $\rho \leq \rho_s^*$, then with very high probability

$$(A_0, E_0) = \arg \min \|A\|_* + \frac{1}{\sqrt{m}} \|E\|_1 \quad \text{subj} \quad A + E = A_0 + E_0,$$

and the minimizer is unique.

GREAT NEWS: “Convex optimization recovers almost any matrix of rank $O\left(\frac{m}{\log^2 n}\right)$ from errors corrupting $O(mn)$ of the observations!”

MAIN THEORY – *Corrupted, Incomplete Matrix*

$$D = \mathcal{P}_\Omega [\ A_0 + E_0 \], \quad \Omega \sim \text{uni}\left(\binom{[m] \times [n]}{mn}\right)$$

Theorem 2 (Matrix Completion and Recovery). *If $A_0, E_0 \in \mathbb{R}^{m \times n}$, $m \geq n$, with*

$$\text{rank}(A_0) \leq C \frac{n}{\mu \log^2(m)}, \quad \text{and} \quad \|E_0\|_0 \leq \rho^\star mn,$$

and we observe only a random subset of size

$$|\Omega| = mn/10$$

entries, then with very high probability, solving the convex program

$$\min \|A\|_* + \frac{1}{\sqrt{m}} \|E\|_1 \quad \text{subj} \quad P_\Omega[A + E] = D,$$

uniquely recovers (A_0, E_0) .

MAIN THEORY – *With Dense Errors and Noise*

Theorem 3 (Dense Error Correction). If A_0 has rank $r \leq \rho_r \frac{m}{\mu^2 \log^2(n)}$ and E_0 has random signs and Bernoulli support with error probability $\rho < 1$, then with very high probability

$$(A_0, E_0) = \arg \min \|A\|_* + \lambda \|E\|_1 \quad \text{subj} \quad A + E = A_0 + E_0,$$

and the minimizer is unique.

Theorem 4 (Robust PCA with Noise). Given $D = A_0 + E_0 + Z$ for any $\|Z\|_F \leq \eta$, if A_0 has rank $r \leq \rho_r \frac{m}{\mu^2 \log^2(n)}$ and E_0 has Bernoulli support with error probability $\rho \leq \rho_s^*$, then with very high probability

$$(\hat{A}, \hat{E}) = \arg \min \|A\|_* + \frac{1}{\sqrt{m}} \|E\|_1 \quad \text{subj} \quad \|D - A - E\| \leq \eta,$$

satisfies $\|(\hat{A}, \hat{E}) - (A_0, E_0)\| \leq C\eta$ for some constant $C > 0$.

FIRST RESULTS OF THIS TYPE

Example: for $D = A_0 + E_0$,

Previous Best Result [Chandrasekharan, Parrilo, Wilsky'11]:

Deterministic error models, success when $\|E\|_0 \leq Cm^{1.5}/r^{.5} \log m$.

Does not guarantee to correct nonzero fractions of errors, even with $r = 1$.

FIRST RESULTS OF THIS TYPE

Example: for $D = A_0 + E_0$,

Previous Best Result [Chandrasekharan et. al.]:

Success when $\|E\|_0 \leq Cm^{1.5}/r^{.5} \log m$.

Does not guarantee to correct nonzero fractions of errors, even with $r = 1$.

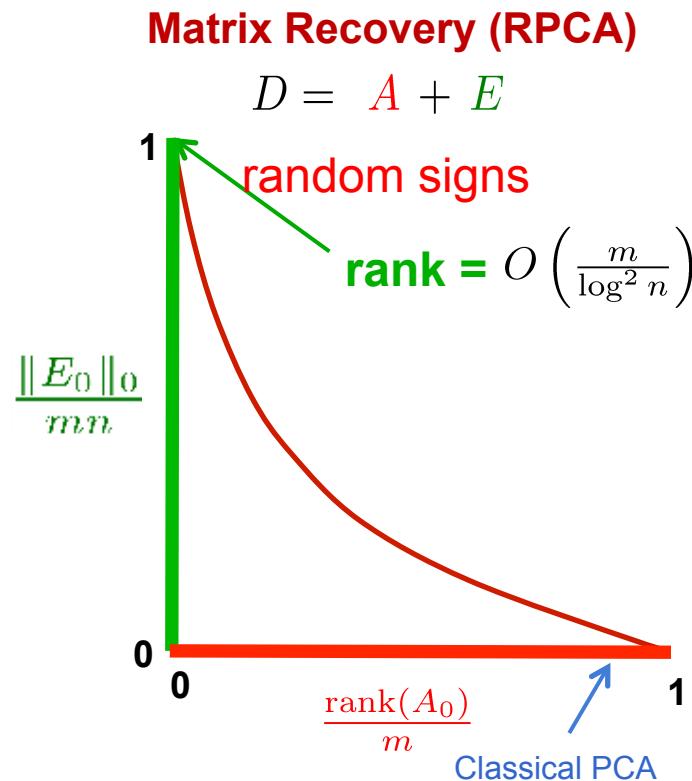
Our results:

Corrects nonzero fractions of errors, even with $r = O(m/\log^2 n)$,

Considers **corruption, missing elements and noise**: $\mathcal{P}_\Omega[A_0 + E_0 + Z]$

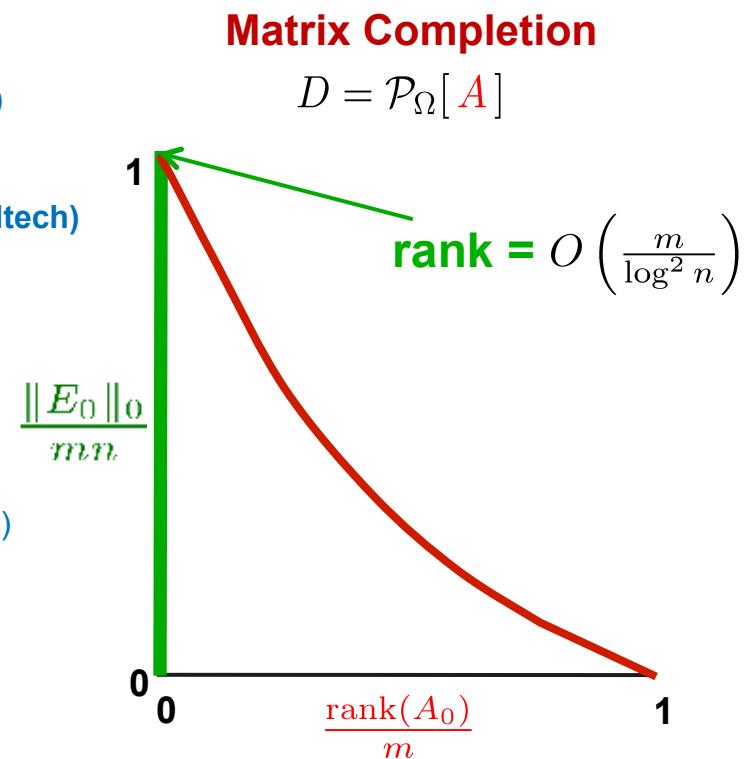
BIG PICTURE – *Landscape of Theoretical Guarantees*

What people have known so far in the past 3-4 years:



D. Gross
E. Candes (Stanford)
B. Recht (UC Berkeley)
J. Wright (Columbia)
J. Tropp (Caltech)
Chandrasekharan (Caltech)

B. Hassibi (Caltech)
P. Parrilo (MIT)
A. Willsky (MIT)
B. Hastie (Stanford)
C. Montanari (Stanford)
M. Jordan (Berkeley)
M. Wainwright (Berkeley)
B. Yu (Berkeley)
A. Singer (Princeton)
T. Tao (UCLA)
S. Osher (UCLA)
O. Milenkovic (UIUC)
Y. Bresler (UIUC)
Y. Ma (UIUC)
M. Fazel (U Wash.)

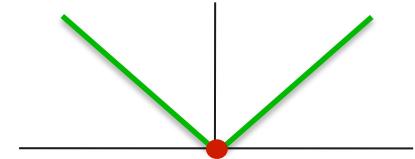


This phase transition landscape has been precisely understood! (Tropp et. al.)

ALGORITHMS – Are scalable solutions possible?

Seemingly BAD NEWS: Our optimization problem

$$\min \|A\|_* + \lambda \|E\|_1 \text{ subj } A + E = D.$$



is high-dimensional and non-smooth.

Convergence rate of solving a generic convex program: $\min_x f(x)$

Second-order Newton method, # of iterations: $O(\log(1-\epsilon))$, but not scalable!
First-order methods depend strongly on the smoothness of f :

Function class \mathcal{F}	Suboptimality $f(\mathbf{x}_k) - f(\mathbf{x}^*)$
<i>smooth</i> f convex, differentiable $\ \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}')\ \leq L\ \mathbf{x} - \mathbf{x}'\ $	$\frac{CL\ \mathbf{x}_0 - \mathbf{x}^*\ ^2}{k^2} = \Theta\left(\frac{1}{k^2}\right)$
<i>smooth + structured nonsmooth:</i> $+$ f, g convex, $\ \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}')\ \leq L\ \mathbf{x} - \mathbf{x}'\ $	$\frac{CL\ \mathbf{x}_0 - \mathbf{x}^*\ ^2}{k^2} = \Theta\left(\frac{1}{k^2}\right)$
<i>nonsmooth</i> f convex $ f(\mathbf{x}) - f(\mathbf{x}') \leq M\ \mathbf{x} - \mathbf{x}'\ $	$\frac{CM\ \mathbf{x}_0 - \mathbf{x}^*\ }{\sqrt{k}} = \Theta\left(\frac{1}{\sqrt{k}}\right)$

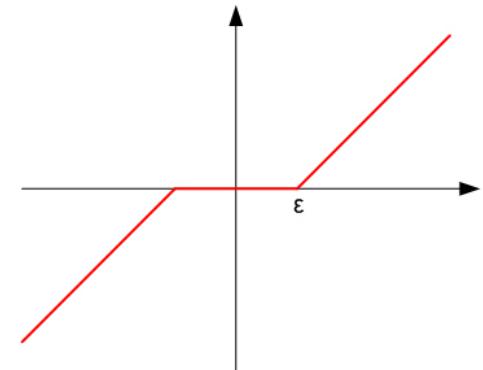
ALGORITHMS – *Why are scalable solutions possible?*

GOOD NEWS: The objective function has **special structures**

$$\min \|A\|_* + \lambda \|E\|_1 \text{ subj } A + E = D.$$

KEY OBSERVATION: **Simple solutions** for the proximal operations, given by soft-thresholding the entries or singular values of the matrix, respectively.

$$\begin{aligned} S_\varepsilon(Q) &= \operatorname{argmin}_X \varepsilon \|X\|_1 + \frac{1}{2} \|X - Q\|_F^2 \\ D_\varepsilon(Q) &= \operatorname{argmin}_X \varepsilon \|X\|_* + \frac{1}{2} \|X - Q\|_F^2 \end{aligned}$$



For composite functions $F = f + g$, with f smooth,
if g has an efficient proximal operator, we achieve
the same (optimal) rate as if F was smooth.

ALGORITHMS – *Evolution of scalable algorithms*

GOOD NEWS: Scalable first-order gradient-descent algorithms:

- Proximal Gradient [Osher, Mao, Dong, Yin '09, Wright et. al.'09, Cai et. al.'09].
- Accelerated Proximal Gradient [Nesterov '83, Beck and Teboulle '09]:
- Augmented Lagrange Multiplier [Hestenes '69, Powell '69]:
- Alternating Direction Method of Multipliers [Gabay and Mercier '76].

A scalable algorithm: alternating direction method (ADMoM) for ALM:

$$l(A, E, Y) = \|A\|_* + \lambda\|E\|_1 + \langle Y, D - A - E \rangle + \frac{\mu}{2}\|D - A - E\|_F^2$$

repeat

$$\begin{cases} A_{k+1} &= \mathcal{D}_{\mu_k^{-1}}(D - E_k + Y_k/\mu_k), & \text{Shrink singular values} \\ E_{k+1} &= \mathcal{S}_{\lambda\mu_k^{-1}}(D - A_{k+1} + Y_k/\mu_k), & \text{Shrink absolute values} \\ Y_{k+1} &= Y_k + \mu_k(D - A_{k+1} - E_{k+1}). \end{cases}$$

Cost of each iteration is a classical PCA, i.e. a (partial) SVD.

ALGORITHMS – *Evolution of fast algorithms (around 2009)*

For a 1000x1000 matrix of rank 50, with 10% (100,000) entries randomly corrupted: $\min \|A\|_* + \lambda \|E\|_1$ subj $A + E = D$.

Algorithms	Accuracy	Rank	$\ E\ _0$	# iterations	time (sec)
IT	5.99e-006	50	101,268	8,550	119,370.3
DUAL	8.65e-006	50	100,024	822	1,855.4
APG	5.85e-006	50	100,347	134	1,468.9
APG _P	5.91e-006	50	100,347	134	82.7
EALM _P	2.07e-007	50	100,014	34	37.5
IALM _P	3.83e-007	50	99,996	23	11.8

10,000
times
speedup!

Provably Robust PCA at only a constant factor (≈ 20) more computation than conventional PCA!

ALGORITHMS – *Convergence rate with strong convexity*

GREAT NEWS: Geometric convergence for gradient algorithms!

- f restricted strong convex: $O(\log(1/\varepsilon))$ [Agarwal, Negahban, Wainwright, NIPS 2010]
 f smooth, ∇f Lipschitz: $O(\varepsilon^{-1/2})$
 f differentiable: $O(\varepsilon^{-1})$
 f non-smooth: $O(\varepsilon^{-2})$

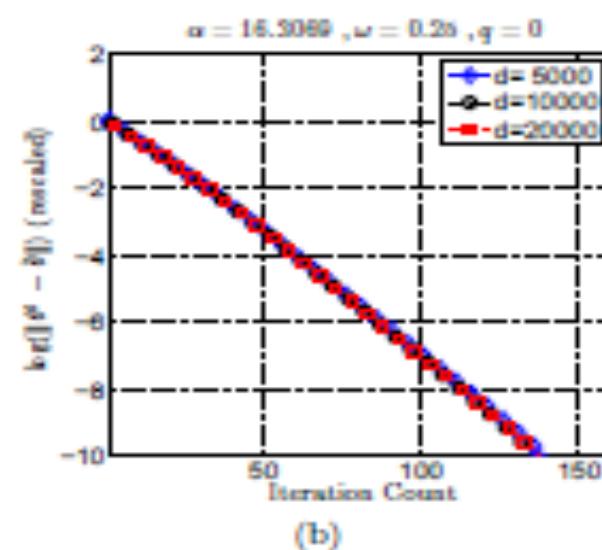
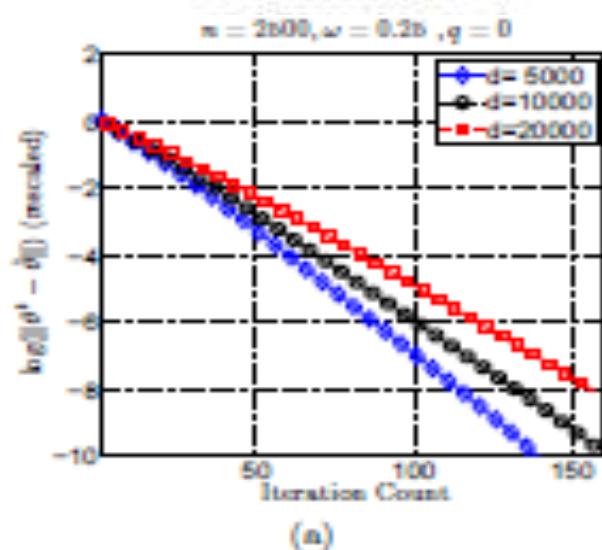


Figure 1. Convergence rates of projected gradient descent in application to Lasso programs (ℓ_1 -constrained least-squares). Each panel shows the log optimization error $\log \|\theta^t - \hat{\theta}\|$ versus the iteration number t . Panel (a) shows three curves, corresponding to dimensions $d \in \{5000, 10000, 20000\}$, sparsity $s = \lceil \sqrt{d} \rceil$, and all with the same sample size $n = 2500$. All cases show geometric convergence, but the rate for larger problems becomes progressively slower. (b) For an appropriately rescaled sample size ($\alpha = \frac{n}{d \log d}$), all three convergence rates should be roughly the same, as predicted by the theory.

ALGORITHMS – *Recap and Conclusions*

Key challenges of **nonsmoothness** and **scale** can be mitigated by using **special structure** in sparse and low-rank optimization problems:

Efficient proximity operators \Rightarrow proximal gradient methods

Separable objectives \Rightarrow alternating directions methods

Efficient **moderate-accuracy solutions** for **very large problems**.

Special tricks can further improve specific cases (factorization for low-rank)

Techniques in this literature apply quite broadly.

Extremely useful tools for creative problem formulation / solution.

Fundamental **theory** guiding engineering **practice**:

What are the basic principles and limitations?

What specific structure in my problem can allow me to do better?

APPLICATIONS

□ Repairing Images and Videos

- Image Repairing, Background Extraction, Street Panorama

□ Reconstructing 3D Geometry

- Shape from Texture, Featureless 3D Reconstruction

□ Registering Multiple Images

- Multiple Image Alignment, Video Stabilization

□ Recognizing Objects

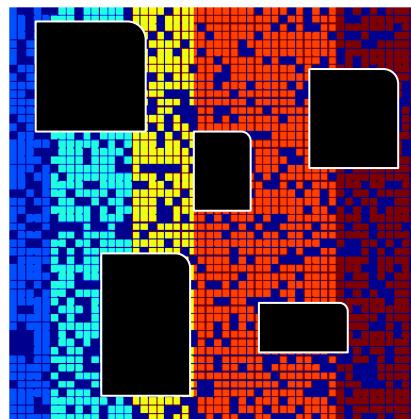
- Faces, Texts, etc

□ Other Data and Applications

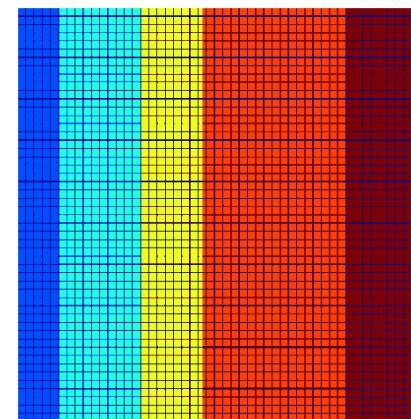
Implications: Highly Compressive Sensing of Structured Information!

Recover low-dimensional structures from a fraction of missing measurements with structured support.

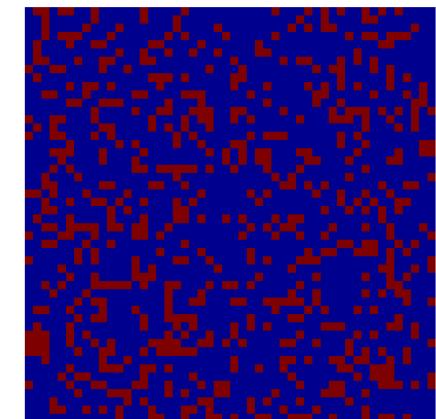
compressive samples



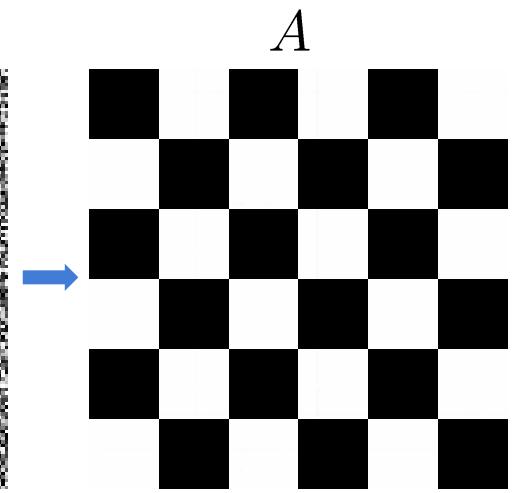
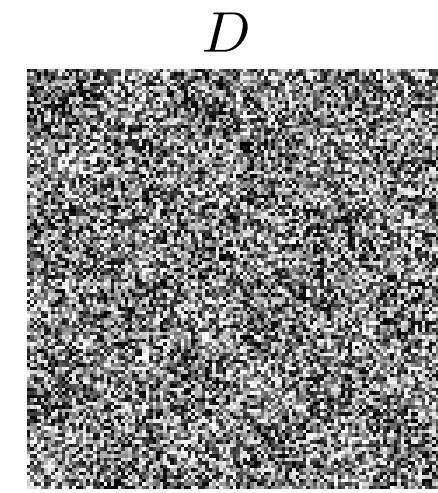
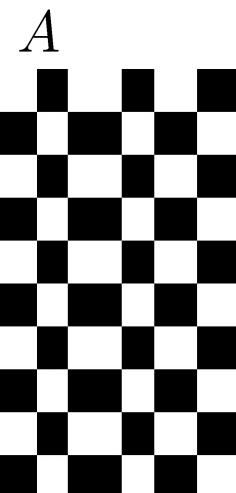
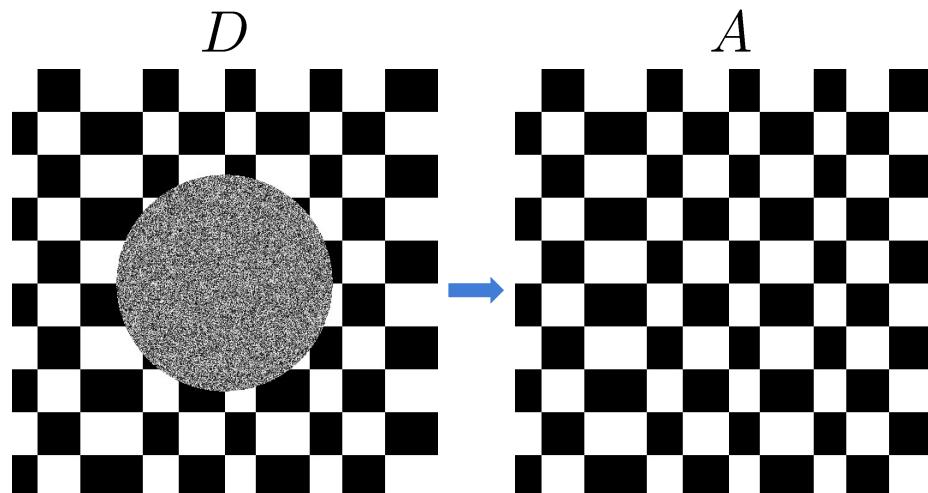
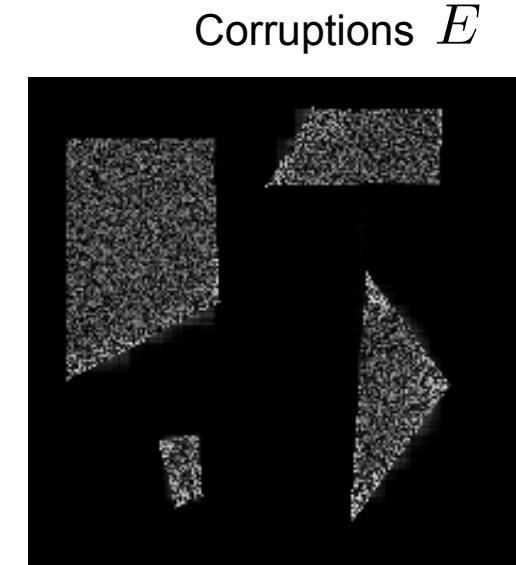
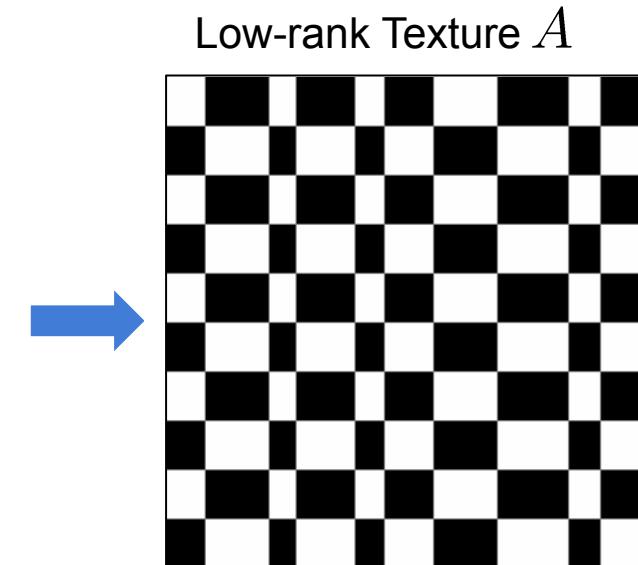
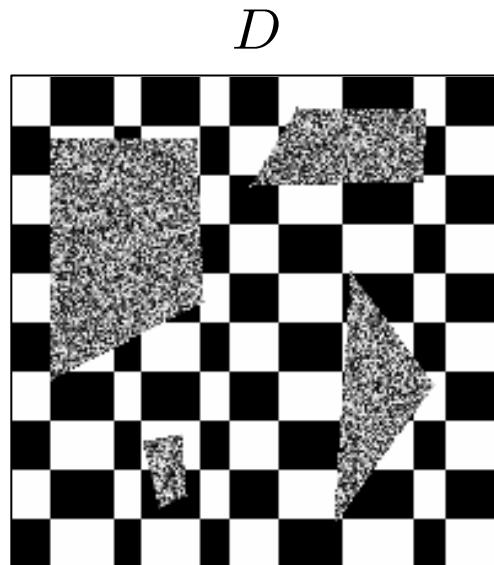
Low-rank Structures



Sparse Structures



Repairing Images: Highly Robust Repairing of Low-rank Textures!



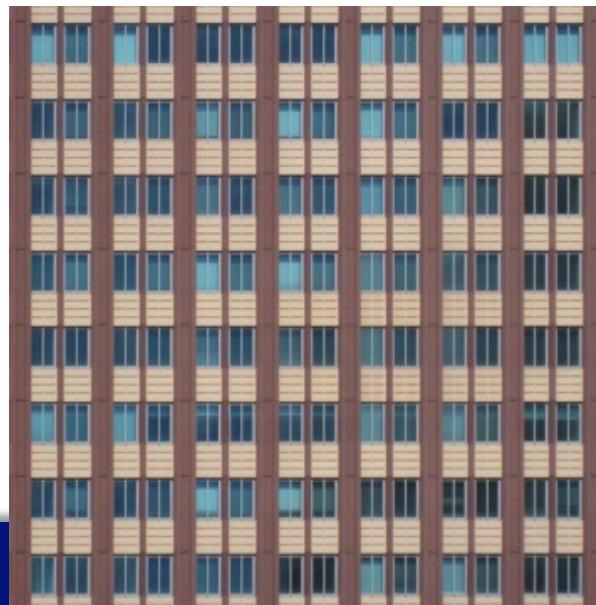
Repairing Low-rank Textures

Low-rank Method

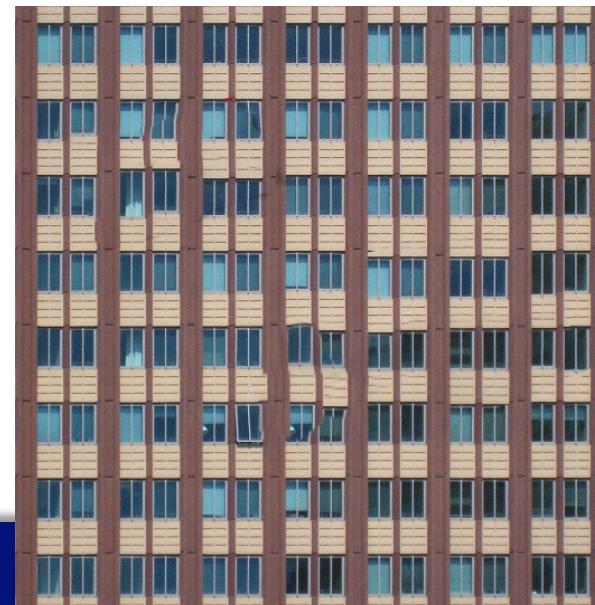
Input



Output



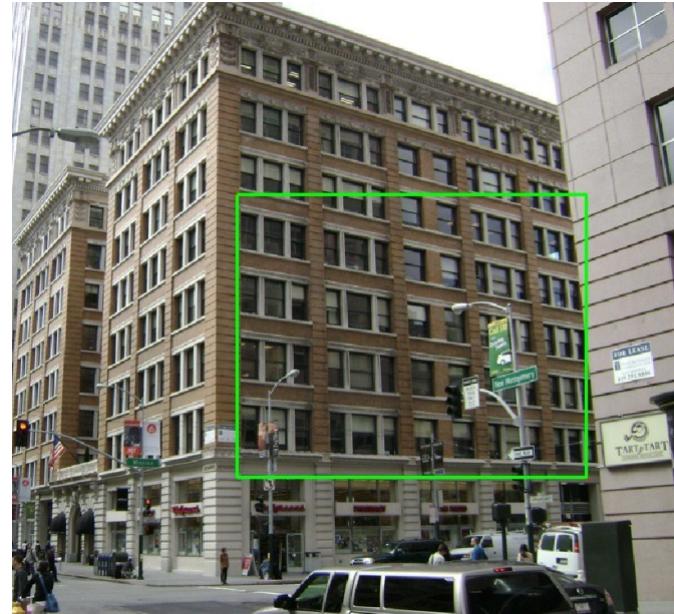
Photoshop



Repairing (Distorted) Low-rank Textures

Low-rank Method

Input



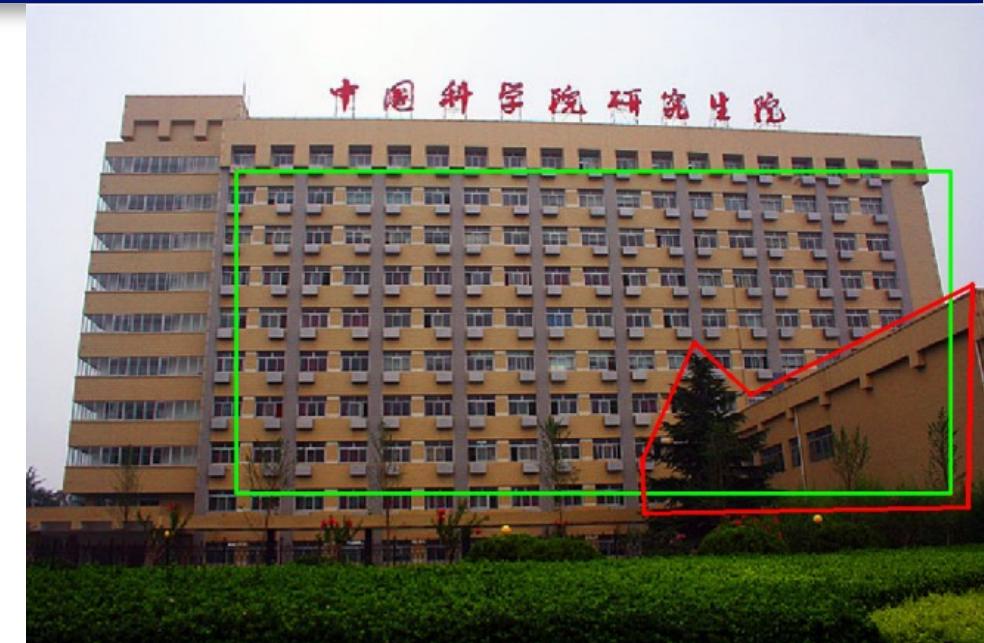
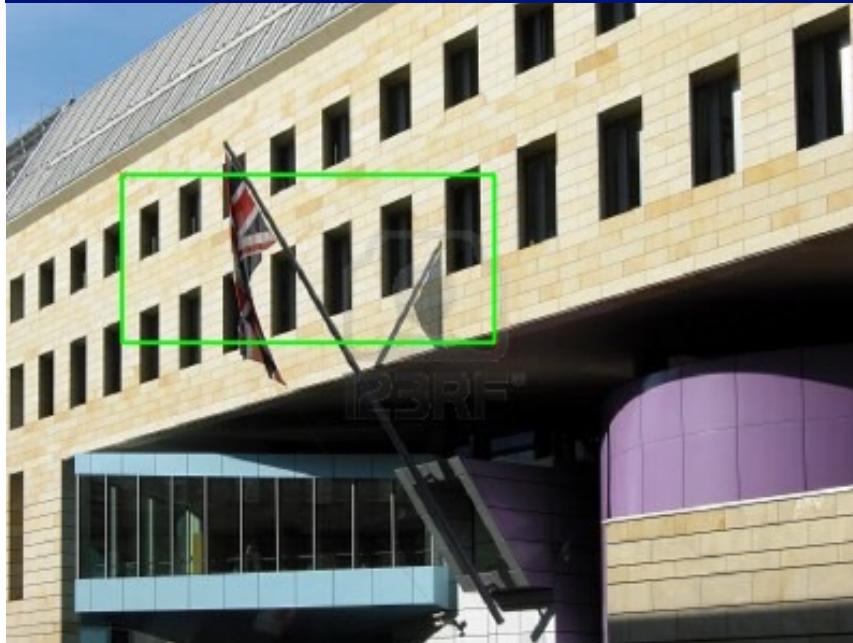
Photoshop



Output

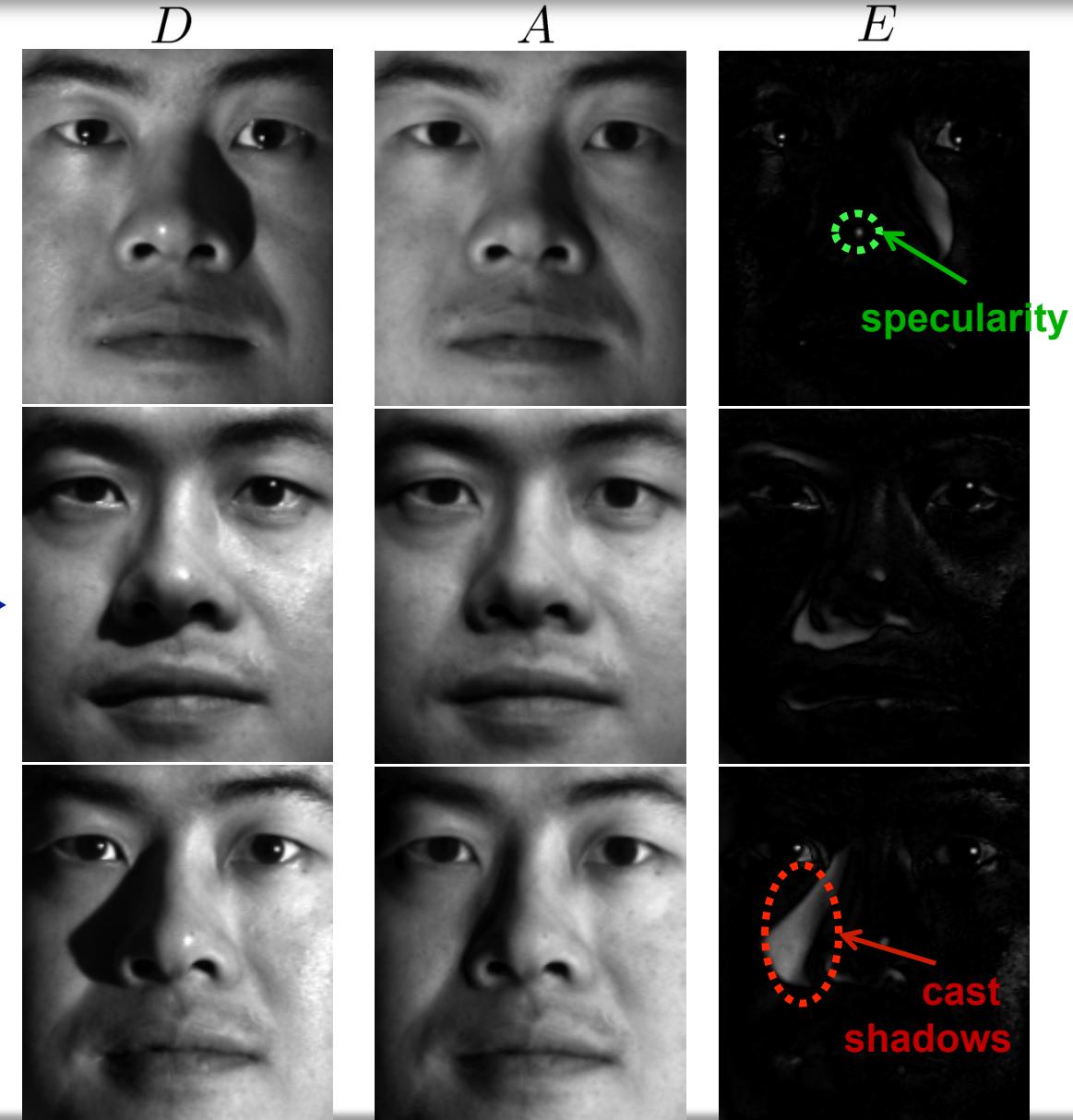
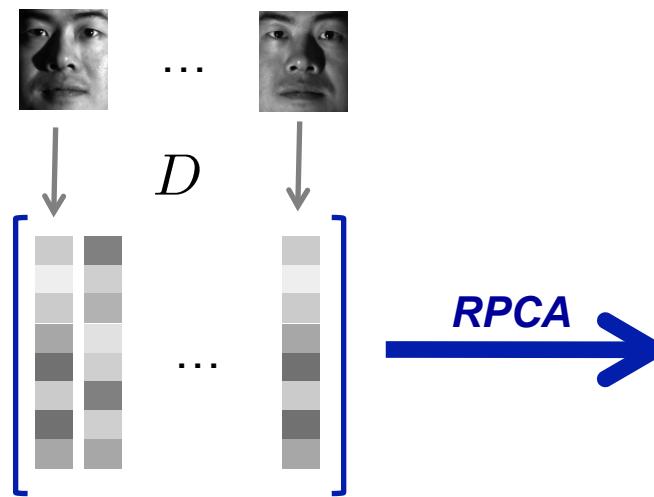


Structured Texture Completion and Repairing



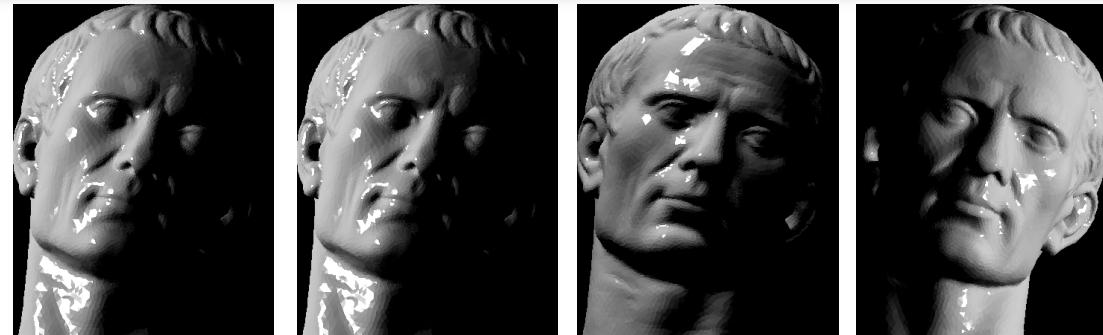
Repairing Multiple Correlated Images

58 images of one person under varying lighting:

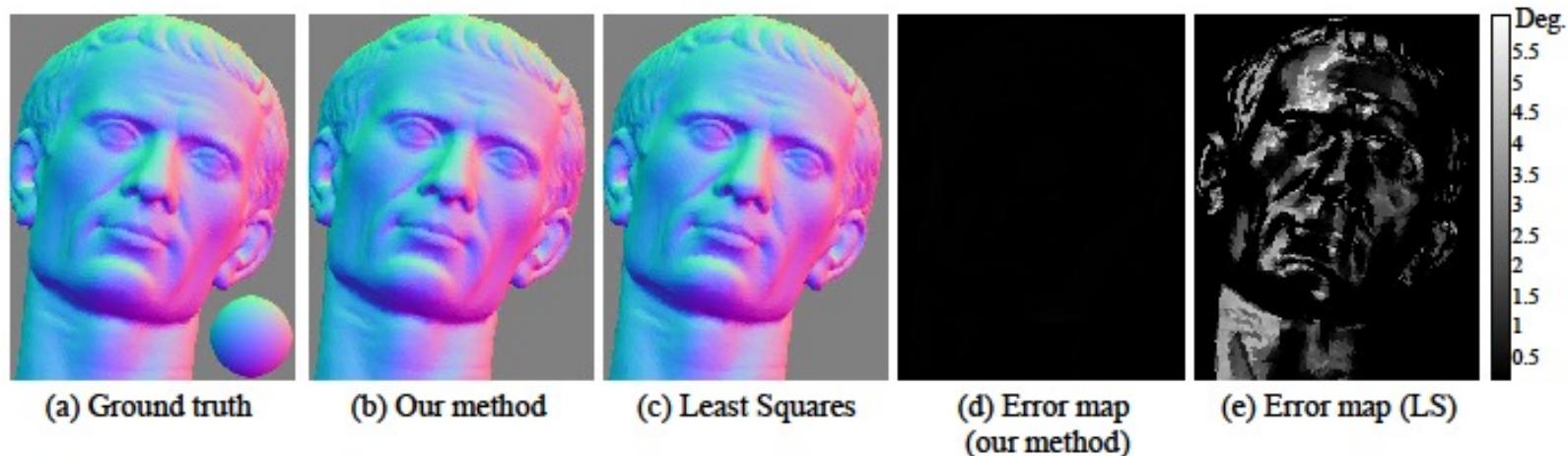


Repairing Images: robust photometric stereo

Input images



$$\min \|A\|_* + \lambda \|E\|_1 \text{ subj } D = \mathcal{P}_\Omega(A + E). \quad \begin{aligned} \Omega^c &\sim \text{shadow}(20.7\%) \\ E &\sim \text{specularities}(13.6\%) \end{aligned}$$



Mean error	0.014°	0.96°
Max error	0.20°	8.0°

Repairing Video Frames: *background modeling from video*

Surveillance video

200 frames,
144 x 172 pixels,

Significant foreground
motion



D

...



\downarrow
 $RPCA$
 \rightarrow



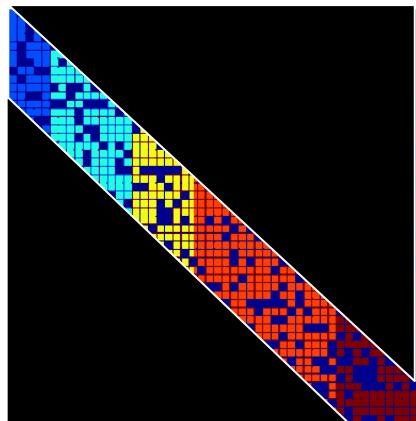
$$\text{Video } D = \text{Low-rank appx. } A + \text{Sparse error } E$$



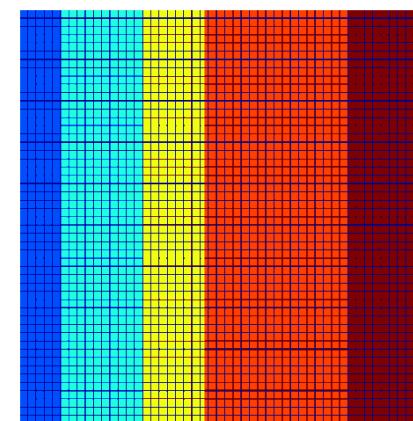
Implications: Highly Compressive Sensing of Structured Information!

Recover low-dimensional structures from diminishing fraction of corrupted measurements.

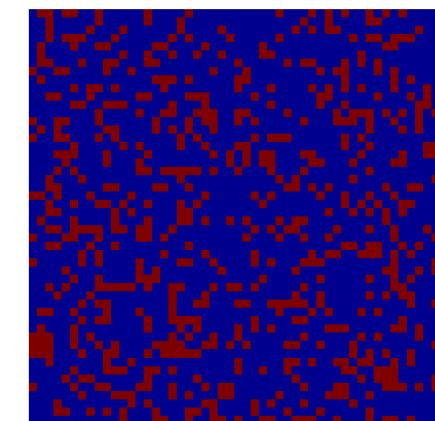
compressive samples



Low-rank Structures

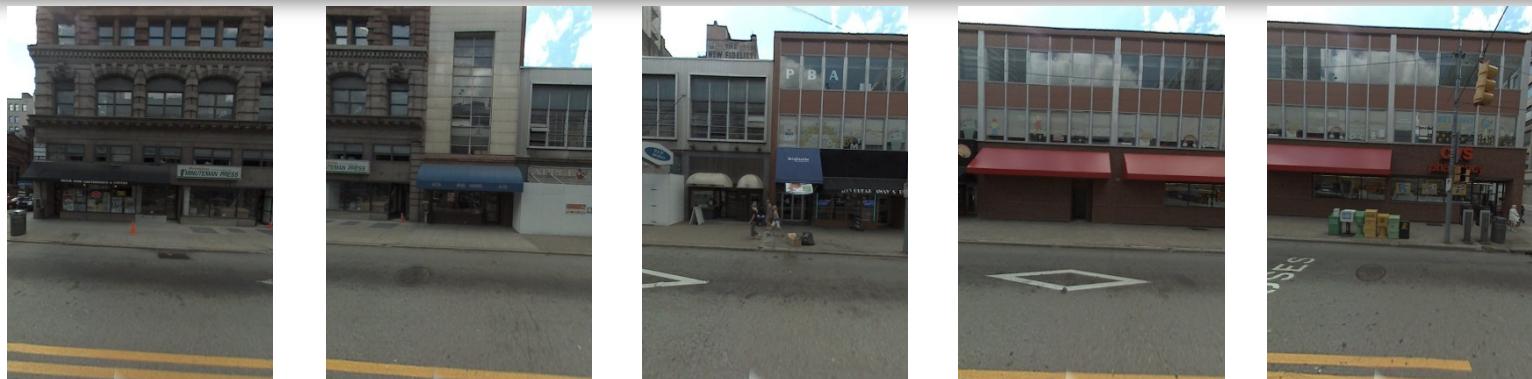


Sparse Structures



Repairing Video Frames: Street Panorama

D



A



E

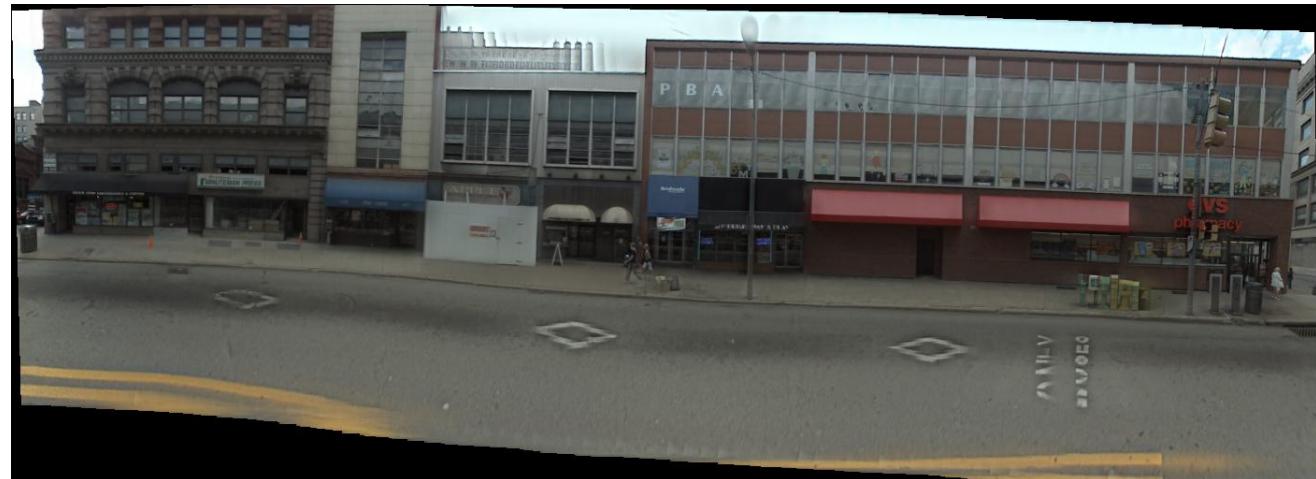


Repairing Video Frames: Street Panorama

Low-rank



AutoStitch



Photoshop



Repairing Video Frames: Street Panorama

Low-rank



AutoStitch

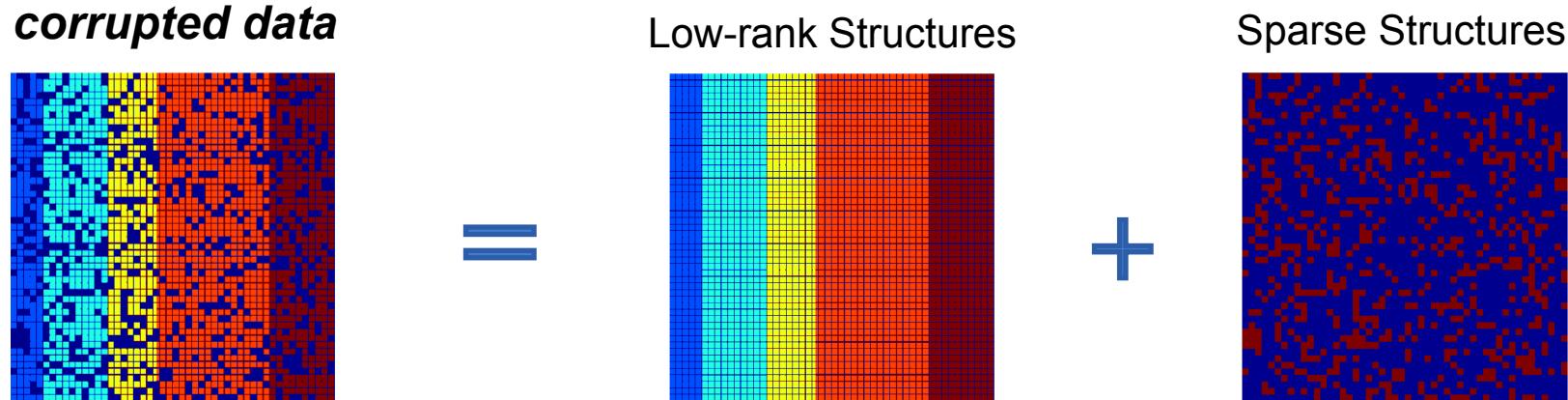


Photoshop

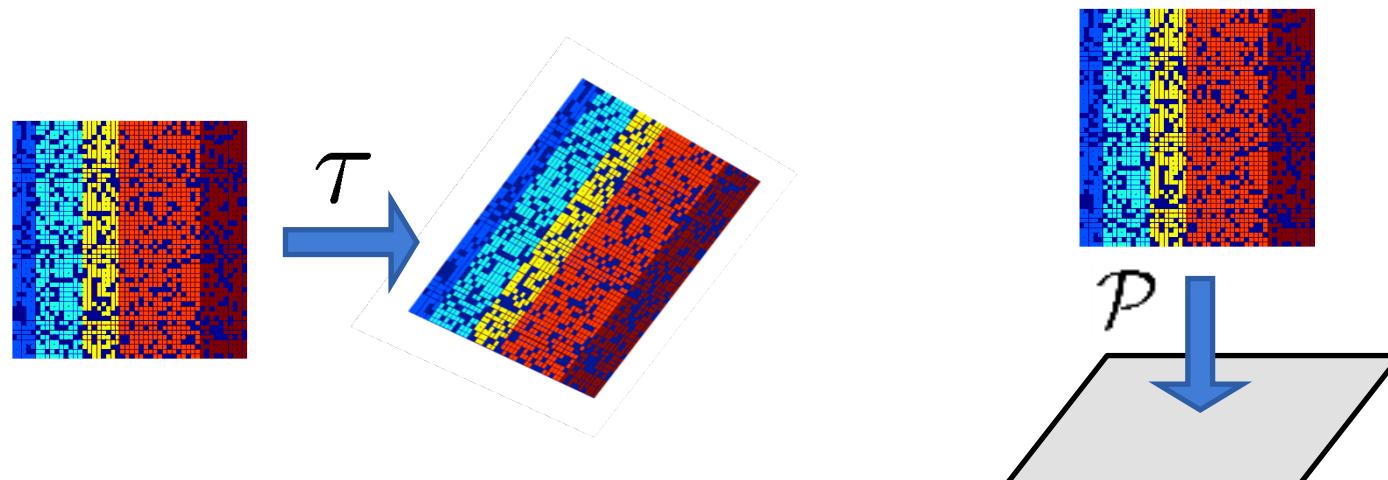


Sensing or Imaging of Low-rank and Sparse Structures

Fundamental Problem: *How to recover low-rank and sparse structures from corrupted data*

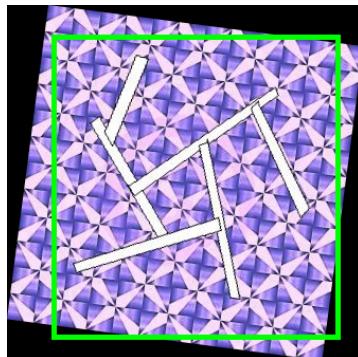


subject to either nonlinear deformation τ or linear compressive sampling \mathcal{P} ?

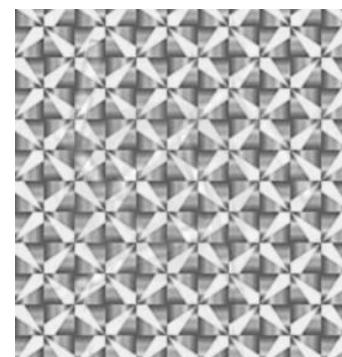


Reconstructing 3D Geometry and Structures

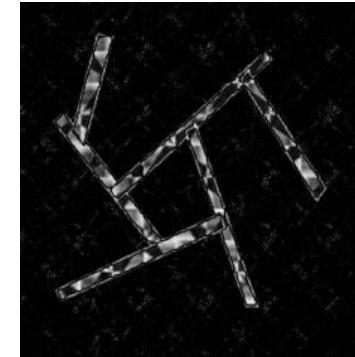
D – deformed observation



A – low-rank structures



E – sparse errors



$$D \circ \tau =$$

+

Problem: Given $D \circ \tau = A_0 + E_0$, recover τ , A_0 and E_0 simultaneously.

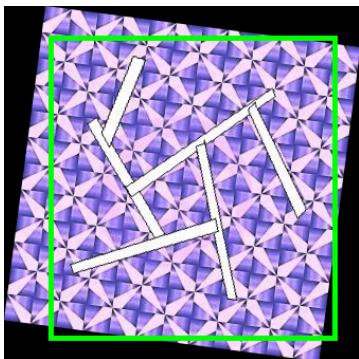
Low-rank component
(regular patterns...)

Sparse component
(occlusion, corruption, foreground...)

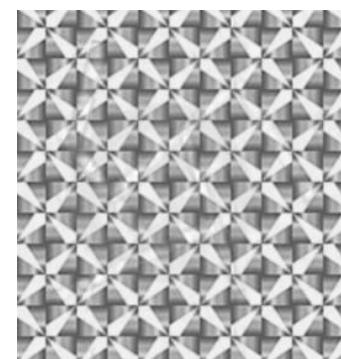
Parametric deformations
(affine, projective, radial distortion, 3D shape...)

Transform Invariant Low-rank Textures (TILT)

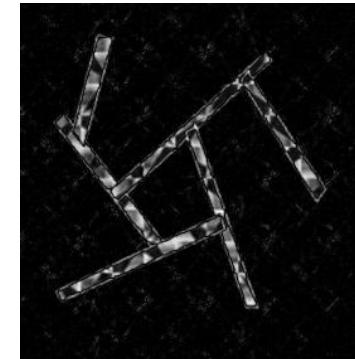
D – deformed observation



A – low-rank structures



E – sparse errors



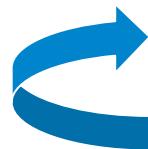
$$\circ \tau =$$

+

Objective: *Transformed Principal Component Pursuit*:

$$\min \|A\|_* + \lambda \|E\|_1 \quad \text{subj } A + E = D \circ \tau$$

Solution: Iteratively solving the linearized convex program:



$$\min \|A\|_* + \lambda \|E\|_1 \quad \text{subj } A + E = D \circ \tau_k + J \cdot \Delta \tau$$



Or reduced version: subj $\mathcal{P}_Q[A + E] = \mathcal{P}_Q[D \circ \tau_k], \mathcal{P}_Q[J] = 0$

THEORY – Compressive Robust PCA

Theorem 5 (Compressive Principal Component Pursuit). Let $\mathbf{A}_0 \in \mathbb{R}^{m \times n}$, $m \geq n$ have rank $r \leq \rho_r \frac{m}{\mu^2 \log^2(n)}$, and \mathbf{E}_0 have a Bernoulli support with error probability $\rho < \rho^*$. Let Q^\perp be a random subspace of $\mathbb{R}^{m \times n}$ of dimension

$$\dim(Q) \geq C_Q(\rho mn + mr) \cdot \log^2 m,$$

distributed according to the Haar measure, independent of the support of E_0 . Then with very high probability

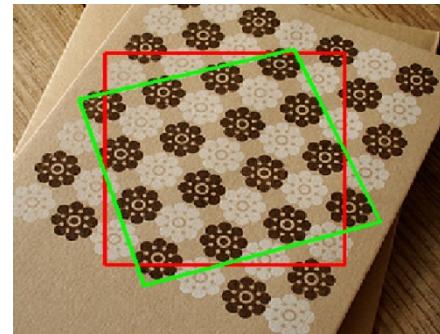
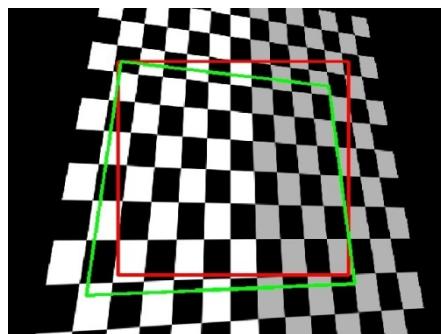
$$(\mathbf{A}_0, \mathbf{E}_0) = \arg \min \|\mathbf{A}\|_* + \frac{1}{\sqrt{m}} \|\mathbf{E}\|_1 \quad \text{subj } \mathcal{P}_Q[\mathbf{A} + \mathbf{E}] = \mathcal{P}_Q[\mathbf{A}_0 + \mathbf{E}_0],$$

for some numerical constant ρ_r , C_p and ρ^* , and the minimizer is unique.

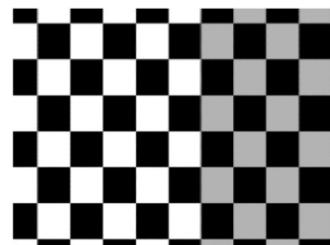
A nearly optimal lower bound on minimum # of measurements!

TILT: *Shape from texture*

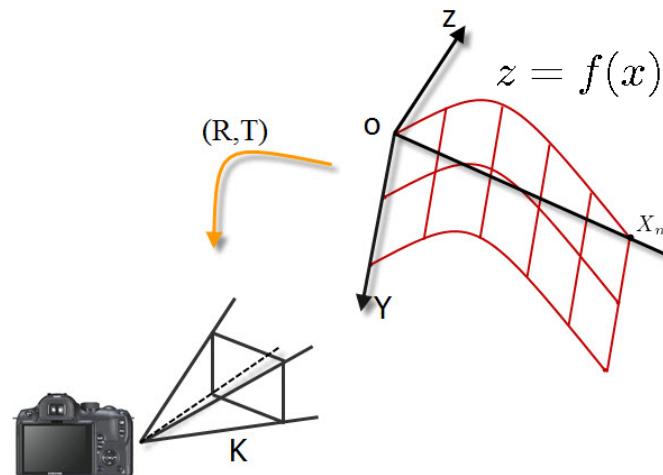
Input (red window D)



Output (rectified green window A)

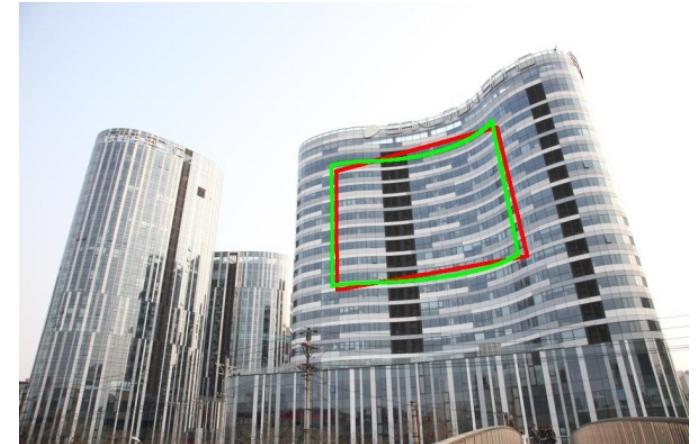


TILT: Shape and geometry from textures



$$D \circ \tau = A + E$$

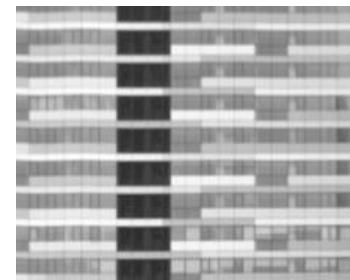
$$\tau = (K, R, T, \{a_i\})$$



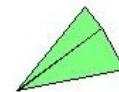
D



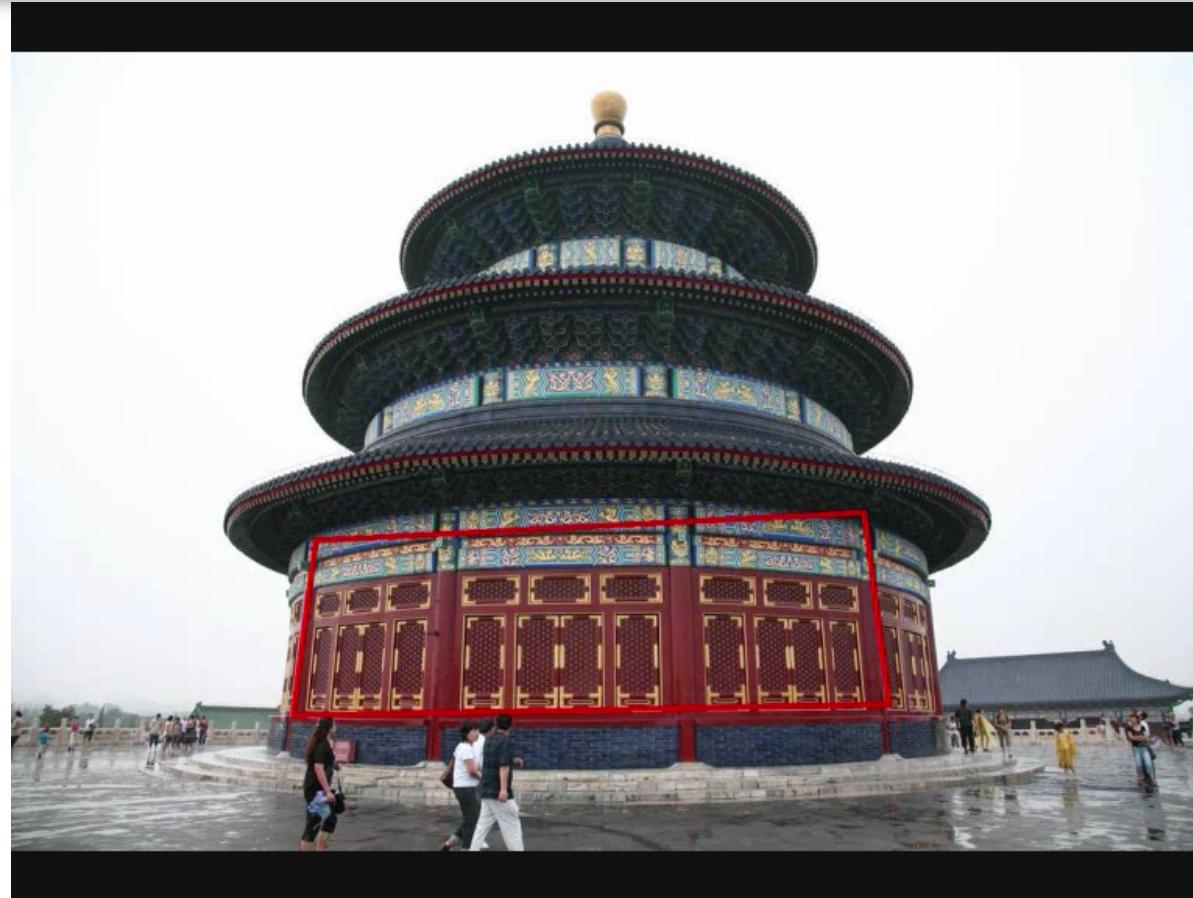
A



E



TILT: *Shape and geometry from textures*

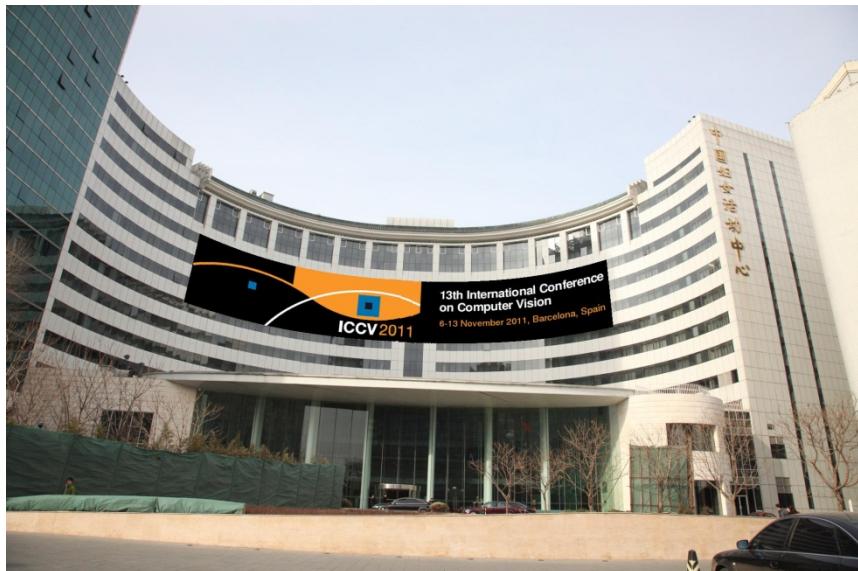
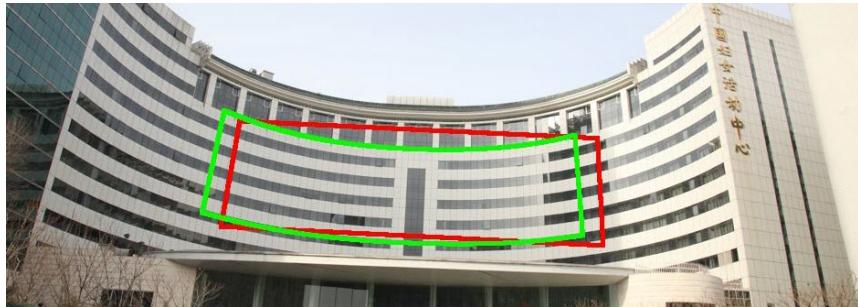


360° panorama



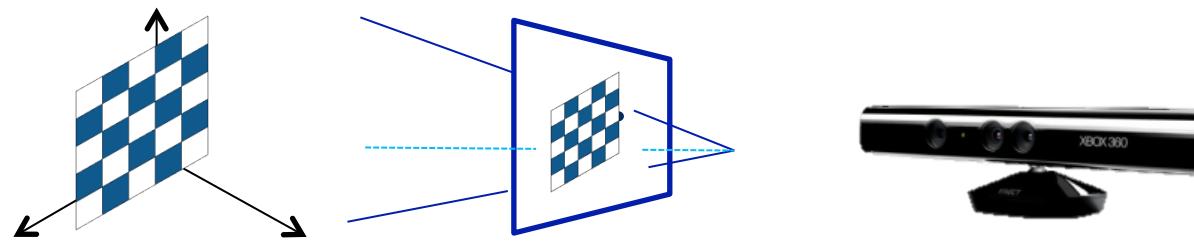
Zhang, Liang, and Ma, in ICCV 2011

TILT: Virtual reality



Zhang, Liang, and Ma, in ICCV 2011

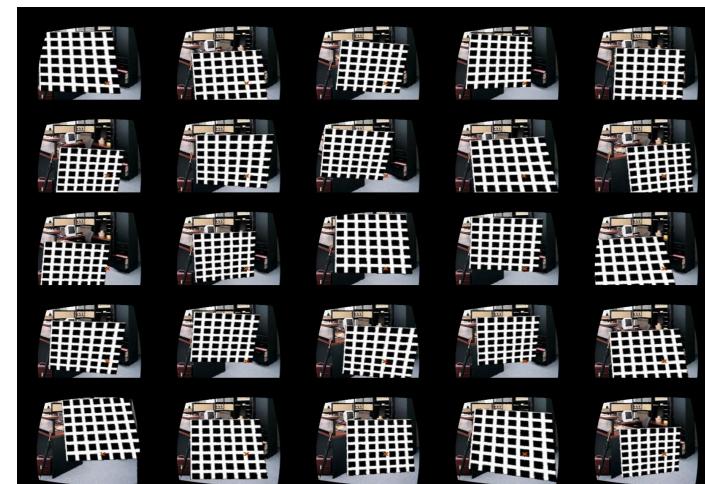
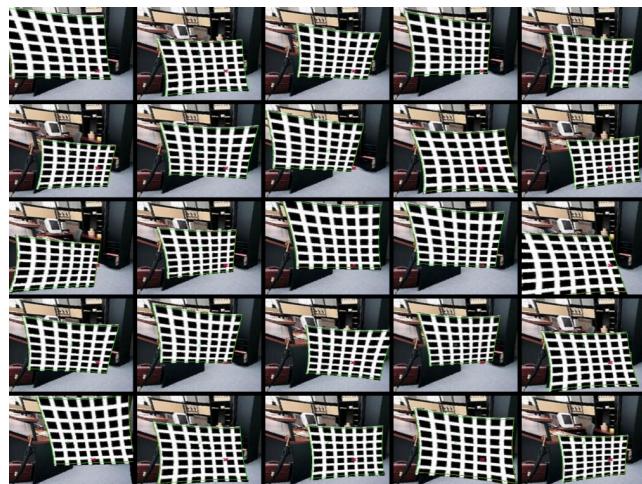
TILT: Camera Calibration with Radial Distortion



$$r = \sqrt{x_0^2 + y_0^2}, f(r) = 1 + kc(1)r^2 + kc(2)r^4 + kc(5)r^6$$

$$\begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} f(r)x_0 + 2kc(3)x_0y_0 + kc(4)(r^2 + 2x_0^2) \\ f(r)y_0 + 2kc(4)x_0y_0 + kc(3)(r^2 + 2y_0^2) \end{pmatrix}$$

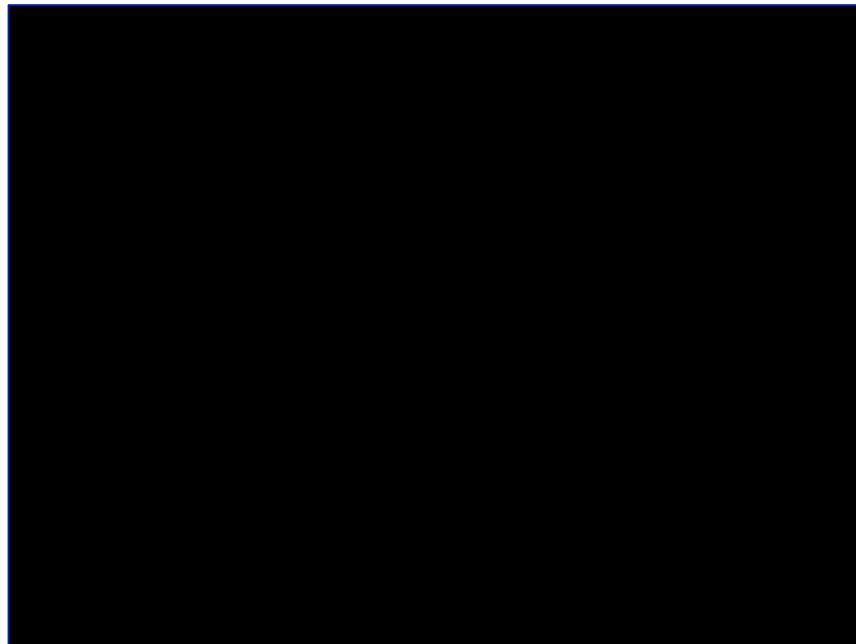
$$K = \begin{bmatrix} f_x & \theta & o_x \\ 0 & f_y & o_y \\ 0 & 0 & 1 \end{bmatrix}$$



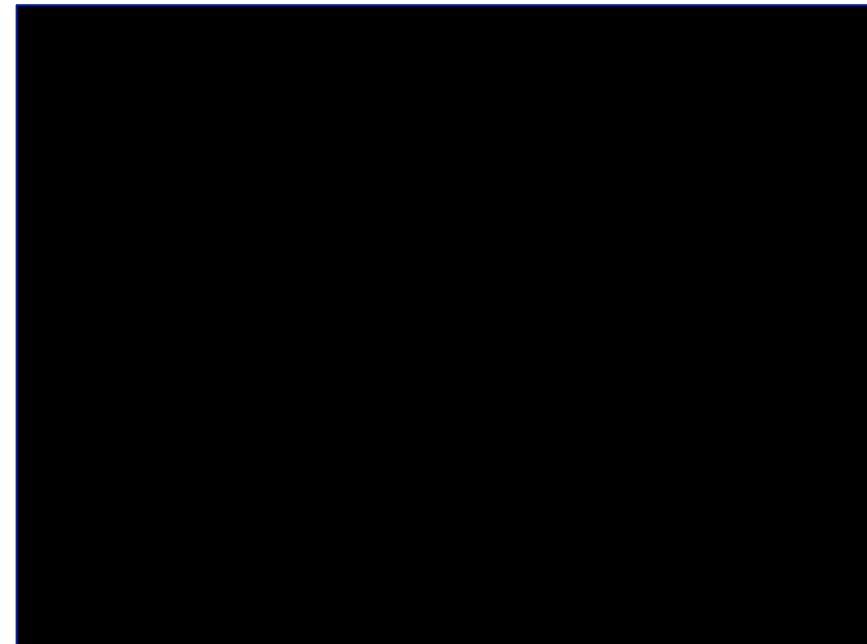
TILT: *Camera Calibration with Radial Distortion*

$$\begin{aligned} \min \sum_{i=1}^N \|A_i\|_* + \lambda \|E_i\|_1 \quad & \text{subj } A_i + E_i = D \circ (\tau_0, \tau_i) \\ \tau_0 = (K, K_c), \quad & \tau_i = (R_i, T_i). \end{aligned}$$

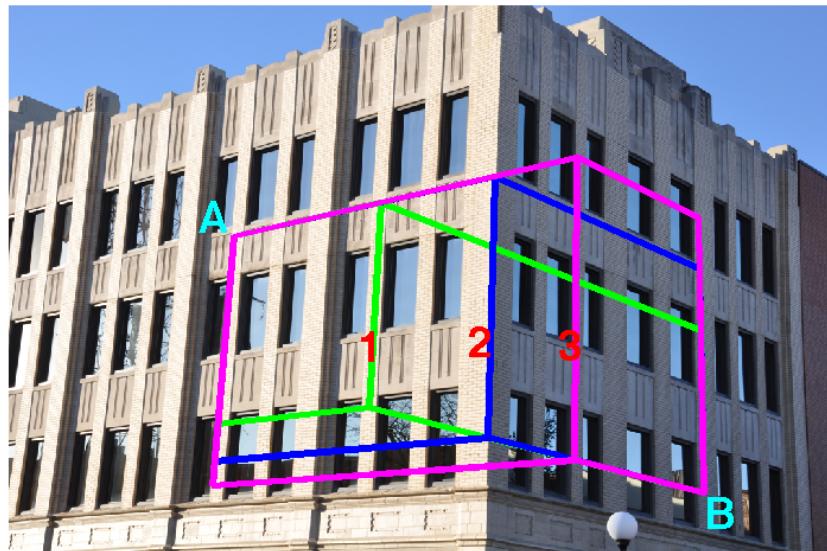
Previous approach



Low-rank method

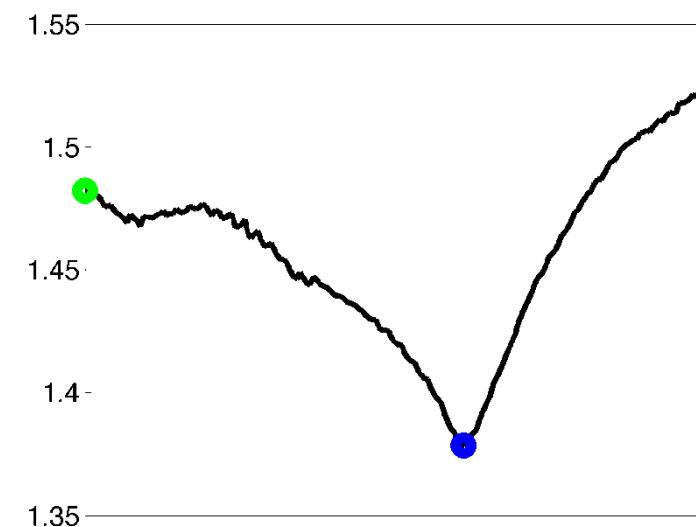


TILT: Holistic 3D Reconstruction of Urban Scenes



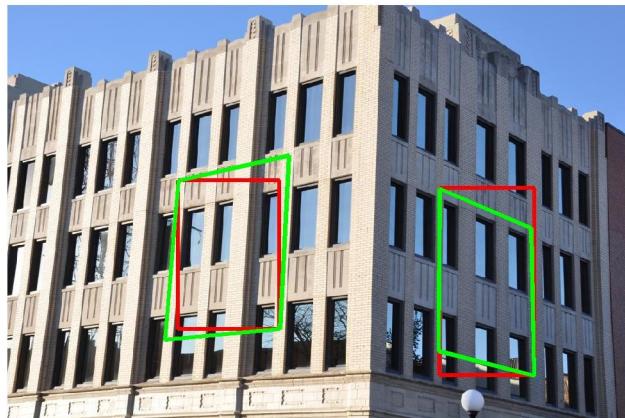
$$\min \|A\|_* + \|E\|_1 \quad \text{s.t.}$$

$$A + E = [D_1 \circ \tau_1, D_2 \circ \tau_2]$$

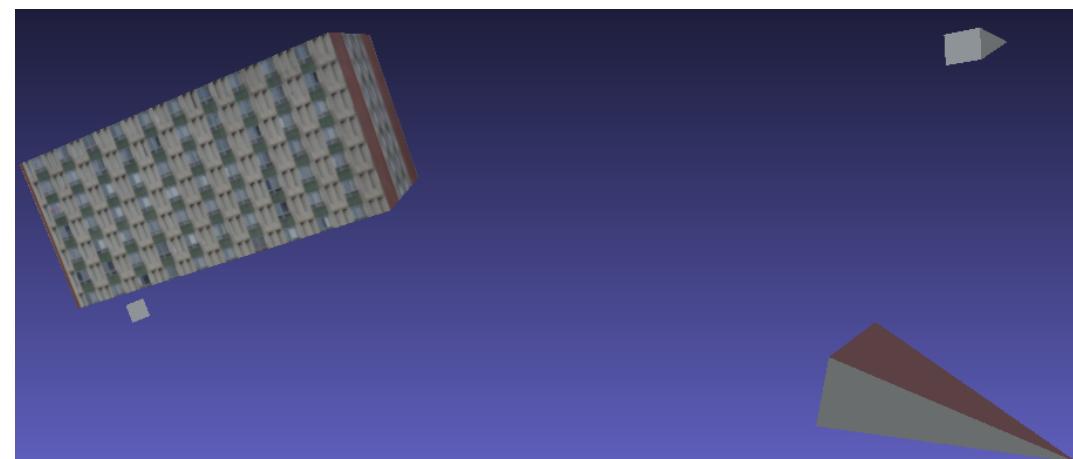
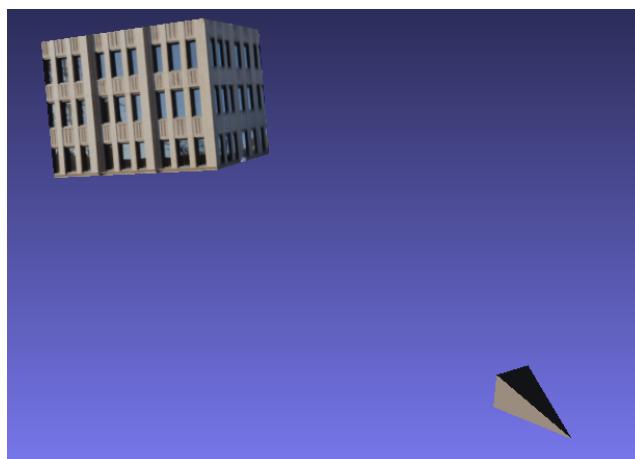


TILT: *Holistic 3D Reconstruction of Urban Scenes*

From one input image

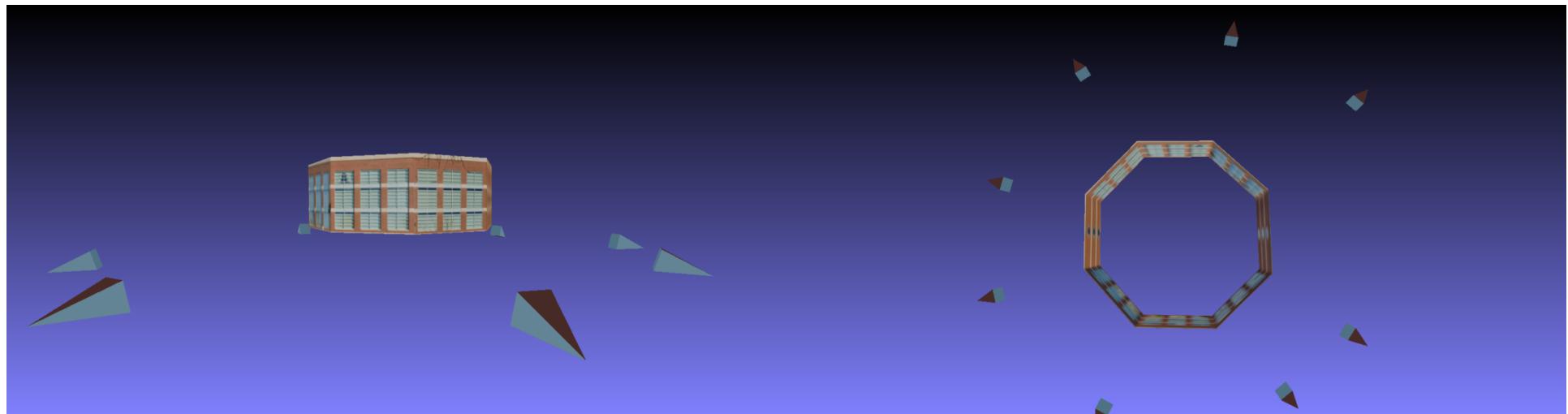
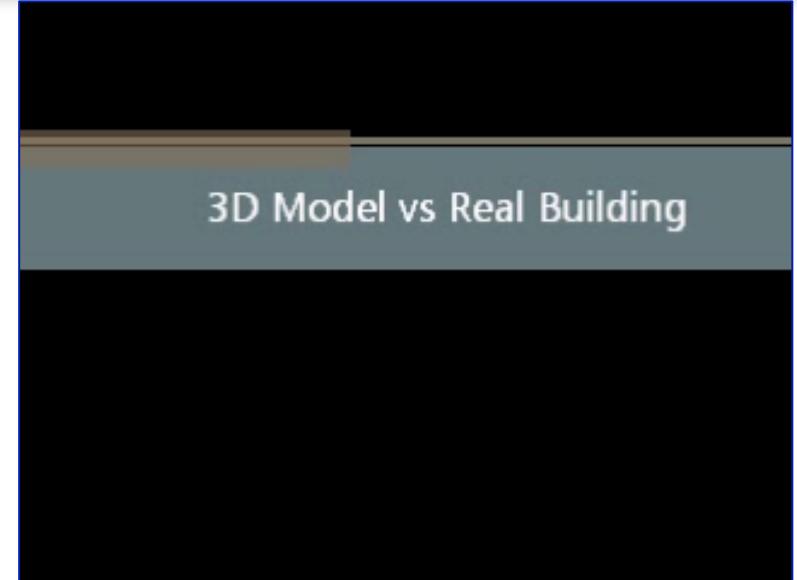


From four input images



TILT: *Holistic 3D Reconstruction of Urban Scenes*

From eight input images



Mobahi, Zhou, and Ma, in ICCV 2011

Virtual reality in urban scenes

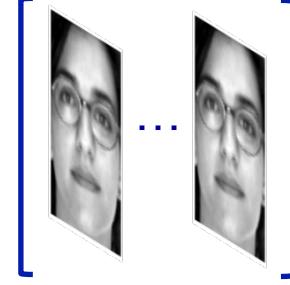


Registering Multiple Images: Robust Alignment

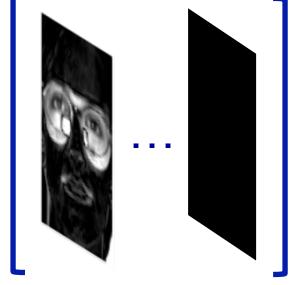
D – corrupted & misaligned observation



A – aligned low-rank signals



E – sparse errors



$\circ \tau =$

$+$

Problem: Given $D \circ \tau = A_0 + E_0$, recover τ , A_0 and E_0 .

Parametric deformations
(rigid, affine, projective...)

Low-rank component

Sparse component

Solution: Robust Alignment via Low-rank and Sparse (**RASL**) Decomposition

Iteratively solving the linearized convex program:



$$\min \|A\|_* + \lambda \|E\|_1 \text{ subj } A + E = D \circ \tau_k + J \Delta \tau \\ (\text{or } Q(A + E) = QD \circ \tau_k, QJ = 0)$$

RASL: Aligning Face Images from the Internet



*48 images collected from internet

Peng, Ganesh, Wright, Ma, CVPR'10, TPAMI'11

RASL: *Faces Detected*

Input: faces detected by a face detector (D)



Average



RASL: *Faces Aligned*

Output: aligned faces ($D \circ \tau$)



Average



RASL: *Faces Repaired and Cleaned*

Output: clean low-rank faces (A)

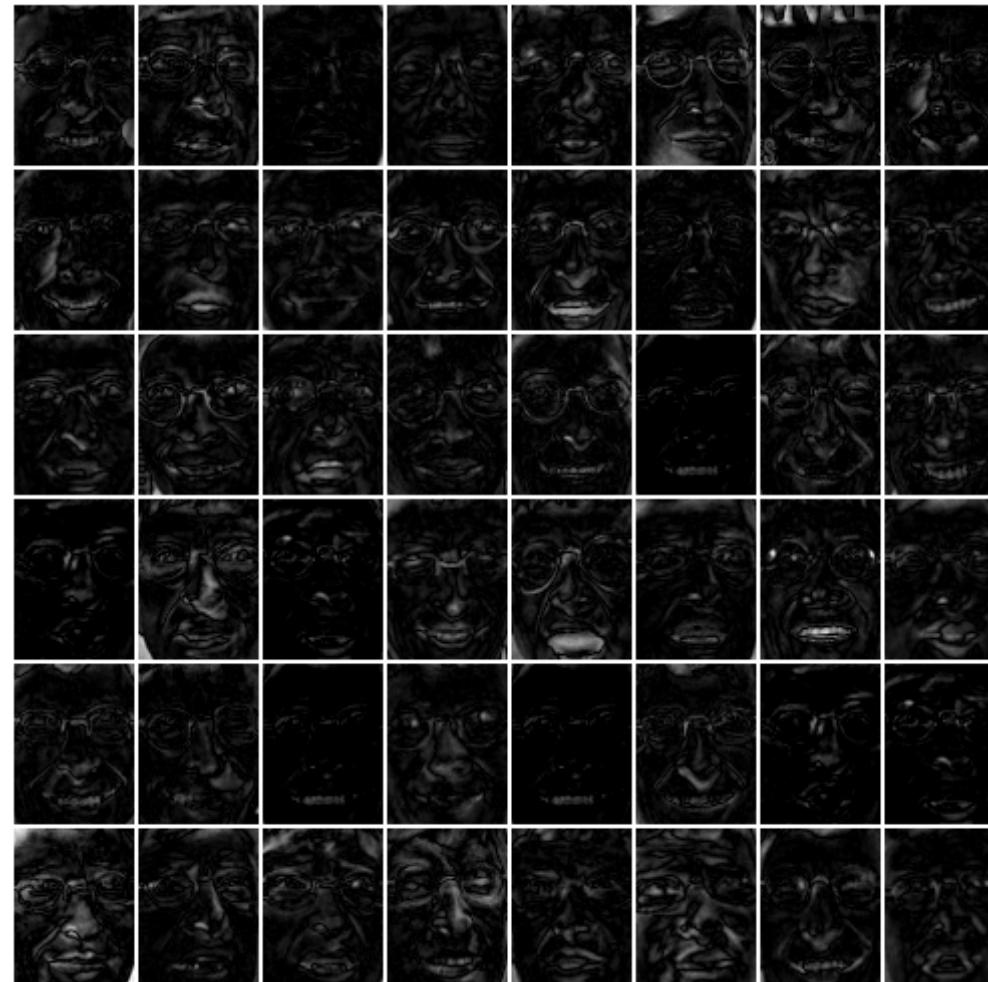


Average



RASL: *Sparse Errors of the Face Images*

Output: sparse error images (E)



RASL: *Video Stabilization and Enhancement*

Original video (D) Aligned video ($D \circ \tau$) Low-rank part (A) Sparse part (E)



RASL: Aligning Handwritten Digits

D



Learned-Miller PAMI'06



Vedaldi CVPR'08



$D \circ \tau$



A

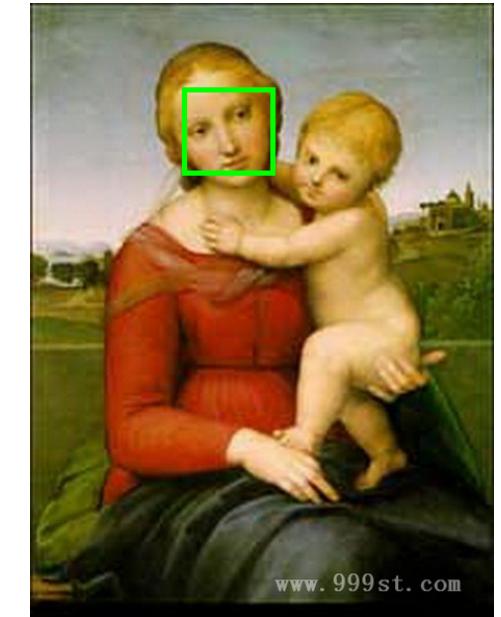
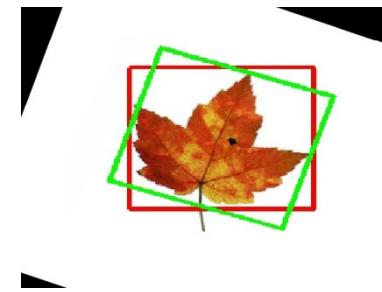
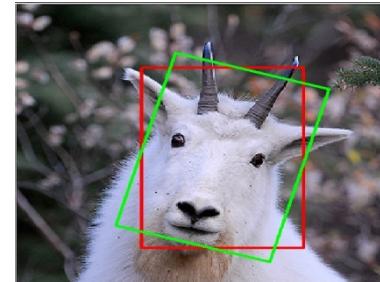
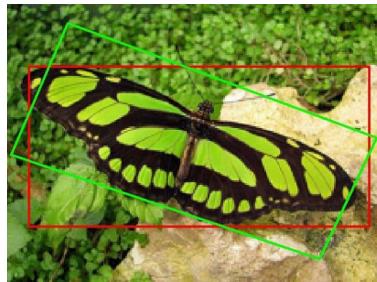


E



Object Recognition: *Rectifying Pose of Objects*

Input (red window D)

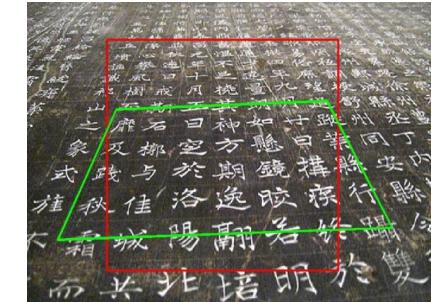
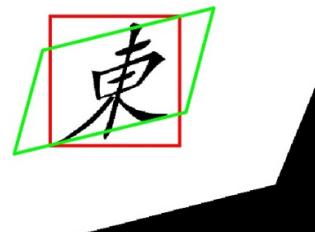


Output (rectified green window A)

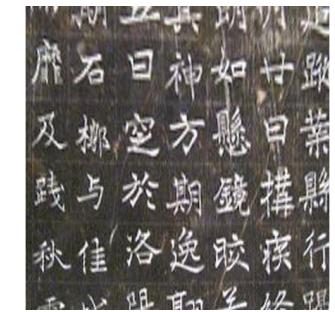


Object Recognition: *Regularity of Texts at All Scales!*

Input (red window D)



Output (rectified green window A)



Recognition: Character/Text Rectification



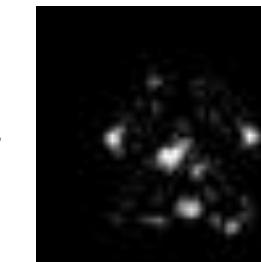
D



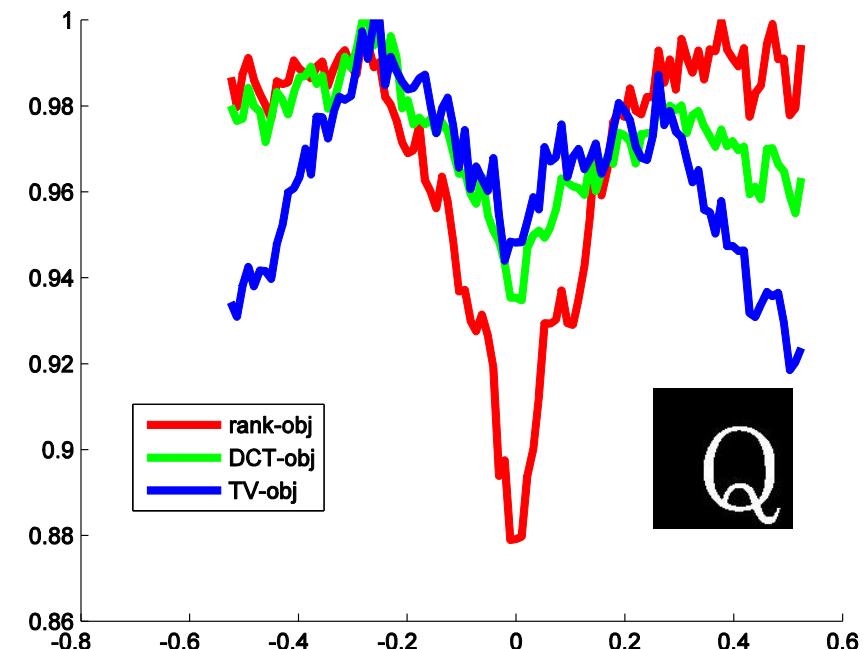
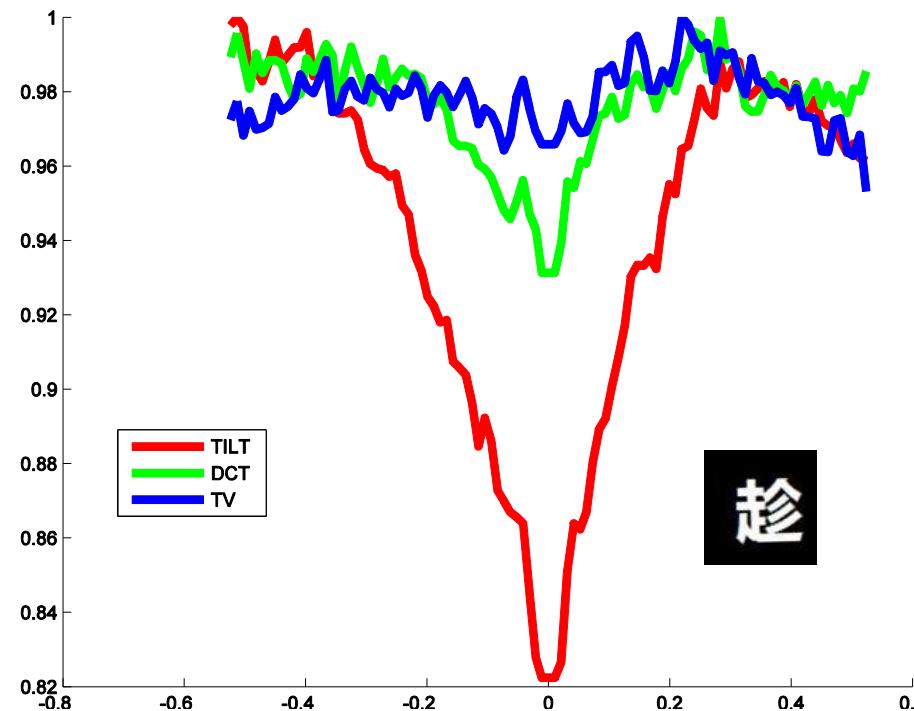
$D \circ \tau$



A



E

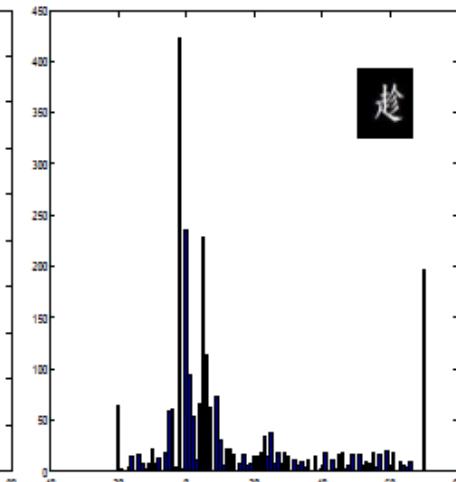
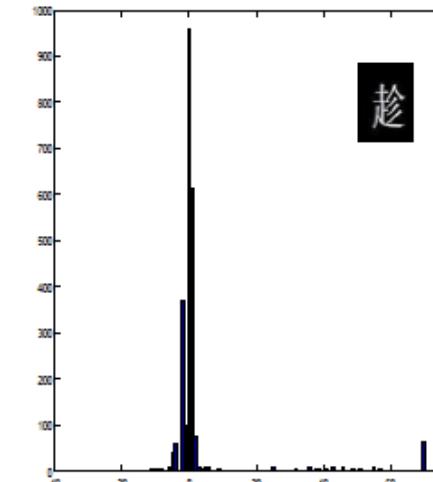
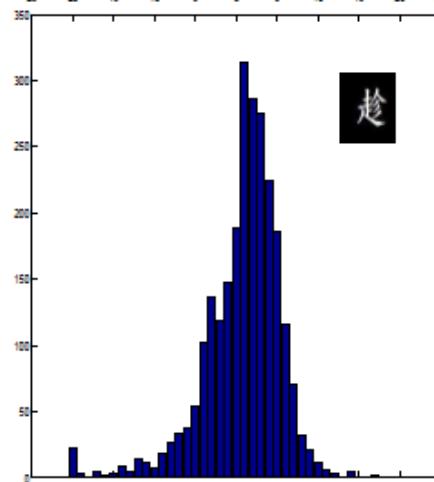
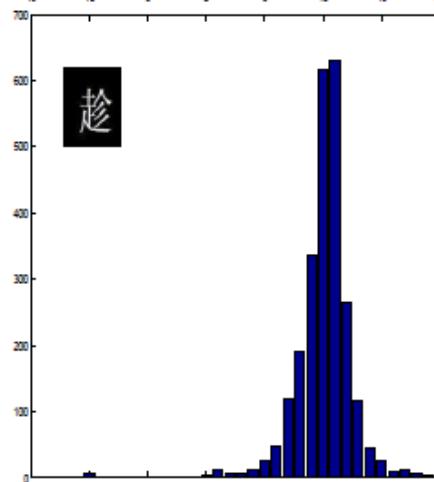
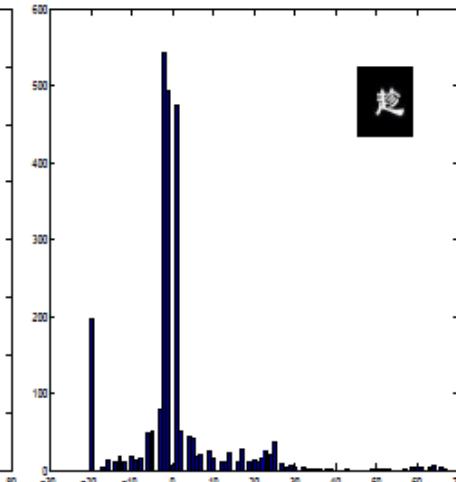
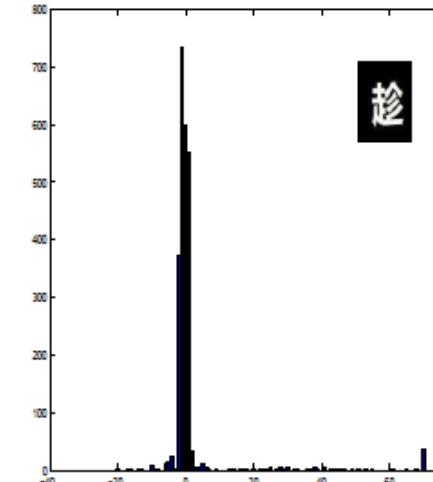
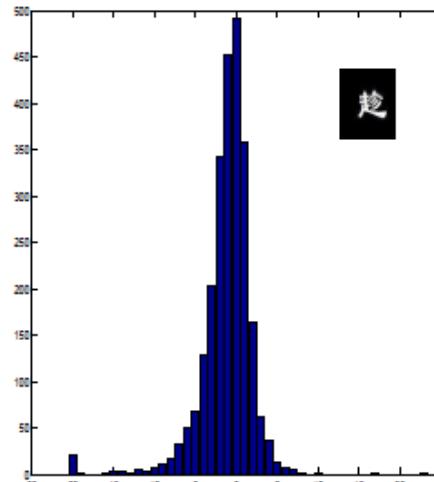
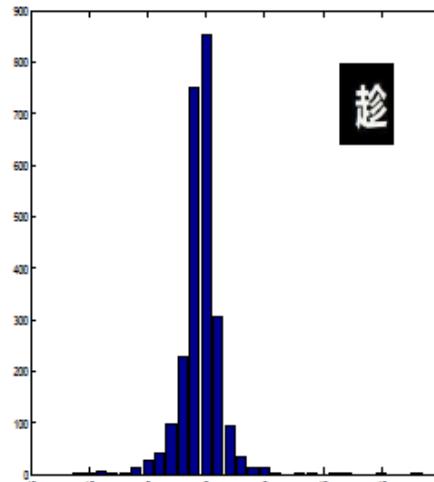


Recognition: *Character/Text Rectification*

TILT

versus

Hough Transform



Recognition: Street Sign Rectification



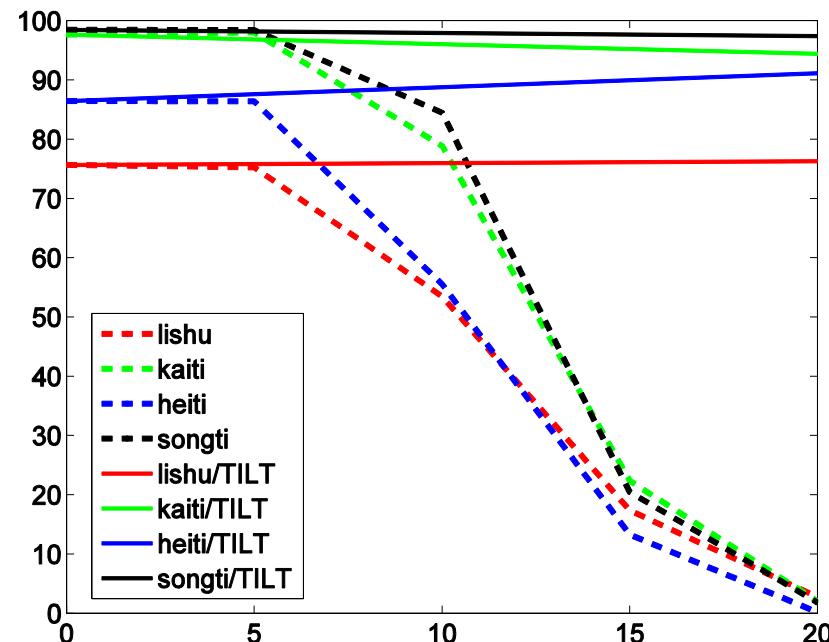
会堂西路
A₁ | A₂ | A₃ | A₄

$$\min \sum_{i=1}^4 \|A_i\|_* + \lambda \|E_i\|_1$$

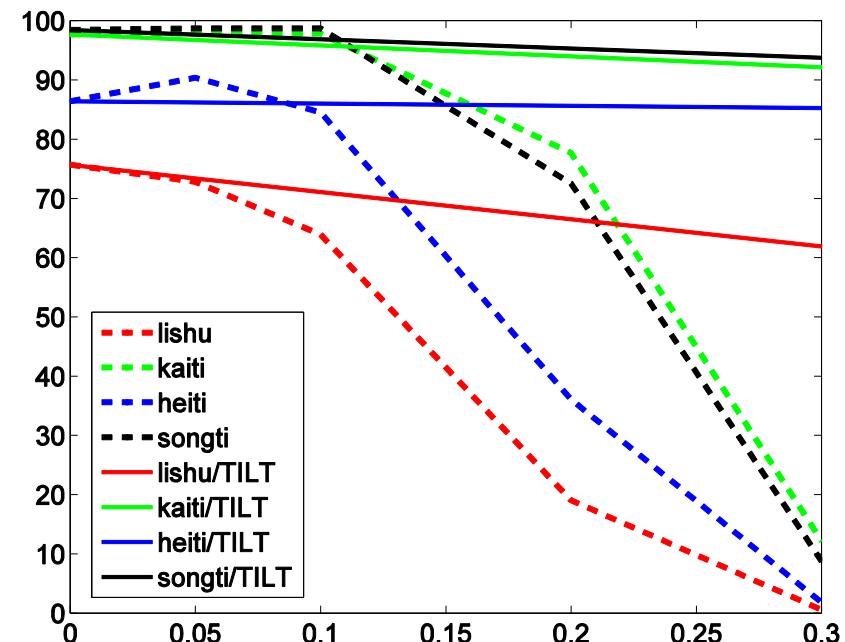
$$\text{subj } D \circ \tau = [A_1 \cdots A_4] + [E_1 \cdots E_4].$$

Recognition: Character Rectification and Recognition

Microsoft OCR for rotated characters
(2,500 common Chinese characters)



Microsoft OCR for skewed characters
(2,500 common Chinese characters)



Take-home Messages for Visual Data Processing:

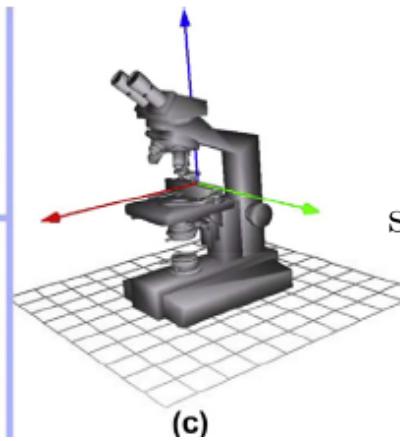
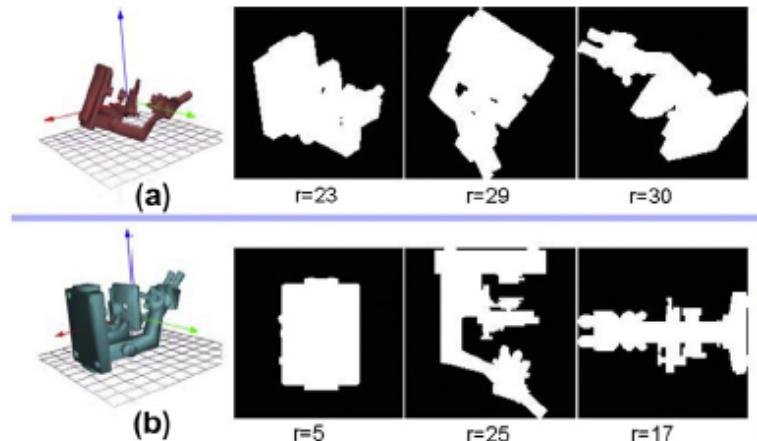
1. (Transformed) **low-rank and sparse** structures are central to visual data modeling, processing, and analyzing;
2. Such structures can now be extracted **correctly, robustly, and efficiently**, from raw image pixels (or high-dim features);
3. These new algorithms **unleash tremendous local or global information** from single or multiple images, emulating or surpassing human capability;
4. These algorithms start to exert significant impact on **image/video processing, 3D reconstruction, and object recognition**.

....

But try not to abuse or misuse them...

Other Applications: *Upright orientation of man-made objects*

TILT for 3D: Unsupervised upright orientation of man-made 3D objects



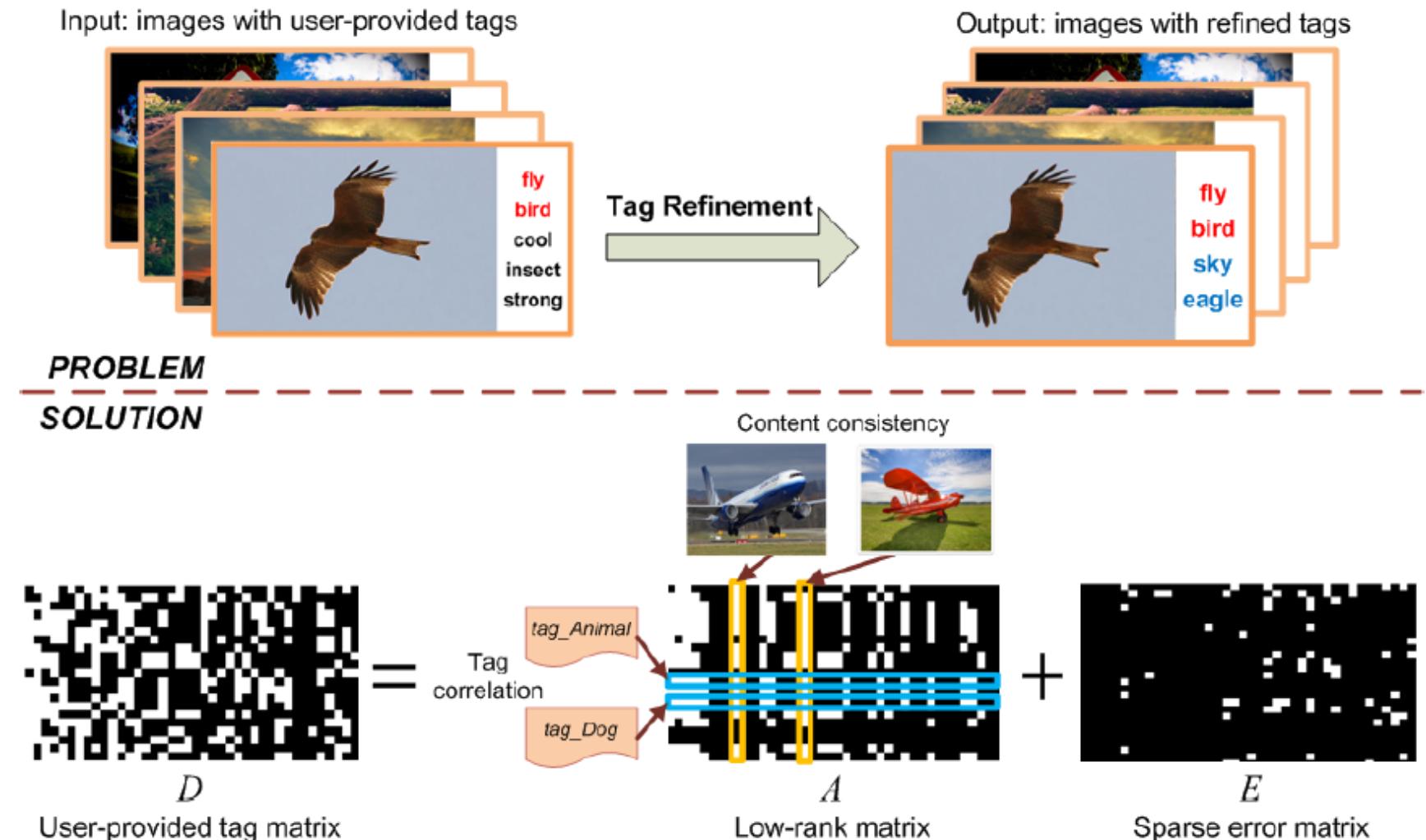
$$\min \sum_{i=1}^3 \|A_i\|_* + \lambda \|E_i\|_1$$

$$\text{st } D \circ \tau = [A_1, A_2, A_3] + [E_1, E_2, E_3].$$



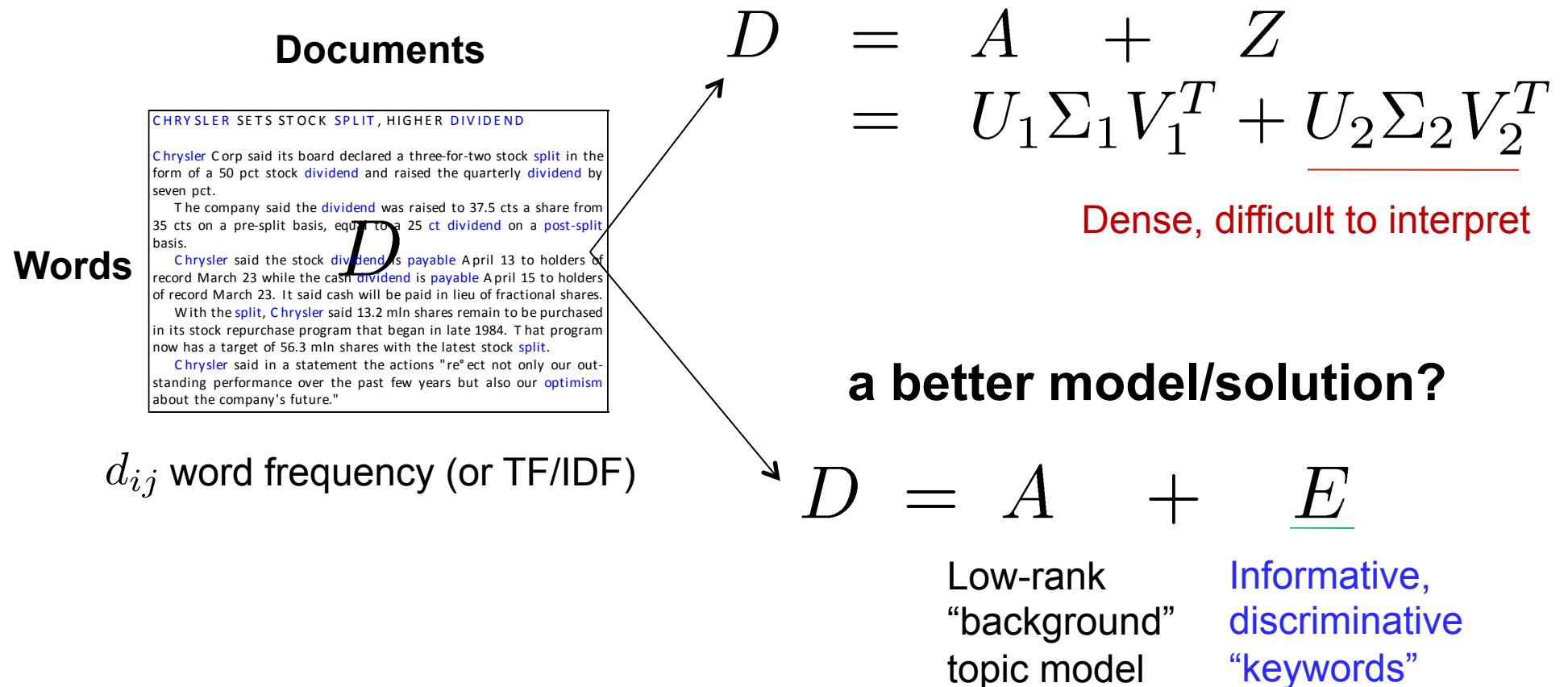
Fig. 10. More models which have been successfully tested through our algorithm.

Other Data/Applications: Web Image/Tag Refinement



Other Data/Applications: Web Document Corpus Analysis

Latent Semantic Indexing: the classical solution (PCA)



Other Data/Applications: Sparse Keywords Extracted

Reuters-21578 dataset: 1,000 longest documents; 3,000 most frequent words

CHRYSLER SETS STOCK SPLIT, HIGHER DIVIDEND

Chrysler Corp said its board declared a three-for-two stock **split** in the form of a 50 pct stock **dividend** and raised the quarterly **dividend** by seven pct.

The company said the **dividend** was raised to 37.5 cts a share from 35 cts on a pre-split basis, equal to a 25 **ct dividend** on a **post-split** basis.

Chrysler said the stock **dividend** is **payable** April 13 to holders of record March 23 while the cash **dividend** is **payable** April 15 to holders of record March 23. It said cash will be paid in lieu of fractional shares.

With the **split**, Chrysler said 13.2 mln shares remain to be purchased in its stock repurchase program that began in late 1984. That program now has a target of 56.3 mln shares with the latest stock **split**.

Chrysler said in a statement the actions "reflect not only our outstanding performance over the past few years but also our **optimism** about the company's future."

Other Data/Applications: Protein-Gene Correlation

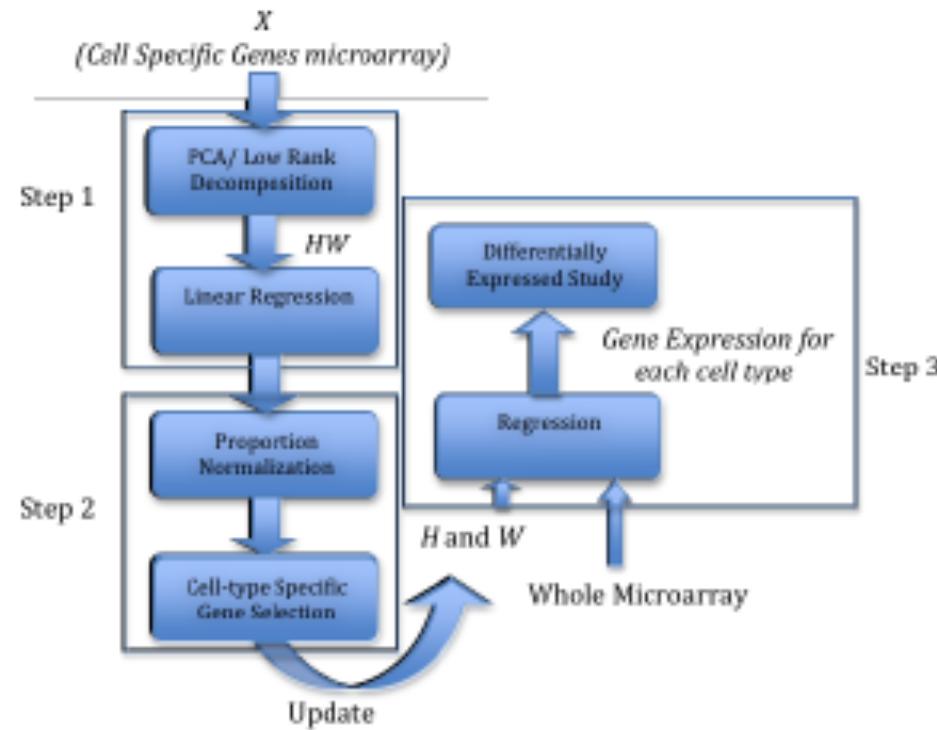


Fig. 1. The diagram of the workflow of the method presented in this paper.

Microarray data

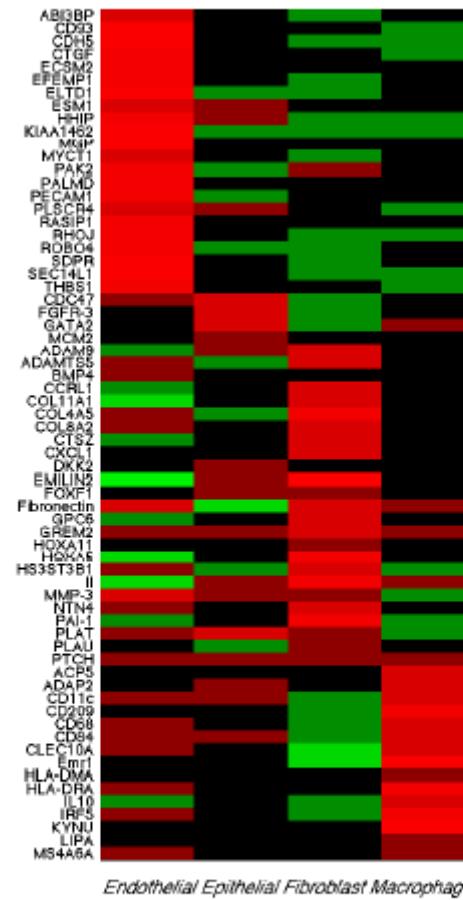


Fig. 6. HeatMap of estimated gene signatures for the sorted cell specific genes after adjustments based on fold changes. RPCA is used in the first step. It is clear that this matrix is close to a block diagonal structure.

Other Data: Time Series Gene Expressions

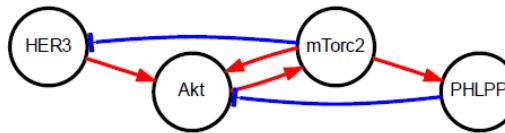


Figure S3. Abstract HER2 overexpressed breast cancer model by Dr. Moasser.

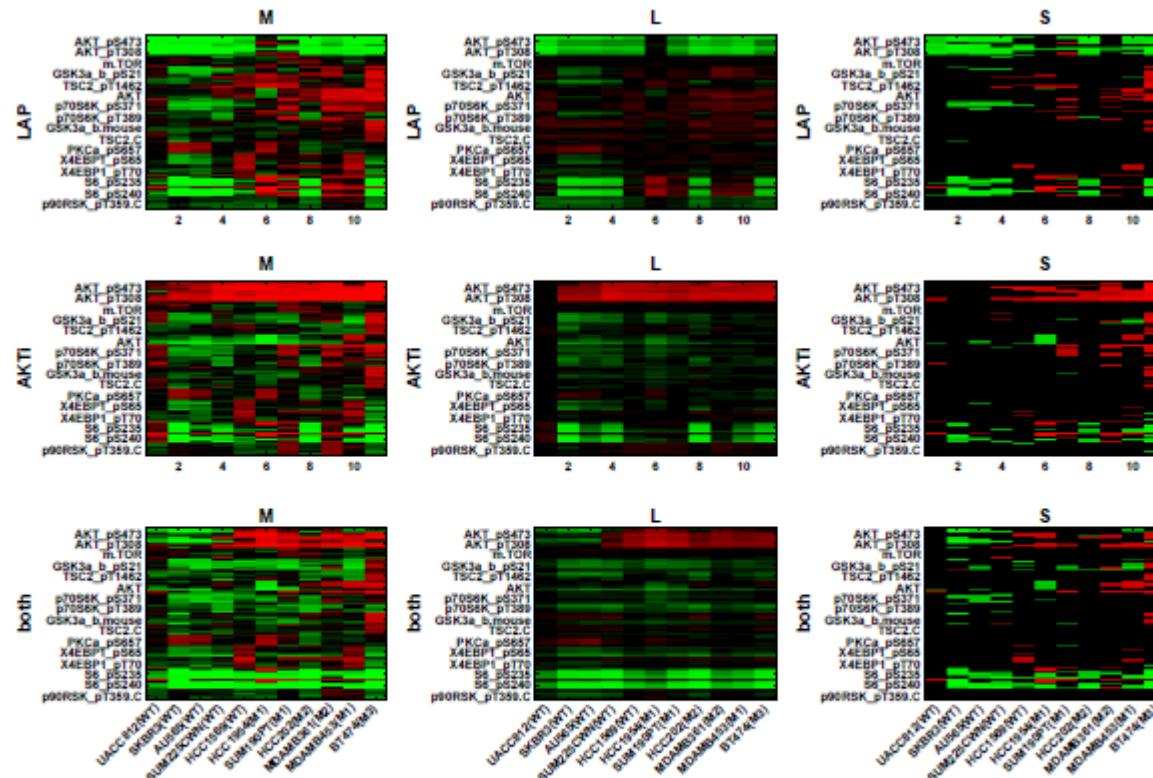
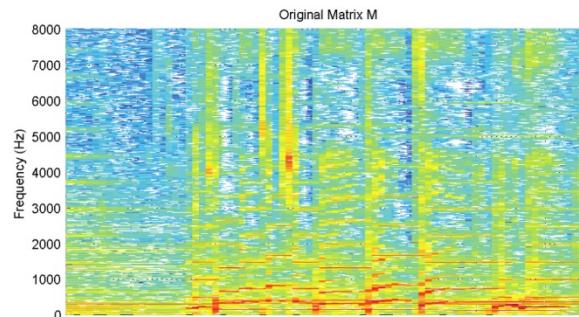


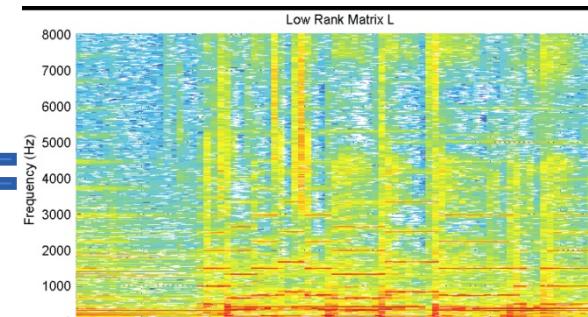
Figure S4. Separation result: (1_{st} column) raw data (2_{nd} column) low-rank component and (3_{rd} column) highly corrupted sparse component using threshold (M1: H1047R (kinase domain mutation) M2: E545K (helical domain mutation), and M3: K111N mutation in PIK3CA).

Other Data/Applications: Lyrics and Music Separation

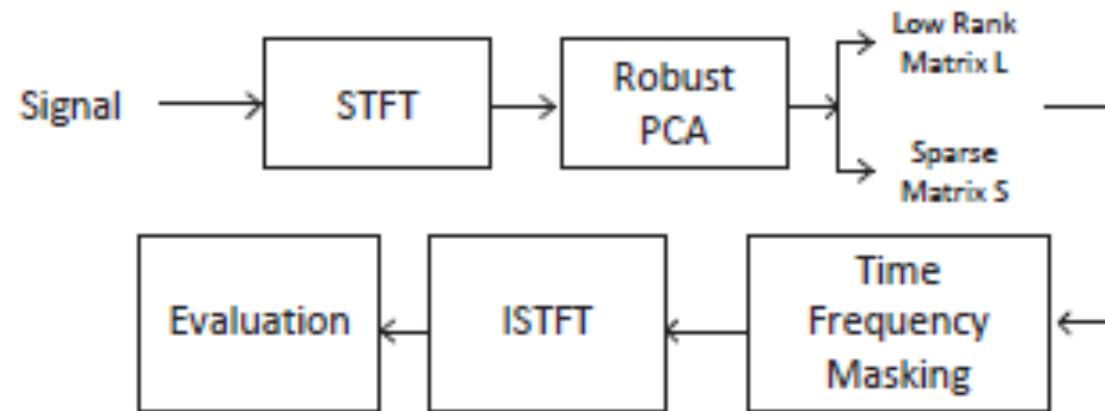
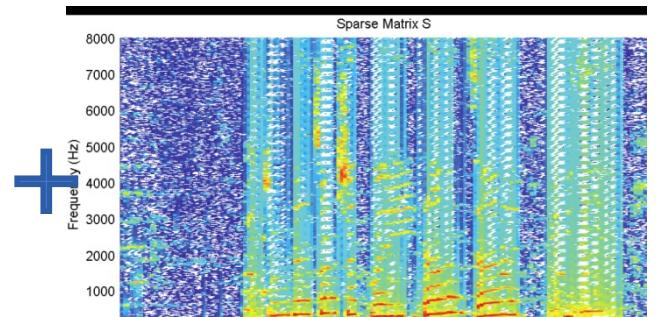
Songs (STFT)



Low-rank (music)



Sparse (voices)



Other Data/Applications: Internet Traffic Anomalies

Network Traffic = Normal Traffic + Sparse Anomalies + Noise

$$D = L + RS + N$$

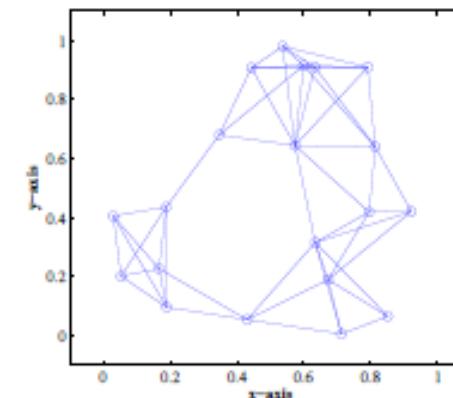
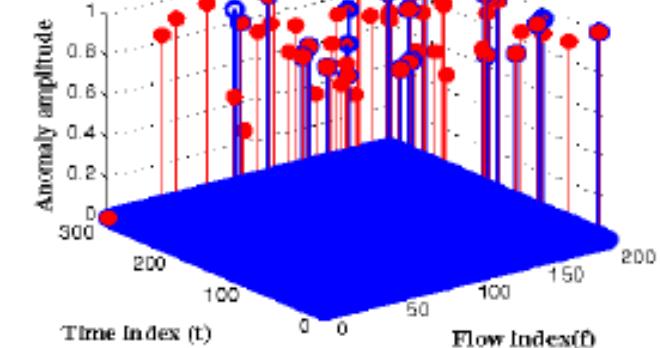
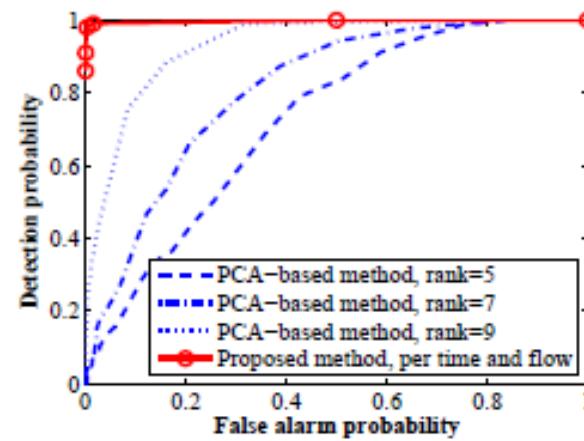


Fig. 2. Network topology graph.



Other Data/Applications: Robust Filtering and System ID



GPS on a Car:

$$\begin{cases} \dot{x} = Ax + Bu, & A \in \Re^{r \times r} \\ y = Cx + z + e \end{cases}$$

gross sparse errors
(due to buildings, trees...)

Robust Kalman Filter: $\hat{x}_{t+1} = Ax_t + K(y_t - C\hat{x}_t)$

Robust System ID:

$$\left[\begin{array}{ccccc} y_n & y_{n-1} & y_{n-2} & \cdots & y_0 \\ y_{n-1} & y_{n-2} & \cdots & \ddots & y_{-1} \\ y_{n-2} & \cdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & y_{-n+2} \\ y_0 & y_{-1} & \cdots & y_{-n+2} & y_{-n+1} \end{array} \right] = O_{n \times r} X_{r \times n} + S$$

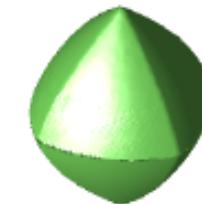
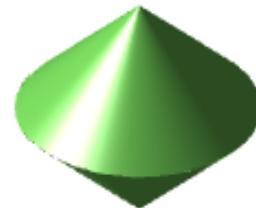
Hankel matrix

CONCLUSIONS – A Unified Theory for Sparsity and Low-Rank

	Sparse Vector	Low-Rank Matrix
Low-dimensionality of	individual signal	correlated signals
Measure	L_0 norm $\ x\ _0$	$\text{rank}(X)$
Convex Surrogate	L_1 norm $\ x\ _1$	Nuclear norm $\ X\ _*$
Compressed Sensing	$y = Ax$	$Y = A(X)$
Error Correction	$y = Ax + e$	$Y = A(X) + E$
Domain Transform	$y \circ \tau = Ax + e$	$Y \circ \tau = A(X) + E$
Mixed Structures	$Y = A(X) + B(E) + Z$	

Compressive Sensing of Low-Dimensional Structures

$$L \quad x + e$$



A norm $\|\cdot\|$ is said to be **decomposable** at \mathbf{X} if there exists a subspace T and a matrix \mathbf{S} such that

$$\partial\|\cdot\|(\mathbf{X}) = \{\Lambda \mid \mathcal{P}_T(\Lambda) = \mathbf{S}, \|P_{T^\perp}(\Lambda)\|^* \leq 1\},$$

where $\|\cdot\|^*$ is the dual norm of $\|\cdot\|$, and \mathcal{P}_{T^\perp} is nonexpansive w.r.t. $\|\cdot\|^*$.

Theorem [Candes, Recht'11] Any low-complexity signal \mathbf{X}^0 can be exactly recovered from high compressive measurements via convex optimization:

$$\|\mathbf{X}\|_\diamond \quad \text{subject to} \quad \mathcal{P}_Q(\mathbf{X}) = \mathcal{P}_Q(\mathbf{X}^0),$$

for a decomposable norm $\|\cdot\|_\diamond$.

Compressive Sensing and Separation of Low-dim Structures

Suppose $(\mathbf{X}_1^0, \dots, \mathbf{X}_k^0) = \arg \min \sum_{i=1}^k \lambda_i \|\mathbf{X}_i\|_{(i)}$ subj $\sum_{i=1}^k \mathbf{X}_i = \sum_{i=1}^k \mathbf{X}_i^0$,
for decomposable norms $\|\cdot\|_{(i)}$ that majorize the Frobenius norm.

Theorem 6 (Compressive Sensing of Mixed Low-Comp. Structures).
Let Q^\perp be a random subspace of $\mathbb{R}^{m \times n}$ of dimension

$$\dim(Q) \geq O(\log^2 m) \times \text{intrinsic degrees of freedom of } (\mathbf{X}_1, \dots, \mathbf{X}_k),$$

distributed according to the Haar measure, independent of \mathbf{X}_i . Then with very high probability

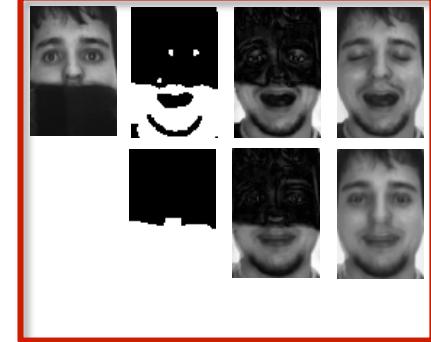
$$(\mathbf{X}_1^0, \dots, \mathbf{X}_k^0) = \arg \min \sum_{i=1}^k \lambda_i \|\mathbf{X}_i\|_{(i)} \quad \text{subj} \quad \mathcal{P}_Q \left[\sum_{i=1}^k \mathbf{X}_i \right] = \mathcal{P}_Q \left[\sum_{i=1}^k \mathbf{X}_i^0 \right],$$

and the minimizer is unique.

Extensions – A Suite of Powerful Regularizers

For compressive robust recovery of a family of low-dimensional structures:

- [Zhou et. al. '09] Spatially contiguous sparse errors via MRF
- [Bach '10] – relaxations from submodular functions
- [Negahban+Yu+Wainwright '10] – geometric analysis of recovery
- [Becker+Candès+Grant '10] – algorithmic templates
- [Xu+Caramanis+Sanghavi '11] column sparse errors $L_{2,1}$ norm
- [Recht+Parillo+Chandrasekaran+Wilsky '11'12] – compressive sensing of various structures
- [Candes+Recht '11] – **compressive sensing of decomposable structures**



$$X^0 = \arg \min \|X\|_\diamond \quad \text{s.t.} \quad \mathcal{P}_Q(X) = \mathcal{P}_Q(X^0)$$

- [McCoy+Tropp'11, Amenlunxen+McCoy+Tropp'13] – **phase transition for recovery and decomposition of structures**

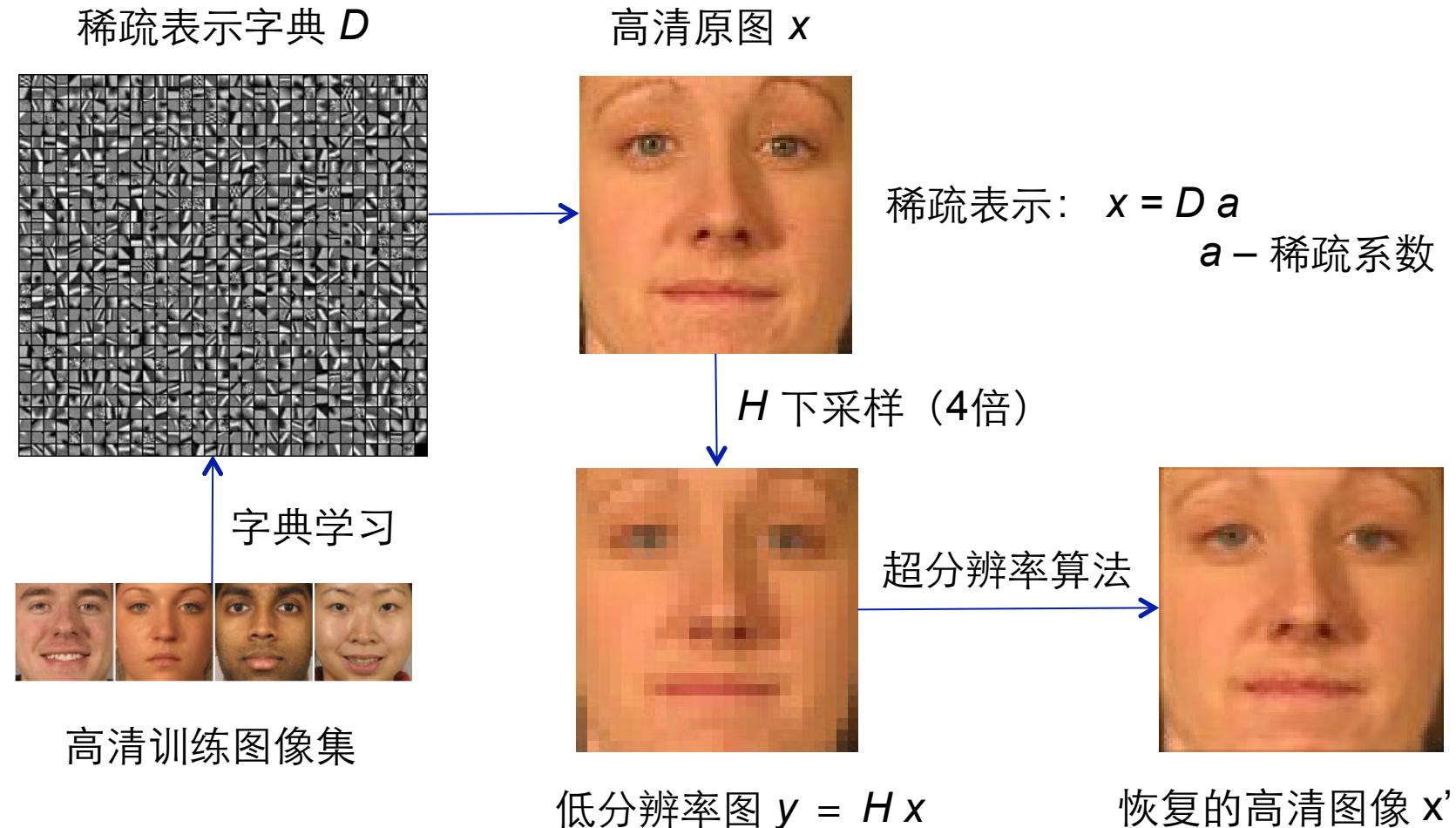
$$(X_1^0, X_2^0) = \arg \min \|X_1\|_{(1)} + \lambda \|X_2\|_{(2)} \quad \text{s.t.} \quad X_1 + X_2 = X_1^0 + X_2^0$$

- [Wright+Ganesh+Min+Ma, ISIT'12, I&I'13] – **compressive superposition of decomposable structures**

$$(X_1^0, \dots, X_k^0) = \arg \min \sum \lambda_i \|X_i\|_{(i)} \quad \text{s.t.} \quad \mathcal{P}_Q(\sum_i X_i) = \mathcal{P}_Q(\sum_i X_i^0)$$

*Take home message: **Let the data and application tell you the structure...***

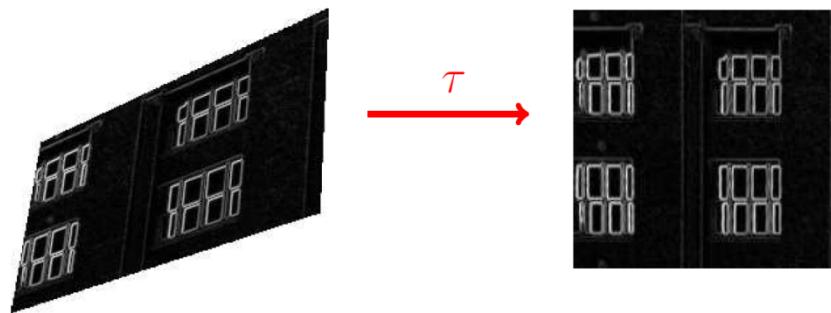
Super Resolution via Transform Invariant Group Sparsity



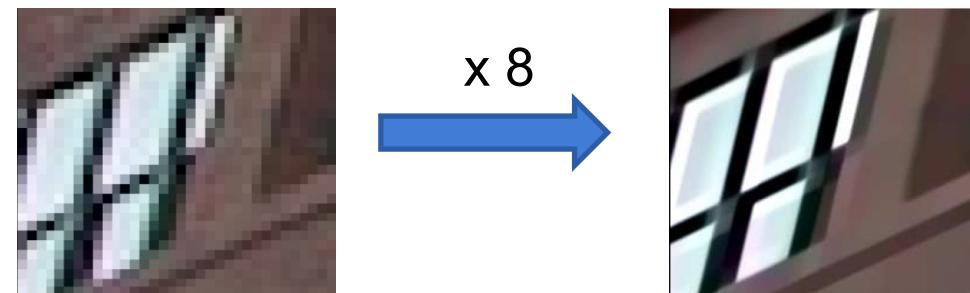
Super Resolution via Transform Invariant Group Sparsity

Aim: Exploiting non-local structures to perform super-resolution at large upsampling factors by

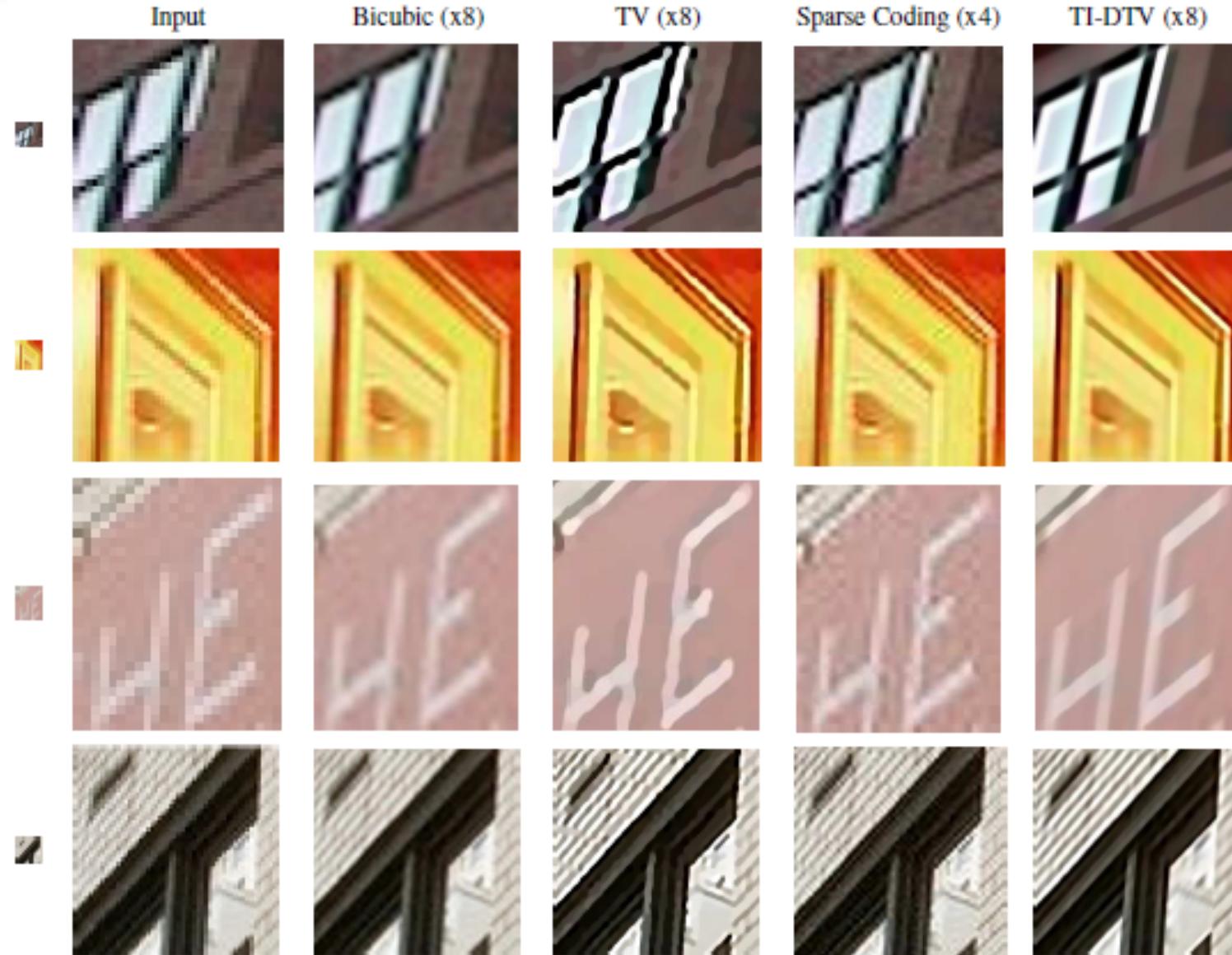
1. Learning the transformation that reveals the group-sparse structure of the image gradient (via TILT)



2. Enforcing this structure through **group-sparse regularizers (DTV)** that incorporates the transform and is consequently invariant to the transform



Super Resolution via Transform Invariant Group Sparsity



Carlos Fernandez and Emmanuel Candes of Stanford, ICCV2013

Image Com. Sensing via Nonlocal Low-rank Regularization

Aim: Exploiting **sparse and non-local structures** to recover natural images from compressive measurements

1. Sparsity (in some domain):

$$\mathbf{x} = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x}\|_0, \text{ s.t. } \mathbf{y} = \Phi \mathbf{x},$$

2. Non-local similarity of patches $\mathbf{X}_i = [\mathbf{x}_{i_0}, \mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{m-1}}]$

$$\min_{\mathbf{L}_i} \|\mathbf{X}_i - \mathbf{L}_i\|_F^2 + \frac{\lambda}{\eta} \sum_{j=1}^{n_0} \log(\sigma_j(\mathbf{L}_i) + \varepsilon).$$

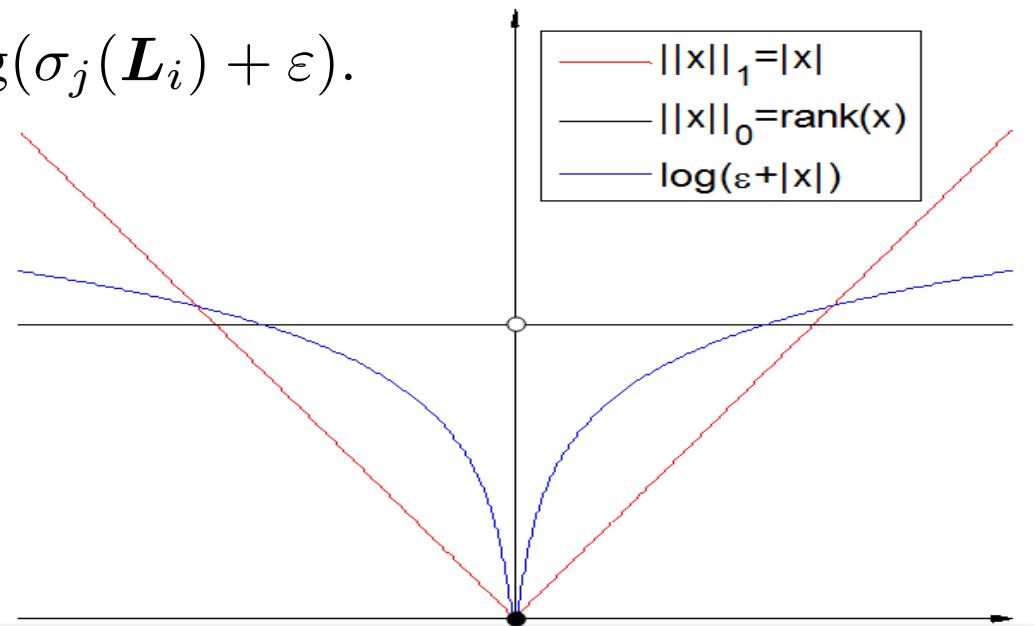
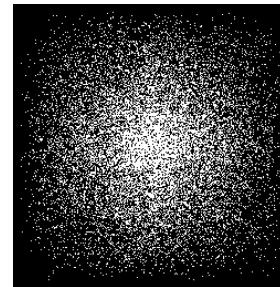
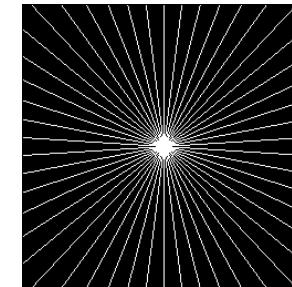


Image Com. Sensing via Nonlocal Low-rank Regularization

Recovery from sampling of Fourier transform coefficients



(a)



(b)



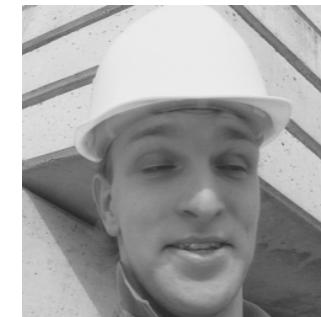
(a) Barbara



(b) boats



(c) Girl



(d) foreman



(e) house



(f) lena256



(g) Monarch



(h) Parrots

Image Com. Sensing via Nonlocal Low-rank Regularization

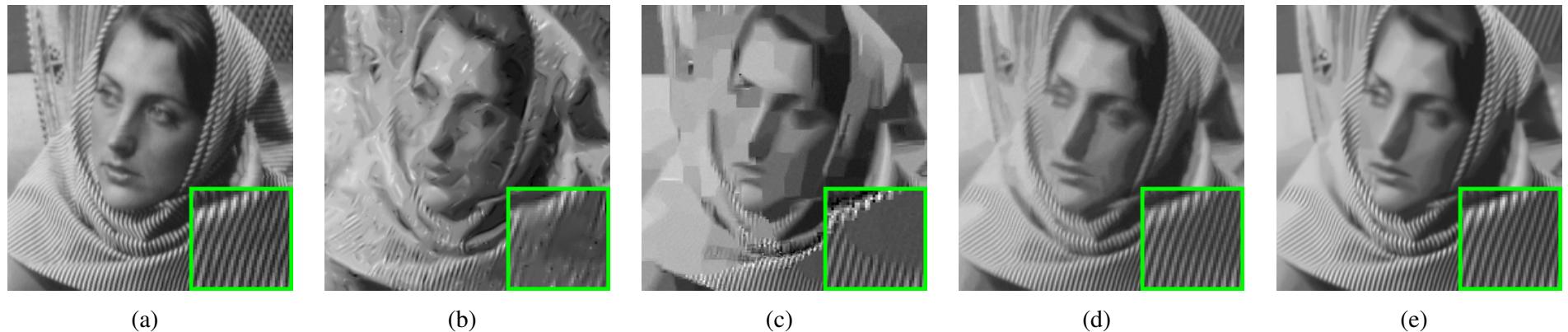


Fig. 4. CS recovered *Barbara* images with $0.05N$ measurements (random sampling). (a) Original image; (b) MARX-PC recovery [27] (24.11 dB); (c) BM3D-CS recovery [25] (24.34 dB); (d) Proposed **NLR-CS-baseline** recovery (27.59 dB) (e) Proposed **NLR-CS** recovery (**29.79 dB**).



Fig. 7. CS recovered *Barbara* images with pseudo radial subsampling (35 radial lines, i.e., $0.13N$ measurements). (a) Original image; (b) MARX-PC recovery [27] (22.99 dB); (c) BM3D-CS recovery [25] (24.38 dB); (d) Proposed **NLR-CS-baseline** recovery (26.99 dB); (e) Proposed **NLR-CS** recovery (**28.07 dB**).

Image Com. Sensing via Nonlocal Low-rank Regularization

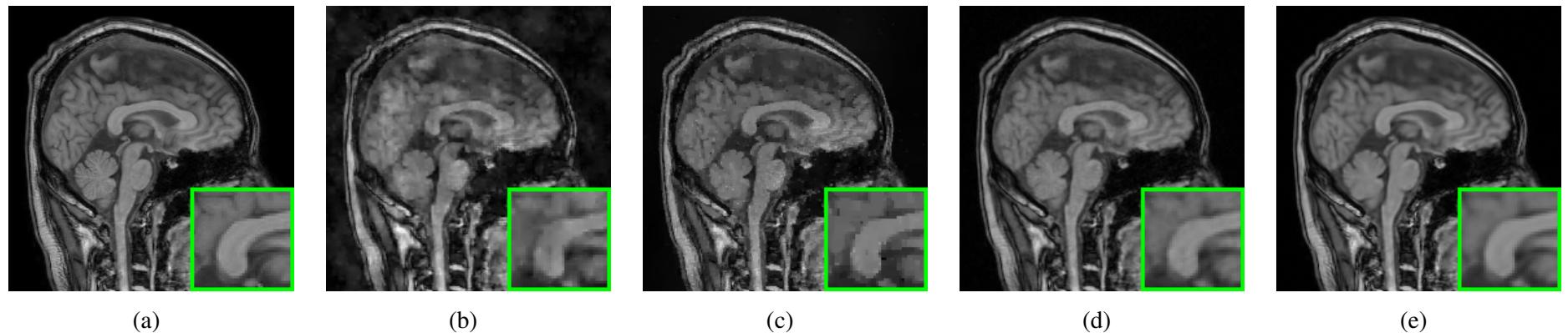


Fig. 9. Reconstructed *Head* MR images using $0.2N$ k -space samples (5 fold under-sampling, random sampling). (a) original MR image (magnitude); (b) SparseMRI method [32] (22.45 dB); (c) ReTV method [28] (27.31 dB); (d) proposed **NLR-CS-baseline** method (30.84 dB); (e) proposed **NLR-CS** method (**33.31** dB).

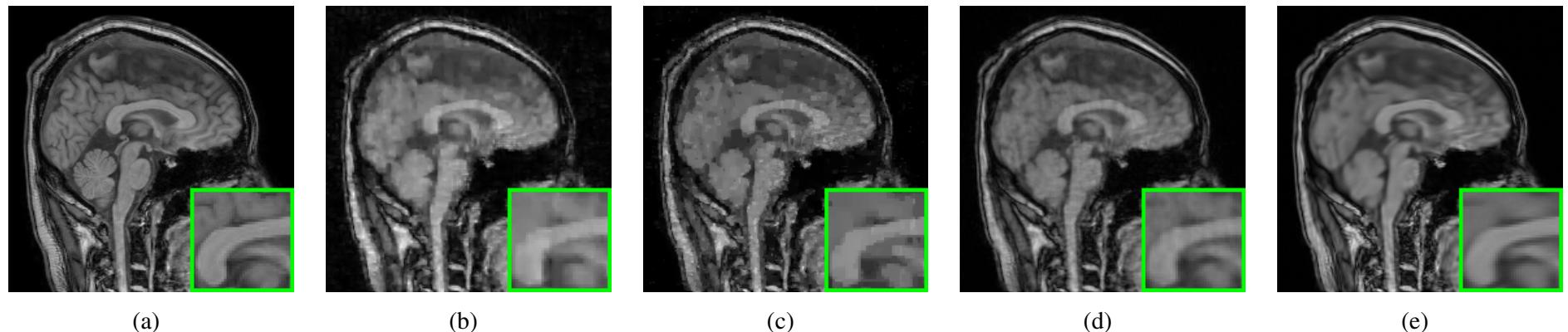


Fig. 11. Reconstruction of *Head* MR images from 45 radial lines (pseudo-radial sampling, 6.0 fold under-sampling). (a) original MR image (magnitude); (b) SparseMRI method [32] (24.02 dB); (c) ReTV method [28] (27.09 dB); (d) proposed **NLR-CS-baseline** (27.97 dB); (e) proposed **NLR-CS** method (**29.67** dB).

Image Com. Sensing via Nonlocal Low-rank Regularization

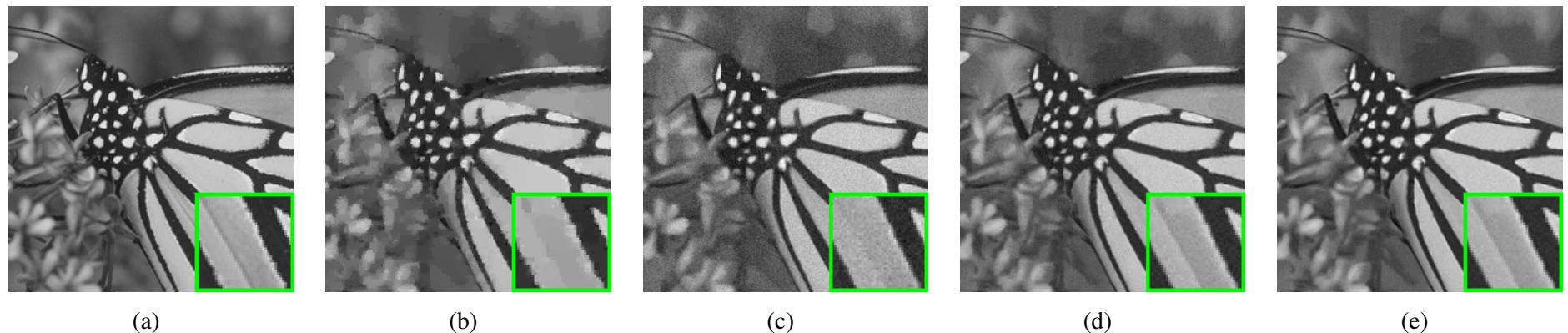


Fig. 15. Reconstruction of *Monarch* images from $0.2N$ noisy CS measurements (SNR=17.5 dB). (a) original image; (b) ReTV recovery [28] (26.73 dB); (c) BM3D-CS recovery [25] (26.00 dB); (c) proposed **NLR-CS-baseline** recovery (27.99 dB); (d) proposed **NLR-CS** recovery (**28.70 dB**).

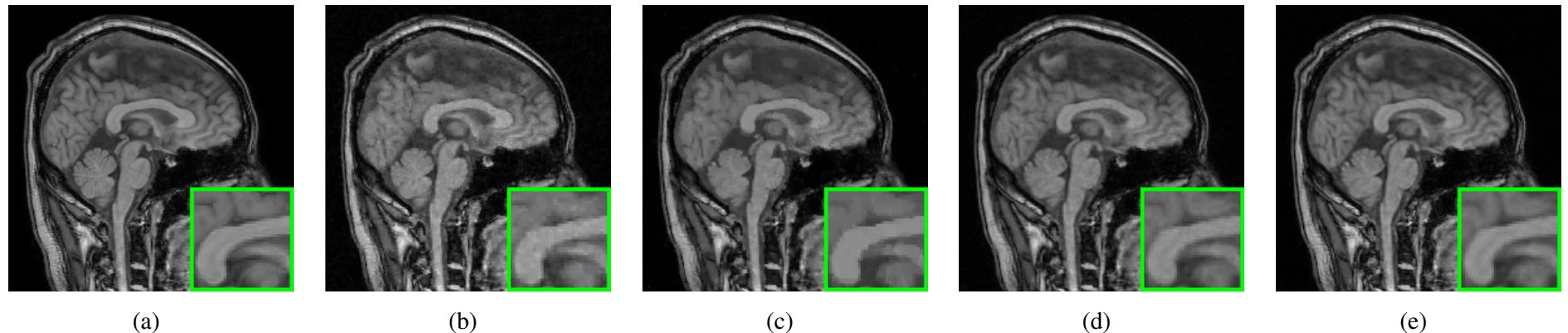
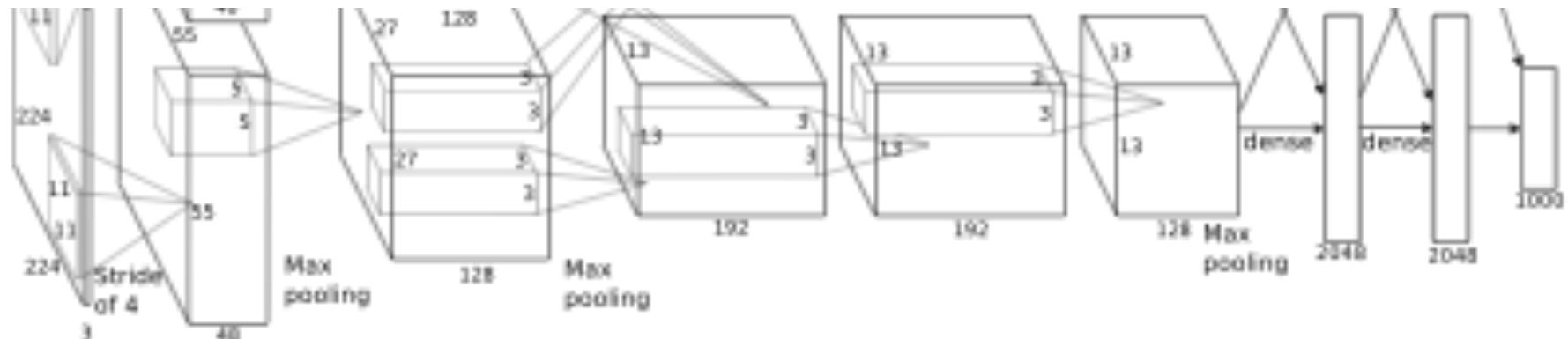


Fig. 17. Reconstruction of *Head* MR images from 65 noisy radial lines (SNR=15.42 dB). (a) original MR image (magnitude); (b) SparseMRI method [32] (28.73 dB); (c) ReTV method [28] (30.88 dB) (d) proposed **NLR-CS-baseline** (31.34 dB) (e) proposed **NLR-CS** method (**32.41 dB**).

Structured Matrix Factorization and Deep Learning

- Deep learning is a cascaded matrix factorization



$$\Phi(X^1, \dots, X^K) = \psi_K(\dots \psi_2(\psi_1(VX^1)X^2)\dots X^K)$$

nonlinearity features weights

$$\min_{X^1, \dots, X^K} \ell(Y, \Phi(X^1, \dots, X^K)) + \lambda \Theta(X^1, \dots, X^K)$$

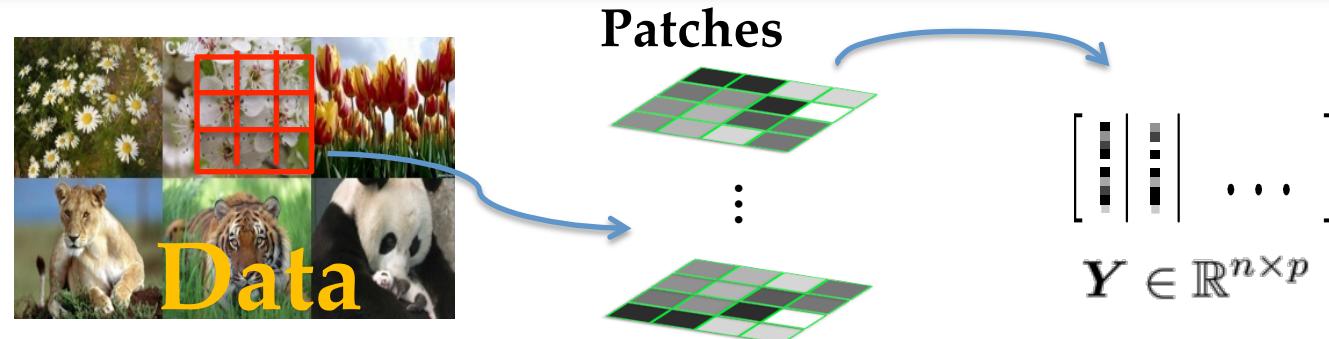
loss labels regularizer

Structured Matrix Factorization and Deep Learning

$$\min_{X^1, \dots, X^K} \ell(Y, \Phi(X^1, \dots, X^K)) + \lambda \Theta(X^1, \dots, X^K)$$

- **Theorem:** If the functions Φ and Θ are sums of positively homogeneous functions, then any local minimizer such that for some i and all k $X_i^k = 0$ gives a global minimizer
- Examples of positively homogeneous compositions Φ
 - Matrix multiplication: matrix factorization
 - CANDECOMP/PARAFAC decompositions: tensor factorization
 - Rectified linear units + max pooling: deep learning
- Examples of positively homogeneous regularizers Θ
 - Sums of products of norms (L1, L2, TV, etc.): structured factorizations

Sparsifying Dictionary Learning



$$\mathbf{Y} \approx \mathbf{Q}\mathbf{X} \quad \mathbf{Q} \in O(n) \quad \mathbf{X} = \boldsymbol{\Omega} \odot \mathbf{V}, \quad \boldsymbol{\Omega} \sim \text{Ber}(\theta), \quad \mathbf{V} \sim \mathcal{N}(0, 1).$$

Efficient algorithms with performance guarantees:

[S] For **dictionary learning**, $\theta n = O(\sqrt{n})$ nonzeros per column of \mathbf{X} is an upper bound for known efficient algorithms.

[A] If \mathbf{Q} is known, **compressed sensing** results imply that we can efficiently recover \mathbf{X} with $O(n/\log(n/m))$ nonzeros per column.

[E] Can we break the $\theta n = O(\sqrt{n})$ barrier?

Other theoretical work on local geometry. [Gribonval+Schnass 11], [Geng, w. 11], [Schnass 14].

Globally Optimal Dictionary Learning: Trust Region Method

Algorithm 1 Trust Region Method

Input: Initialization $\mathbf{q}^{(0)} \in \mathbb{S}^{n-1}$, data matrix $\mathbf{Y} \in \mathbb{R}^{n \times p}$, smoothing parameter μ , initial trust region size Δ , parameters $\Delta_{\min}, \Delta_{\max}, \rho_{\min}$.

Output: $\widehat{\mathbf{q}} \in \mathbb{S}^{n-1}$

- 1: Set $k = 1$,
- 2: **while** not converged **do**
- 3: Set $\mathbf{U} \in \mathbb{R}^{n \times n-1}$ to be an orthonormal basis for $\mathbf{q}^{(k-1)^\perp}$
- 4: **repeat**
- 5: Solve the trust region subproblem

$$\min_{\|\zeta\|_2 \leq \Delta} \widehat{f}(\mathbf{q}^{(k-1)}, \mathbf{U}\zeta) \quad (4.7)$$

- 6: Set

$$\widehat{\mathbf{q}} = \mathbf{q}^{(k-1)} \cos(\|\widehat{\boldsymbol{\delta}}\|_2) + \frac{\widehat{\boldsymbol{\delta}}}{\|\widehat{\boldsymbol{\delta}}\|_2} \sin(\|\widehat{\boldsymbol{\delta}}\|_2). \quad (4.8)$$

- 7: Set

$$\rho = \frac{f(\mathbf{q}^{(k-1)}) - f(\widehat{\mathbf{q}})}{f(\mathbf{q}^{(k-1)}) - \widehat{f}(\mathbf{q}^{(k-1)}, \widehat{\boldsymbol{\delta}})} \quad (4.9)$$

- 8: **if** $\rho \leq 1/4$ **then**
 - 9: $\Delta = \Delta/4$.
 - 10: **end if**
 - 11: **if** $\rho \geq 3/4$ and $\|\widehat{\boldsymbol{\delta}}\|_2 = \Delta$ **then**
 - 12: $\Delta = \min\{2\Delta, \Delta_{\max}\}$.
 - 13: **end if**
 - 14: **until** $\rho \geq \rho_{\min}$ or $\Delta \leq \Delta_{\min}$
 - 15: Set $\mathbf{q}^{(k)} = \widehat{\mathbf{q}}$.
 - 16: Set $k = k + 1$.
 - 17: **end while**
-

Globally Optimal Dictionary Learning

Theorem: Suppose that $\theta \in \left(0, \frac{1}{3}\right)$, $\mu \leq c \min\left\{\theta n^{-1}, n^{-5/4}\right\}$,

$\mathbf{Y} = \mathbf{Q}\mathbf{X}$, $\mathbf{Q} \in O(n)$ we observe $p \geq \frac{C}{\mu^2 \theta^2} n^3 \log \frac{n}{\mu \theta}$ samples.

Apply the TRM with $\Delta = \frac{1}{\text{poly}(n, \mu^{-1}, \theta^{-1})}$. W.p. $\geq 1 - cp^{-6}$, in

$T = \text{poly}(n, \mu^{-1}, \theta^{-1})$ iterations, the TRM produces $\hat{\mathbf{q}}$ such that

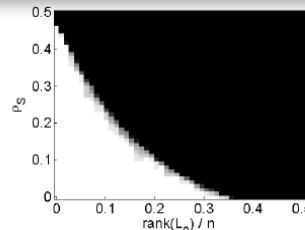
$$\|\hat{\mathbf{q}} - \mathbf{q}_*\|_2 \leq \mu/7,$$

for some target solution \mathbf{q}_* satisfying $\mathbf{q}_*^* \mathbf{Y} = \pm e_i^* \mathbf{X}$.

Using **LP rounding + deflation**, can recover all of \mathbf{X} , and subsequently all of \mathbf{Q} .

If \mathbf{Q} is not an orthobasis, **precondition**, but need $p \geq \text{poly}(n, \mu^{-1}, \theta^{-1}, \kappa(\mathbf{Q}))$.

A Perfect Storm...



(a) Robust PCA, Random Signs

BIG DATA
(images, videos,
voices, texts,
biomedical, geospatial,
consumer data...)



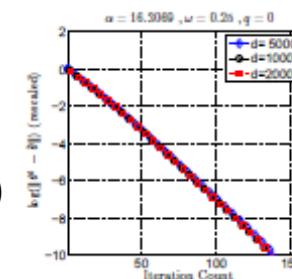
Mathematical Theory
(high-dimensional statistics, convex geometry
measure concentration, combinatorics...)

Cloud Computing
(parallel, distributed,
scalable platforms)



Applications & Services
(data processing,
analysis, compression,
knowledge discovery,
search, recognition...)

Computational Methods
(convex optimization, first-order algorithms,
random sampling, approximate solutions...)



A Perfect Storm...



**Dr. Arvind Ganesh, vision architect of Baarzo.com
web video analysis
purchased by Google in June, 2014**



**Kerui Min, CTO of Bosonnlp.com
web document analysis,
found in Shanghai, 2013**

CHRYSLER SETS STOCK SPLIT, HIGHER DIVIDEND
Chrysler Corp said its board declared a three-for-two stock **split** in the form of a 50 pct stock **dividend** and raised the quarterly **dividend** by seven pct.
The company said the **dividend** was raised to 37.5 cts a share from 35 cts on a pre-split basis, equal to a 25 ct **dividend** on a **post-split** basis.
Chrysler said the stock **dividend** is **payable** April 13 to holders of record March 23 while the cash **dividend** is **payable** April 15 to holders of record March 23. It said cash will be paid in lieu of fractional shares.
With the **split**, Chrysler said 13.2 mln shares remain to be purchased in its stock repurchase program that began in late 1984. That program now has a target of 56.3 mln shares with the latest stock **split**.
Chrysler said in a statement the actions "reflect not only our outstanding performance over the past few years but also our **optimism** about the company's future."



**Dr. Allen Yang, CTO of Atheerlabs.com
stereo gargle, object & gesture recognition,
found on Google campus, 2012**



REFERENCES + ACKNOWLEDGEMENT

Core References:

- *Robust Principal Component Analysis?* Candes, Li, Ma, Wright, Journal of the ACM, 2011.
- *TILT: Transform Invariant Low-rank Textures*, Zhang, Liang, Ganesh, and Ma, IJCV 2012.
- *Compressive Principal Component Pursuit*, Wright, Ganesh, Min, and Ma, IMA I&I 2013.

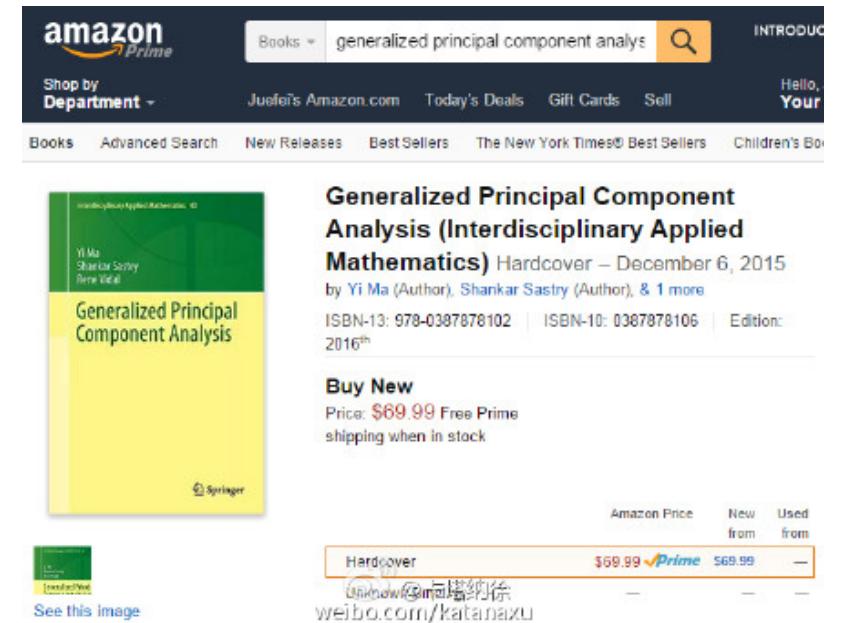
Website (codes, applications, & references):

<http://perception.csl.illinois.edu/matrix-rank/home.html>

A New Textbook:

- **Generalized Principal Component Analysis**

R. Vidal, Yi Ma, S. Sastry, Springer 2015



ACKNOWLEDGEMENT

Colleagues:

- Prof. Emmanuel Candes (Stanford)
- Prof. John Wright (Columbia)
- Prof. Zhouchen Lin (Peking Univ.)
- Dr. Yasuyuki Matsushita (Osaka Univ.)
- Dr. Arvind Ganesh (Google)
- Prof. Shuicheng Yan (Na. Univ. Singapore)
- Prof. Lei Zhang (HK Polytech Univ.)
- Prof. Liangshen Zhuang (USTC)
- Prof. Weisheng Dong (Xidian Univ.)

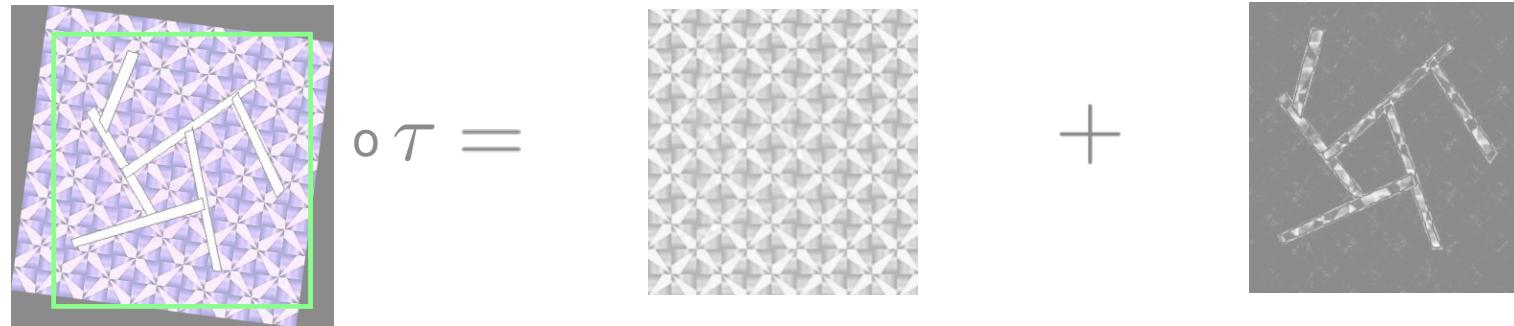
Students:

- Zhengdong Zhang (MSRA, now MIT)
- Xiao Liang (MSRA, Tsinghua University)
- Xin Zhang (MSRA, Tsinghua University)
- Kerui Min (UIUC)
- Zhihan Zhou (UIUC, now PennState)
- Hossein Mobahi (UIUC, now MIT)
- Guangcan Liu (UIUC, now UPenn)
- Xiaodong Li (Stanford)
- Carlos Fernandez (Stanford, MSRA)



THANK YOU!

Questions, please?



$$D \circ \tau = A + E \quad \min \|A\|_* + \lambda \|E\|_1$$