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

## Yi Ma

**School of Information Science & Technology**

**ShanghaiTech University, China**

**Images**

⇓➢ **1M pixels**

*Compression*
*De-noising*
*Super-resolution*
*Recognition…*

**Videos**

⇓➢ **1B voxels**

*Streaming*
*Tracking*
*Stabilization…*

**User data**

⇓➢ **1B users**

*Clustering*
*Classification*
*Collaborative filtering…*

**Web data**

⇓➢ **100B webpages**

*Indexing*
*Ranking*
*Search…*

**U.S. COMMERCE'S ORTNER SAYS YEN UNDERVALUED**

Commerce Dept. undersecretary of economic a®airs 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 di®erentiations.
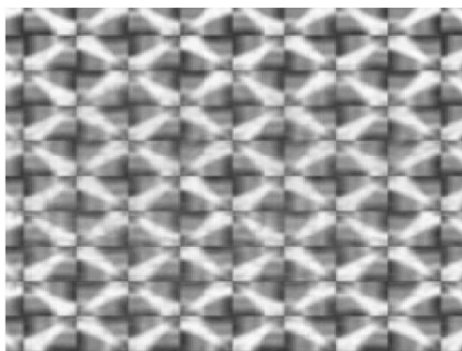
Ortner said there had been little impact on U.S. trade de¯cit by the decline of the dollar because at the time of the Plaza Accord, the dollar was extremely overvalued and that the ¯rst 15 pct decline had little impact.

He said there were indications now that the trade de¯cit was beginning to level o®.

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.

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

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

If we view the data (image) as a matrix

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

then

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



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

- Optimal estimate of $A$ under iid Gaussian noise $D = A + Z$

- 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…***

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

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

A. Legendre. *Nouvelles methodes pour la determination des orbites des cometes*, 1806
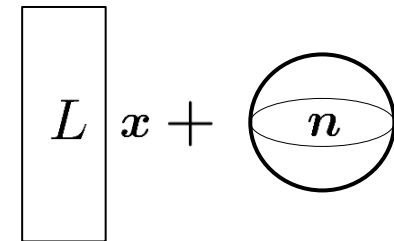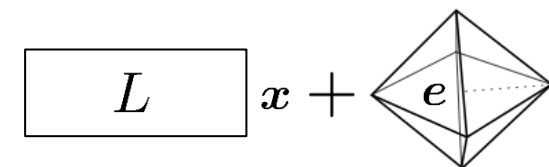
C. Gauss. *Theory of motion of heavenly bodies, 1809*

A. Beurling. *Sur les integrales de Fourier absolument convergentes et leur application a une transformation functionelle*, 1938

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

$$L \, x + \bigcirc n$$

over-determined
+ dense, Gaussian

$$L \quad x + \diamondsuit e$$

underdetermined
+ sparse, Laplacian

*Today, robust estimation in high dimension is more **urgent** and increasingly **better understood**.*

**Theory –** high-dimensional geometry & statistics, measure concentration, combinatorics, coding theory…

**Algorithms –** large scale convex optimization, geometric convergence rate, parallel and distributed computing …

**Applications –** big data driven methods, sensing and hashing, denoising, superresolution, MRI, bioinformatics, image classification, recognition …

Tukey, Bickel, Huber, Hampel, Tibishirani, Donoho, … Candes and Tao 2004 …

and many more I will mention later…

$$L \quad x + \diamond e$$
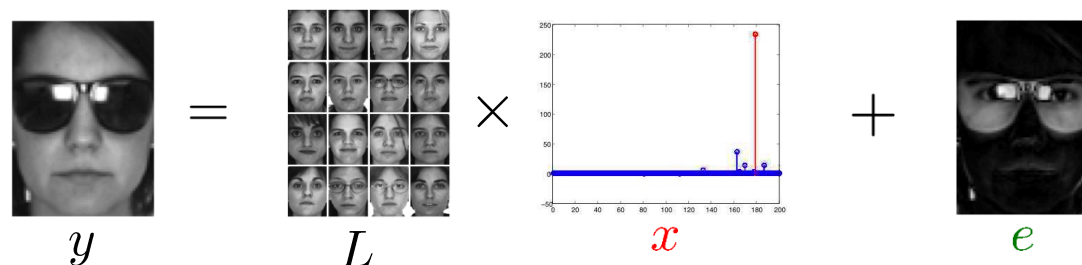
underdetermined
+ sparse, Laplacian

$$\min \|x\|_1 + \|e\|_1$$

**Sparse recovery:** Given $y = Lx_0$, $L \in \mathbb{R}^{m \times n}$, $m \ll n$, recover $x_0$.

$$y \in \mathbb{R}^m \qquad = \qquad L \in \mathbb{R}^{m \times n} \qquad x \in \mathbb{R}^n$$

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

**Well-posed** if $x_0$ is structured *(sparse),* but still **NP-hard**

**Tractable** via convex optimization: $\min \|x\|_1$ s.t. $y = Lx$

    … 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*

**Robust recovery:** Given $y = Lx_0 + e_0$, $L \in \mathbb{R}^{m \times n}$, $m \ll n$, recover $x_0$ and $e_0$.



$$y \qquad L \qquad x \qquad e$$

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

**Well-posed** if $x_0$ is *sparse,* errors $e_0$ not too dense, but still **NP-hard**

**Tractable:** via convex optimization: $\min \|x\|_1 + \|e\|_1$ s.t. $y = Lx + e$

     … 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…*

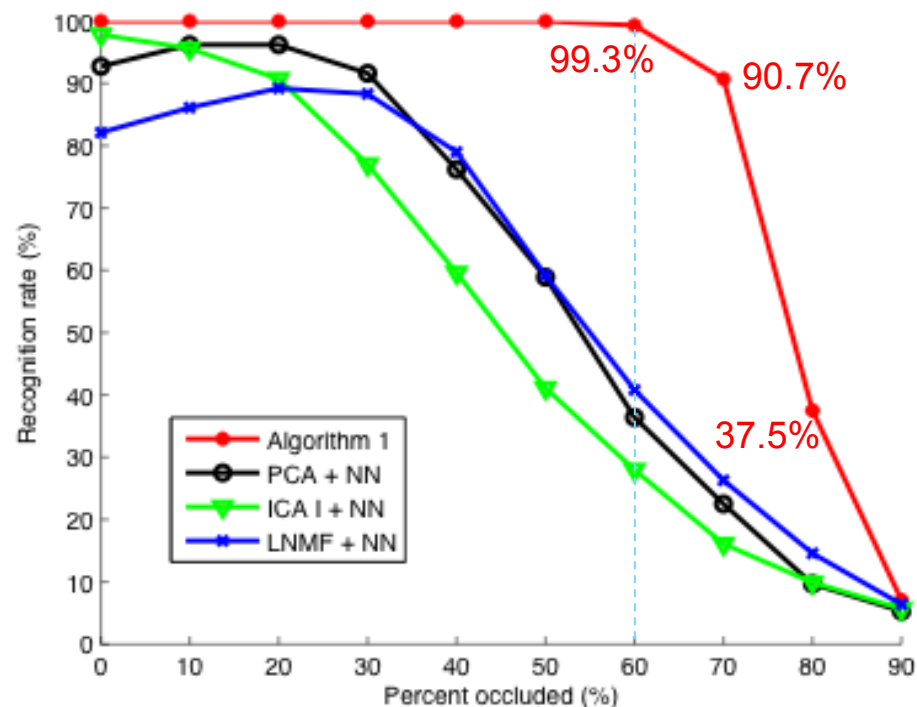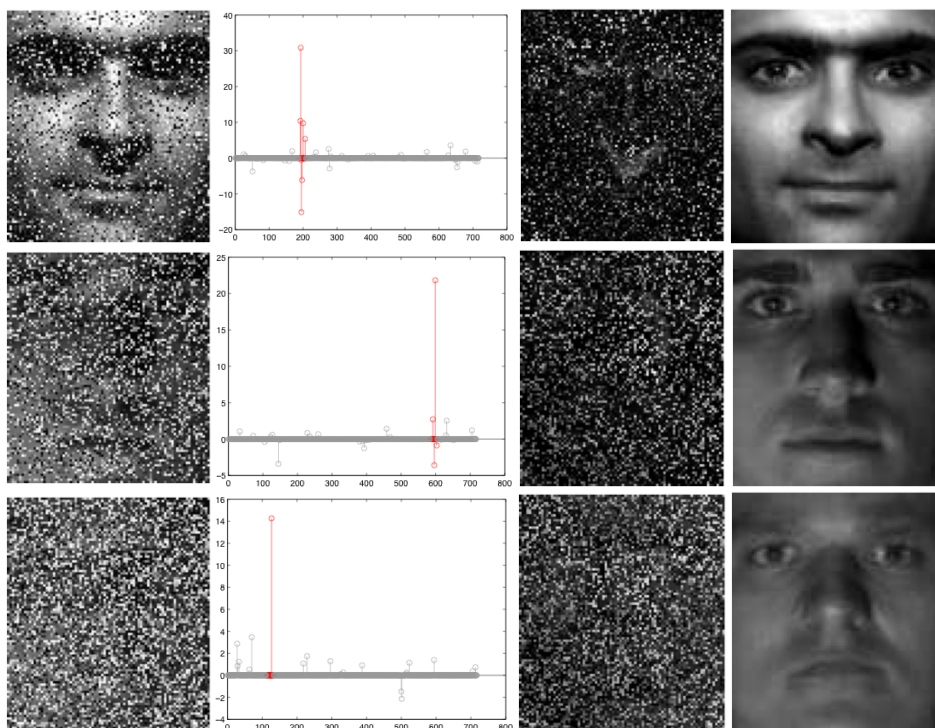Extended Yale B Database
(38 subjects)

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



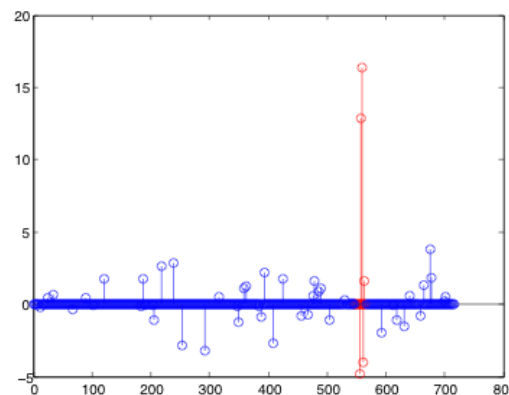$y$   $\widehat{x}_1$   $\widehat{e}_1$   $\widehat{y}_0 = A\widehat{x}_1$

99.3%   90.7%   37.5%

$$\widehat{x}_1 = \arg\min \; \|x\|_1 + \|e\|_1.$$

Input: $y \in \mathbb{R}^D$

$$y = Ax + e$$

$$\widehat{x}_2 = \arg\min_x \; \|y - Ax\|_2.$$

$$\widehat{e}_1 \qquad \widehat{y}_0 = A\widehat{x}_1$$

$$\widehat{e}_2 \qquad \widehat{y}_0 = A\widehat{x}_2$$

$$y = Lx + Ab + e$$

*A*: a common dictionary for intraclass variabilities: illumination, expression, and pose.

$x, b, e$ are sparse

# FERET Dataset

General training: 1,002 images of 429 people
Gallery training: 1,196 images of 1,196 people

Probe sets:

*fb (1,195, expression), fc (194, lighting),*
*dup1 (722, different time), dup2 (234, a year)*

TABLE 3
Comparative Recognition Rates of SRC and ESRC on the FERET Database Using the FERET'96 Testing Protocol

| Probe set | Feature / Dim | Dsampled Image 24×24 | Pixel-Rfaces 540 | Pixel 16384 | Gabor-Rfaces 540 | Gabor 10240 | LBP-Rfaces 540 | LBP 15104 |
|---|---|---|---|---|---|---|---|---|
| fb | SRC | 86.4 | 82.4 | 85.3 | 89.5 | 92.8 | 91.5 | 96.7 |
|    | ESRC | 94.8(+8.4) | 91.5(+9.1) | 92.8(+7.5) | 94.1(+4.6) | **97.3**(+4.5) | 95.2(+3.7) | **97.3**(+0.6) |
| fc | SRC | 69.6 | 75.8 | 76.3 | 96.4 | 97.4 | 72.7 | 93.3 |
|    | ESRC | 67.5(−2.1) | 78.9(+3.1) | 79.4(+3.1) | 96.9(+0.5) | **99.0**(+1.6) | 71.1(−1.6) | 95.4(+2.1) |
| dup1 | SRC | 62.7 | 60.9 | 63.7 | 63.0 | 72.7 | 75.2 | 87.7 |
|    | ESRC | 75.6(+12.9) | 73.1(+12.2) | 77.0(+13.3) | 73.5(+10.5) | 85.0(+12.3) | 81.0(+5.8) | **93.8**(+6.1) |
| dup2 | SRC | 52.6 | 53.0 | 55.6 | 70.1 | 76.5 | 69.7 | 83.8 |
|    | ESRC | 62.4(+9.8) | 59.8(+6.8) | 66.2(+10.6) | 72.6(+2.5) | 85.9(+9.4) | 71.4(+1.7) | **92.3**(+8.5) |

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



$$y \in \mathbb{R}^p \left[\;\right] = \underbrace{\left\langle \quad , \quad \underset{A \in \mathbb{R}^{m \times n}}{} \right\rangle}_{\mathcal{L}}{}_{i = 1 \ldots p}$$

**Impossible** in general ($p \ll mn$)

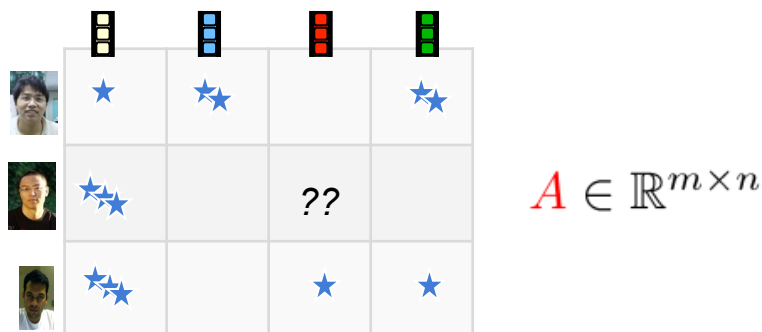**Well-posed** if $A_0$ is structured *(low-rank)*, but still **NP-hard**

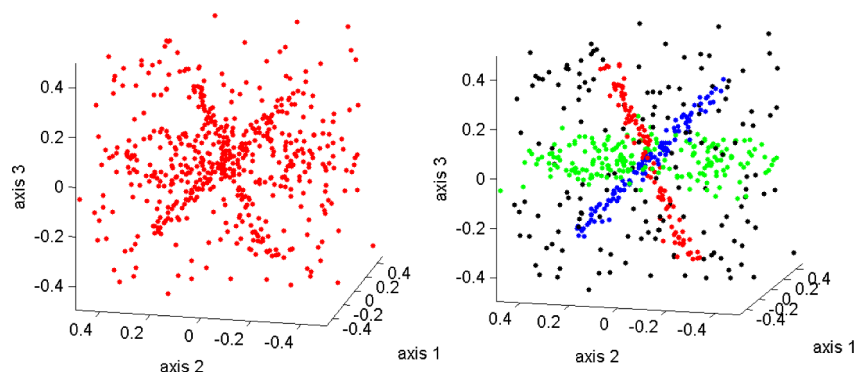**Tractable** via convex optimization: $\min \|A\|_*$ s.t. $y = \mathcal{L}(A)$

   … 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 …*

**Matrix completion:** Given $y = \mathcal{P}_\Omega[A_0]$, $\Omega \subset [m] \times [n]$, recover $A_0$.



$A \in \mathbb{R}^{m \times n}$

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

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

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

… if $\Omega$ is "nice" *(random subset) ...*

… and $A_0$ interacts "nicely" with $\mathcal{P}_\Omega$ ( $A_0$ *incoherent – not "spiky").*

*Hugely active area: Candès+Recht '08, Keshevan+Oh+Montonari '09, Candès+Tao '09, Gross '10, Recht '10, Negahban+Wainwright '10*

**Subspace Clustering:** Given $Y : [y_1, \ldots, y_n] \subset S_1, \ldots, S_k$, recover the subspaces.



$Y$ (with outliers)

**Impossible** in general (solutions highly ambiguous)

**Well-posed** if $\{S_i\}$ are few and structured *(low-dim),* but still **combinatorial**

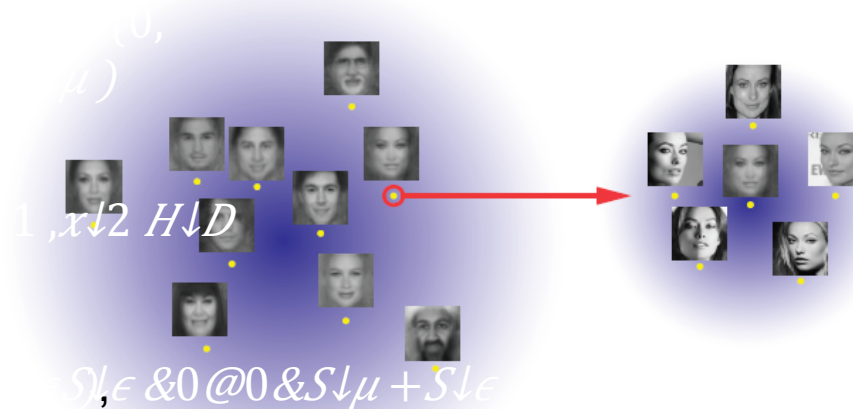**Tractable** via convex optimization: $\min \|X\|_\diamond + \|E\|_1 \text{ s.t. } Y = YX + E.$

    … for random samples $Y$

    … $X$ and *outliers* $E$ *are sparse (or low-rank, column-wise sparse).*

*Hugely active area: Rao, Tron, Ma, Vidal'08, Elhamifar and Vidal'2010, Liu, Lin, Sun, Yan, Ma et. al.' 2011, Soltanolkotabi and Candes' 2011*

**Bayesian Face Verification:**



**Impossible** to learn the covariance matrices in general case.

**Well-posed** if they are structured *(low-dim),* but still **high-dimensional**

**Tractable** via **rank-regularized** optimization:

*Hugely active area:* non-convex, Bayesian sparsity or low-rank regularization, *Wipf '2004, 2011, 2012…*

# CONTEXT – *Recent related progress*

- LFW dataset: 13,000 images, 2,000+ subjects
- Training and testing using the same LFW unconstraint protocol
- Using the same open source feature*



| Methods | Accuracy |
|---|---|
| **Bayesian (MSRA)** | **87.5%** |
| PLDA(2012) | 86.2% |
| LDML(2009) | 83.2% |
| DML-eig(2012) | 81.3% |

*http://lear.inrialpes.fr/people/guillaumin/data.php

Prince, S., Li, P., Fu, Y., Mohammed, U., Elder, J.: **Probabilistic models for inference about identity**. PAMI **34** (2012) 144–157

## MSRA WDRef

- 99,773 images
- 2,995 subjects
- Wide & Deep



| Methods | Accuracy |
|---|---|
| **Bayesian (MSRA)** | **92.4%** |
| face.com (2011) | **91.3%** |
| combined PLDA, funneled & aligned(2012) | 90.07% |
| Associate-Predict(2011)<br>*our previous work* | 90.57% |
| Combined multishot, aligned(2010) | 89.50% |
| LDML-MkNN, funneled(2009) | 87.50% |
| Attribute and Simile classifiers(2009) | 85.29% |

Billions of data

3D face model

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

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

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

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



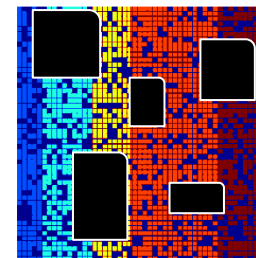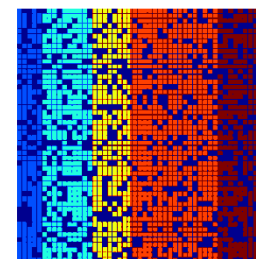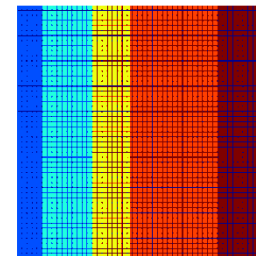*… but some of the observations are **grossly corrupted**:*

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

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

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

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

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

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

$E_{ij}$ 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.
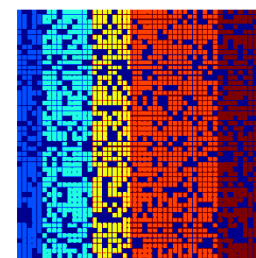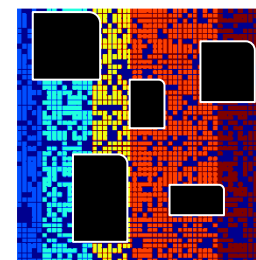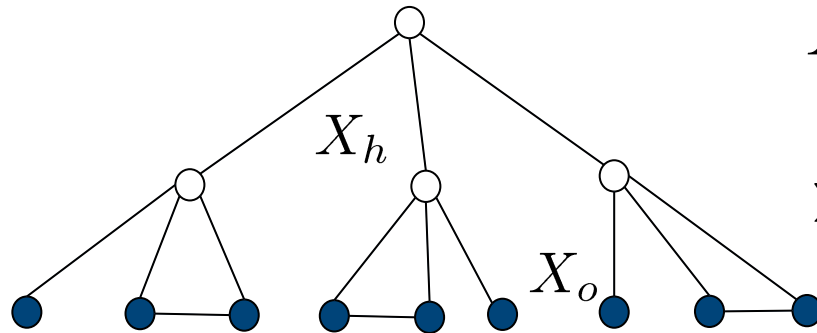
# CONTEXT – *Low-dimensional Models*

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

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

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

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

$E_{ij}$ 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:*

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

# CONTEXT: Learning Graphical Models



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

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

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

Separation Principle:

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

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

*Given observations* $D = \mathcal{P}_Q[A + E + Z]$*, with*
  $A$ *low-rank,*
  $E$ *sparse,*
  $Z$ *small, dense noise,*
*recover a good estimate of* $A$ *and* $E$*.*

❑ **Theory and Algorithm**

 • **Provably Correct and Tractable Solution**

 • **Provably Optimal and Efficient Algorithms**
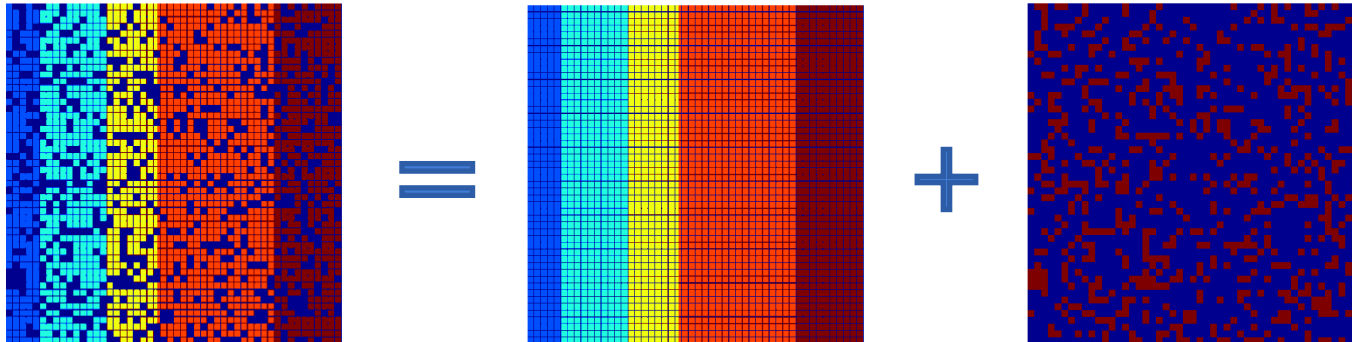
❑ **Potential Applications**

 • **Visual Data (Restoration, Reconstruction, Recognition)**

 • **Other Data**

❑ **Conclusions**

# ROBUST PCA – *Problem Formulation*

$D$ - observation          $A_0$ – low-rank          $E_0$ – sparse



$=$ ... $+$

**Problem**: Given $D = A_0 + E_0$, recover $A_0$ and $E_0$.

**Low-rank component**          **Sparse component (gross errors)**

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?*

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

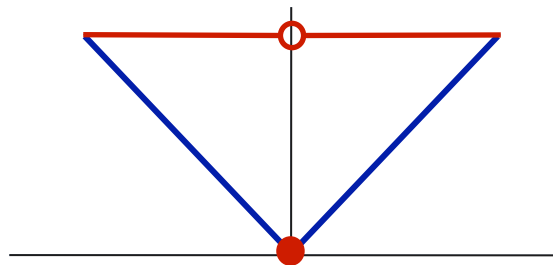$$\min \ \mathrm{rank}(A) \ + \ \gamma\|E\|_0 \ \ \mathrm{subj} \ \ A + E = D.$$

**But INTRACTABLE!** Relax with convex surrogates:

$$\|E\|_0 = \#\{E_{ij} \neq 0\} \quad \rightarrow \quad \|E\|_1 = \sum_{ij} |E_{ij}|. \qquad \text{L}_1 \text{ norm}$$
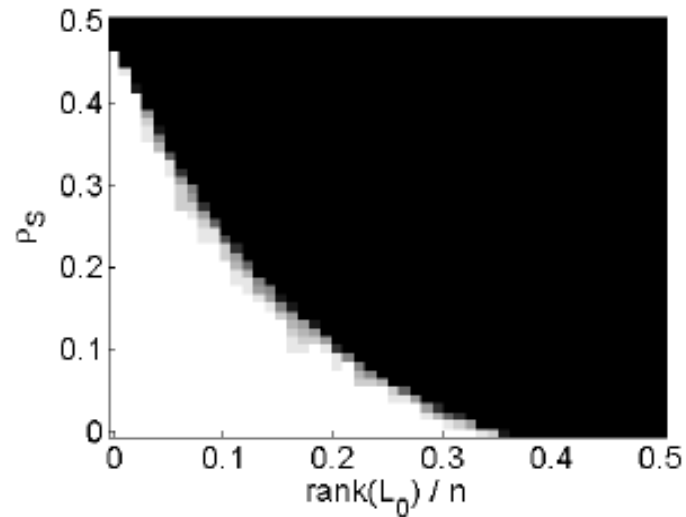
$$\mathrm{rank}(A) = \#\{\sigma_i(A) \neq 0\} \quad \rightarrow \quad \|A\|_* = \sum_i \sigma_i(A). \qquad \text{Nuclear norm}$$

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

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

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

**But INTRACTABLE!** Relax with convex surrogates:

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

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

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

*Semidefinite program, solvable in polynomial time*

$$D = A + E$$

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



Robust PCA, Random Signs



Matrix Completion

White regions are instances with perfect recovery.

Correct recovery when $A$ is indeed **low-rank** and $E$ is indeed **sparse**?

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

$m$

**Non-adaptive weight factor**

*and $E_0$ has Bernoulli support with error probability $\rho \leq \rho_s^\star$, 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.*

GREE 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!"*

$$D = \mathcal{P}_\Omega[\ A_0\ +\ E_0\ ], \qquad \Omega \sim \mathrm{uni}\binom{[m]\times[n]}{mn}$$

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

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

*and we observe only a random subset of size*

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

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

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

*uniquely recovers* $(A_0, E_0)$.

**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 $\rho \leq \rho_s^\star$, then with very high probability*

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

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

**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*.

**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 $r = O\left(m/\log^2 n\right)$,

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:**

**Matrix Recovery (RPCA)**

$$D = A + E$$

random signs

$$\textbf{rank} = O\left(\frac{m}{\log^2 n}\right)$$

$\frac{\|E_0\|_0}{mn}$

$\frac{\text{rank}(A_0)}{m}$

Classical PCA

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.)

**Matrix Completion**

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

$$\textbf{rank} = O\left(\frac{m}{\log^2 n}\right)$$

$\frac{\|E_0\|_0}{mn}$

$\frac{\text{rank}(A_0)}{m}$

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

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: $\min_{\mathbf{x}} f(x)$

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

| Function class $\mathcal{F}$ | Suboptimality $f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*)$ |
|---|---|
| *smooth* $\quad f$ convex, differentiable $\quad \|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{x}')\| \le L\|\boldsymbol{x} - \boldsymbol{x}'\|$ | $\dfrac{CL\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2}{k^2} = \Theta\left(\frac{1}{k^2}\right)$ |
| *smooth + structured nonsmooth:* $\quad F = f + g$ $\quad f, g$ convex, $\quad \|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{x}')\| \le L\|\boldsymbol{x} - \boldsymbol{x}'\|$ | $\dfrac{CL\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|^2}{k^2} = \Theta\left(\frac{1}{k^2}\right)$ |
| *nonsmooth* $\quad f$ convex $\quad |f(\boldsymbol{x}) - f(\boldsymbol{x}')| \le M\|\boldsymbol{x} - \boldsymbol{x}'\|$ | $\dfrac{CM\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|}{\sqrt{k}} = \Theta\left(\frac{1}{\sqrt{k}}\right)$ |

**Y. Nesterov, Introductory Lectures on Convex Optimization: A Basic Course, 2003.**

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.

$$\mathcal{D}_\varepsilon(Q) \ = \ \text{argmin}_X \ \varepsilon\|X\|_* + \frac{1}{2}\|X - Q\|_F^2$$

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.

**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:

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

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

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

| 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!**

GREAT NEWS: Geometric convergence for gradient algorithms!

$f$ restricted strong convex:     $O(\log(1/\varepsilon))$ [Agarwal, Negahban, Wainwright, NIPS 2010]

$f$ smooth, $\nabla f$ 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 ($\ell_1$-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 ($\alpha = \frac{n}{s \log d}$), 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**

*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!



$D$     Low-rank Texture $A$     Corruptions $E$

$D$     $A$     $D$     $A$

Liang, Ren, Zhang, and Ma, in ECCV 2012.

# Repairing Low-rank Textures



Low-rank Method

Photoshop

Input

Output

# Repairing (Distorted) Low-rank Textures



Low-rank Method

Photoshop

Input

Output

# Repairing Multiple Correlated Images

58 images of one person under varying lighting:



$D$   $A$   $E$

specularity

cast shadows

$D$

RPCA

Candes, Li, Ma, and Wright, Journal of the ACM, May 2011.

# Repairing Images: robust photometric stereo

Input images

$$\min \ \|A\|_* + \lambda\|E\|_1 \ \ \text{subj} \ \ D = \mathcal{P}_\Omega(A + E).$$

$$\Omega^c \sim \text{shadow}(20.7\%)$$
$$E \sim \text{specularities}(13.6\%)$$

(a) Ground truth     (b) Our method     (c) Least Squares     (d) Error map (our method)     (e) Error map (LS)

| | (d) Error map (our method) | (e) Error map (LS) |
|---|---|---|
| Mean error | **0.014°** | 0.96° |
| Max error | **0.20°** | 8.0° |

Wu, Ganesh, Li, Matsushita, and Ma, in ACCV 2010.

# Repairing Video Frames: *background modeling from video*

Surveillance video

200 frames,
144 x 172 pixels,

Significant foreground motion



$D$

RPCA

Video $D$ = Low-rank appx. $A$ + 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

D

A

E

# Repairing Video Frames: Street Panorama

Low-rank

AutoStitch

Photoshop

Low-rank

AutoStitch

Photoshop

# Street Panorama: Highly Compressive Sensing of Low-dim Structures!

nips12_video.mp4

# Sensing or Imaging of Low-rank and Sparse Structures

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



Low-rank Structures          Sparse Structures

*subject to either nonlinear deformation $\tau$ or linear compressive sampling $\mathcal{P}$?*

# Reconstructing 3D Geometry and Structures

$D$ – deformed observation    $A$ – low-rank structures    $E$ – sparse errors



$\circ \, \tau \; = \qquad\qquad\qquad\qquad +$

**Problem**: Given $D \circ \tau \; = \; A_0 \, + \, E_0$, recover $\tau$, $A_0$ 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



$$\circ \tau \;=\; <image> \;+\; <image>$$

**Objective:** *Transformed* Principal Component Pursuit::

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

**Solution:** *Iteratively solving the linearized convex program:*:

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

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

**Theorem 5 (Compressive Principal Component Pursuit).** *Let $A_0 \in \mathbb{R}^{m \times n}$, $m \geq n$ have rank $r \leq \rho_r \frac{m}{\mu^2 \log^2(n)}$, and $E_0$ have a Bernoulli support with error probability $\rho < \rho^\star$. Let $Q^\perp$ be a random subspac 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*

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

*for some numerical constant $\rho_r$, $C_p$ and $\rho^\star$, 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 $A$)

$$z = f(x) : x(x - x_m) \sum_{i=1}^{d} a_i x^i$$

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

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

$$D \qquad A \qquad E$$

360° panorama
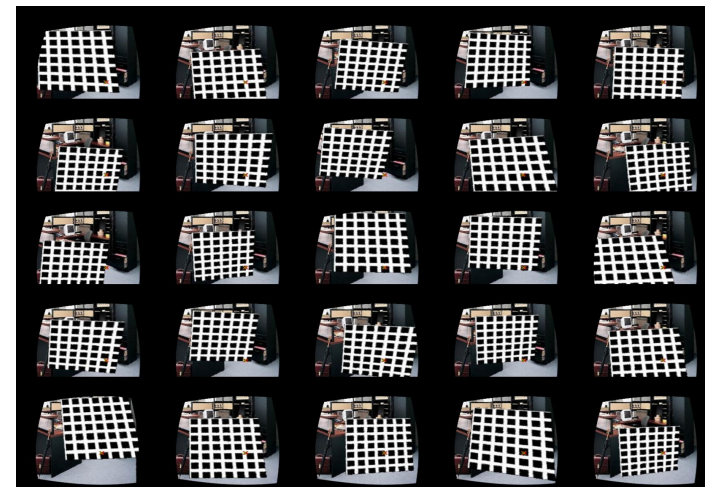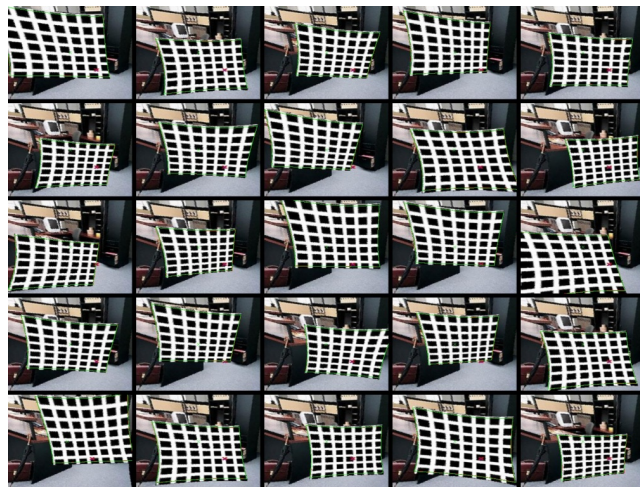
# TILT: *Virtual reality*

# TILT: *Camera Calibration with Radial Distortion*



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

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

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



**Zhang, Matsushita, and Ma, in CVPR 2011**

# TILT: *Camera Calibration with Radial Distortion*

$$\min \sum_{i=1}^{N} \|A_i\|_* + \lambda\|E_i\|_1 \quad \text{subj} \quad A_i + E_i = D \circ (\tau_0, \tau_i)$$

$$\tau_0 = (K, K_c), \quad \tau_i = (R_i, T_i).$$

Previous approach

Low-rank method

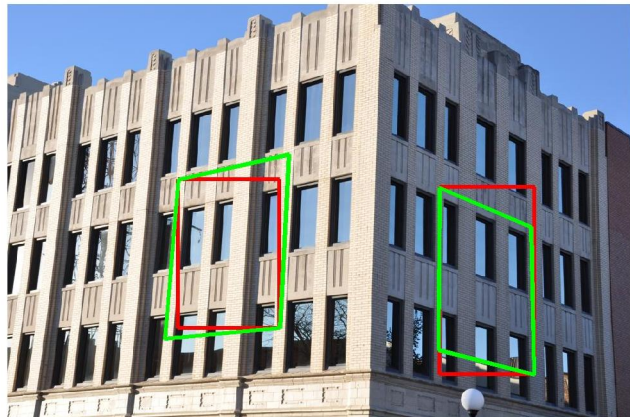# TILT: *Holistic 3D Reconstruction of Urban Scenes*



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

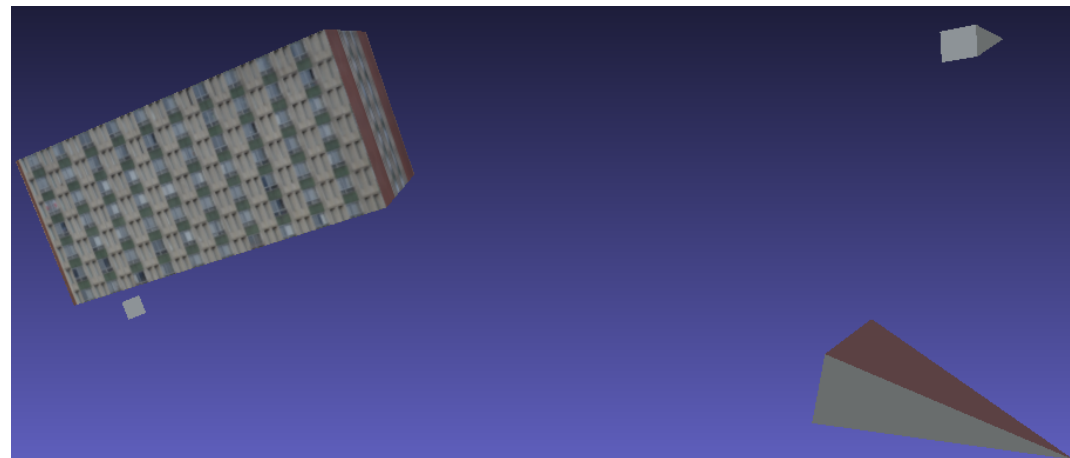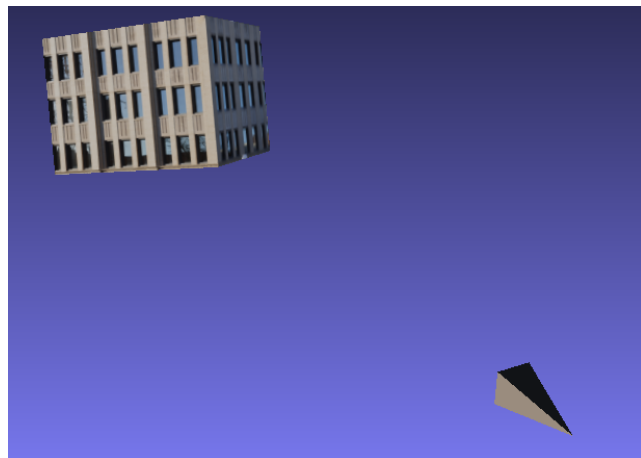$$A + E = [D_1 \circ \tau_1, D_2 \circ \tau_2]$$

**Mobahi, Zhou, and Ma, in ICCV 2011**

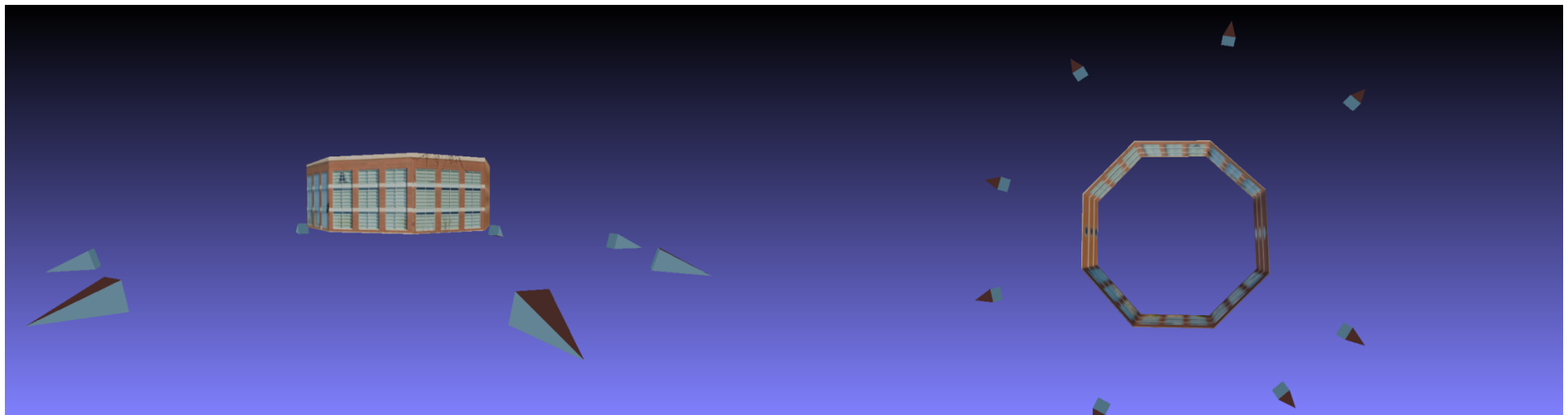# TILT: *Holistic 3D Reconstruction of Urban Scenes*

From one input image

From four input images

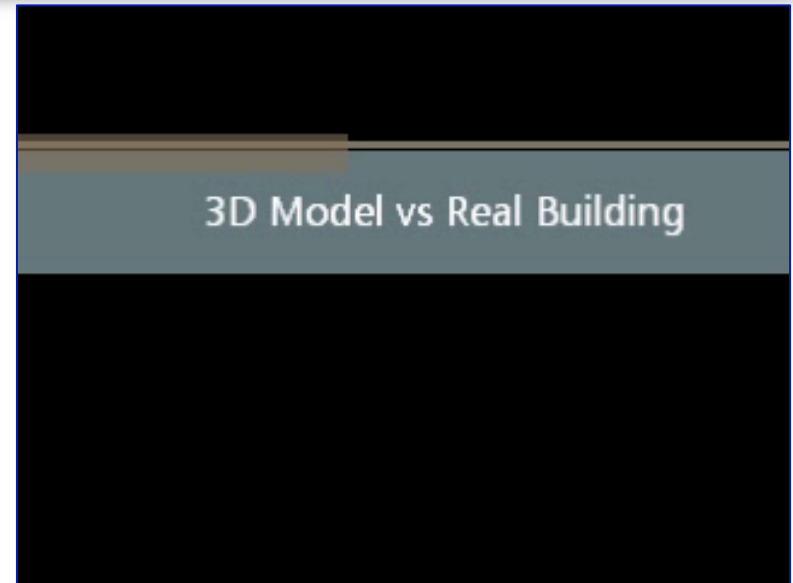# TILT: *Holistic 3D Reconstruction of Urban Scenes*
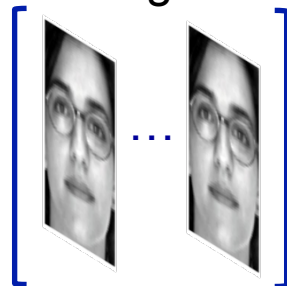
From eight input images



3D Model vs Real Building

# Registering Multiple Images: Robust Alignment
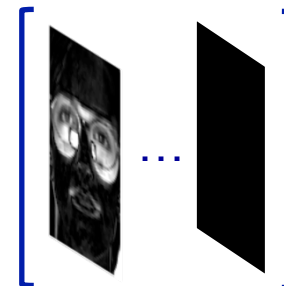
$D$ – corrupted & misaligned observation

$A$ – aligned low-rank signals

$E$ – sparse errors

$$\begin{bmatrix} \ldots \end{bmatrix} \circ \tau = \begin{bmatrix} \ldots \end{bmatrix} + \begin{bmatrix} \ldots \end{bmatrix}$$

**Problem**: Given $D \circ \tau = A_0 + E_0$, recover $\tau$, $A_0$ 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:*

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

# RASL: *Aligning Face Images from the Internet*

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

**Input**: faces detected by a face detector ($D$)



Average

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



Average

**Output**: clean low-rank faces ( $A$ )



Average

**Output**: sparse error images ($E$)

# RASL: *Video Stabilization and Enhancement*

Original video ( $D$ )   Aligned video ( $D \circ \tau$ )   Low-rank part ( $A$ )   Sparse part ( $E$ )
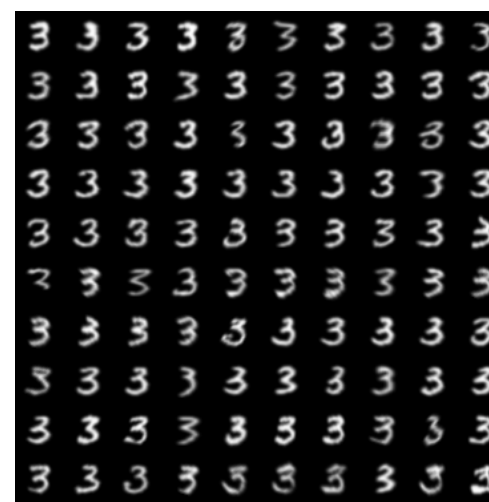
# RASL: *Aligning Handwritten Digits*

$D$

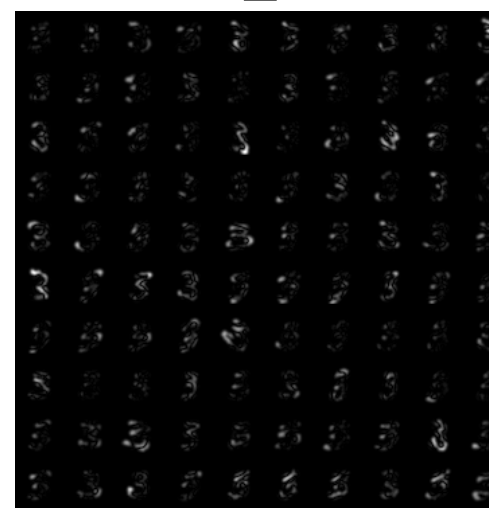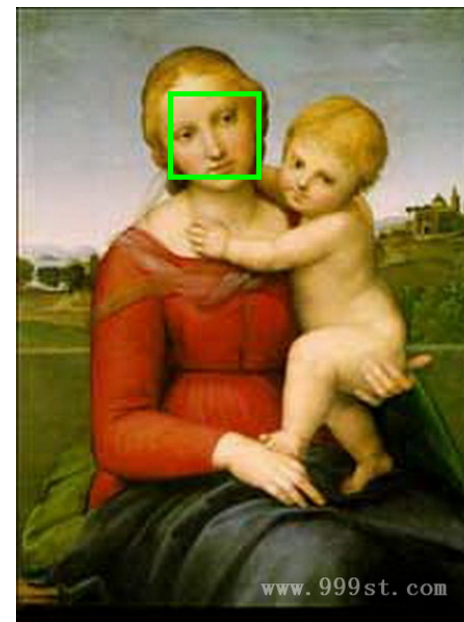Learned-Miller PAMI'06

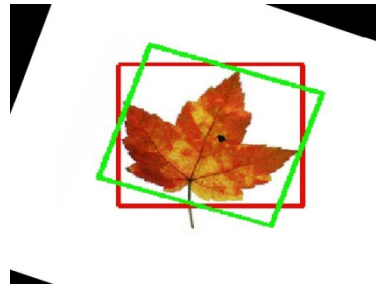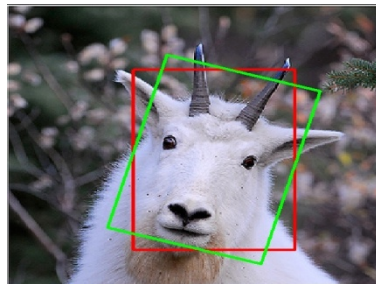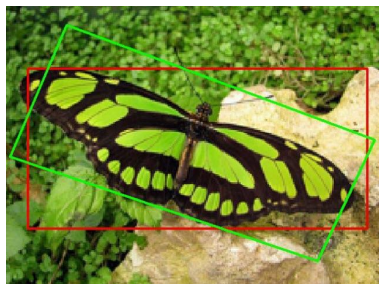Vedaldi CVPR'08

$D \circ \tau$

$A$

$E$

# Object Recognition: *Rectifying Pose of Objects*
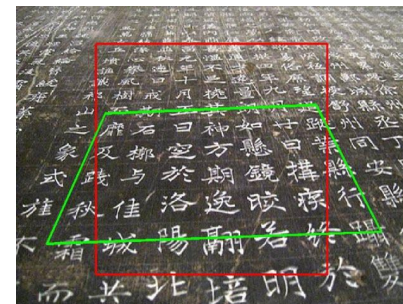
Input (red window $D$ )



Output (rectified green window $A$ )

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

Input (red window $D$ )



Output (rectified green window $A$ )

$$D \qquad D \circ \tau \qquad A \qquad E$$

# Recognition: *Character/Text Rectification*

**TILT**      **versus**      **Hough Transform**

# Recognition: Street Sign Rectification



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

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

**Xin Zhang, Zhouchen Lin, and Ma, ICDAR 2013**

# 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)

**Xin Zhang, Zhouchen Lin, and Ma, ICDAR 2013**

# 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…*

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



Input: images with user-provided tags

fly
bird
cool
insect
strong

Tag Refinement

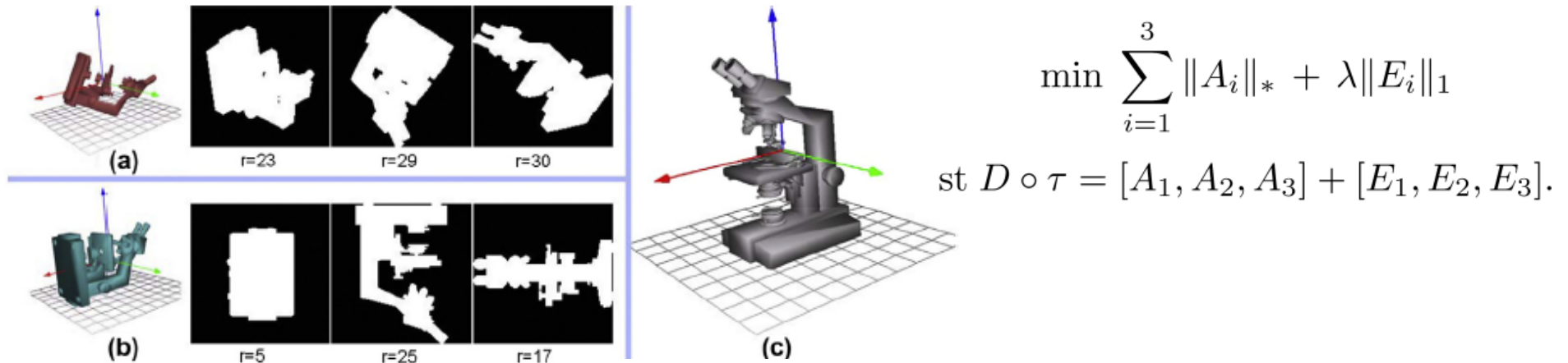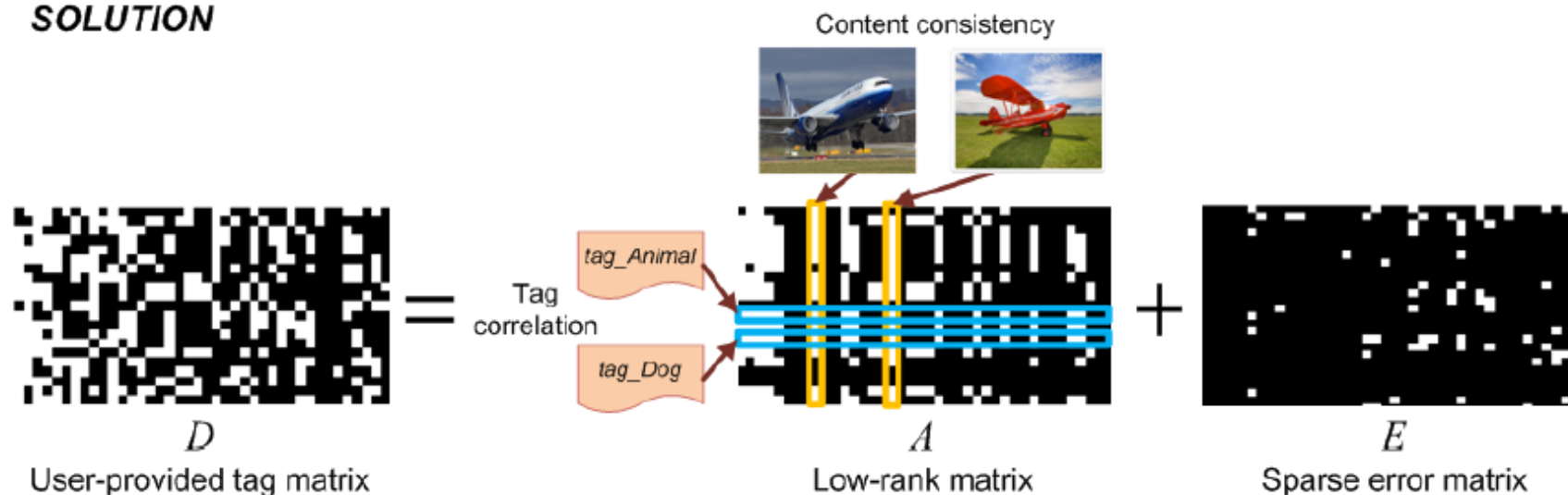Output: images with refined tags

fly
bird
sky
eagle

PROBLEM

SOLUTION

Content consistency

$D$
User-provided tag matrix

$=$

Tag correlation

tag_Animal

tag_Dog

$A$
Low-rank matrix

$+$

$E$
Sparse error matrix

## Latent Semantic Indexing:   the classical solution (PCA)

**Documents**

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."

**Words**

$$D$$

$d_{ij}$ word frequency (or TF/IDF)

$$
\begin{aligned}
D &= A + Z \\
&= U_1 \Sigma_1 V_1^T + U_2 \Sigma_2 V_2^T
\end{aligned}
$$

Dense, difficult to interpret

## a better model/solution?

$$D = A + E$$

Low-rank "background" topic model

Informative, discriminative "keywords"

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."

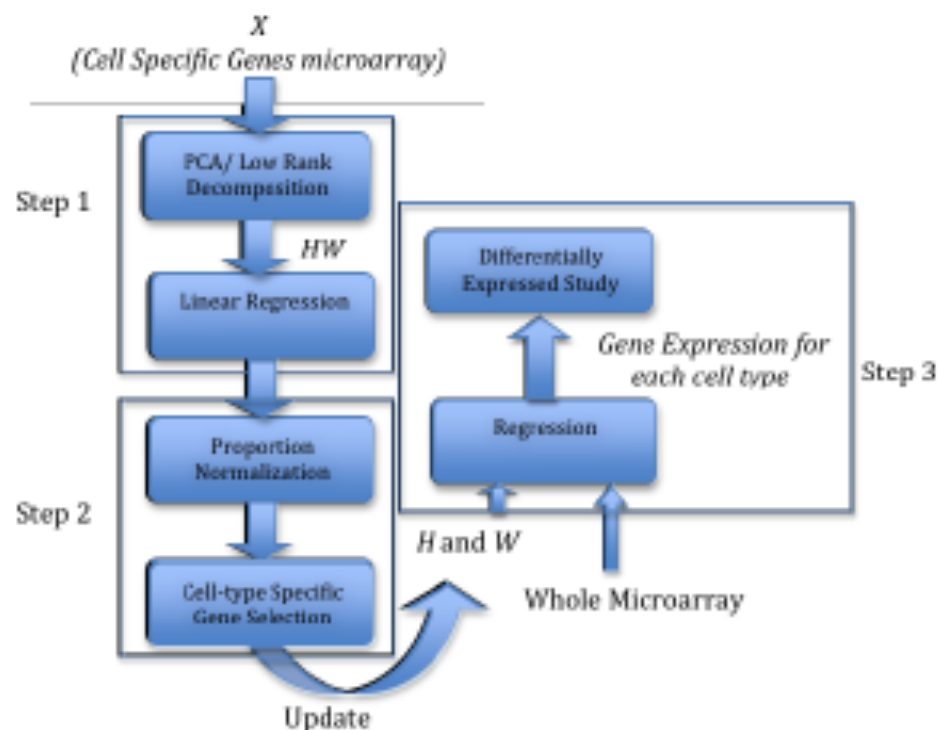# Other Data/Applications: Protein-Gene Correlation

Microarray data



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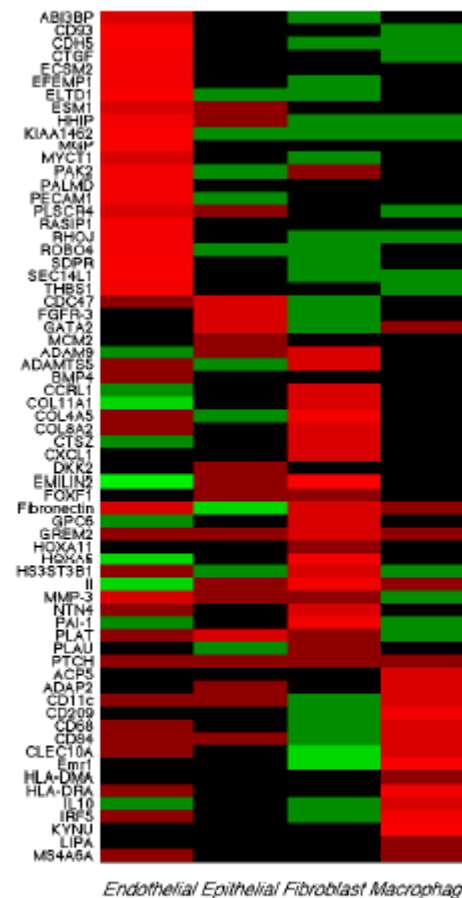


*Endothelial Epithelial Fibroblast Macrophage*

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.
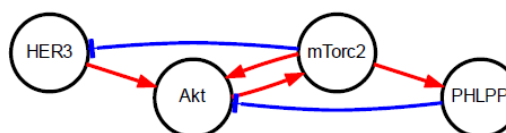
# Other Data: Time Series Gene Expressions



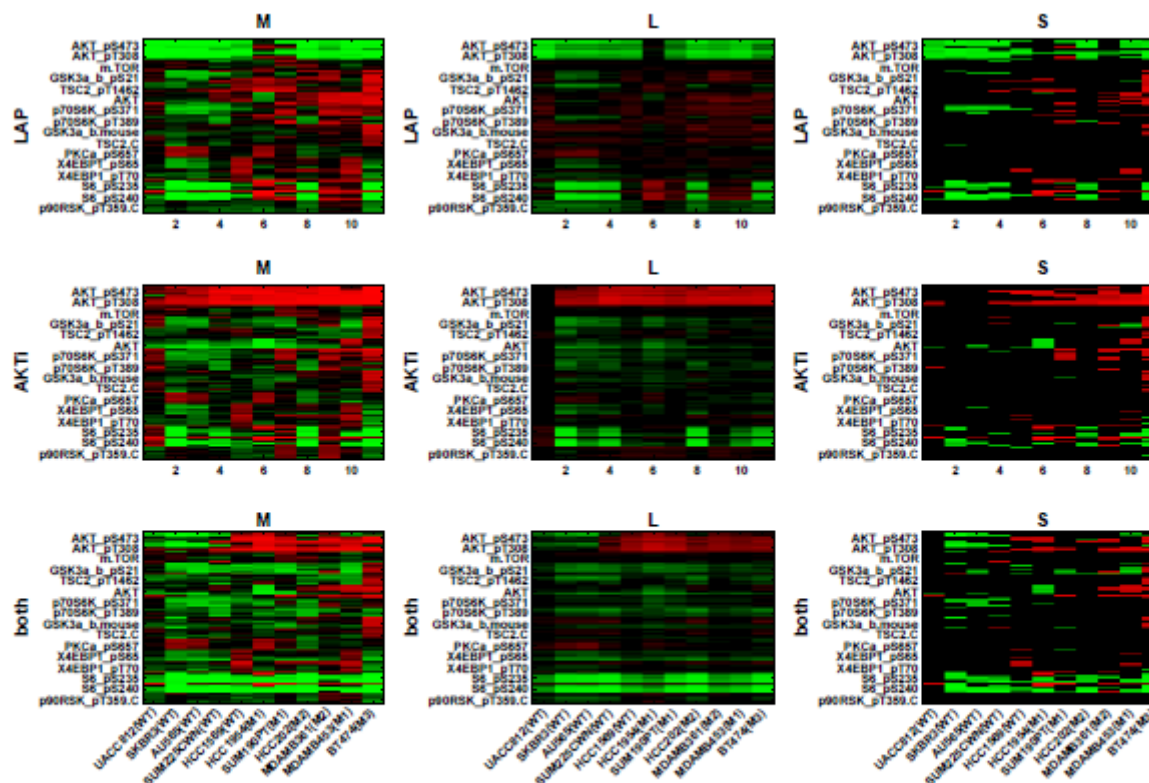Figure S3. Abstract HER2 overexpressed breast cancer model by Dr. Moasser.
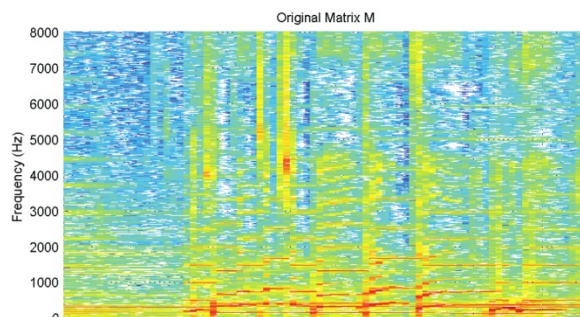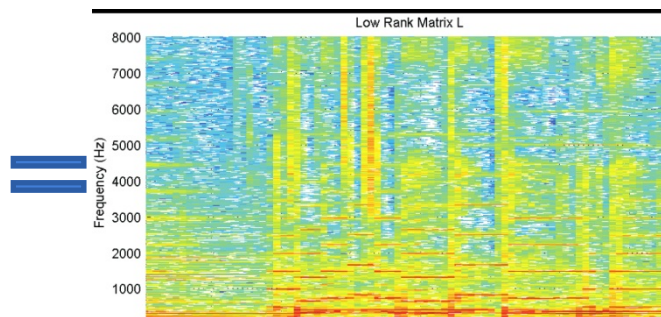
Figure S4. Separation result: ($1_{st}$ column) raw data ($2_{nd}$ column) low-rank component and ($3_{rd}$ column) highly corrupted sparse component using threshold (M1: H1047R (kinase domain mutation) M2: E545K (helical domain mutation), and M3: K111N mutation in PIK3CA).

Chang, Korkola, Amin, Tomlin of Berkeley, BiorXiv, 2014.

Songs (STFT)   Low-rank (music)   Sparse (voices)

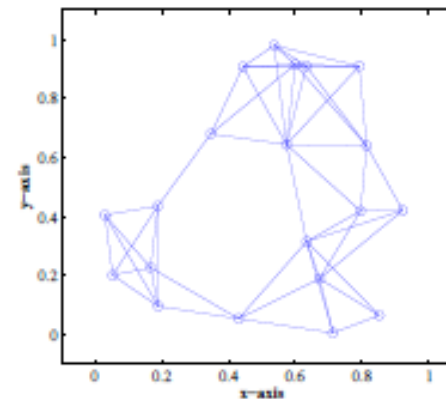Network Traffic = Normal Traffic + Sparse Anomalies + Noise

$$D = L + RS + N$$



Fig. 2. Network topology graph.

# Other Data/pplications: View-Invariant Gait Recognition

**Same gait from different views**

**Perspective distortion rectified**

**GPS on a Car:**

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

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:
$$\begin{bmatrix} y_n & y_{n-1} & y_{n-2} & \cdots & y_0 \\ y_{n-1} & y_{n-2} & \cdots & \ddots & y_{-1} \\ y_{n-2} & \cdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & y_{-n+2} \\ y_0 & y_{-1} & \cdots & y_{-n+2} & y_{-n+1} \end{bmatrix} = \mathcal{O}_{n \times r} X_{r \times n} + S$$

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 $\|x\|_0$ | $\mathrm{rank}(X)$ |
| 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



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

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

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 $X^0$ can be exactly recovered from high compressive measurements via convex optimization:

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

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

Suppose $(\boldsymbol{X}_1^0, \ldots, \boldsymbol{X}_k^0) = \arg\min \sum_{i=1}^k \lambda_i \|\boldsymbol{X}_i\|_{(i)}$ subj $\sum_{i=1}^k \boldsymbol{X}_i = \sum_{i=1}^k \boldsymbol{X}_i^0$, 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 subspac of $\mathbb{R}^{m \times n}$ of dimension*

$$\dim(Q) \geq O(\log^2 m) \times \text{intrinsic degrees of freedomof } (\boldsymbol{X}_1, \ldots, \boldsymbol{X}_k),$$

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

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

*and the minimizer is unique.*

**Compressive Sensing:**

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

**Multiple-Structure Decomposition:**

$$\min \sum_i \lambda_i \|\boldsymbol{X}_i\|_{\diamond_i} \quad \text{s.t.} \quad \sum_i \boldsymbol{X}_i = \boldsymbol{D}$$
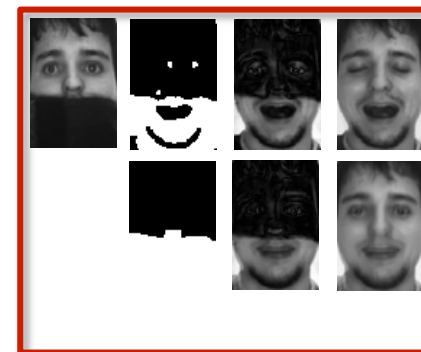
**Compressive Multiple-Structure Decomposition:**

$$\min \sum_i \lambda_i \|\boldsymbol{X}_i\|_{\diamond_i} \quad \text{s.t.} \quad \mathcal{P}_Q[\sum_i \boldsymbol{X}_i] = \mathcal{P}_Q[\boldsymbol{D}]$$

Examples: **PCP** [CLMW'11], **outlier pursuit** [Xu+Caramanis+Sanghavi], **morphological component analysis** [Bobin et. al.], many more …

# A Unified THEORY – *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

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

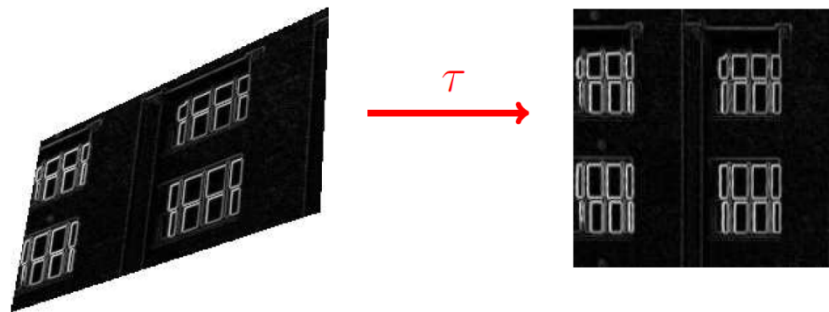- [Wright+Ganesh+Min+Ma, ISIT'12,I&I'13] – compressive superposition of decomposable structures

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

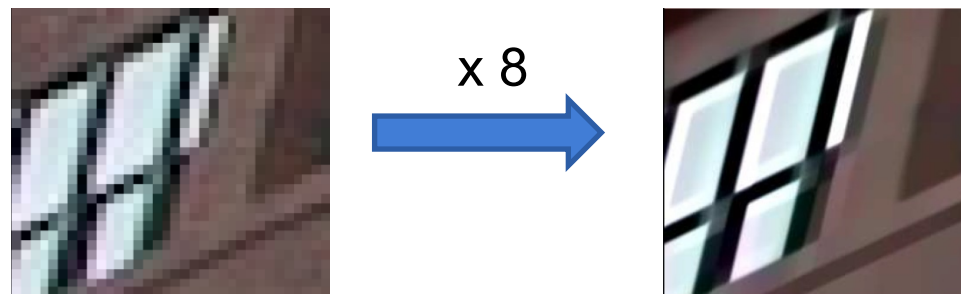*Take home message:* **Let the data and application tell you the structure…**

# Super Resolution via Transform Invariant Group Sparsity

# *A Perfect Storm…*



(a) Robust PCA, Random Signs

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

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

**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**
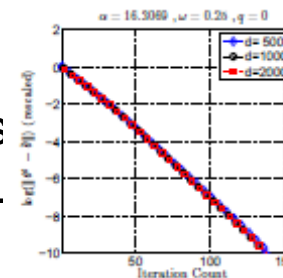


CHRYSLER SETS STOCK SPLIT, HIGHER DIVIDEND

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.

**More references, codes, and applications on the website:**

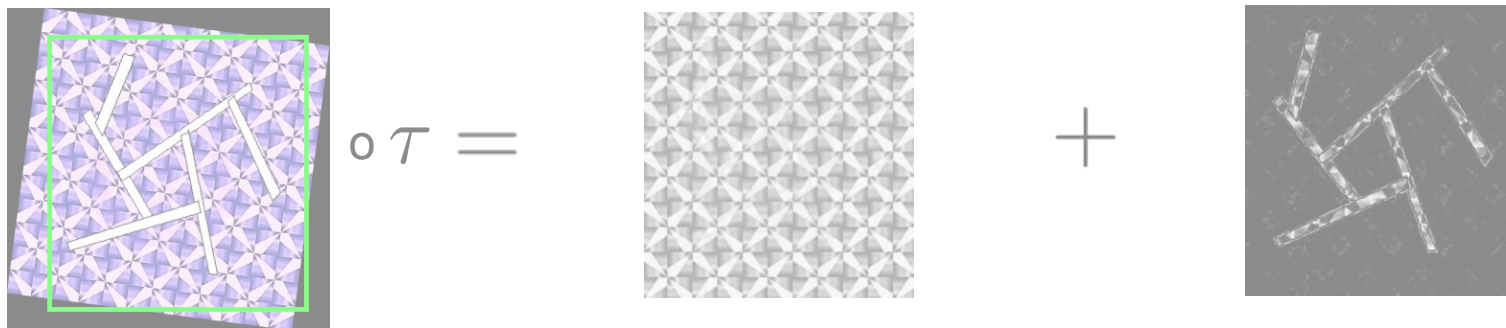http://perception.csl.illinois.edu/matrix-rank/home.html

**Colleagues:**
- Prof. Emmanuel Candes (Stanford)
- Prof. John Wright (Columbia)
- Prof. Zhouchen Lin (Peking University)
- Dr. Yasuyuki Matsushita (MSRA)
- Dr. Arvind Ganesh (IBM Research, India)
- Prof. Shuicheng Yan (Na. Univ. Singapore)
- Prof. Jian Zhang (Sydney Tech. Univ.)
- Prof. Lei Zhang (HK Polytech Univ.)
- Prof. Liangshen Zhuang (USTC)

**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)

# Questions, please?



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