

# Tutorial: Sparse Recovery Using Sparse Matrices

Piotr Indyk  
MIT

# Problem Formulation

(approximation theory, learning Fourier coeffs, linear sketching, finite rate of innovation, **compressed sensing...**)

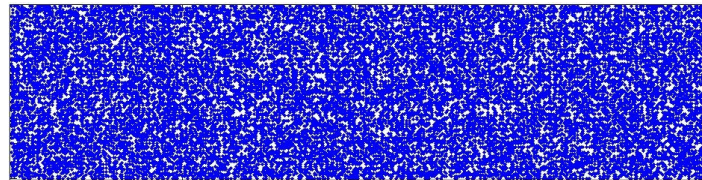
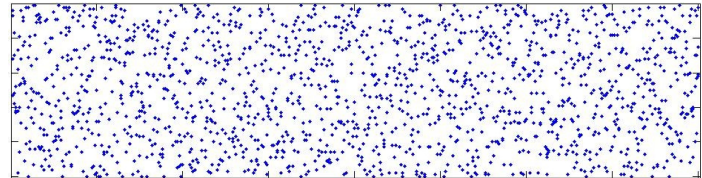
- Setup:
  - Data/signal in  $n$ -dimensional space :  $x$   
E.g.,  $x$  is an  $256 \times 256$  image  $\Rightarrow n=65536$
  - Goal: compress  $x$  into a “sketch”  $Ax$ ,  
where  $A$  is a  $m \times n$  “sketch matrix”,  $m \ll n$
- Requirements:
  - Plan A: want to recover  $x$  from  $Ax$ 
    - Impossible: underdetermined system of equations
  - Plan B: want to recover an “approximation”  $x^*$  of  $x$ 
    - Sparsity parameter  $k$
    - Informally: want to recover largest  $k$  coordinates of  $x$
    - Formally: want  $x^*$  such that
$$\|x^* - x\|_p \leq C(k) \min_{x'} \|x' - x\|_q$$
over all  $x'$  that are  $k$ -sparse (at most  $k$  non-zero entries)
- Want:
  - Good compression (small  $m=m(k,n)$ )
  - Efficient algorithms for encoding and recovery
- Why **linear** compression ?
  - Broader functionality!
  - Useful for compressed signal acquisition, streaming algorithms, etc  
(see Appendix for more info)



$$\begin{pmatrix} A \end{pmatrix} \begin{pmatrix} x \end{pmatrix} = \begin{pmatrix} Ax \end{pmatrix}$$

# Constructing matrix $A$

- “Most” matrices  $A$  work
  - Sparse matrices:
    - Data stream algorithms
    - Coding theory (LDPCs)
  - Dense matrices:
    - Compressed sensing
    - Complexity/learning theory (Fourier matrices)
- “Traditional” tradeoffs:
  - Sparse: computationally more efficient, explicit
  - Dense: shorter sketches
- Recent results: the “best of both worlds”



# Prior and New Results

Paper	Rand. / Det.	Sketch length	Encode time	Column sparsity	Recovery time	Approx
-------	-----------------	------------------	----------------	--------------------	---------------	--------

Scale: Excellent Very Good Good Fair

# Results

Paper	R/D	Sketch length	Encode time	Column sparsity	Recovery time	Approx
[CCF'02], [CM'06]	R	$k \log n$	$n \log n$	$\log n$	$n \log n$	I2 / I2
	R	$k \log^c n$	$n \log^c n$	$\log^c n$	$k \log^c n$	I2 / I2
[CM'04]	R	$k \log n$	$n \log n$	$\log n$	$n \log n$	I1 / I1
	R	$k \log^c n$	$n \log^c n$	$\log^c n$	$k \log^c n$	I1 / I1
[CRT'04] [RV'05]	D	$k \log(n/k)$	$nk \log(n/k)$	$k \log(n/k)$	$n^c$	I2 / I1
	D	$k \log^c n$	$n \log n$	$k \log^c n$	$n^c$	I2 / I1
[GSTV'06] [GSTV'07]	D	$k \log^c n$	$n \log^c n$	$\log^c n$	$k \log^c n$	I1 / I1
	D	$k \log^c n$	$n \log^c n$	$k \log^c n$	$k^2 \log^c n$	I2 / I1
[BGIKS'08]	D	$k \log(n/k)$	$n \log(n/k)$	$\log(n/k)$	$n^c$	I1 / I1
[GLR'08]	D	$k \log n^{\log \log \log n}$	$kn^{1-a}$	$n^{1-a}$	$n^c$	I2 / I1
[NV'07], [DM'08], [NT'08], [BD'08], [GK'09], ...	D	$k \log(n/k)$	$nk \log(n/k)$	$k \log(n/k)$	$nk \log(n/k) * \log$	I2 / I1
	D	$k \log^c n$	$n \log n$	$k \log^c n$	$n \log n * \log$	I2 / I1
[IR'08], [BIR'08],[BI'09]	D	$k \log(n/k)$	$n \log(n/k)$	$\log(n/k)$	$n \log(n/k) * \log$	I1 / I1
[GLSP'09]	R	$k \log(n/k)$	$n \log^c n$	$\log^c n$	$k \log^c n$	I2 / I1

# Part I

Paper	R/D	Sketch length	Encode time	Column sparsity	Recovery time	Approx
[CM'04]	R	$k \log n$	$n \log n$	$\log n$	$n \log n$	$l_1 / l_1$

Theorem: There is a distribution over  $m \times n$  matrices  $A$ ,  $m = O(k \log n)$ , such that for any  $x$ , given  $Ax$ , we can recover  $x^*$  such that

$$\|x - x^*\|_1 \leq C \text{Err}_1, \text{ where } \text{Err}_1 = \min_{k\text{-sparse } x'} \|x - x'\|_1$$

with probability  $1 - 1/n$ .

The recovery algorithm runs in  $O(n \log n)$  time.

This talk:

- Assume  $x \geq 0$  – this simplifies the algorithm and analysis; see the original paper for the general case
- Prove the following  $l_\infty / l_1$  guarantee:  $\|x - x^*\|_\infty \leq C \text{Err}_1 / k$

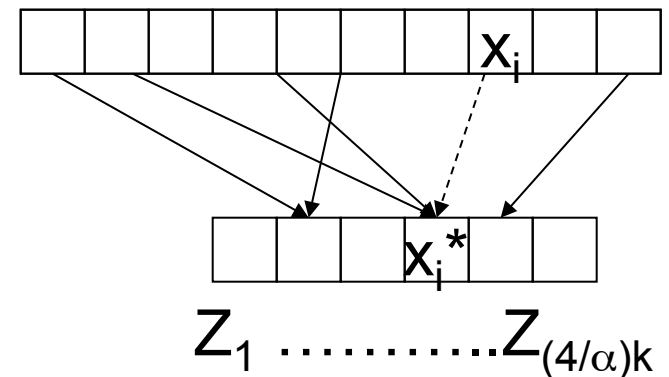
This is actually stronger than the  $l_1 / l_1$  guarantee (cf. [CM'06], see also the Appendix)

Note: [CM'04] originally proved a weaker statement where  $\|x - x^*\|_\infty \leq C \|x\|_1 / k$ . The stronger guarantee follows from the analysis of [CCF'02] (cf. [GGIKMS'02]) who applied it to  $\text{Err}_2$

# First attempt

- Matrix view:
  - A 0-1  $w \times n$  matrix  $A$ , with one 1 per column
  - The  $i$ -th column has 1 at position  $h(i)$ , where  $h(i)$  be chosen uniformly at random from  $\{1 \dots w\}$
- Hashing view:
  - $Z = Ax$
  - $h$  hashes coordinates into “buckets”  $Z_1 \dots Z_w$
- Estimator:  $x_i^* = Z_{h(i)}$

0	0	1	0	0	1	0
1	0	0	0	1	0	0
0	1	0	0	0	0	1
0	0	0	1	0	0	0



Closely related: [Estan-Varghese'03], “counting” Bloom filters

# Analysis

- We show

$$x_i^* \leq x_i \pm \alpha \text{Err}/k$$

with probability  $> 1/2$

- Assume

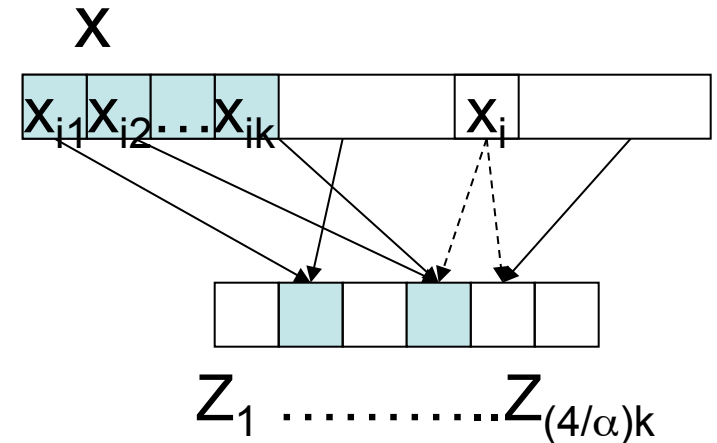
$$|x_{i_1}| \geq \dots \geq |x_{i_m}|$$

and let  $S = \{i_1 \dots i_k\}$  (“elephants”)

- When is  $x_i^* > x_i \pm \alpha \text{Err}/k$  ?

- **Event 1:**  $S$  and  $i$  collide, i.e.,  $h(i) \in h(S - \{i\})$   
Probability: at most  $k/(4/\alpha)k = \alpha/4 < 1/4$  (if  $\alpha < 1$ )
- **Event 2:** many “mice” collide with  $i$ , i.e.,  
$$\sum_{t \text{ not in } S \cup \{i\}, h(i)=h(t)} x_t > \alpha \text{Err}/k$$
  
Probability: at most  $1/4$  (Markov inequality)

- Total probability of “bad” events  $< 1/2$





# Second try

- Algorithm:

- Maintain  $d$  functions  $h_1 \dots h_d$  and vectors  $Z^1 \dots Z^d$
- Estimator:

$$X_i^* = \min_t Z_{ht(i)}^t$$

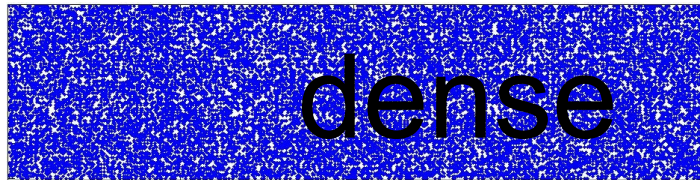
- Analysis:

- $\Pr[|x_i^* - x_i| \geq \alpha \text{Err}/k] \leq 1/2^d$
- Setting  $d = O(\log n)$  (and thus  $m = O(k \log n)$ ) ensures that w.h.p

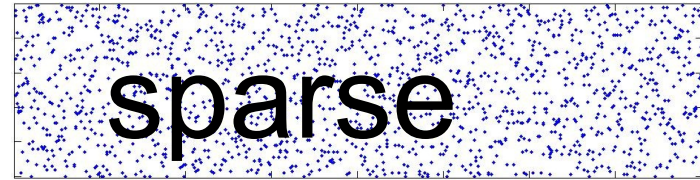
$$|x_i^* - x_i| < \alpha \text{Err}/k$$

# Part II

Paper	R/ D	Sketch length	Encode time	Column sparsity	Recovery time	Approx
[BGIKS'08]	D	$k \log(n/k)$	$n \log(n/k)$	$\log(n/k)$	$n^c$	11 / 11
[IR'08], [BIR'08],[BI'09]	D	$k \log(n/k)$	$n \log(n/k)$	$\log(n/k)$	$n \log(n/k) * \log$	11 / 11



vs.



- Restricted Isometry Property (RIP) [Candes-Tao'04]

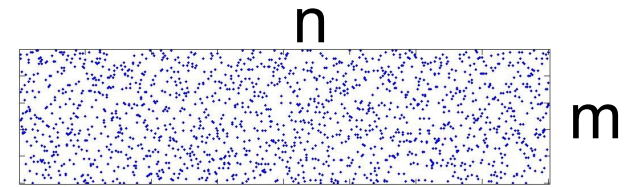
$$\Delta \text{ is } k\text{-sparse} \Rightarrow \|\Delta\|_2 \leq \|A\Delta\|_2 \leq C \|\Delta\|_2$$

- Holds w.h.p. for:
  - Random Gaussian/Bernoulli:  $m = O(k \log(n/k))$
  - Random Fourier:  $m = O(k \log^{O(1)} n)$
- Consider  $m \times n$  0-1 matrices with  $d$  ones per column
- Do they satisfy RIP ?
  - No, unless  $m = \Omega(k^2)$  [Chandar'07]
- However, they can satisfy the following RIP-1 property [Berinde-Gilbert-Indyk-Karloff-Strauss'08]:

$$\Delta \text{ is } k\text{-sparse} \Rightarrow d(1-\epsilon) \|\Delta\|_1 \leq \|A\Delta\|_1 \leq d \|\Delta\|_1$$

- Sufficient (and necessary) condition: the underlying graph is a  $(k, d(1-\epsilon/2))$ -expander

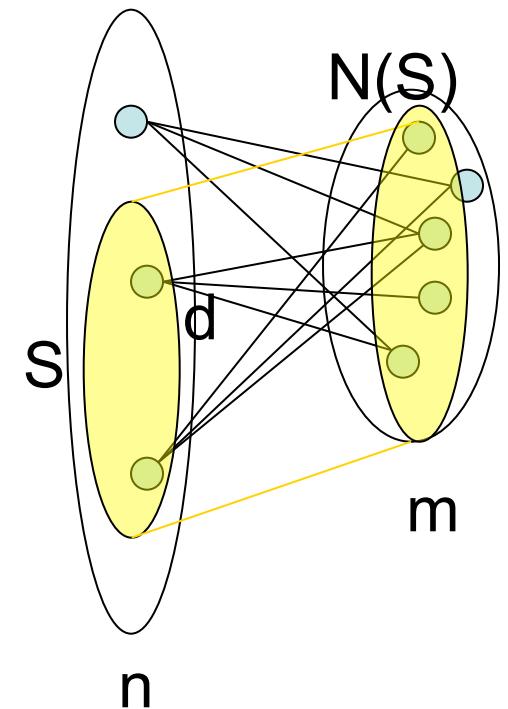
# Expanders



- A bipartite graph is a  $(k, d(1-\epsilon))$ -**expander** if for any left set  $S$ ,  $|S| \leq k$ , we have  $|N(S)| \geq (1-\epsilon)d |S|$
- Objects well-studied in theoretical computer science and coding theory
- Constructions:
  - Probabilistic:  $m = O(k \log(n/k))$
  - Explicit:  $m = k \text{ quasipolylog } n$
- High expansion implies RIP-1:

$$\Delta \text{ is } k\text{-sparse} \Rightarrow d(1-\epsilon) \|\Delta\|_1 \leq \|A\Delta\|_1 \leq d\|\Delta\|_1$$

[Berinde-Gilbert-Indyk-Karloff-Strauss'08]



# Proof: $d(1-\varepsilon/2)$ -expansion $\Rightarrow$ RIP-1

- Want to show that for any  $k$ -sparse  $\Delta$  we have

$$d(1-\varepsilon) \|\Delta\|_1 \leq \|A\Delta\|_1 \leq d\|\Delta\|_1$$

- RHS inequality holds for **any**  $\Delta$

- LHS inequality:

- W.l.o.g. assume

$$|\Delta_1| \geq \dots \geq |\Delta_k| \geq |\Delta_{k+1}| = \dots = |\Delta_n| = 0$$

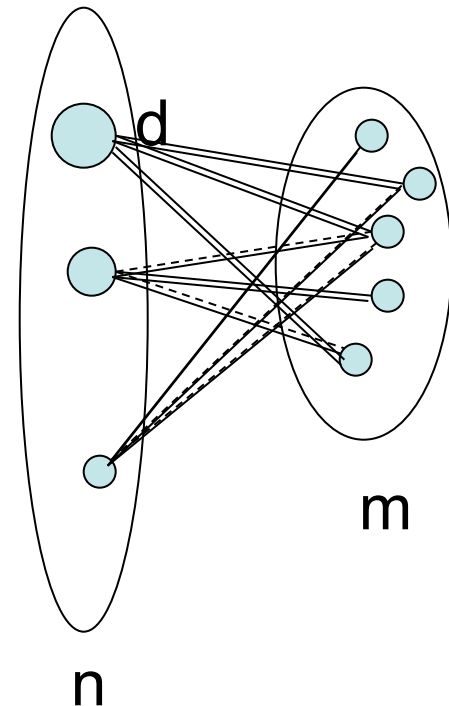
- Consider the edges  $e=(i,j)$  in a lexicographic order

- For each edge  $e=(i,j)$  define  $r(e)$  s.t.

- $r(e)=-1$  if there exists an edge  $(i',j) < (i,j)$
- $r(e)=1$  if there is no such edge

- Claim 1:  $\|A\Delta\|_1 \geq \sum_{e=(i,j)} |\Delta_i| r_e$

- Claim 2:  $\sum_{e=(i,j)} |\Delta_i| r_e \geq (1-\varepsilon) d\|\Delta\|_1$



Recovery: algorithms

# Matching Pursuit(s)

$$\left( \begin{array}{c} A \\ \color{red}{\square} \\ \color{blue}{\square} \\ \color{green}{\square} \\ \color{grey}{\square} \\ i \end{array} \right) i \begin{pmatrix} x^* - x \\ \color{red}{\square} \end{pmatrix} = \begin{pmatrix} \color{red}{\square} \\ \color{blue}{\square} \\ \color{green}{\square} \\ \color{grey}{\square} \end{pmatrix} \begin{array}{c} Ax - Ax^* \end{array}$$

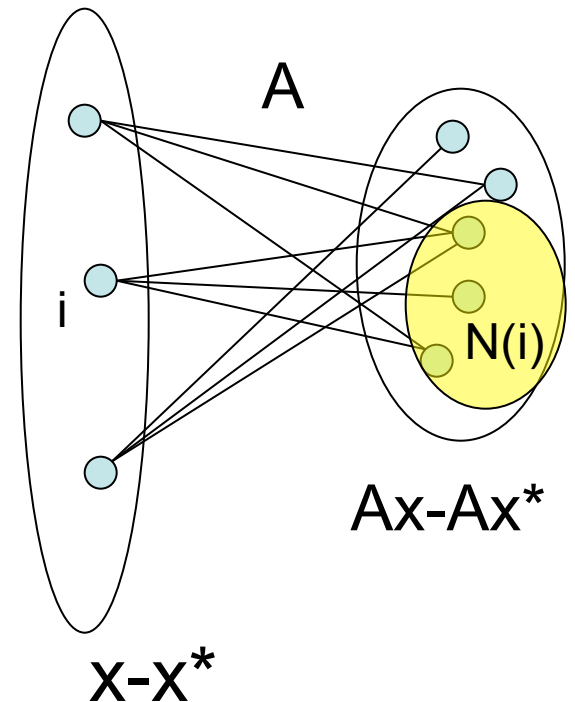
- Iterative algorithm: given current approximation  $x^*$  :
  - Find (possibly several)  $i$  s. t.  $A_i$  “correlates” with  $Ax - Ax^*$  . This yields  $i$  and  $z$  s. t.

$$\|x^* + ze_i - x\|_p \ll \|x^* - x\|_p$$

- Update  $x^*$
- Sparsify  $x^*$  (keep only  $k$  largest entries)
- Repeat
- Norms:
  - $p=2$  : CoSaMP, SP, IHT etc (RIP)
  - $p=1$  : SMP, SSMP (RIP-1)
  - $p=0$  : LDPC bit flipping (sparse matrices)

# Sequential Sparse Matching Pursuit

- Algorithm:
  - $x^*=0$
  - Repeat  $T$  times
    - Repeat  $S=O(k)$  times
      - Find  $i$  and  $z$  that minimize\*  $\|A(x^*+ze_i)-Ax\|_1$
      - $x^* = x^*+ze_i$
    - Sparsify  $x^*$   
(set all but  $k$  largest entries of  $x^*$  to 0)
- Similar to SMP, but updates done sequentially

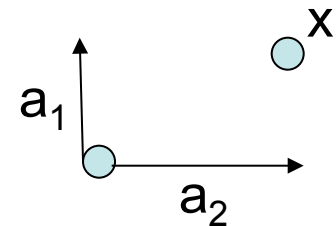
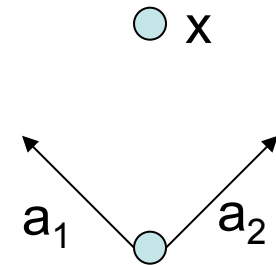


\* Set  $z=\text{median}[(Ax^*-Ax)_{N(i)}]$ . Instead, one could first optimize (gradient)  $i$  and then  $z$  [ Fuchs'09]



# SSMP: Approximation guarantee

- Want to find  $k$ -sparse  $x^*$  that minimizes  $\|x-x^*\|_1$
- By RIP1, this is approximately the same as minimizing  $\|Ax-Ax^*\|_1$
- Need to show we can do it *greedily*



Supports of  $a_1$  and  $a_2$  have small overlap (typically)

# Conclusions

- Sparse approximation using sparse matrices
  - State of the art: deterministically can do 2 out of 3:
    - Near-linear encoding/decoding
    - $O(k \log(n/k))$  measurements
    - Approximation guarantee with respect to L2/L1 norm
- } This talk
- Open problems:
    - 3 out of 3 ?
    - Explicit constructions ?
  - For more, see  
A. Gilbert, P. Indyk, “Sparse recovery using sparse matrices”,  
Proceedings of IEEE, June 2010.

# Appendix

# $l_\infty/l_1$ implies $l_1/l_1$

- Algorithm:
  - Assume we have  $x^*$  s.t.  $\|x-x^*\|_\infty \leq C \text{Err}_1 / k$ .
  - Let vector  $x'$  consist of  $k$  largest (in magnitude) elements of  $x^*$

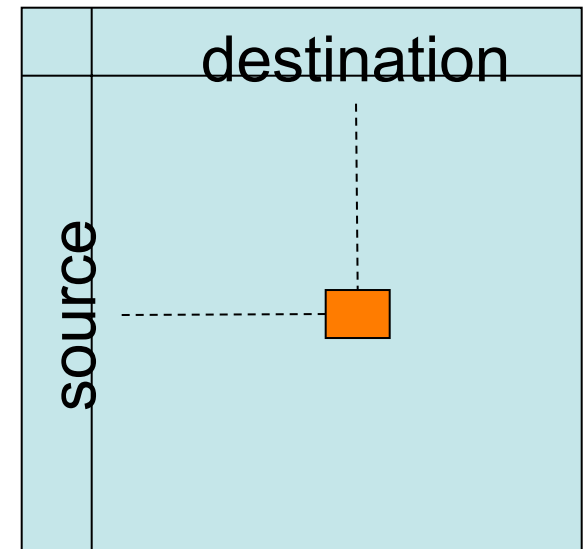
- Analysis

- Let  $S$  (or  $S^*$ ) be the set of  $k$  largest in magnitude coordinates of  $x$  (or  $x^*$ )
- Note that  $\|x^*_S\|_1 \leq \|x^*_{S^*}\|_1$
- We have

$$\begin{aligned}
 \|x-x'\|_1 &\leq \|x\|_1 - \|x_S\|_1 + \|x_{S^*}-x^*_{S^*}\|_1 \\
 &\leq \|x\|_1 - \|x^*_{S^*}\|_1 + 2\|x_{S^*}-x^*_{S^*}\|_1 \\
 &\leq \|x\|_1 - \|x^*_S\|_1 + 2\|x_{S^*}-x^*_{S^*}\|_1 \\
 &\leq \|x\|_1 - \|x_S\|_1 + \|x^*_S-x_S\|_1 + 2\|x_{S^*}-x^*_{S^*}\|_1 \\
 &\leq \text{Err} + 3\alpha/k * k \\
 &\leq (1+3\alpha)\text{Err}
 \end{aligned}$$

# Application I: Monitoring Network Traffic Data Streams

- Router routs packets
  - Where do they come from ?
  - Where do they go to ?
- Ideally, would like to maintain a traffic matrix  $x[.,.]$ 
  - Easy to update: given a (src,dst) packet, increment  $x_{src,dst}$
  - Requires way too much space! ( $2^{32} \times 2^{32}$  entries)
  - Need to **compress**  $x$ , **increment** easily
- Using linear compression we can:
  - Maintain sketch  $Ax$  under increments to  $x$ , since
$$A(x+\Delta) = Ax + A\Delta$$
  - Recover  $x^*$  from  $Ax$

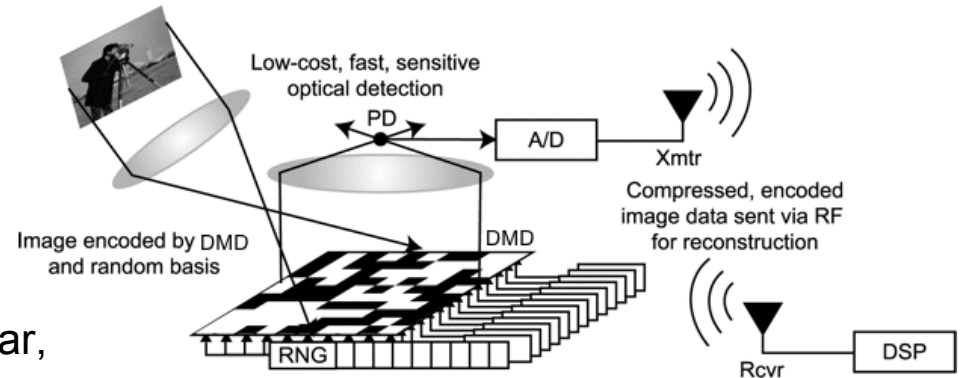


$x$

# Applications, ctd.

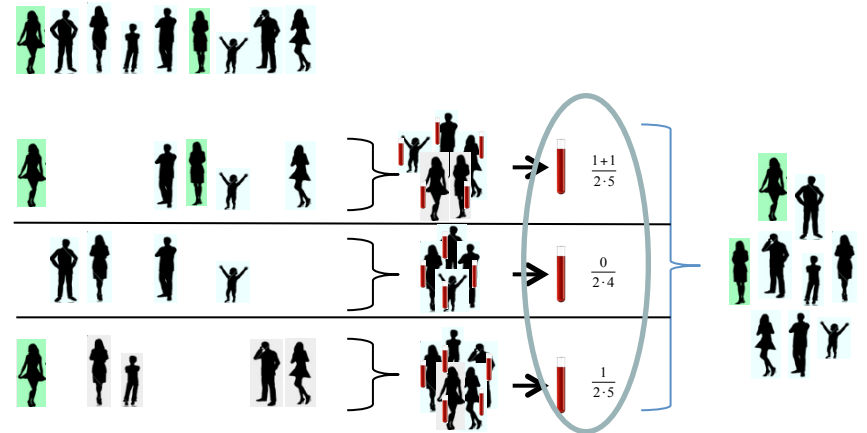
- Single pixel camera

[Wakin, Laska, Duarte, Baron, Sarvotham, Takhar, Kelly, Baraniuk'06]



- Pooling Experiments

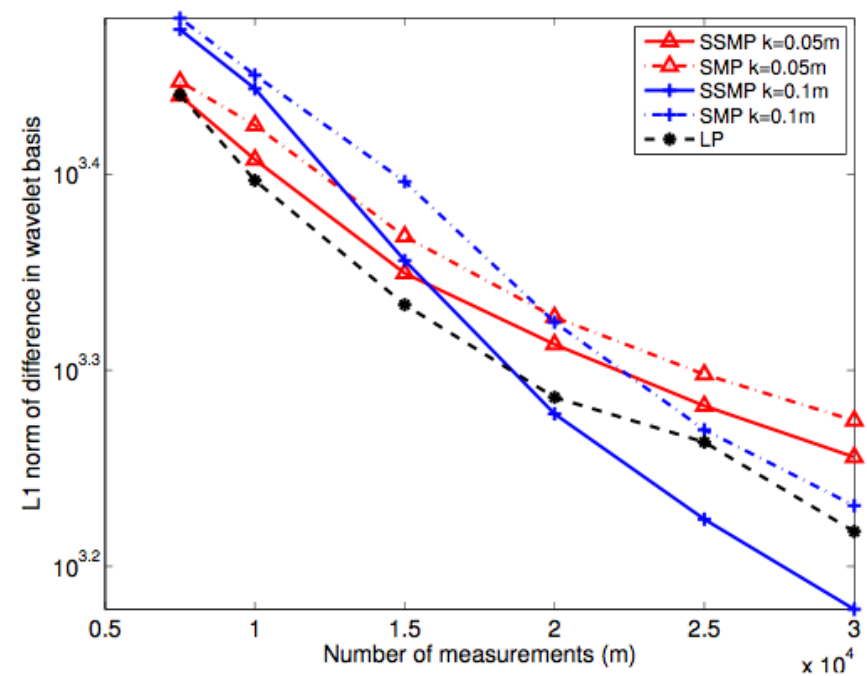
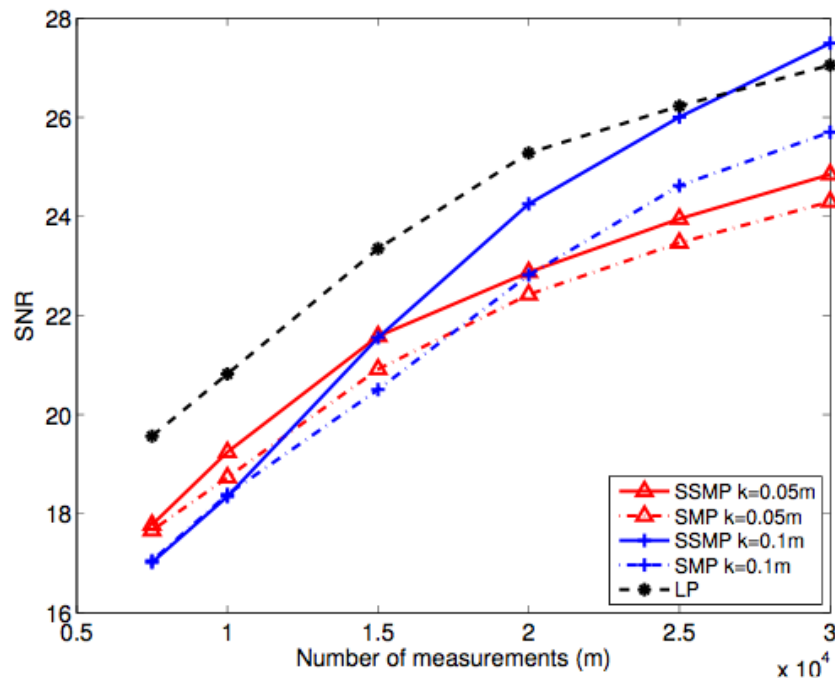
[Kainkaryam, Woolf'08], [Hassibi et al'07], [Dai-Sheikh, Milenkovic, Baraniuk], [Shental-Amir-Zuk'09],[Erlich-Shental-Amir-Zuk'09]



# Experiments



256x256



SSMP is ran with  $S=10000, T=20$ . SMP is ran for 100 iterations. Matrix sparsity is  $d=8$ .

# SSMP: Running time

- Algorithm:
  - $x^*=0$
  - Repeat  $T$  times
    - For each  $i=1\dots n$  compute\*  $z_i$  that achieves
 
$$D_i = \min_z \|A(x^* + ze_i) - b\|_1$$
 and store  $D_i$  in a heap
    - Repeat  $S=O(k)$  times
      - Pick  $i, z$  that yield the best gain
      - Update  $x^* = x^* + ze_i$
      - Recompute and store  $D_{i'}$  for all  $i'$  such that  $N(i)$  and  $N(i')$  intersect
    - Sparsify  $x^*$   
(set all but  $k$  largest entries of  $x^*$  to 0)
- Running time:
 
$$T [ n(d+\log n) + k nd/m*d (d+\log n) ]$$

$$= T [ n(d+\log n) + nd (d+\log n) ] = T [ nd (d+\log n) ]$$

