

# To Adapt or Not To Adapt

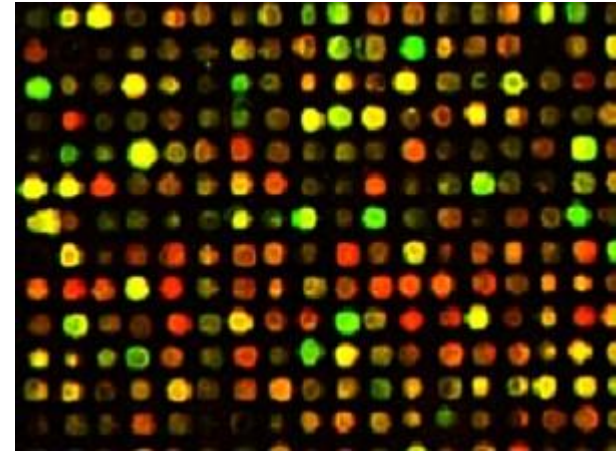
## The Power and Limits of Adaptive Sensing

*Mark A. Davenport*

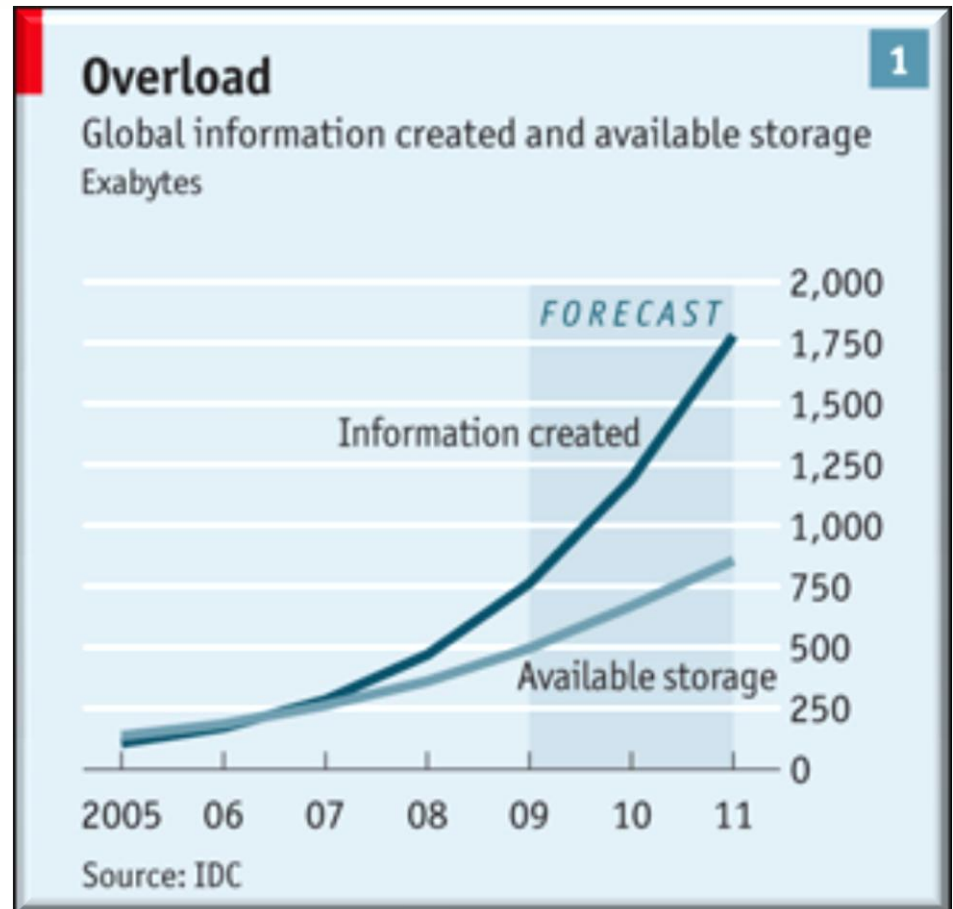
Stanford University  
Department of Statistics



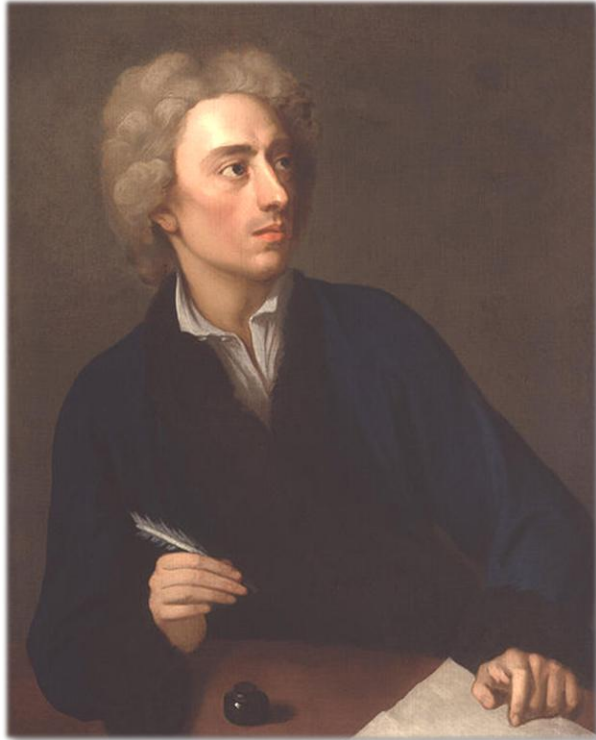
# Sensor Explosion



# Data Deluge



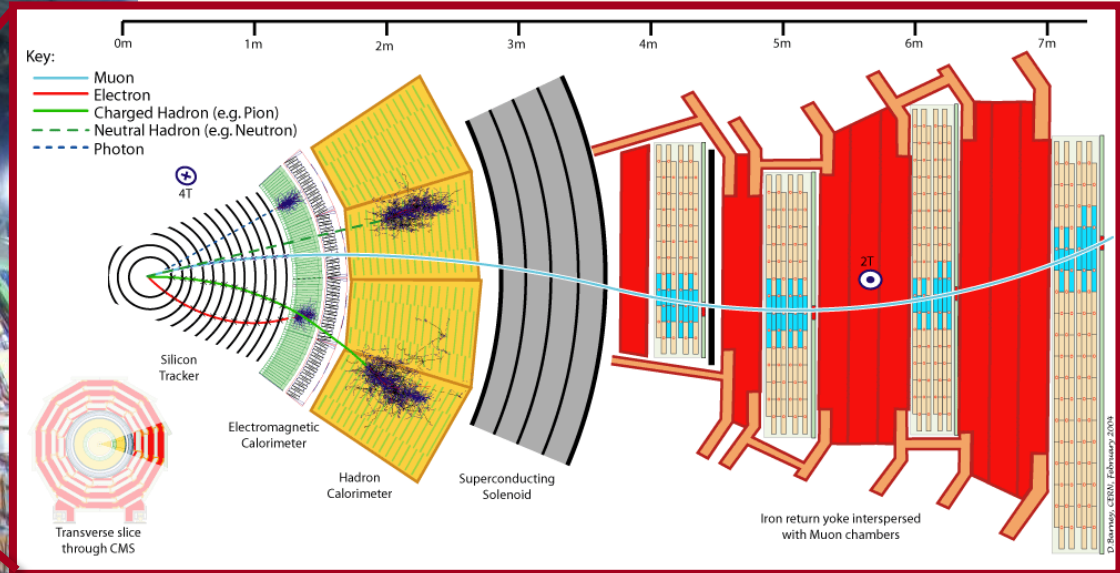
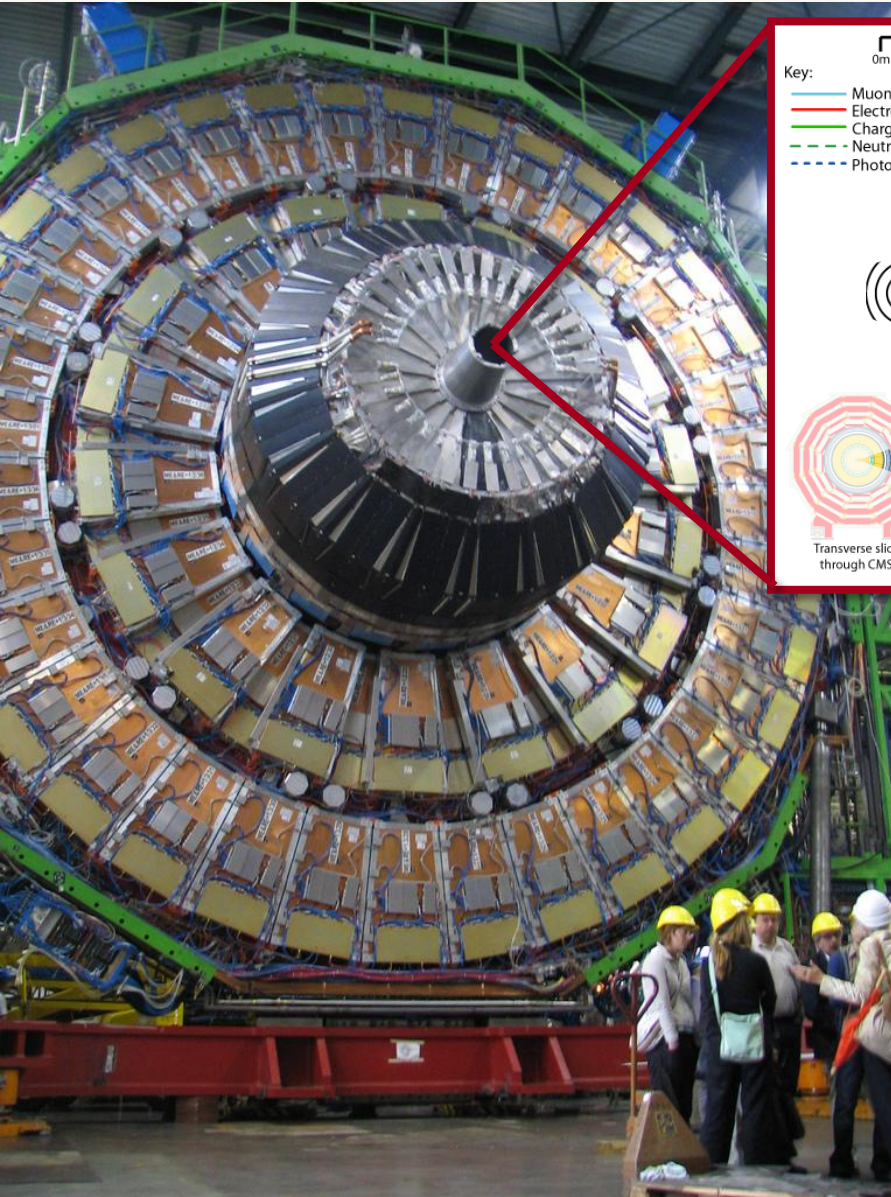
# Ye Olde Data Deluge



“Paper became so cheap, and printers so numerous, that a deluge of authors covered the land”

Alexander Pope, 1728

# Large Hadron Collider at CERN



Compact Muon Solenoid detector

*320 terabits per second* raw data

Stop-gap: perform ad-hoc triage to 800 Gbps, recording only “interesting events”

# Data Deluge Challenges

~~How can we get our hands on as much data as possible?~~

~~How can we extract as much information as possible from a limited amount of data?~~



How can we avoid having to acquire so much data to begin with?



How can we extract any information at all from a massive amount of high-dimensional data?

# Low-Dimensional Structure

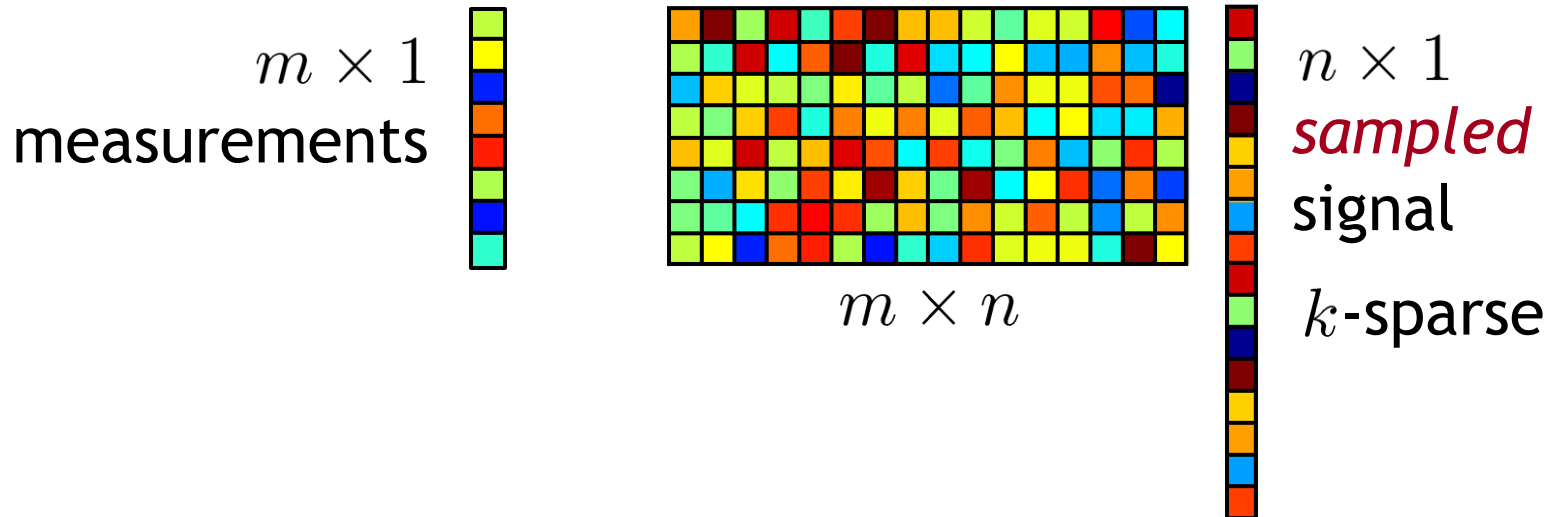
How can we exploit low-dimensional structure to address the challenges posed by the “data deluge”?

- Visualization
- Feature extraction/selection
- Compression
- Regularization of ill-posed inverse problems
- Underpins *compressive sensing*

# Compressive Sensing

Replace samples with general *linear measurements*

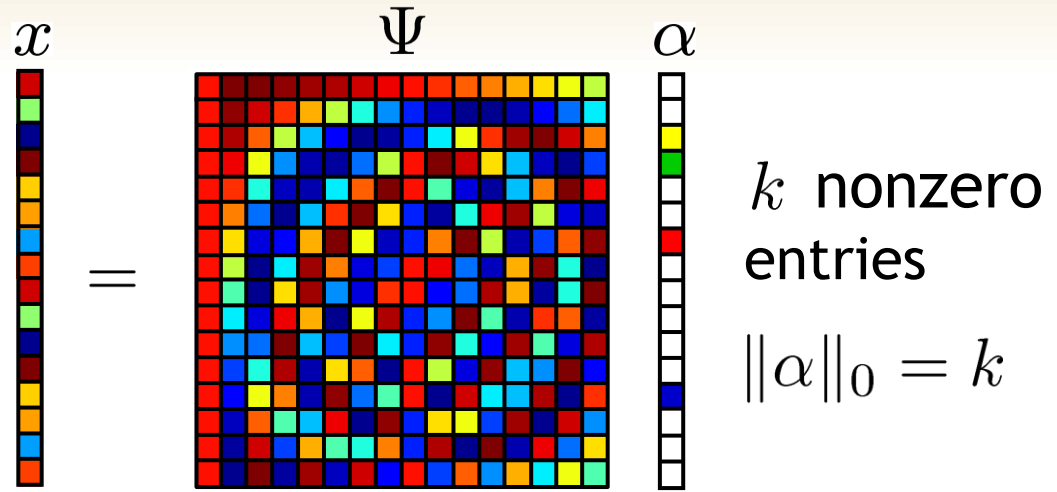
$$y = A x$$



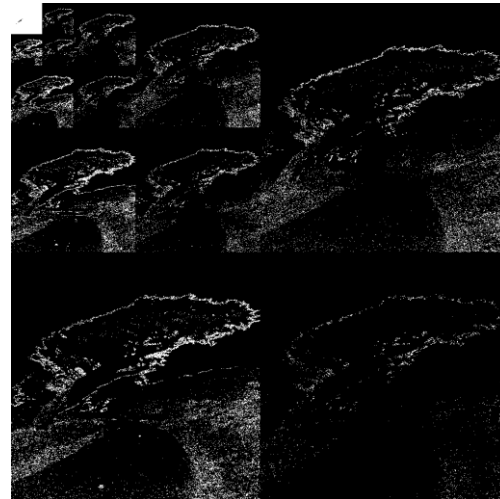


# Sparsity

$$x = \sum_{j=1}^n \alpha_j \psi_j$$
$$= \Psi \alpha$$



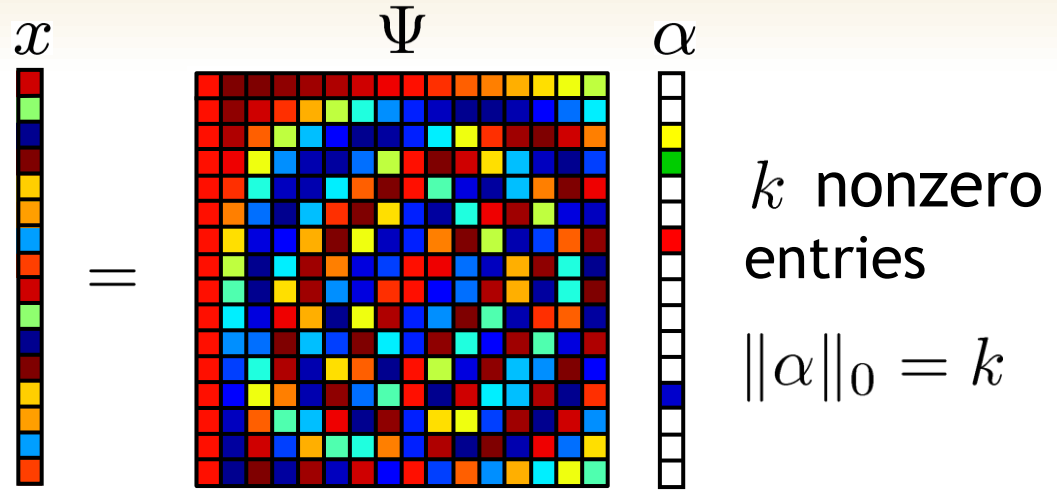
$n$   
pixels



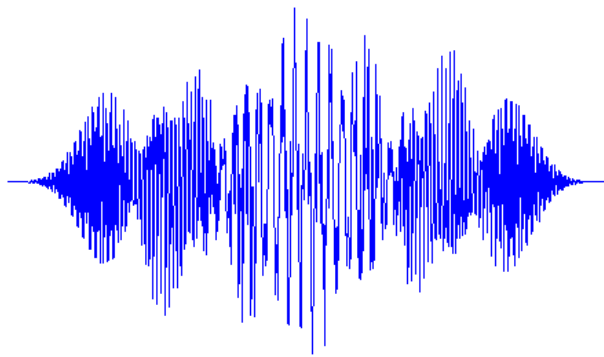
$k \ll n$   
large  
wavelet  
coefficients

# Sparsity

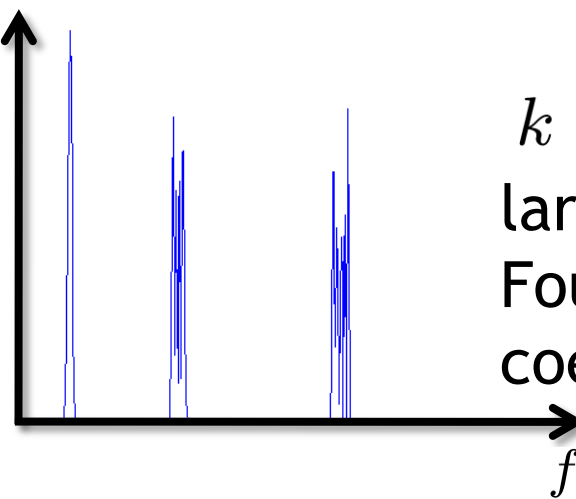
$$x = \sum_{j=1}^n \alpha_j \psi_j$$
$$= \Psi \alpha$$



$n$   
samples



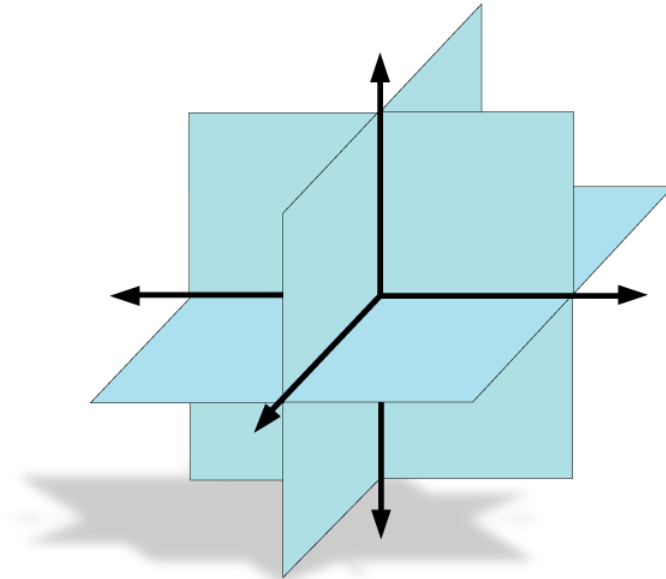
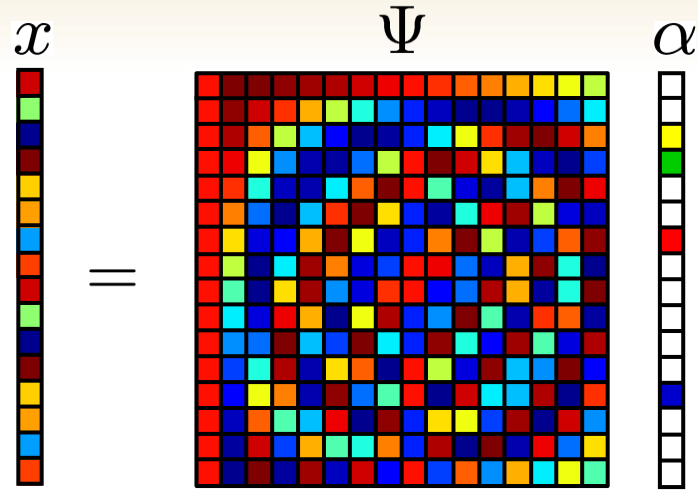
$X(f)$



$k \ll n$   
large  
Fourier  
coefficients

# Sparsity

$$x = \sum_{j=1}^n \alpha_j \psi_j$$
$$= \Psi \alpha$$



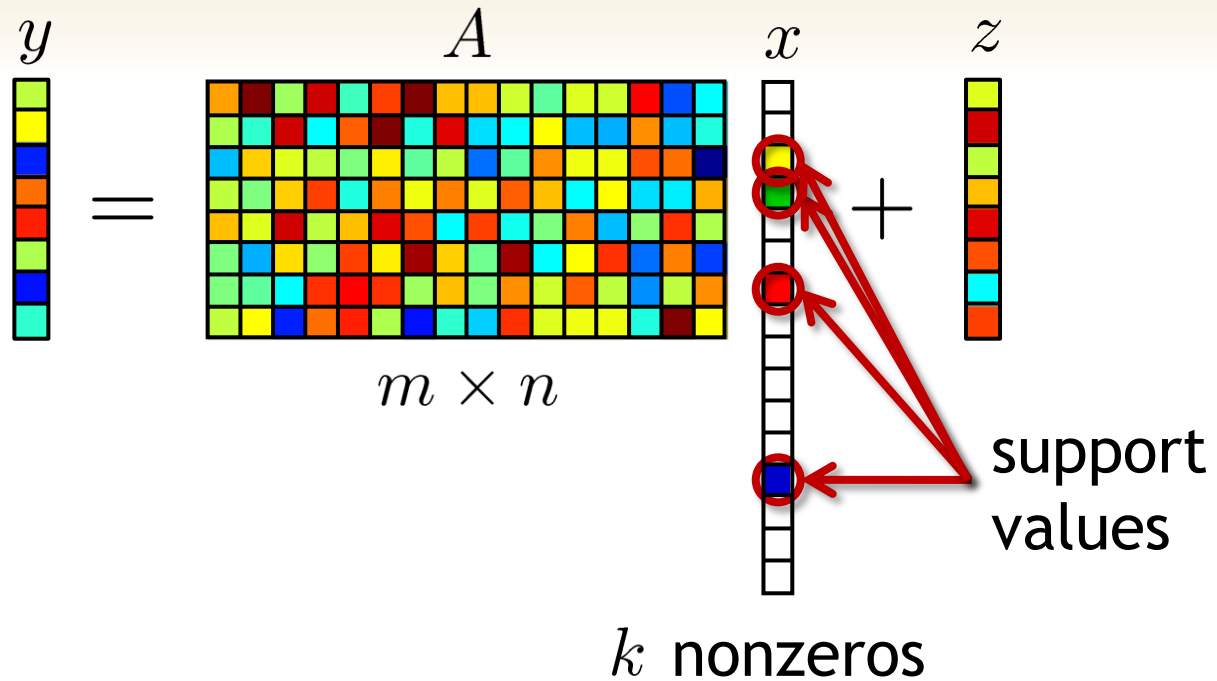
# Compressive Sensing

The diagram illustrates the compressive sensing equation  $y = Ax + z$ . On the left, a vertical vector  $y$  of size  $m \times 1$  is shown. This is equal to the product of a matrix  $A$  of size  $m \times n$  and a vector  $x$  of size  $n \times 1$ . The matrix  $A$  is represented as a grid of colored squares. Below the matrix, the dimensions are given as  $m \times n$  and  $m \ll n$ . The vector  $x$  is shown as a vertical column of colored squares, with the label  $n \times 1$  and  $k$ -sparse below it. A plus sign indicates the addition of a noise vector  $z$ , which is a vertical column of colored squares of size  $m \times 1$ .

When (and how well) can we estimate  $x$  from the measurements  $y$ ?

# **Review of Nonadaptive Compressive Sensing**

# Compressive Sensing



- How should we design  $A$  to ensure that  $y$  contains as much information about  $x$  as possible?
- What algorithms do we have for recovering  $x$  from  $y$ ?

# How To Design $A$ ?

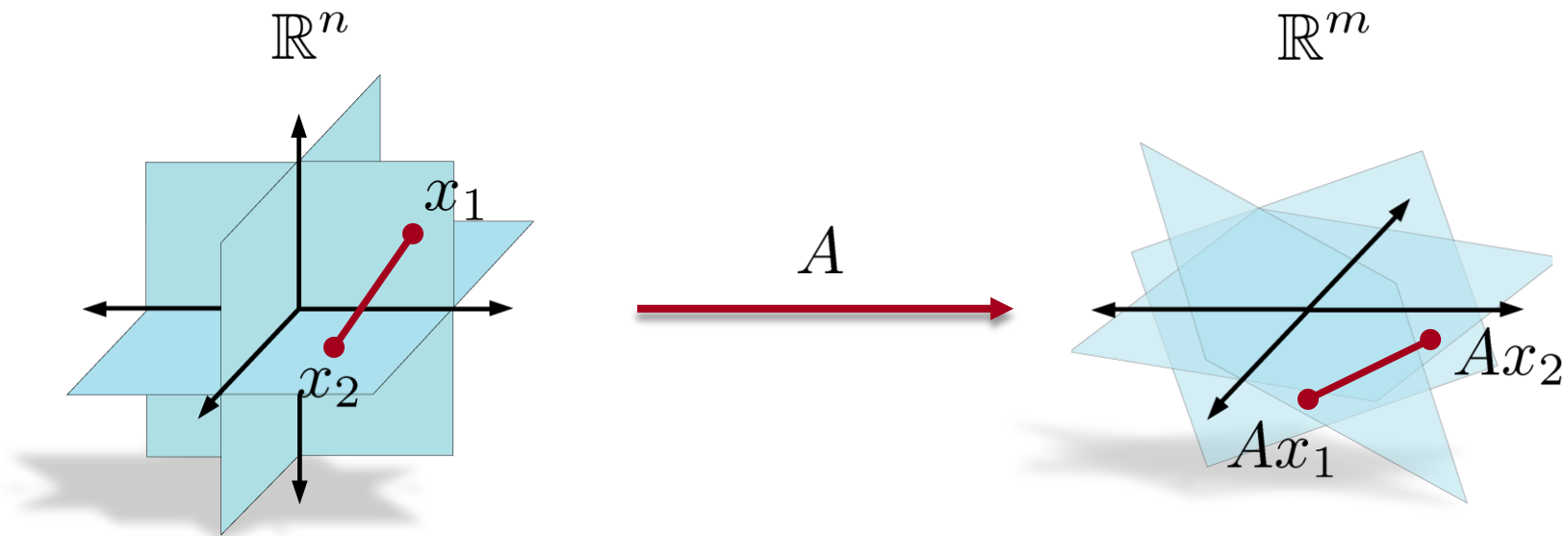
Prototypical sensing model:

$$y = Ax + z \quad z \sim \mathcal{N}(0, \sigma^2 I)$$

- Constrain  $A$  to have unit-norm rows
- Pick  $A$  at *random!*
  - i.i.d. Gaussian entries (with variance  $1/n$ )
  - random rows from a unitary matrix
- As long as  $m = O(k \log(n/k))$ , with high probability a random  $A$  will satisfy the *restricted isometry property*

# Restricted Isometry Property (RIP)

$$\frac{\|Ax_1 - Ax_2\|_2^2}{\|x_1 - x_2\|_2^2} \approx \frac{m}{n} \quad \|x_1\|_0, \|x_2\|_0 \leq k$$





# How To Design $A$ ?

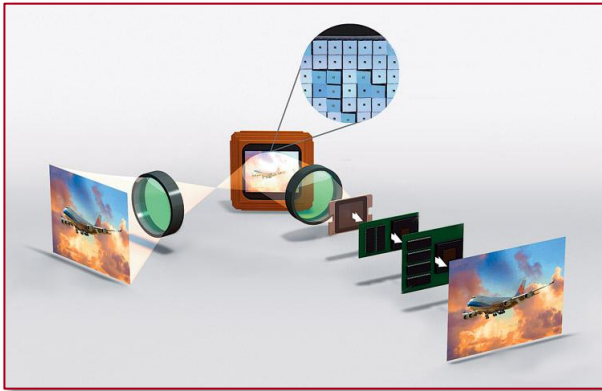
Prototypical sensing model:

$$y = Ax + z \quad z \sim \mathcal{N}(0, \sigma^2 I)$$

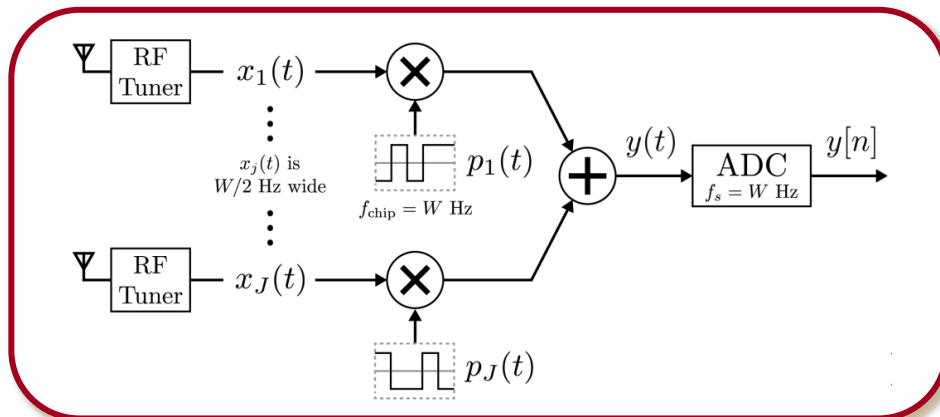
- Constrain  $A$  to have unit-norm rows
- Pick  $A$  at *random!*
  - i.i.d. Gaussian entries (with variance  $1/n$ )
  - random rows from a unitary matrix
- As long as  $m = O(k \log(n/k))$ , with high probability a random  $A$  will satisfy the *restricted isometry property*
- Deep connections with *Johnson-Lindenstrauss Lemma*
  - see Baraniuk, Davenport, DeVore, and Wakin (2008)

# Architectures For “Random Sensing”

Single-Pixel Camera [Duarte, Davenport, et al. - 2008]



Compressive Multiplexor [Slavinsky, Laska, Davenport, and Baraniuk - 2011]



# How To Recover $x$ ?

- Lots and lots of algorithms
  - $\ell_1$ -minimization
  - greedy algorithms (matching pursuit, CoSaMP, IHT)

If  $A$  satisfies the RIP,  $\|x\|_0 \leq k$ , and  $y = Ax + z$  with  $z \sim \mathcal{N}(0, \sigma^2 I)$ , then

$$\hat{x} = \arg \min_{x' \in \mathbb{R}^n} \|x'\|_1$$

$$\text{s.t. } \|A^*(y - Ax')\|_\infty \leq c\sqrt{\log n}\sigma$$

satisfies

$$\mathbb{E} \|\hat{x} - x\|_2^2 \leq C \frac{n}{m} k \sigma^2 \log n.$$

[Candès and Tao - 2005]

# Room For Improvement?

There exists matrices  $A$  such that for *any* (sparse)  $x$  we have

$$\mathbb{E} \|\hat{x} - x\|_2^2 \leq C \frac{n}{m} k \sigma^2 \log n.$$

$$y_i = \langle a_i, x \rangle + z_i$$



$a_i$  and  $x$  are almost orthogonal

- We are using most of our “sensing power” to sense entries that aren’t even there!
- Tremendous loss in signal-to-noise ratio (SNR)
- It’s hard to imagine any way to avoid this...

# Can We Do Better?

## Theorem

For *any* matrix  $A$  (with unit-norm rows) and *any* recovery procedure  $\hat{x}$ , there exists an  $x$  with  $\|x\|_0 \leq k$  such that if  $y = Ax + z$  with  $z \sim \mathcal{N}(0, \sigma^2 I)$ , then

$$\mathbb{E} \|\hat{x}(y) - x\|_2^2 \geq C' \frac{n}{m} k \sigma^2 \log(n/k).$$

Compressive sensing is already operating at the limit

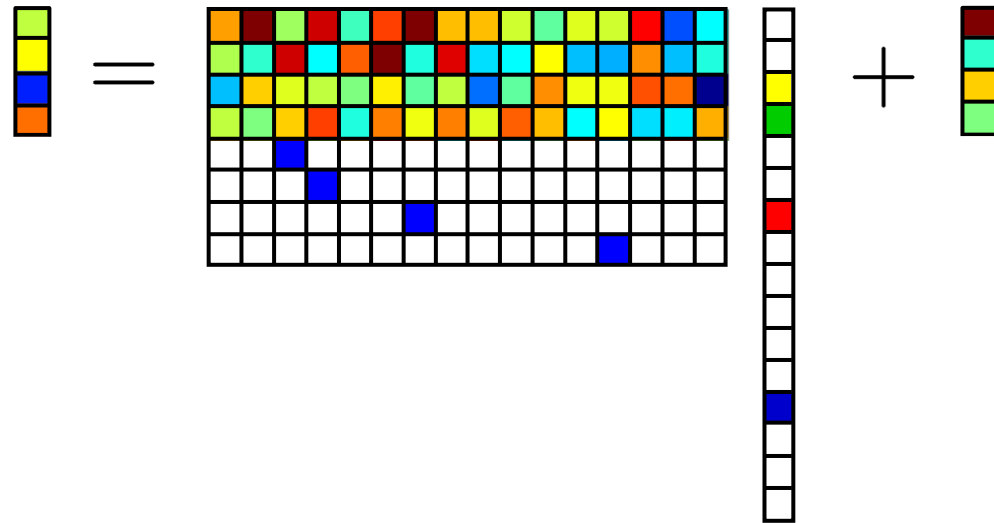
Proof ingredients:

- construct unfavorable prior: *Matrix Bernstein inequality*
- use *Fano's inequality* to show that Bayes risk is large

# **Adaptive Sensing**

# Adaptive Sensing

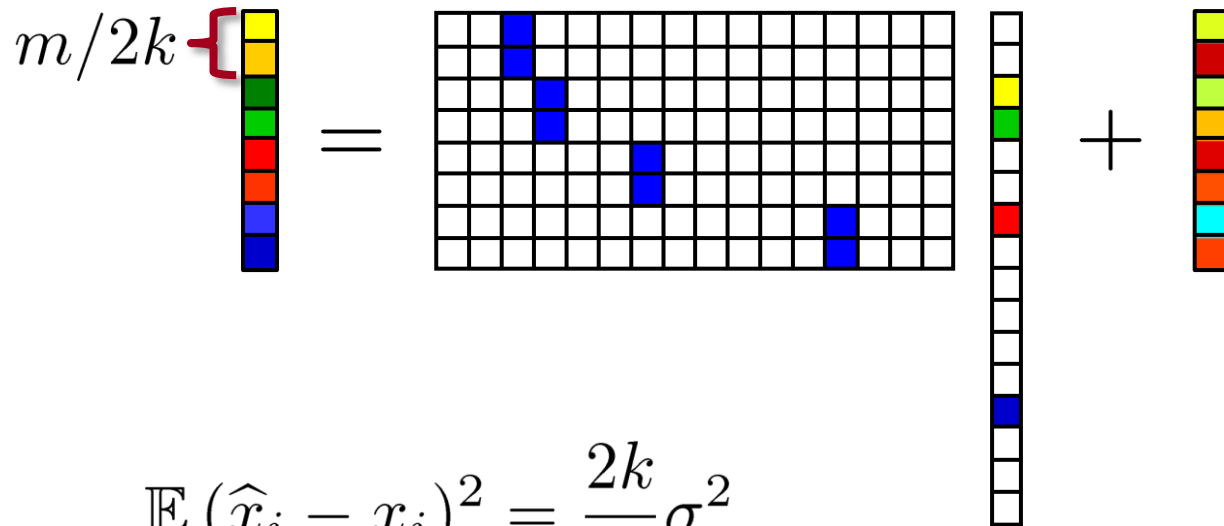
Think of sensing as a game of 20 questions



Simple strategy: Use  $m/2$  measurements to find the support, and the remainder to estimate the values.

# Thought Experiment

Suppose that after  $m/2$  measurements we have perfectly estimated the support.



$$\mathbb{E} (\hat{x}_i - x_i)^2 = \frac{2k}{m} \sigma^2$$

$$\mathbb{E} \|\hat{x} - x\|_2^2 = \frac{2k}{m} k \sigma^2 \ll \frac{n}{m} k \sigma^2 \log n$$



# Does Adaptivity *Really* Help?

Sometimes...

- Noise-free measurements, but non-sparse signal
  - adaptivity doesn't help if you want a uniform guarantee
  - probabilistic adaptive algorithms can reduce the required number of measurements from  $O(k \log(n/k))$  to  $O(k \log \log(n/k))$  [Indyk et al. - 2011]
- Noisy setting
  - distilled sensing [Haupt et al. - 2007, 2010]
  - adaptivity can reduce the estimation error to

$$\mathbb{E} \|\hat{x} - x\|_2^2 = \frac{n}{m} k \sigma^2$$

$$\mathbb{E} \|\hat{x} - x\|_2^2 = \frac{k}{m} k \sigma^2$$

*Which is it?*



# Which Is It?

Suppose we have a budget of  $m$  measurements of the form  $y_i = \langle a_i, x \rangle + z_i$  where  $\|a_i\|_2 = 1$  and  $z_i \sim \mathcal{N}(0, \sigma^2)$

The vector  $a_i$  can have an arbitrary dependence on the measurement history, i.e.,  $(a_1, y_1), \dots, (a_{i-1}, y_{i-1})$

## Theorem

There exist  $x$  with  $\|x\|_0 \leq k$  such that for *any* adaptive measurement strategy and *any* recovery procedure  $\hat{x}$ ,

$$\mathbb{E} \|\hat{x}(y) - x\|_2^2 \geq C \frac{n}{m} k \sigma^2.$$

Thus, in general, adaptivity does *not* significantly help!

# Proof Strategy

**Step 1:** Consider sparse signals with nonzeros of amplitude

$$\mu \approx \sigma \sqrt{n/m}$$

**Step 2:** Show that if given a budget of  $m$  measurements, you cannot detect the support very well

**Step 3:** Immediately translate this into a lower bound on the MSE

To make things simpler, we will consider a Bernoulli prior  $\pi(x)$  instead of a uniform  $k$ -sparse prior:

$$x_j = \begin{cases} 0 & \text{with probability } 1 - k/n \\ \mu > 0 & \text{with probability } k/n \end{cases}$$


# Proof of Main Result

Let  $S = \{j : x_j \neq 0\}$  and set  $\sigma^2 = 1$

For any estimator  $\hat{x}$ , define  $\hat{S} := \{j : |\hat{x}_j| \geq \mu/2\}$

Whenever  $j \in S \setminus \hat{S}$  or  $j \in \hat{S} \setminus S$ ,  $|\hat{x}_j - x_j| \geq \mu/2$

$$\|\hat{x} - x\|_2^2 \geq \frac{\mu^2}{4} |S \setminus \hat{S}| + \frac{\mu^2}{4} |\hat{S} \setminus S| = \frac{\mu^2}{4} |\hat{S} \Delta S|$$

  $\mathbb{E} \|\hat{x} - x\|_2^2 \geq \frac{\mu^2}{4} \mathbb{E} |\hat{S} \Delta S|$

# Proof of Main Result

## Lemma

Under the Bernoulli prior, *any* estimate  $\hat{S}$  satisfies

$$\mathbb{E} |\hat{S} \Delta S| \geq k \left( 1 - \frac{\mu}{2} \sqrt{\frac{m}{n}} \right).$$

Thus, 
$$\begin{aligned} \mathbb{E} \|\hat{x} - x\|_2^2 &\geq \frac{\mu^2}{4} \mathbb{E} |\hat{S} \Delta S| \\ &\geq k \cdot \frac{\mu^2}{4} \left( 1 - \frac{\mu}{2} \sqrt{\frac{m}{n}} \right) \end{aligned}$$

Plug in  $\mu = \frac{8}{3} \sqrt{\frac{n}{m}}$  and this reduces to

$$\mathbb{E} \|\hat{x} - x\|_2^2 \geq \frac{4}{27} \cdot \frac{kn}{m} \geq \frac{1}{7} \cdot \frac{kn}{m}$$

# Key Ideas in Proof of Lemma

$$\mathbb{P}_{0,j}(y_1, \dots, y_m) = \mathbb{P}(y_1, \dots, y_m | x_j = 0)$$

$$\mathbb{P}_{1,j}(y_1, \dots, y_m) = \mathbb{P}(y_1, \dots, y_m | x_j = \mu)$$

$$\begin{aligned} \mathbb{E} |\widehat{S} \Delta S| &\geq \frac{k}{n} \sum_j (1 - \|\mathbb{P}_{1,j} - \mathbb{P}_{0,j}\|_{\text{TV}}) \\ &\geq k - \frac{k}{\sqrt{n}} \sqrt{\sum_j \|\mathbb{P}_{1,j} - \mathbb{P}_{0,j}\|_{\text{TV}}^2} \end{aligned}$$


$$\sum_j \|\mathbb{P}_{1,j} - \mathbb{P}_{0,j}\|_{\text{TV}}^2 \leq \frac{\mu^2}{4} m \quad \longrightarrow \quad \mathbb{E} |\widehat{S} \Delta S| \geq k \left( 1 - \frac{\mu}{2} \sqrt{\frac{m}{n}} \right)$$

# Key Ideas in Proof of Lemma

## Pinsker's Inequality

$$\|\mathbb{P} - \mathbb{Q}\|_{\text{TV}} \leq \sqrt{K(\mathbb{P}, \mathbb{Q})/2}$$

$$\begin{aligned} \|\mathbb{P}_{1,j} - \mathbb{P}_{0,j}\|_{\text{TV}}^2 &\leq \frac{\pi_0}{2} K(\mathbb{P}_{0,j}, \mathbb{P}_{1,j}) + \frac{\pi_1}{2} K(\mathbb{P}_{1,j}, \mathbb{P}_{0,j}) \\ &\leq \frac{\mu^2}{4} \sum_i \mathbb{E} a_{i,j}^2 \end{aligned}$$

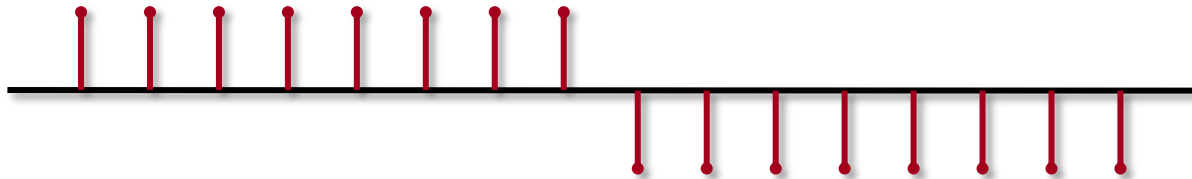

$$\sum_j \|\mathbb{P}_{1,j} - \mathbb{P}_{0,j}\|_{\text{TV}}^2 \leq \frac{\mu^2}{4} \sum_{i,j} \mathbb{E} a_{i,j}^2 = \frac{\mu^2}{4} m$$

# Adaptivity In Practice

Suppose that  $k = 1$  and that  $x_{j^*} = \mu$

Binary Search [Iwen and Tewfik - 2011, Davenport and Arias-Castro - 2012]

- split measurements into  $\log n$  stages
- in each stage, use measurements to decide if the nonzero is in the left or right half of the “active set”
- after subdividing  $\log n$  times, return support



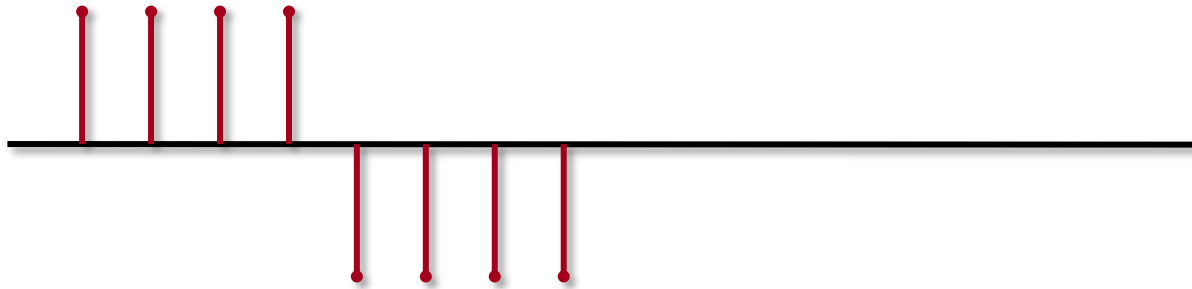


# Adaptivity In Practice

Suppose that  $k = 1$  and that  $x_{j^*} = \mu$

Binary Search [Iwen and Tewfik - 2011, Davenport and Arias-Castro - 2012]

- split measurements into  $\log n$  stages
- in each stage, use measurements to decide if the nonzero is in the left or right half of the “active set”
- after subdividing  $\log n$  times, return support

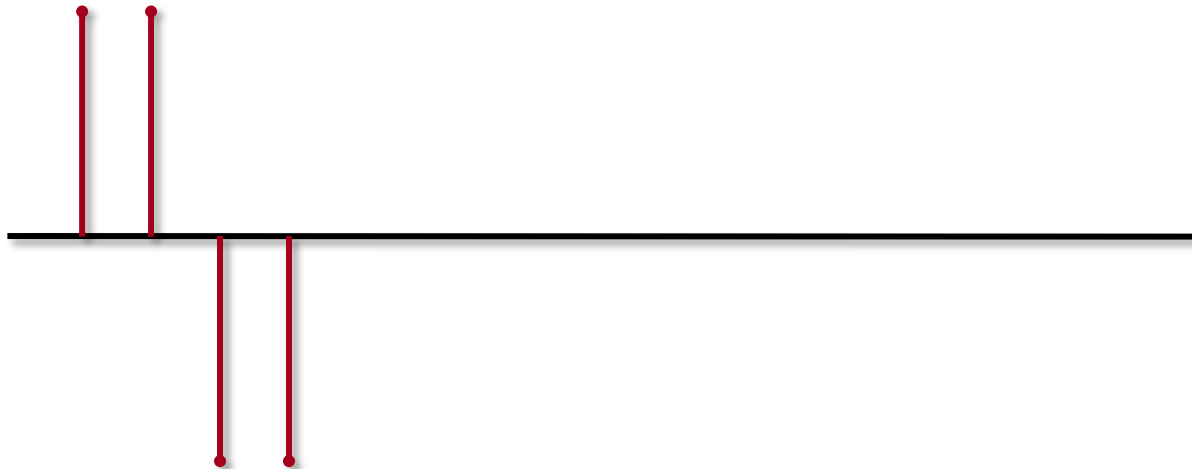


# Adaptivity In Practice

Suppose that  $k = 1$  and that  $x_{j^*} = \mu$

Binary Search [Iwen and Tewfik - 2011, Davenport and Arias-Castro - 2012]

- split measurements into  $\log n$  stages
- in each stage, use measurements to decide if the nonzero is in the left or right half of the “active set”
- after subdividing  $\log n$  times, return support

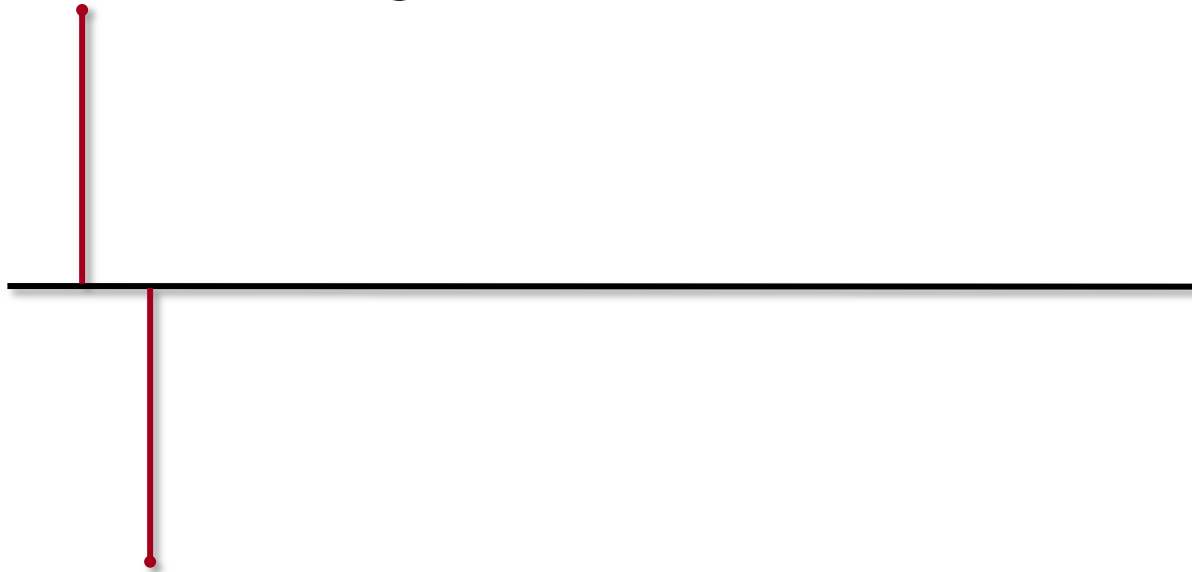


# Adaptivity In Practice

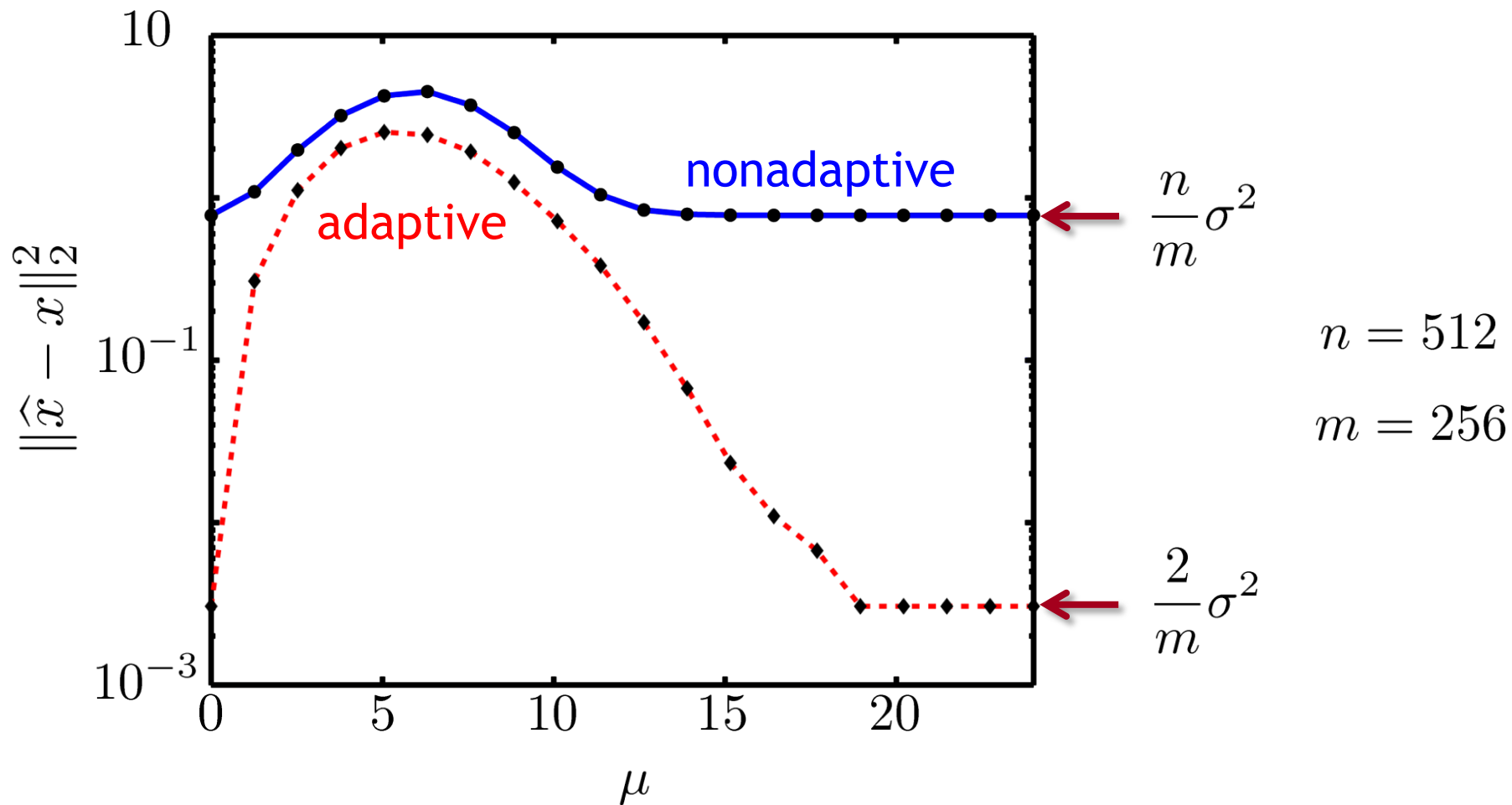
Suppose that  $k = 1$  and that  $x_{j^*} = \mu$

Binary Search [Iwen and Tewfik - 2011, Davenport and Arias-Castro - 2012]

- split measurements into  $\log n$  stages
- in each stage, use measurements to decide if the nonzero is in the left or right half of the “active set”
- after subdividing  $\log n$  times, return support



# Experimental Results



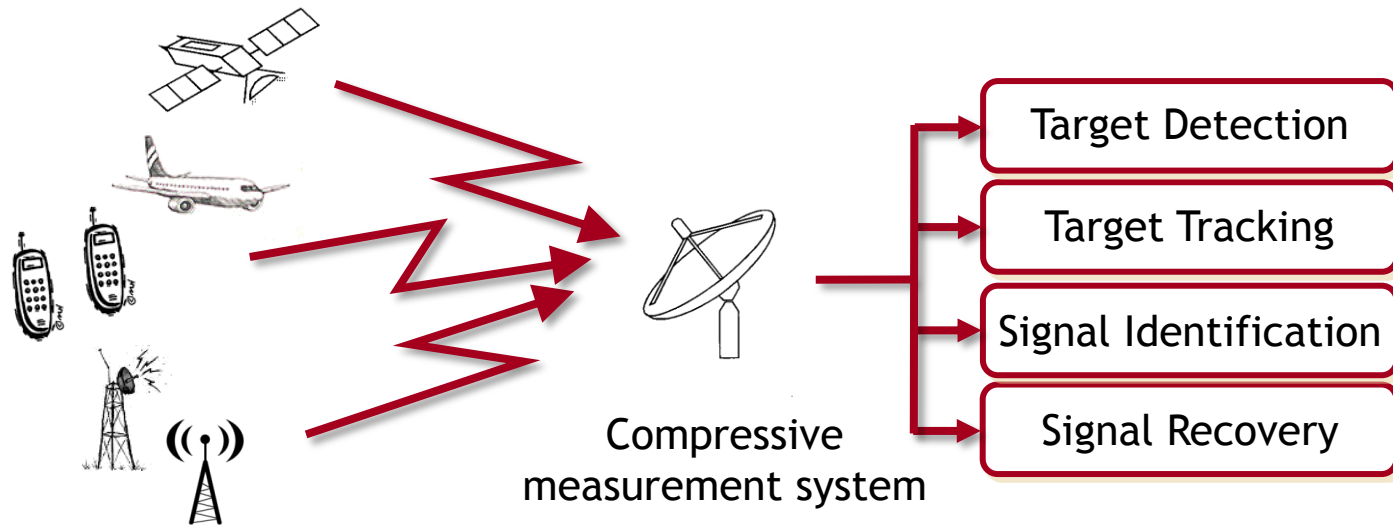
**Looking Forward**

# Adaptivity in Practice

- No method can succeed when  $\frac{\mu}{\sigma} \approx \sqrt{\frac{n}{m}}$ , but the binary search approach succeeds as long as  $\frac{\mu}{\sigma} \geq C \sqrt{\frac{n}{m} \log \log n}$   
[Davenport and Arias-Castro - 2012]
- Practical algorithms that work well for all values of  $\mu$
- New theory for restricted adaptive measurements
  - single-pixel camera: 0/1 measurements
  - magnetic resonance imaging (MRI): Fourier measurements
  - analog-to-digital converters: linear filter measurements
- New sensors and architectures that can actually acquire adaptive measurements

# Beyond Recovery

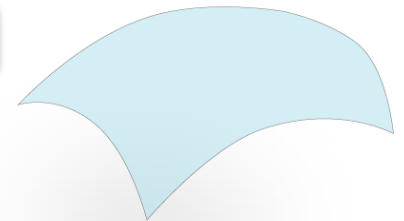
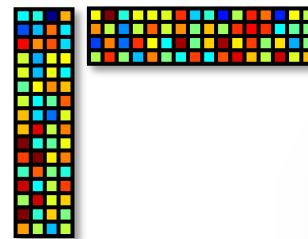
When and how can we directly solve inference problems directly from measurements?



- “Compressive signal processing”
- Links with machine learning
  - Johnson-Lindenstrauss lemma and “compressive learning”
  - quantized compressive sensing and sparse learning

# Beyond Sparsity

- Learned dictionaries, structured sparsity, models for continuous-time signals
- Multi-signal models
  - e.g., sensor networks/arrays, multi-modal data, ...
- Low-rank matrix models
  - matrix completion
- Manifold/parametric models



## Acquisition

- how to design  $A$
- practical devices
- adaptivity

## Recovery

- practical algorithms
- robustness
- quantization

## Inference

- classification
- estimation
- learning



# Acknowledgements

- Richard Baraniuk (Rice)
- Emmanuel Candès (Stanford)
- Ery Arias-Castro (UC San Diego)
- Ronald DeVore (Texas A&M)
- Marco Duarte (UMass Amherst)
- Kevin Kelly (Rice)
- Michael Wakin (Colorado School of Mines)



# More Information

<http://stat.stanford.edu/~markad>

markad@stanford.edu