

# **Modern Quantization Strategies for Compressive Sensing and Acquisition Systems**

---

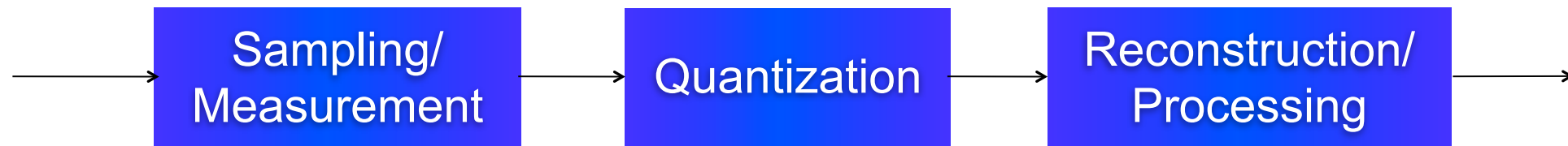


Petros Boufounos  
*petros@boufounos.com*



Laurent Jacques  
*laurent.jacques@uclouvain.be*

# Signal Acquisition Pipeline



- Typically linear (at least today's discussion)
- Can be designed to be invertible (e.g. Nyquist theorem, compressive sensing, signals with finite rate of innovation, etc...)

- Classical: Linear reconstruction
- Modern: Non-linear, heavy computation (e.g., compressive sensing, finite rate of innovation)

- Highly non-linear
- Not invertible  $\Leftrightarrow$  loss of information
- Design to minimize loss

**Today**  
Quantizer Design  
Interaction with measurement system  
Optimal reconstruction



# Today's Topics

---

1. Modern Scalar Quantization
2. Compressive Sensing Overview
3. Compressive Sensing and Quantization
4. 1-bit Compressive Sensing
5. Locality Sensitive Hashing and Universal Quantization

# Today's Topics

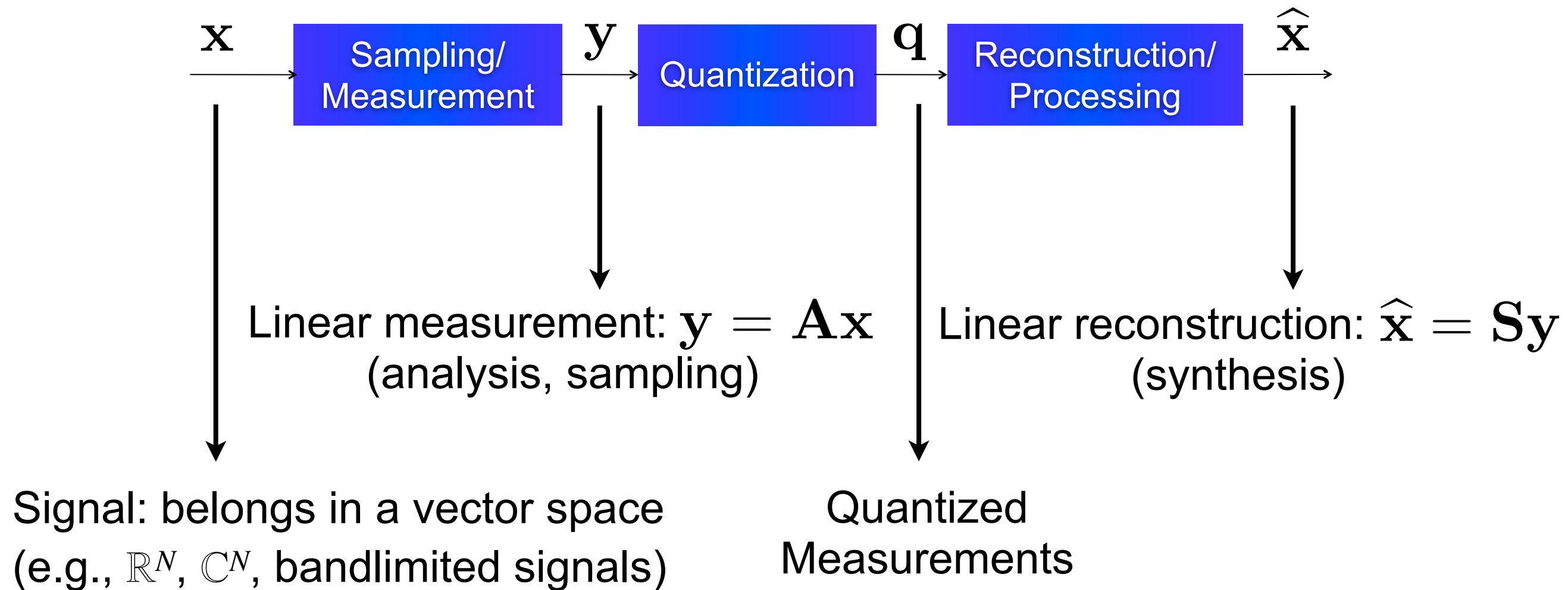
---

1. Modern Scalar Quantization
2. Compressive Sensing Overview
3. Compressive Sensing and Quantization
4. 1-bit Compressive Sensing
5. Locality Sensitive Hashing and Universal Quantization

---

# **SIGNAL REPRESENTATION**

# Linear Measurement and Reconstruction Model



$\mathbf{A}$ : Basis expansion (critically sampled) or frame expansion (oversampled)

In absence of quantization:  $\mathbf{S} = \mathbf{A}^{-1}$  or  $\mathbf{S} = \mathbf{A}^\dagger$

Biorthogonal (dual) basis      Dual frame

# Basis Expansions

Analysis (Measurement)

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

$$y_i = \langle \mathbf{a}_i, \mathbf{x} \rangle$$

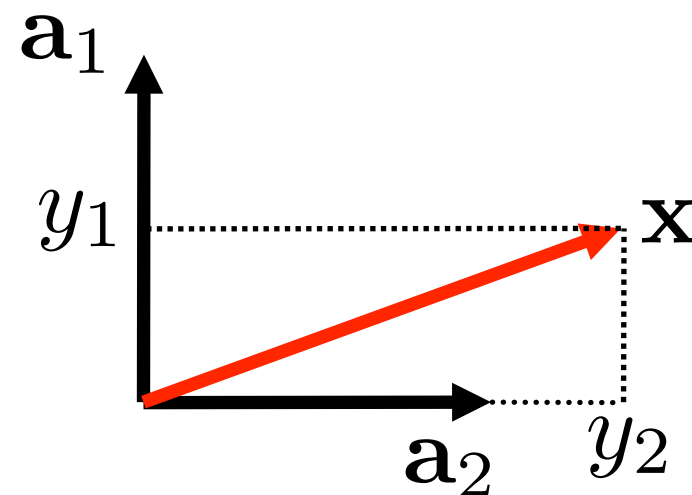
Synthesis (Reconstruction)

$$\hat{\mathbf{x}} = \mathbf{S}\mathbf{y}$$

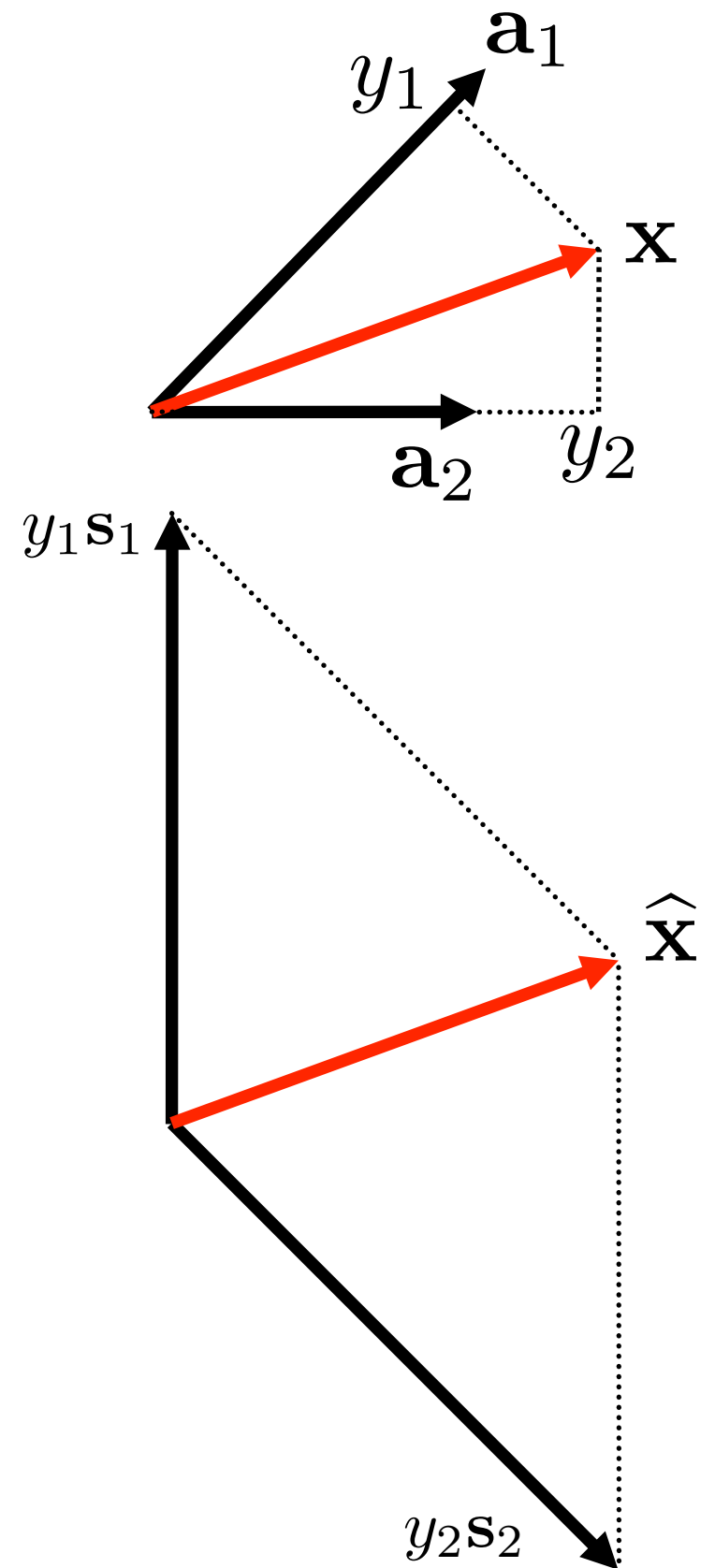
$$\hat{\mathbf{x}} = \sum_i y_i \mathbf{s}_i$$

$$\mathbf{S} = \mathbf{A}^{-1}$$

Orthonormal



Oblique

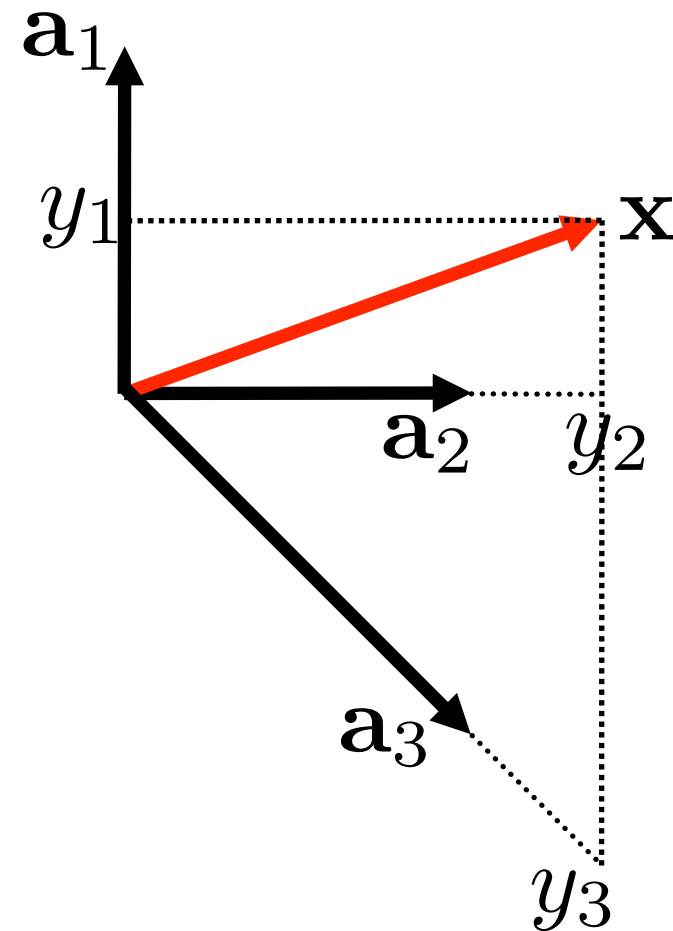


# Frame Representations and Oversampling

Analysis (Measurement)

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

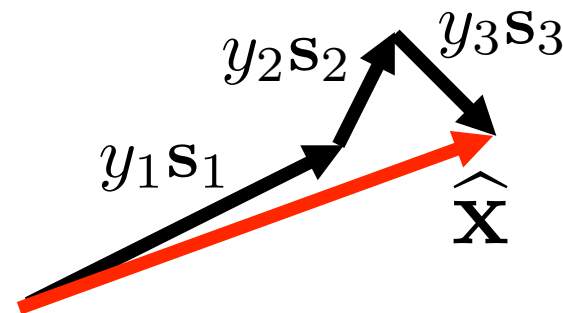
$$y_i = \langle \mathbf{a}_i, \mathbf{x} \rangle$$



Synthesis (Reconstruction)

$$\hat{\mathbf{x}} = \mathbf{S}\mathbf{y}$$

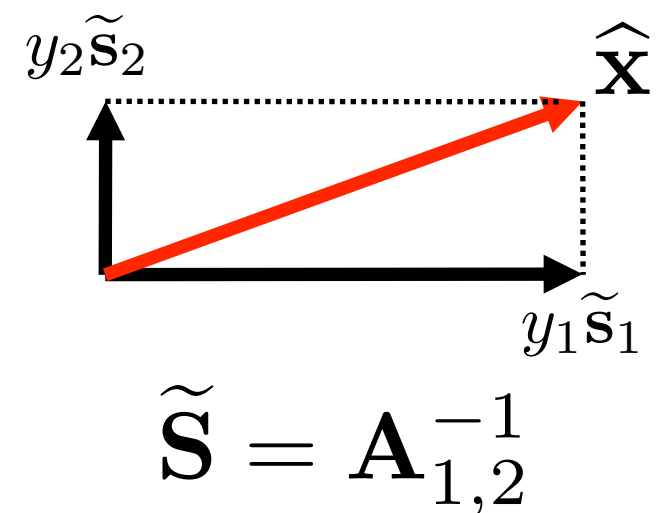
$$\hat{\mathbf{x}} = \sum_i y_i \mathbf{s}_i$$



$$\mathbf{S} = \mathbf{A}^{-1}, \text{ i.e. } \mathbf{S}\mathbf{A} = \mathbf{I}$$

$$\mathbf{S} = \mathbf{A}^\dagger$$

Canonical  
Dual Frame



$$\tilde{\mathbf{S}} = \mathbf{A}_{1,2}^{-1}$$

# Examples of Frames and Frame Expansions

Matrix Operations in  $\mathbb{R}^{M \times N}$

Analysis  
(Measurement)

$$\begin{bmatrix} -\mathbf{a}_1- \\ \vdots \\ -\mathbf{a}_M- \end{bmatrix} \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} \Leftrightarrow y_i = \langle \mathbf{a}_i, \mathbf{x} \rangle$$

Redundancy  
 $r=M/N$

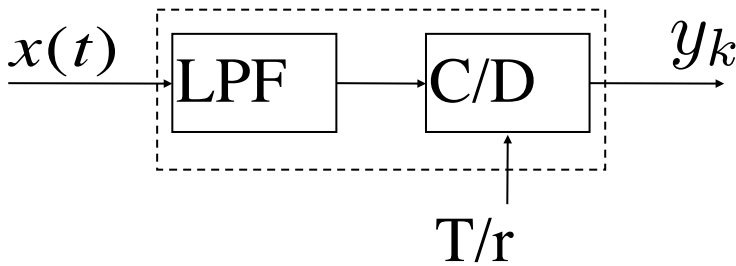
Synthesis  
(Reconstruction)

$$\begin{bmatrix} | & & | \\ \mathbf{s}_1 & \cdots & \mathbf{s}_M \\ | & & | \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} \Leftrightarrow \mathbf{x} = \sum_i y_i \mathbf{s}_i$$

$r$ -times Oversampling:

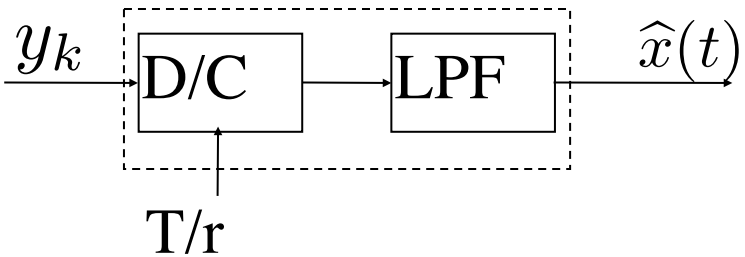
Analysis  
(Measurement)

$$y_i = \int_{-\infty}^{+\infty} x(t) \frac{1}{rT} \text{sinc} \left( \frac{r}{T} t - i \right) dt \Leftrightarrow y_i = \langle \mathbf{a}_i, \mathbf{x} \rangle$$

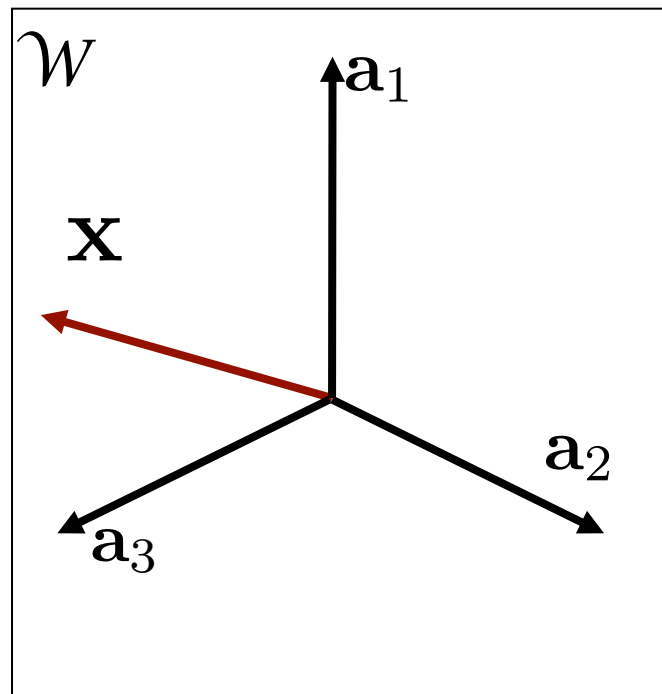


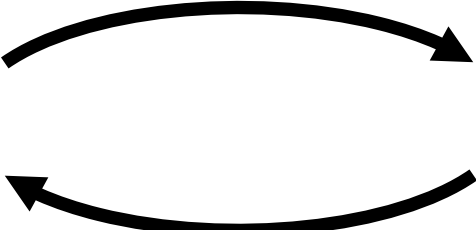
Synthesis  
(Reconstruction)

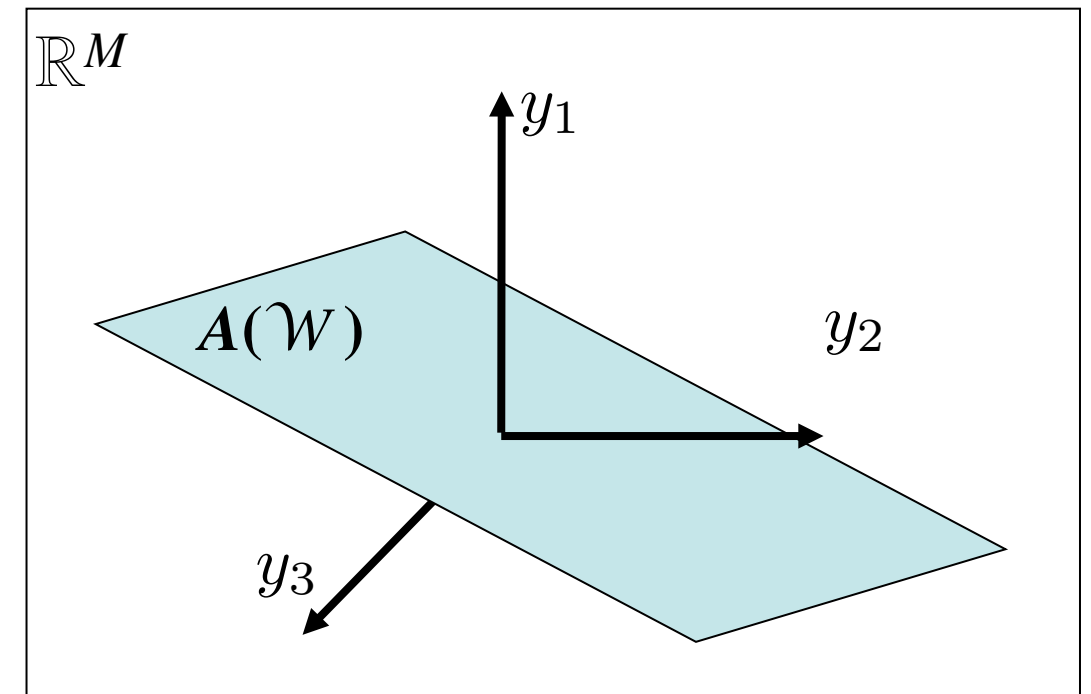
$$x(t) = \sum_i y_i \text{sinc} \left( \frac{r}{T} t - i \right) \Leftrightarrow \mathbf{x} = \sum_i y_i \mathbf{s}_i$$



# Frame Expansion/Oversampling: Subspace Mapping



$$\begin{aligned} \mathbf{y} &= \mathbf{A}\mathbf{x} \\ y_i &= \langle \mathbf{a}_i, \mathbf{x} \rangle \end{aligned}$$

$$\begin{aligned} \hat{\mathbf{x}} &= \mathbf{S}\mathbf{y} \\ \hat{\mathbf{x}} &= \sum_i y_i \mathbf{s}_i \end{aligned}$$



Signal Space  $\mathcal{W}$

Frame:  $\{\mathbf{a}_i, i = 1, \dots, M \mid \mathbf{a} \in \mathcal{W}\}$

$\dim(\mathcal{W}) = N$

Coefficient/Measurement Space  $\mathbb{R}^M$

Image is  $N$ -dimensional

$\dim(\mathbf{A}(\mathcal{W})) \leq \text{rank}(\mathbf{A}) \leq N < M$

Frames provide **redundancy**

Mechanism: nullspace of synthesis operator.

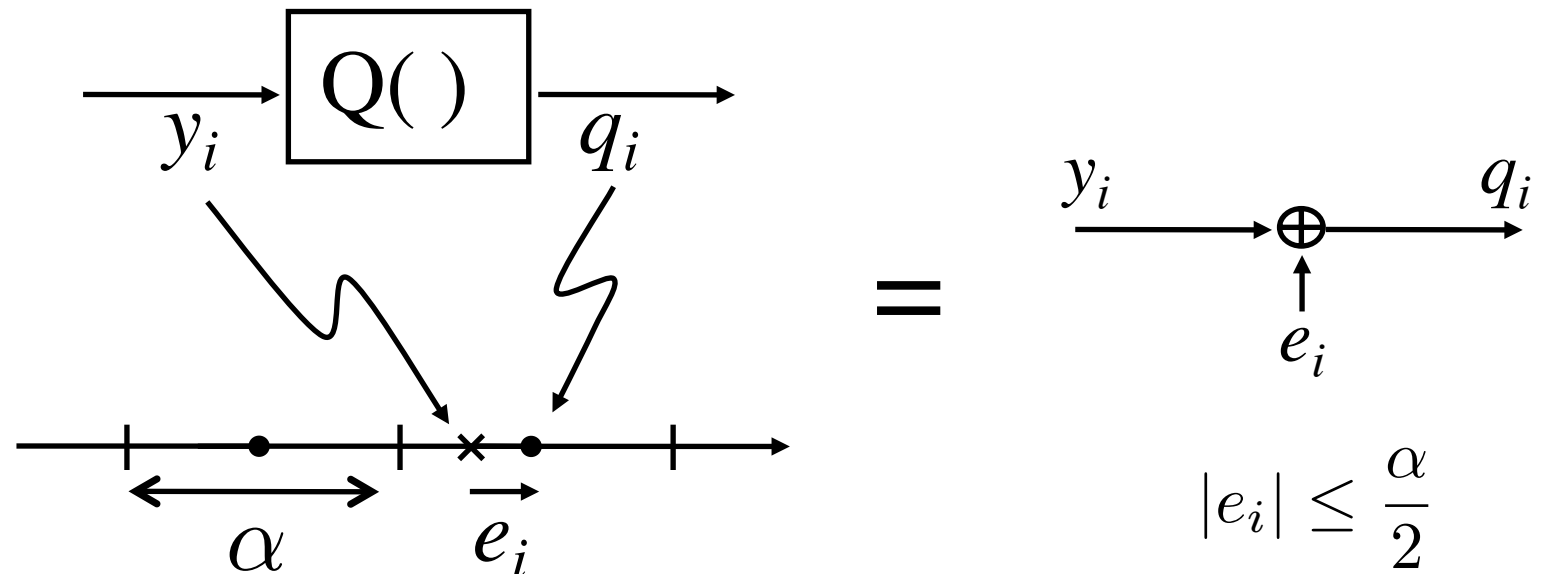
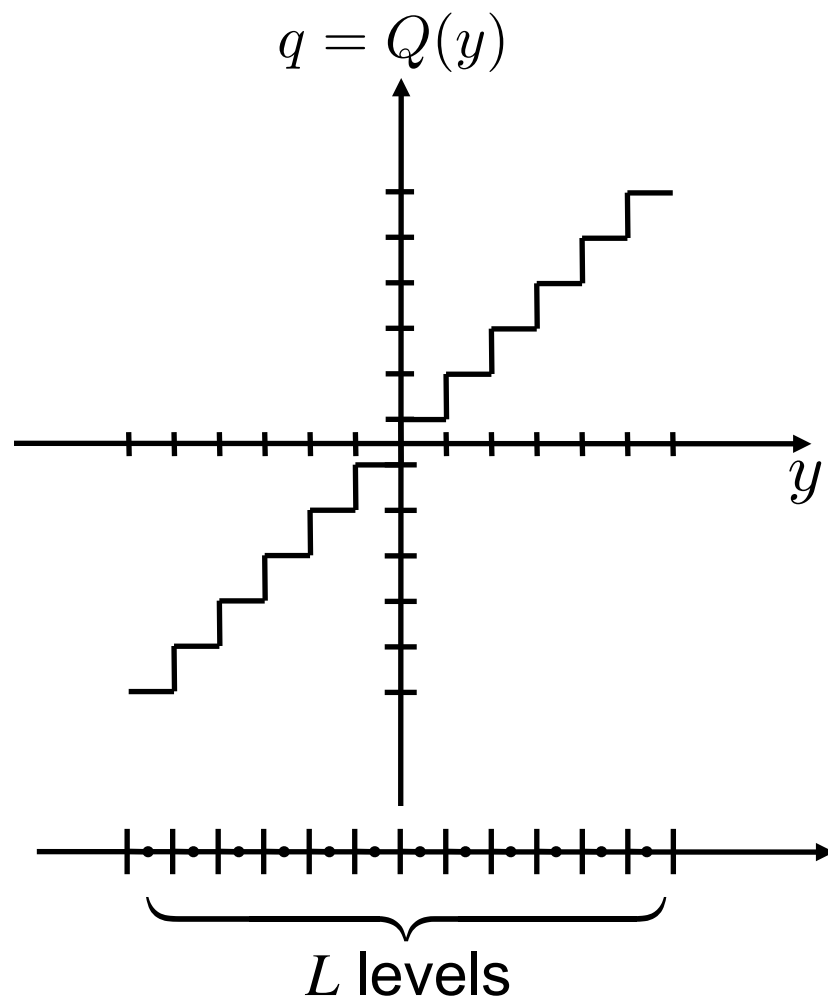
**Redundancy** can be exploited for quantization **robustness**



---

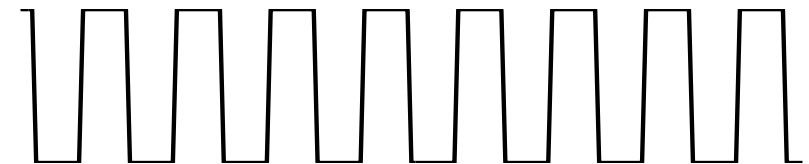
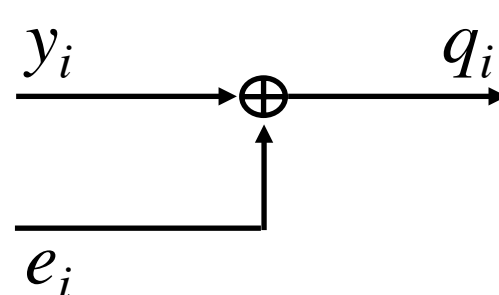
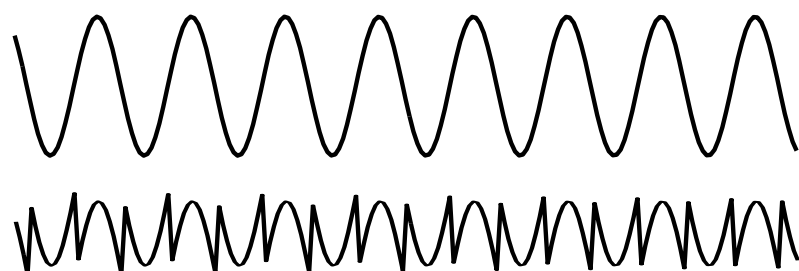
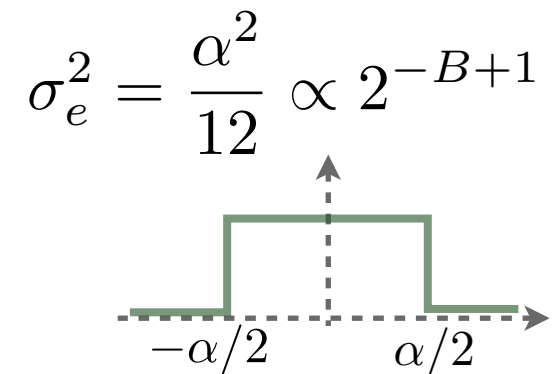
# **SCALAR QUANTIZATION**

# Scalar Quantization

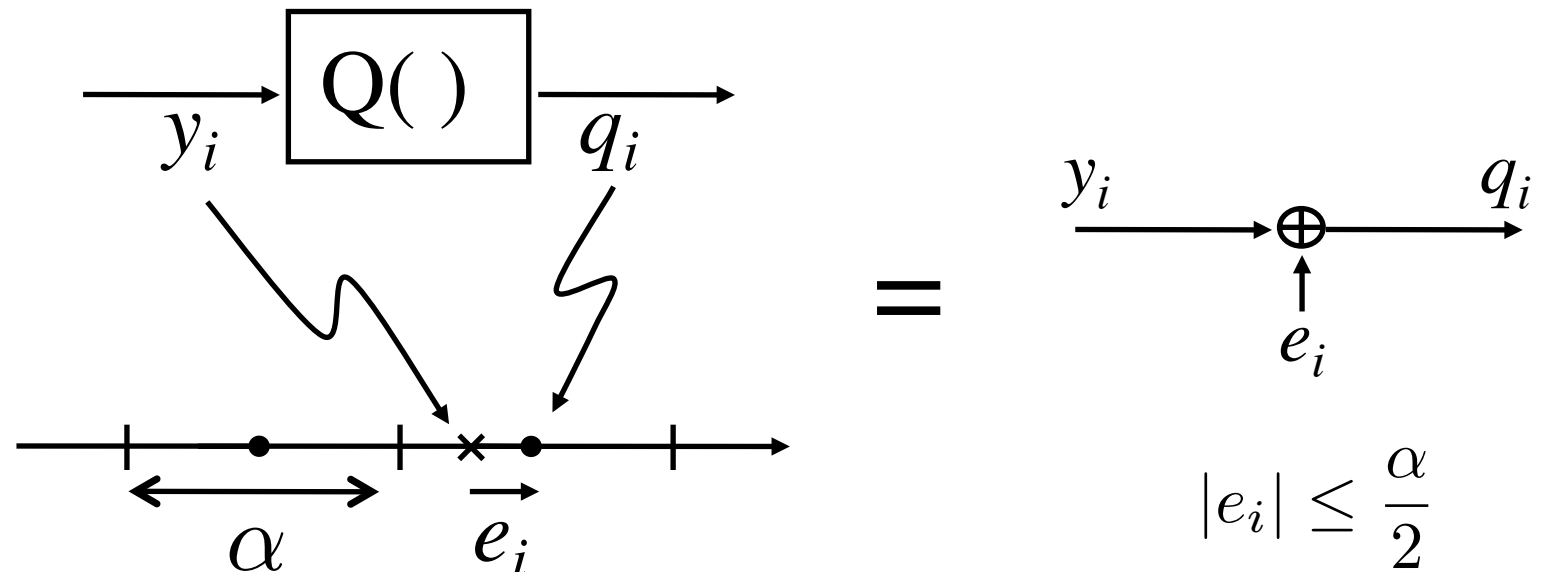
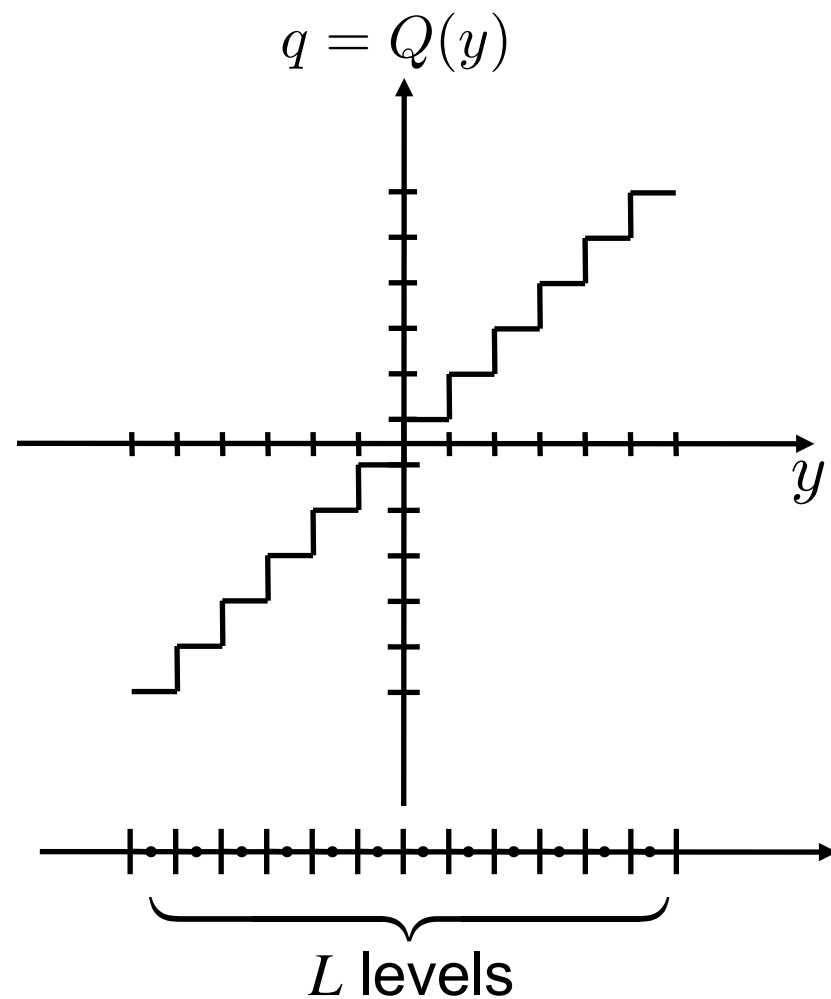


$L$  level quantizer:  $B \sim \log_2(L)$  bits per coefficient

Additive noise model:  $e_i$  uncorrelated, uniform in  $\pm \frac{\alpha}{2}$

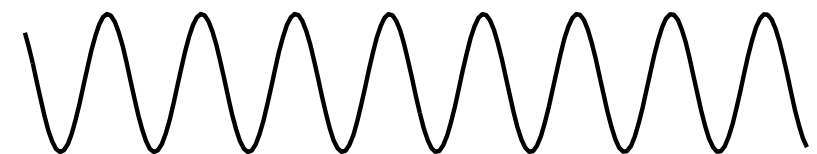
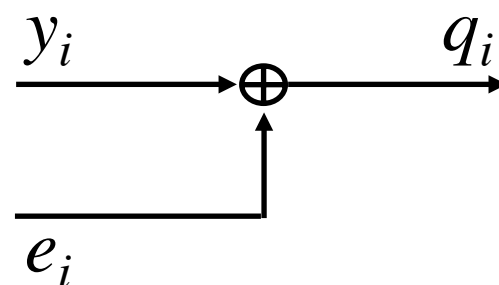
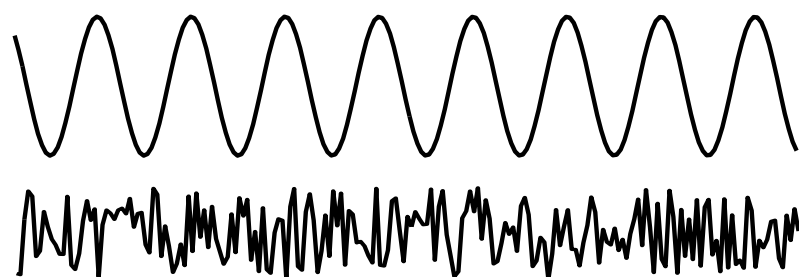
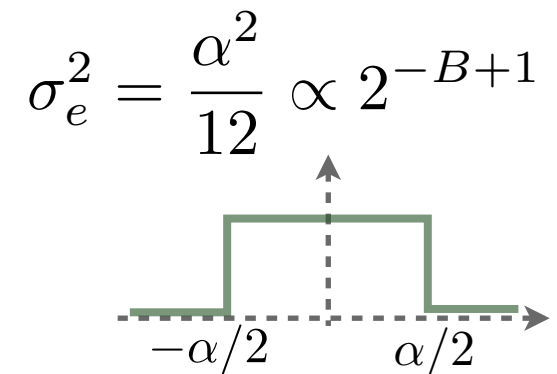


# Scalar Quantization

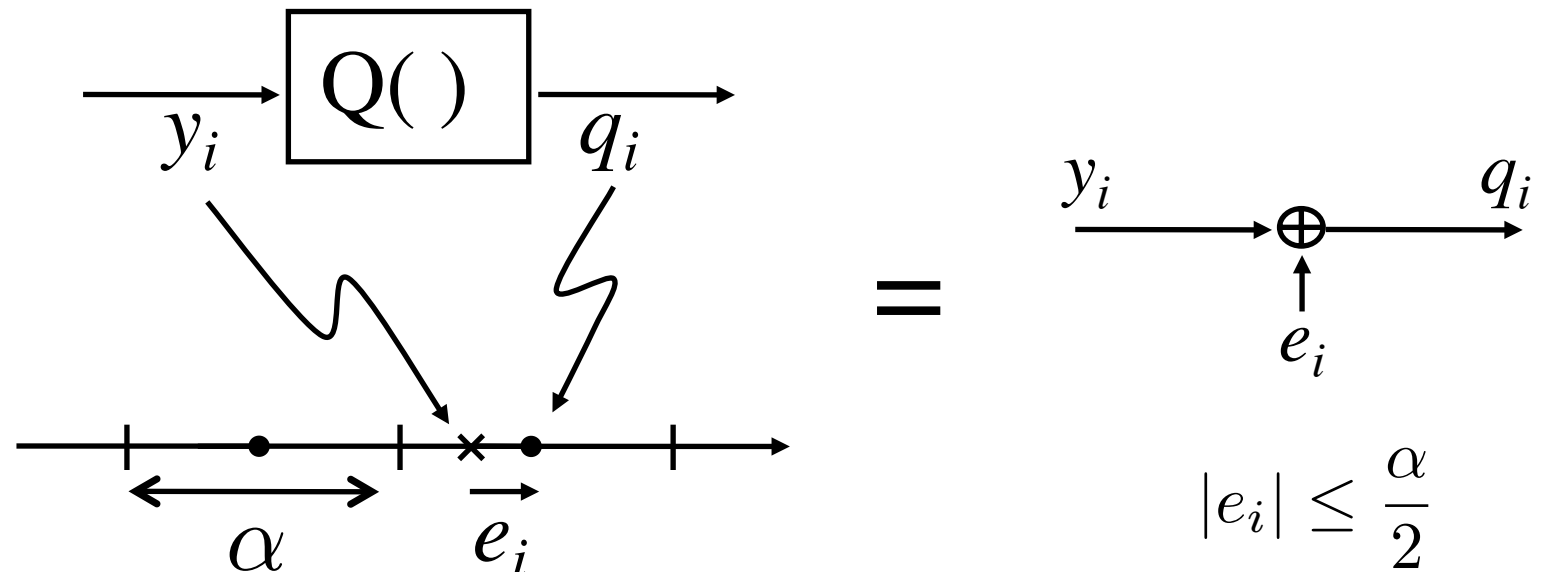
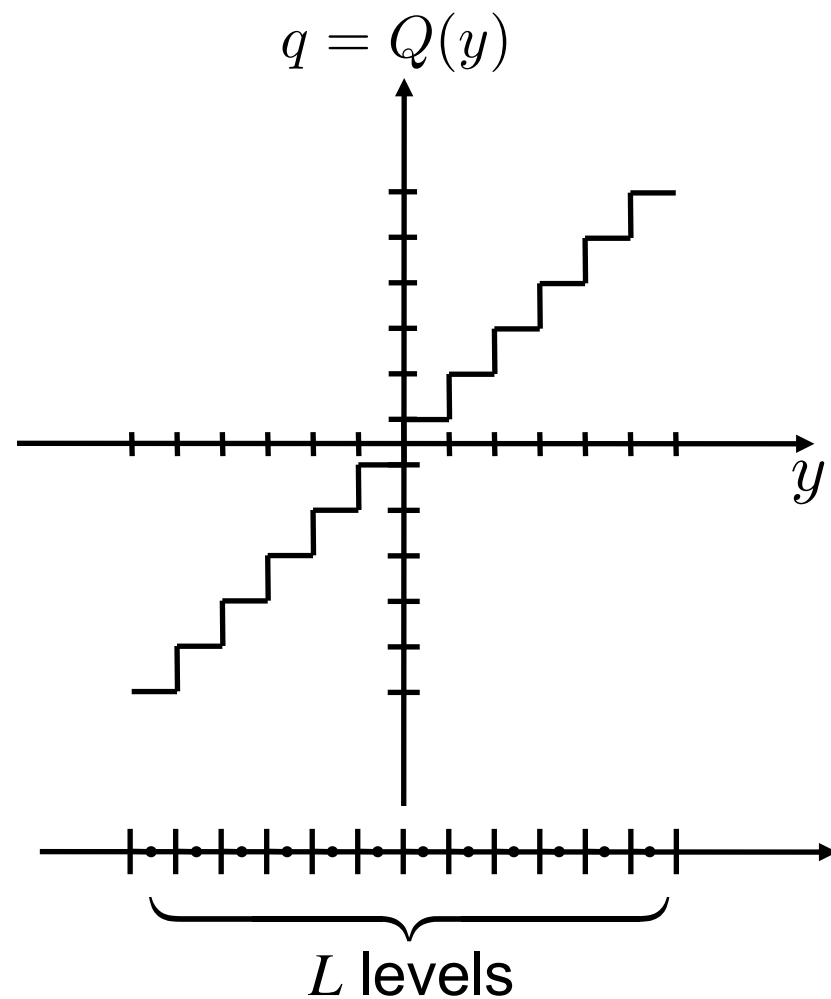


$L$  level quantizer:  $B \sim \log_2(L)$  bits per coefficient

Additive noise model:  $e_i$  uncorrelated, uniform in  $\pm \frac{\alpha}{2}$

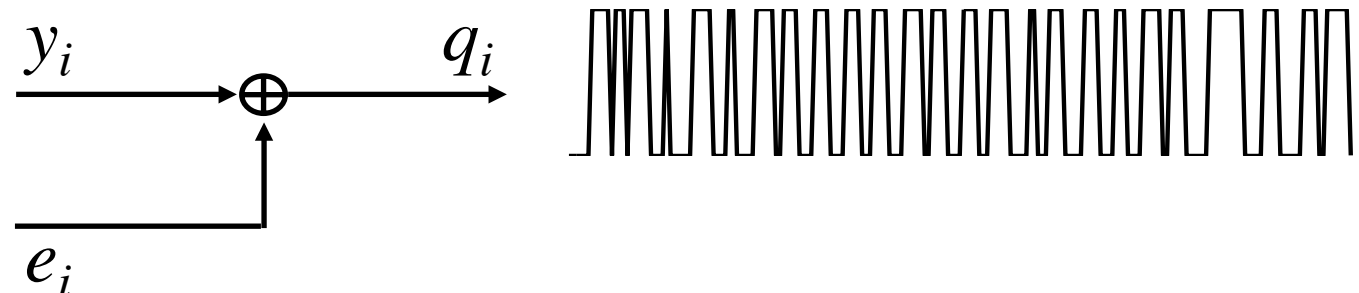
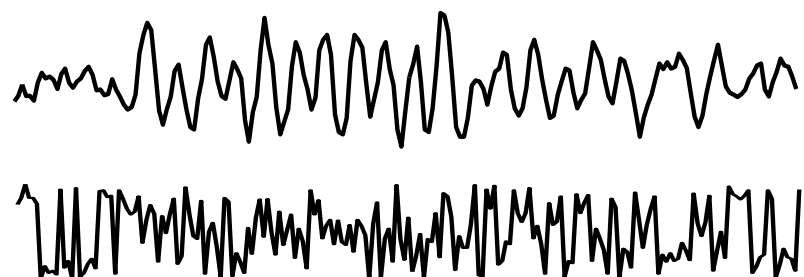
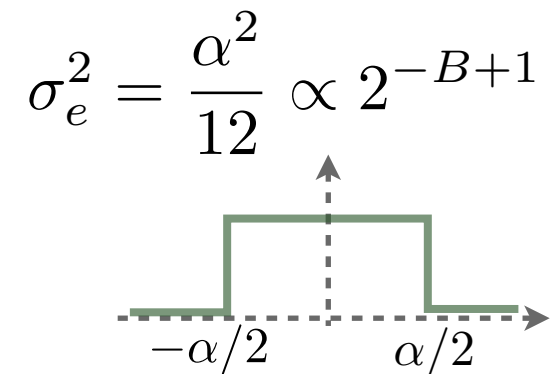


# Scalar Quantization

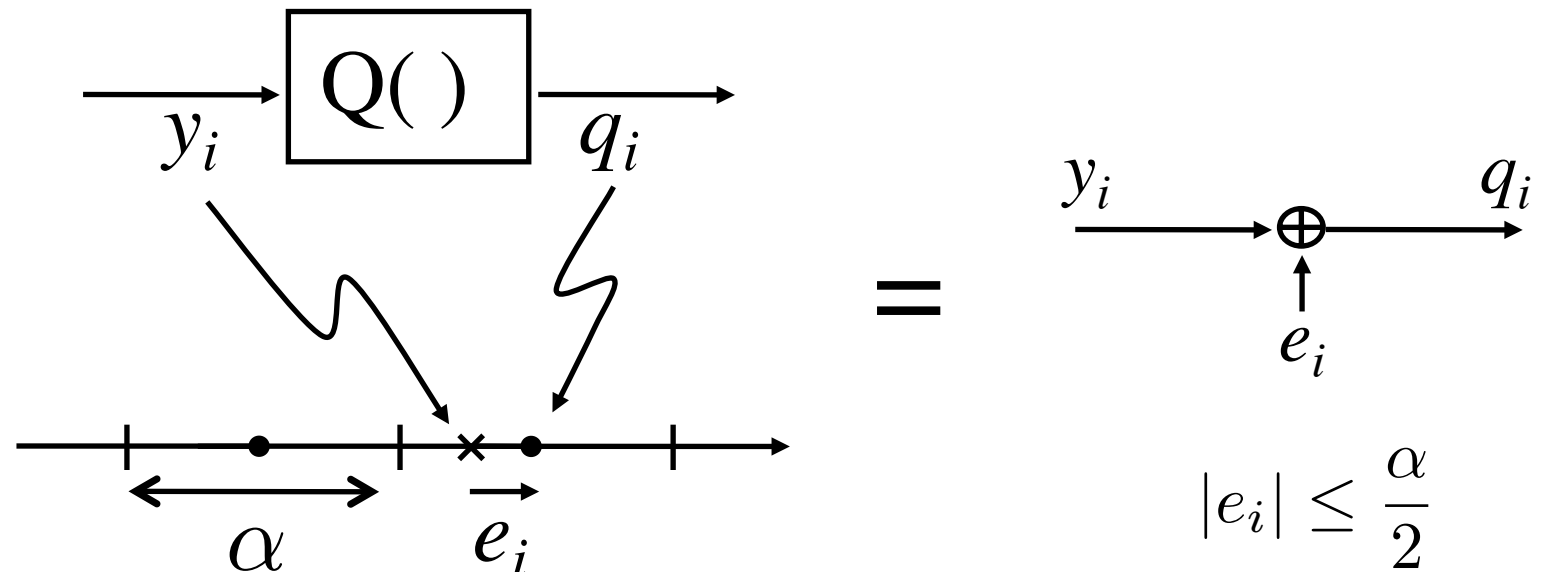
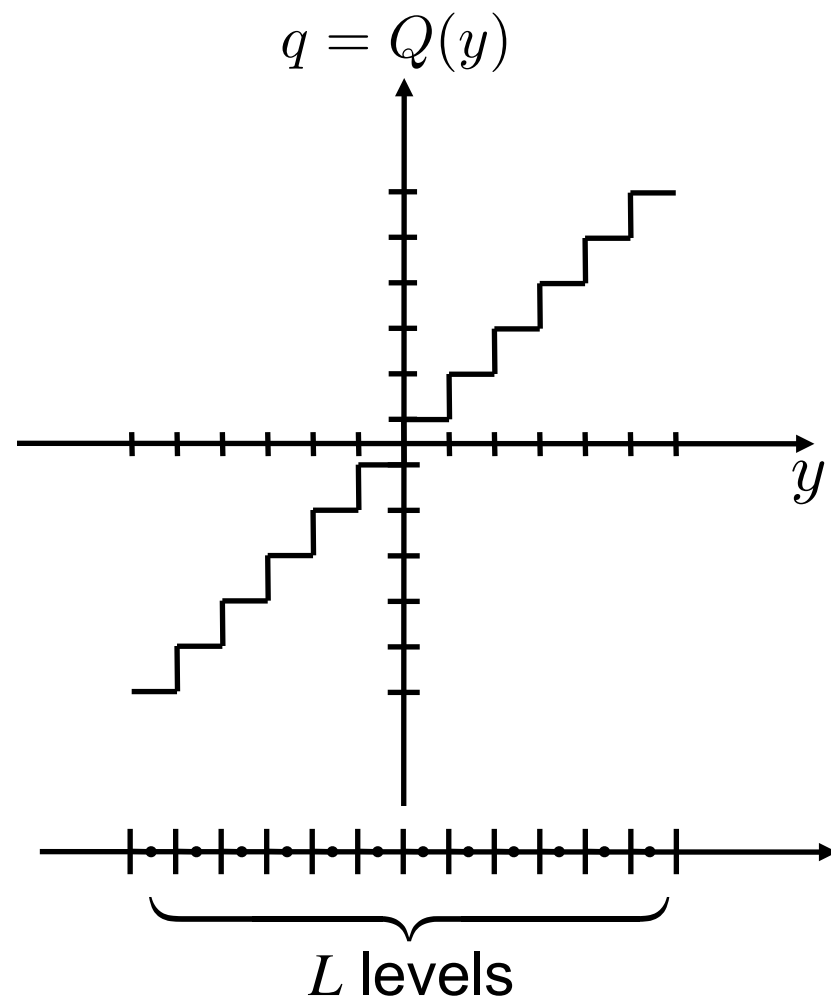


$L$  level quantizer:  $B \sim \log_2(L)$  bits per coefficient

Additive noise model:  $e_i$  uncorrelated, uniform in  $\pm \frac{\alpha}{2}$

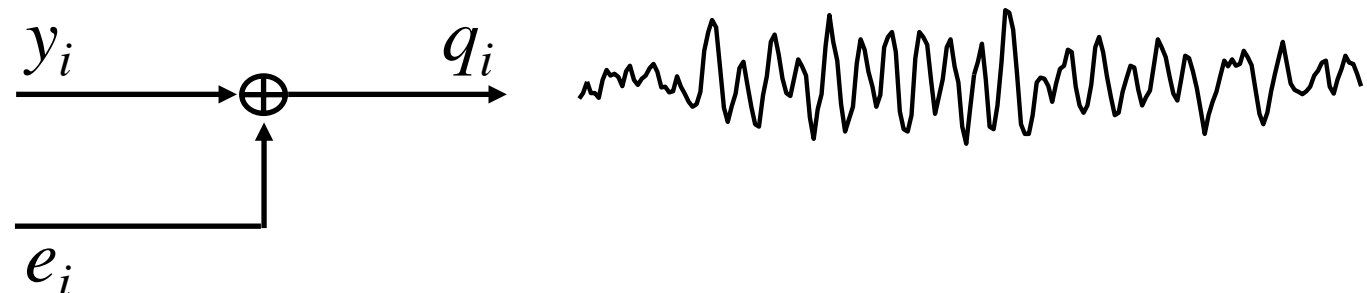
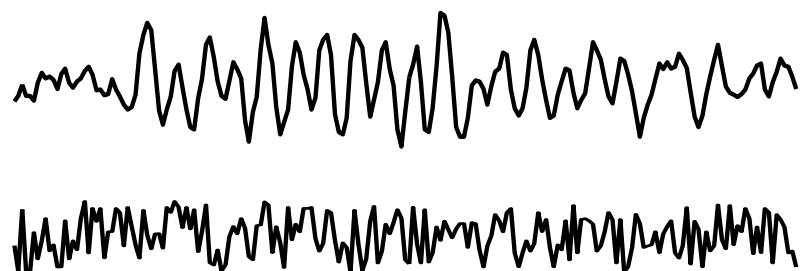
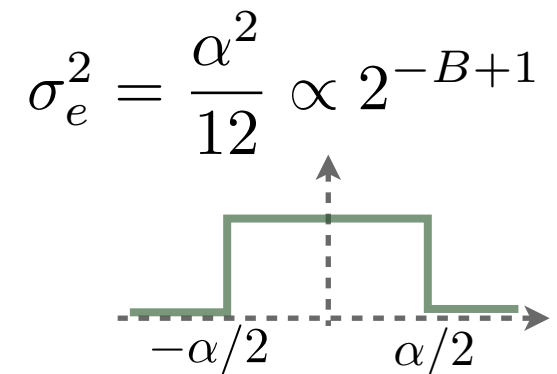


# Scalar Quantization

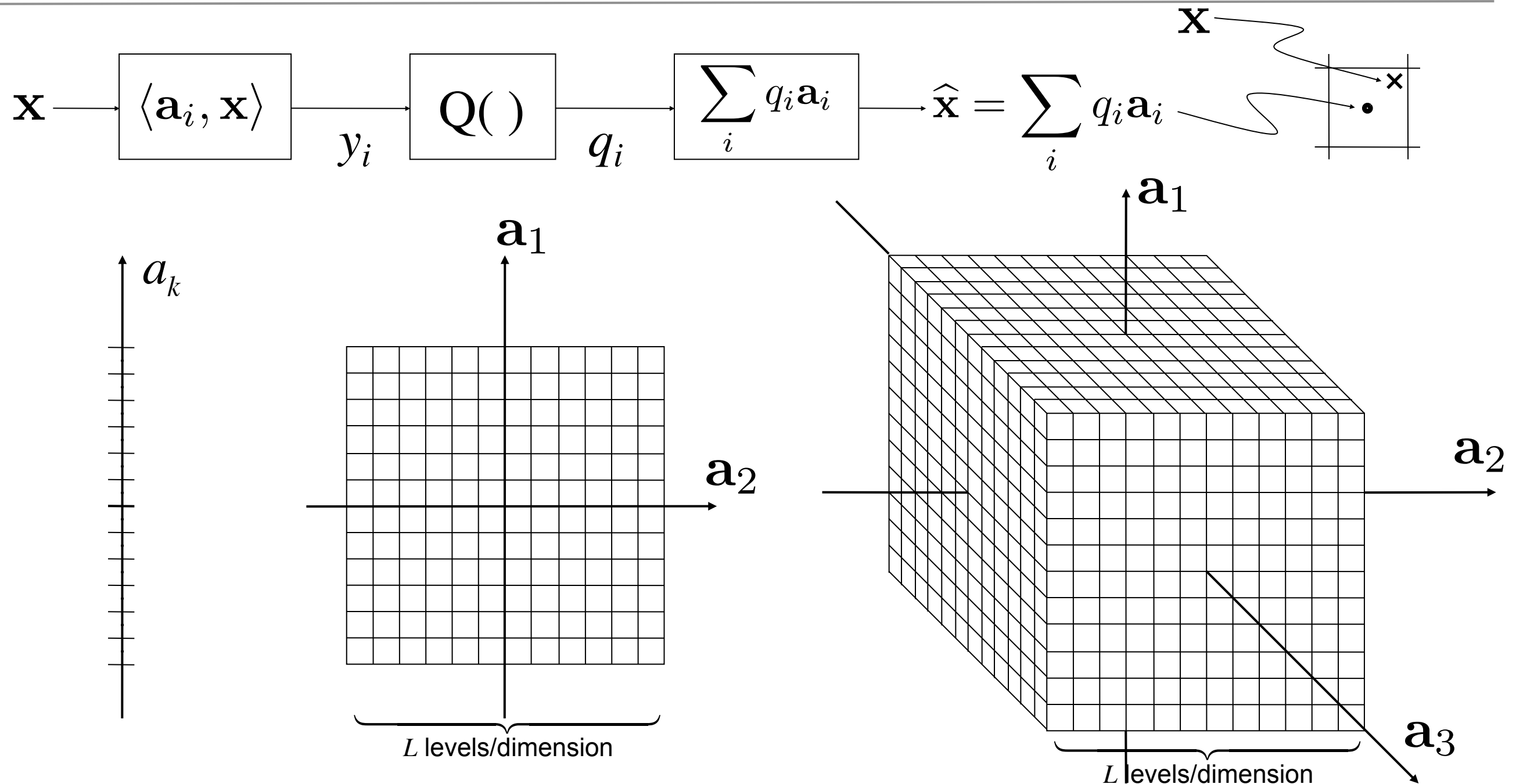


$L$  level quantizer:  $B \sim \log_2(L)$  bits per coefficient

Additive noise model:  $e_i$  uncorrelated, uniform in  $\pm \frac{\alpha}{2}$



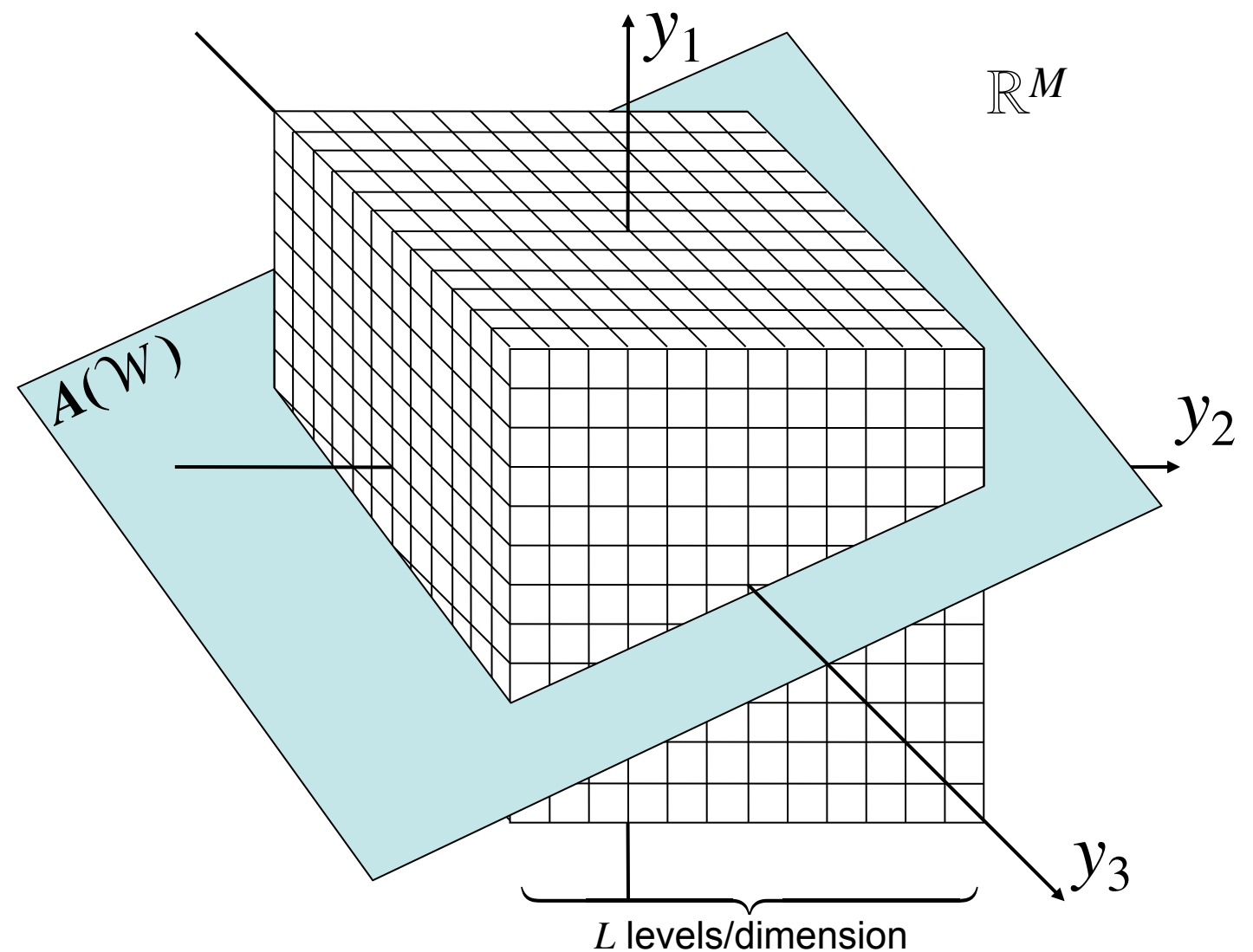
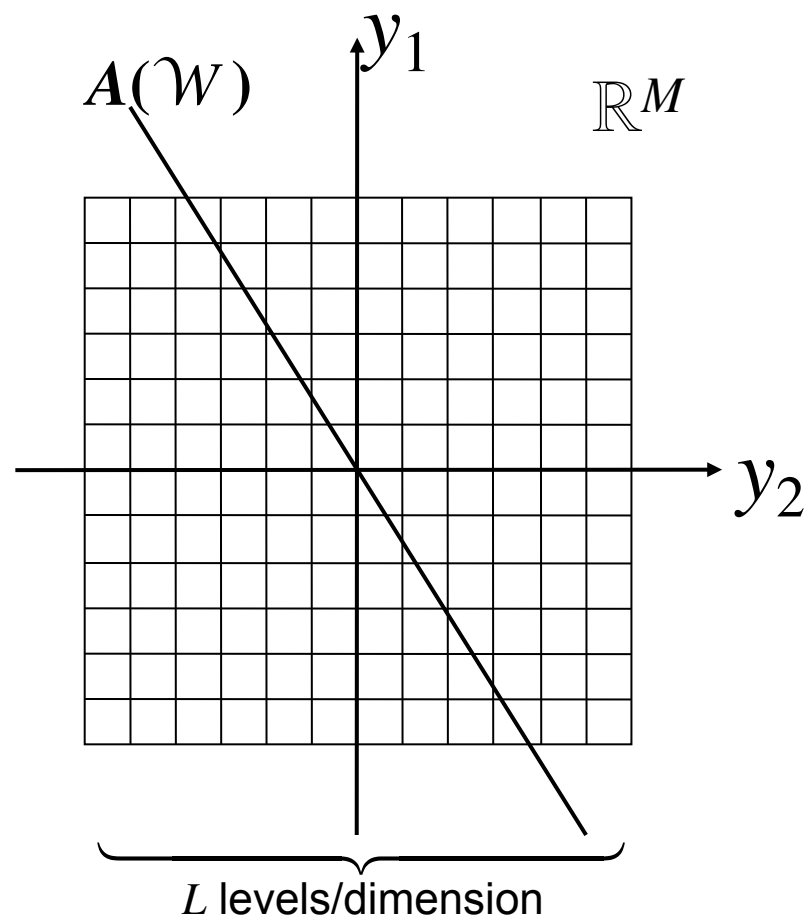
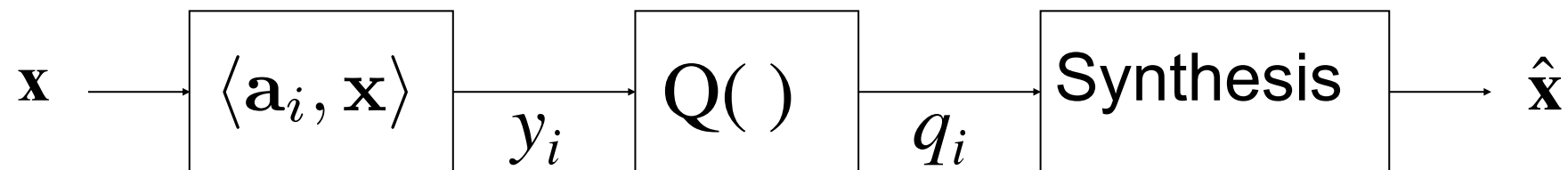
# Quantization of Orthonormal Basis Expansions



$B = \log_2 L$  bits per coefficient  
 $M$  expansion coefficients  $\Rightarrow R = MB = M \log_2 L$  bits used (rate)

**Total Error**  $\varepsilon = O(c^{-R})$

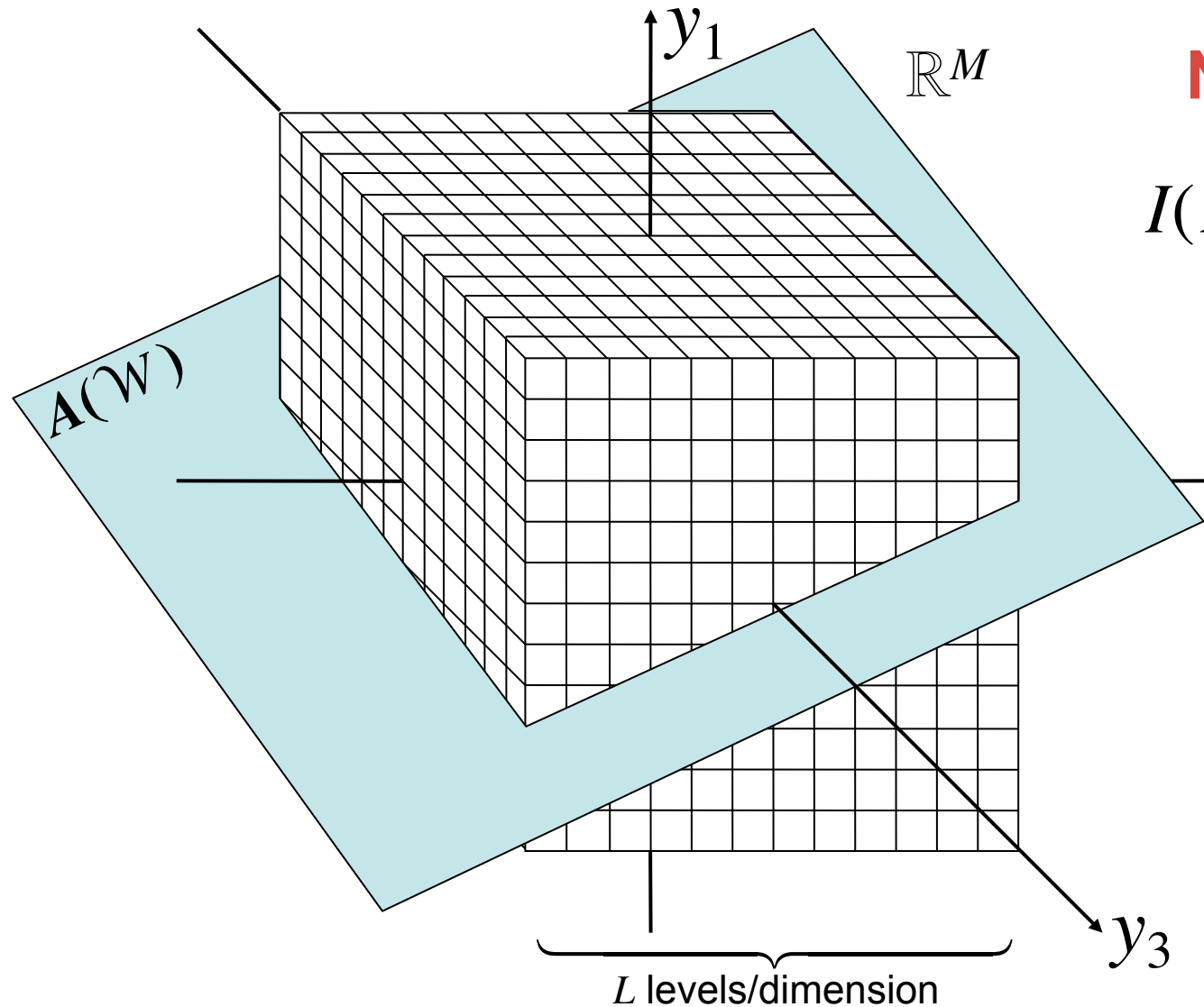
# Quantization of Frame Representations



**Very few quantization cells are intersected!**

**Oversampling** provides **robustness**,  
but also introduces **inefficiency**

# Bounds on Scalar Quantization



**Number of cells intersected  $I(M,N,L)$ :**

$$I(M,N,L) \leq (2L)^N \binom{M}{N} \leq \left( \frac{2LMe}{N} \right)^N = (2Lre)^N$$

**Bit-use efficiency:**

$$\frac{\text{bits necessary}}{\text{bits used}} \leq \frac{\log_2(I(M,N,L))}{M \log_2 L} \leq \frac{\log_2(2Lre)}{r \log_2 L},$$

**Quantization Error Reduction Rate:**

$$\varepsilon^2 = \Omega(r^{-2})$$

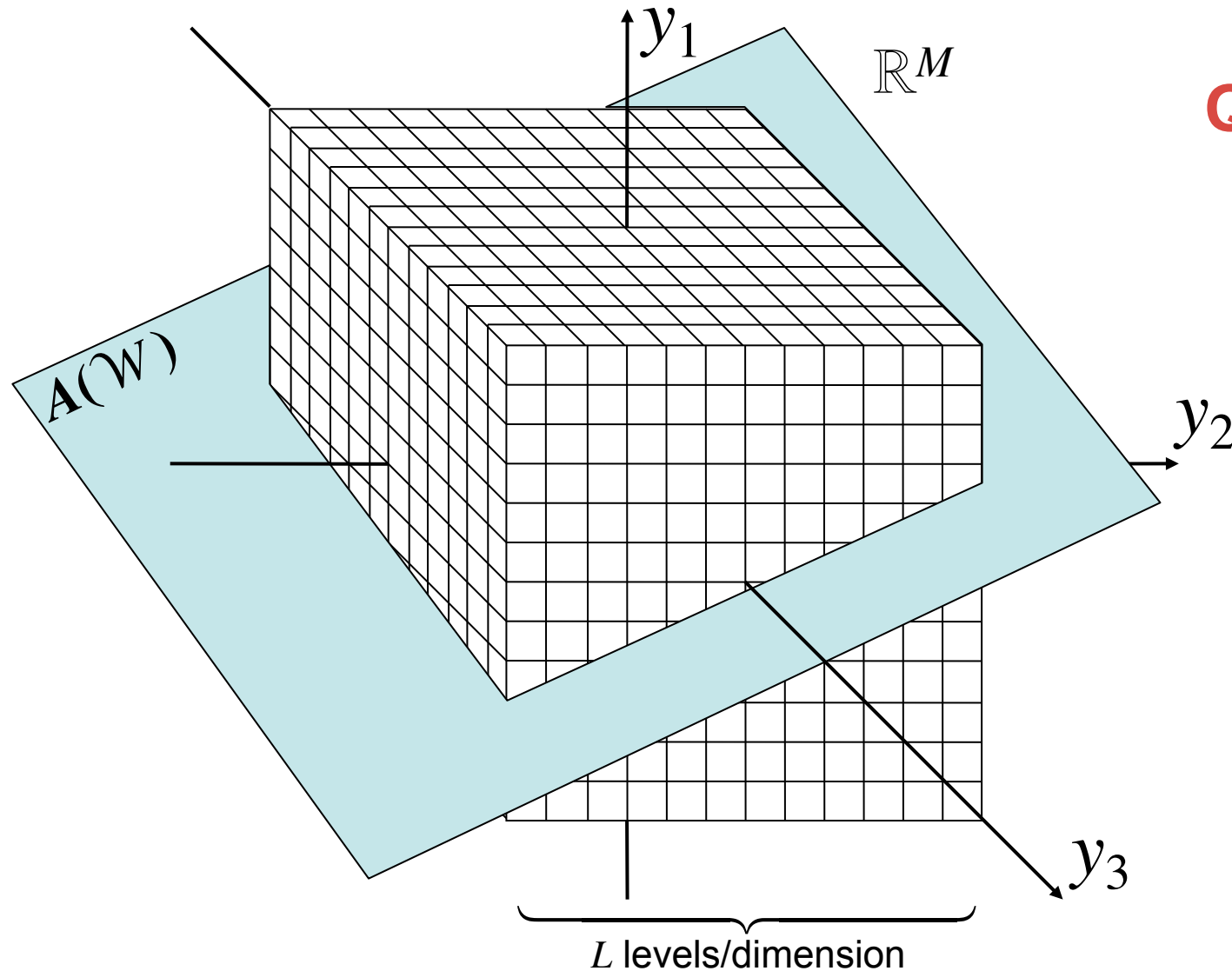
**Total error  $\varepsilon = \Omega(1/R)$   
(vs.  $\varepsilon = O(c^{-R})$  for basis expansions)**

**Oversampled scalar quantization is inefficient!**

- Thao N. T. and Vetterli M., "Lower bound on the mean-squared error in oversampled quantization of periodic signals using vector quantization analysis," *IEEE Trans. Info. Theory*, vol. 42, no. 2, pp. 469–479, Mar. 1996.
- Boufounos P. T., "Quantization and erasures in frame representations," *MIT D.Sc. Thesis*, Cambridge, MA, January 2006.



# Bounds on Scalar Quantization



Quantization Error Reduction Rate:

$$\varepsilon^2 = \Omega(r^{-2})$$

But: Can we achieve it?

Linear reconstruction

$$\mathbf{q} = Q(\mathbf{A}\mathbf{x})$$

$$\Rightarrow \hat{\mathbf{x}} = \mathbf{A}^\dagger \mathbf{q}$$

$$\Rightarrow \varepsilon^2 = \Omega(r^{-1})$$

**Solution: “Consistent reconstruction”**

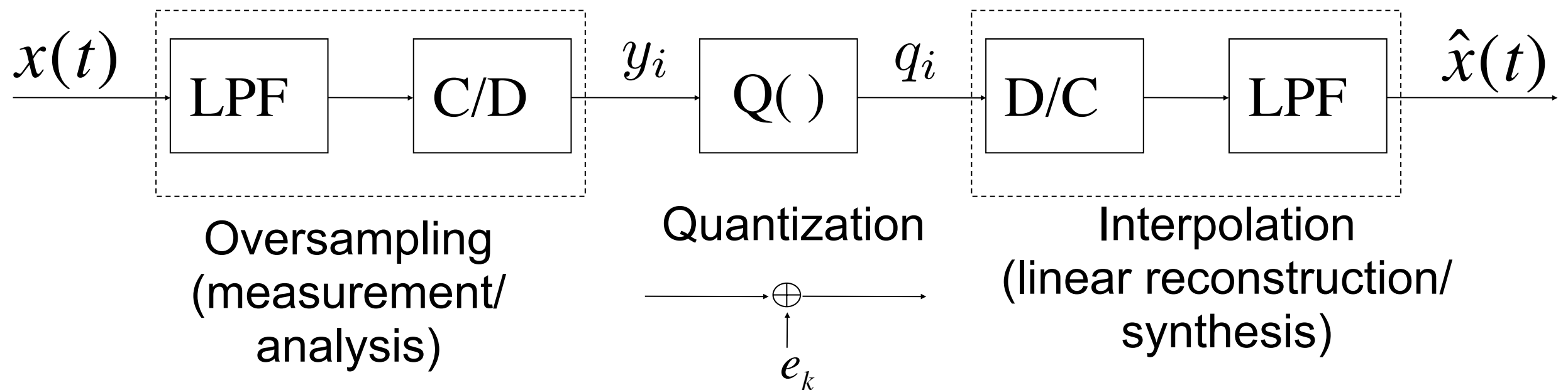
Reconstruct a signal that explains *quantized* measurements

$$\hat{\mathbf{x}} \quad \text{s.t.} \quad \mathbf{q} = Q(\mathbf{A}\hat{\mathbf{x}})$$

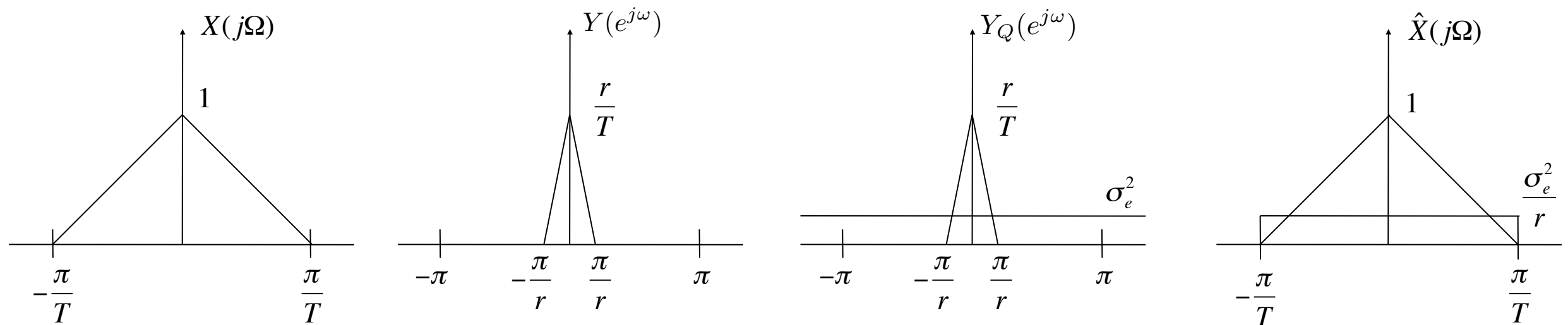
$$\text{i.e.} \quad q_i - \frac{\alpha}{2} \leq \langle \mathbf{a}_i, \hat{\mathbf{x}} \rangle \leq q_i + \frac{\alpha}{2}$$

- Thao N. and Vetterli M., “Reduction of the MSE in R-times oversampled A/D conversion  $O(1/R)$  to  $O(1/R^2)$ ,” *IEEE Trans. Signal Processing*, vol. 42, no. 1, pp. 200–203, Jan 1994.

# Oversampling and Quantization

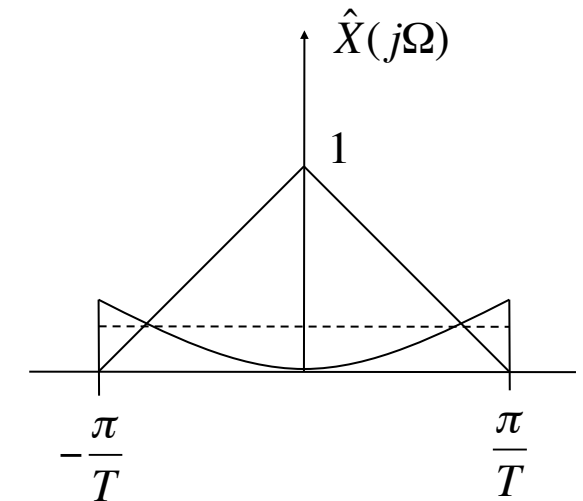
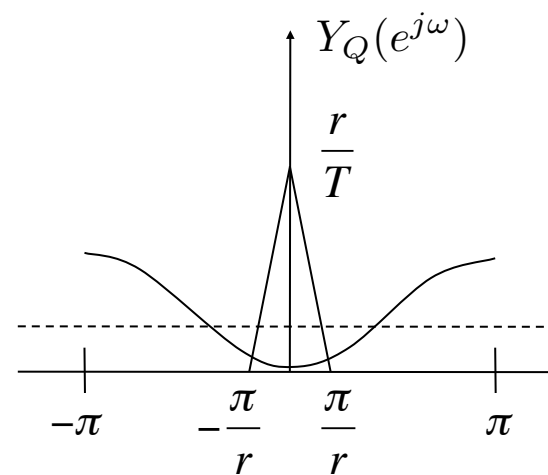
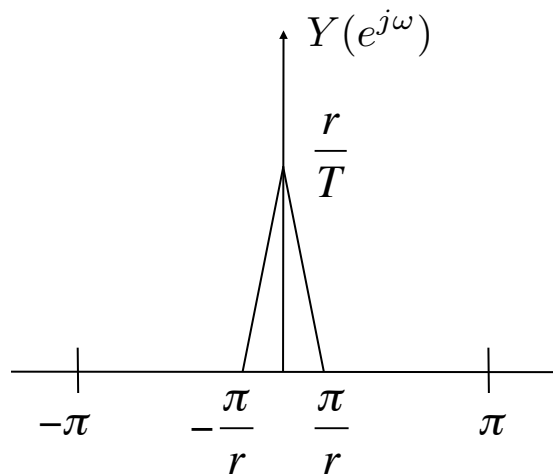
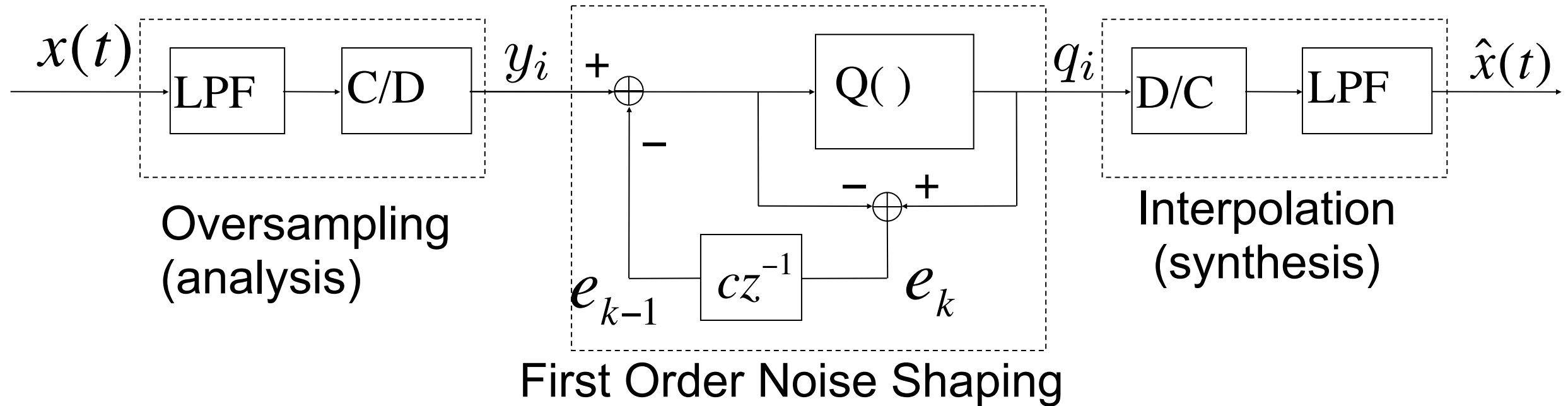


Using the additive noise model,  $e_k$  uncorrelated, uniform in  $\pm \frac{\Delta}{2}$



**Tradeoff:** Gain **1 bit** for each **4 times oversampling**  
**Quantization error**  $\varepsilon^2 \sim \Omega(1/r)$

# First Order Noise Shaping



Optimal choice  $c = \text{sinc}\left(\frac{\pi}{r}\right)$  ( $\approx 1$  for  $r \geq 4$ )

**Can we extend noise shaping to arbitrary frames?**

# Error compensation using projections

$$\mathbf{x} = y_1 \mathbf{s}_1 + y_2 \mathbf{s}_2 + y_3 \mathbf{s}_3$$

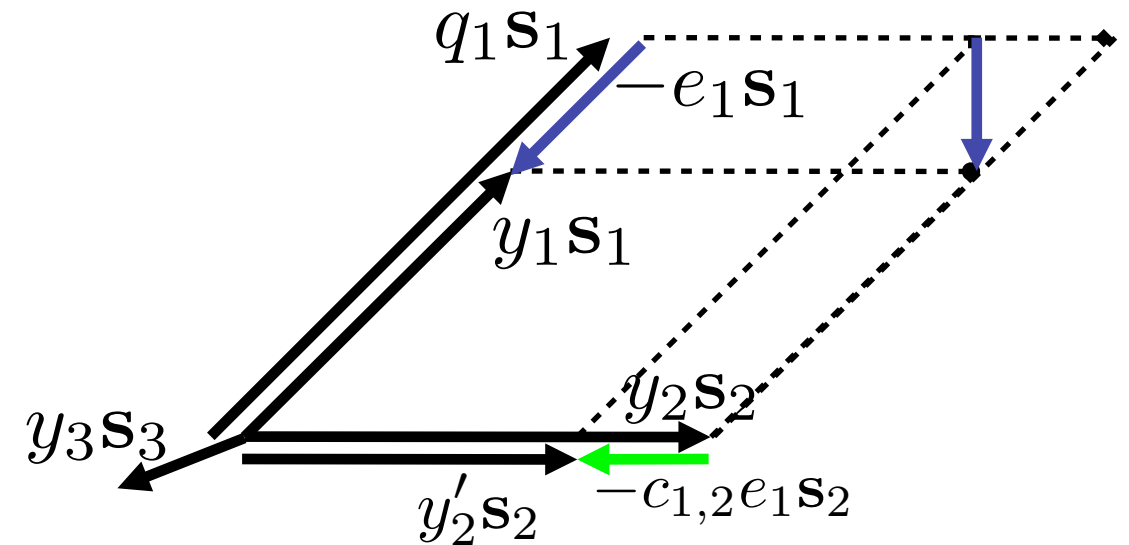
## 1. Quantization

$$\mathbf{x} = q_1 \mathbf{s}_1 + y_2 \mathbf{s}_2 + y_3 \mathbf{s}_3$$

## 2. Compensation using projection

$$y'_2 = y_2 - e_1 c_{1,2}$$

$$\mathbf{x} = q_1 \mathbf{s}_1 + y'_2 \mathbf{s}_2 + y_3 \mathbf{s}_3 - e_1 (\mathbf{s}_1 - c_{1,2} \mathbf{s}_2)$$

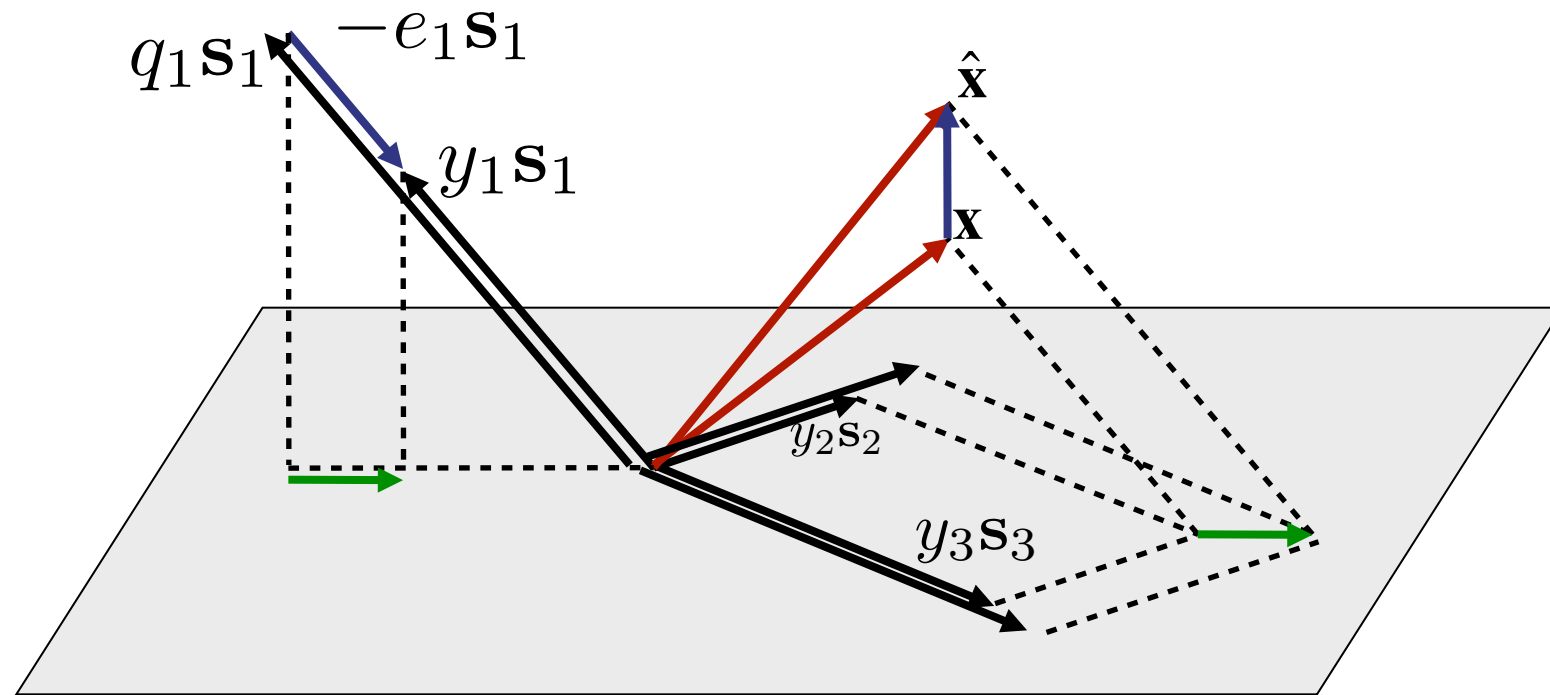


$$\text{Incremental error : } -e_1 (\mathbf{s}_1 - c_{1,2} \mathbf{s}_2) \Rightarrow c_{1,2} = \frac{\langle \mathbf{s}_1, \mathbf{s}_2 \rangle}{\|\mathbf{s}_2\|^2}$$

**Compensation linear in the error.**  
**Coefficients can be pre-computed.**

- Boufounos P. and Oppenheim A. V., “Quantization noise shaping on arbitrary frame expansions,” *EURASIP Journal on Applied Signal Processing, Special issue on Frames and Overcomplete Representations in Signal Processing, Communications, and Information Theory*, vol. 2006, pp. Article ID 53 807, 12 pages, DOI:10.1155/ASP/2006/53 807, 2006.

# Higher Order Projections



$$\mathbf{x} = y_1 \mathbf{s}_1 + y_2 \mathbf{s}_2 + y_3 \mathbf{s}_3$$

## 1. Quantization:

$$q_i = Q(y'_i) = y'_i + e_i$$

## 2. Projection:

$$y'_{i+1} = y_{i+1} - c_{i,i+1} e_i$$

$$\vdots$$

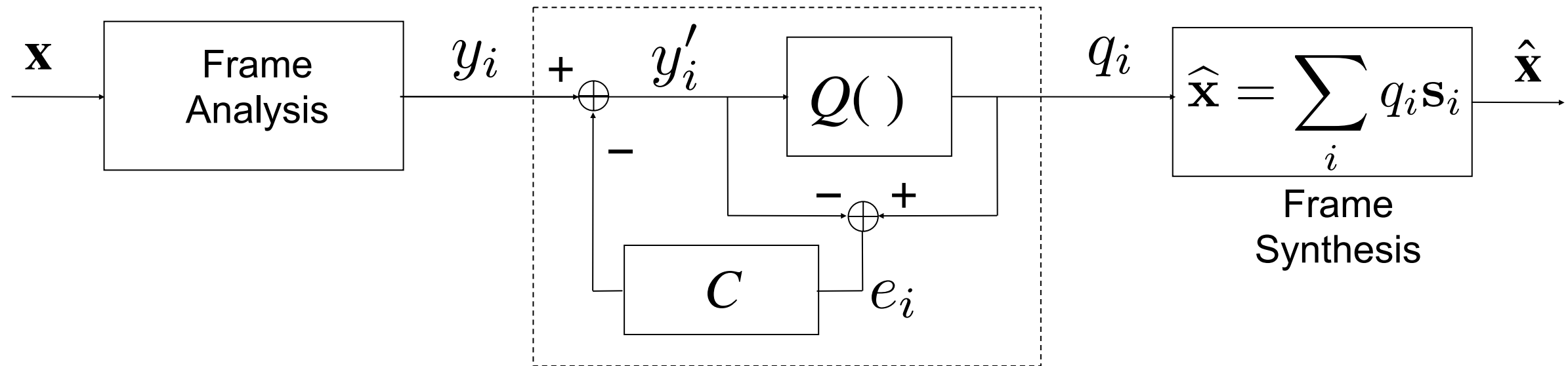
$$y'_{i+p} = y_{i+p} - c_{i,i+p} e_i$$

Projection coefficients  $c_{i,i+k}$  designed to reduce or minimize

$$\left\| \mathbf{s}_i - \sum_{k=1}^p c_{i,i+k} \mathbf{s}_{i+k} \right\|_2$$

- Boufounos P. and Oppenheim A. V., “Quantization noise shaping on arbitrary frame expansions,” *EURASIP Journal on Applied Signal Processing, Special issue on Frames and Overcomplete Representations in Signal Processing, Communications, and Information Theory*, vol. 2006, pp. Article ID 53 807, 12 pages, DOI:10.1155/ASP/2006/53 807, 2006.
- Benedetto J. J., Powell A. M., and Yilmaz O., “Sigma-Delta quantization and finite frames,” *IEEE Trans. Info. Theory*, vol. 52, no. 5, pp. 1990–2005, May 2006.
- Deift, P., Krahmer, F. and Güntürk, C. S. (2011), “An optimal family of exponentially accurate one-bit Sigma-Delta quantization schemes.” *Comm. Pure Appl. Math.*, 64: 883–919. doi: 10.1002/cpa.20367

# System Description



## 1. Quantization:

$$q_i = Q(y'_i) = y'_i + e_i$$

## 2. Projection (Coefficient Update):

$$y'_{i+1} = y_{i+1} - c_{i,i+1} e_i$$

$\vdots$

$$y'_{i+p} = y_{i+p} - c_{i,i+p} e_i$$

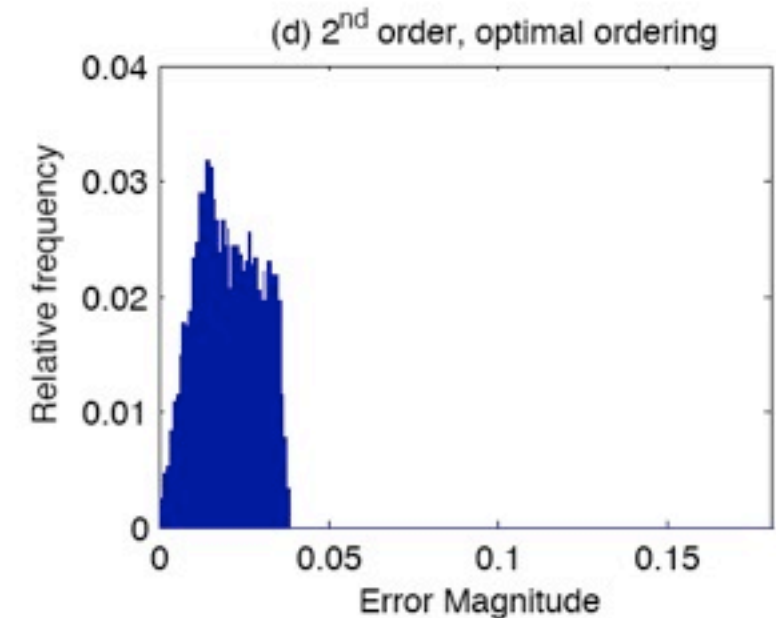
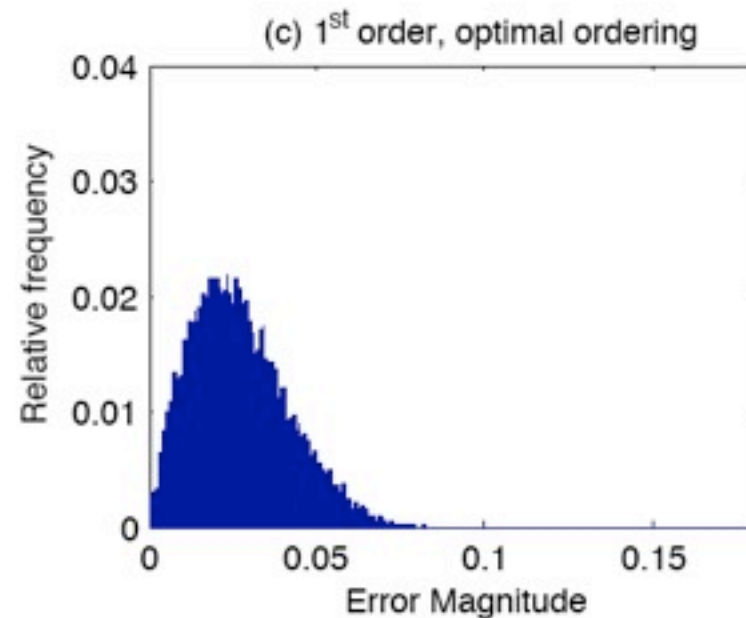
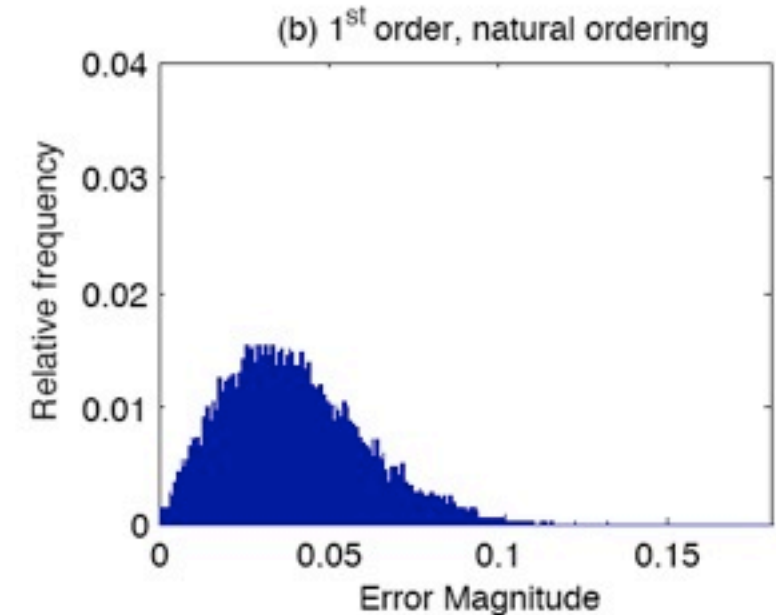
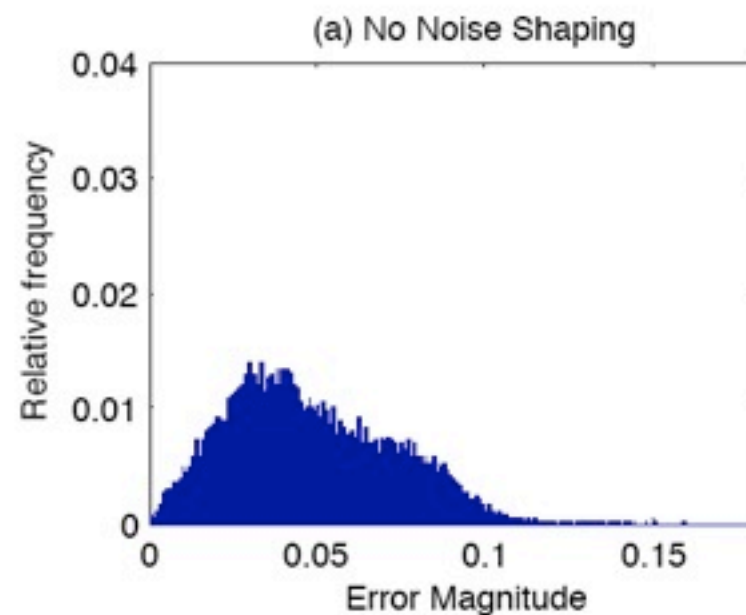
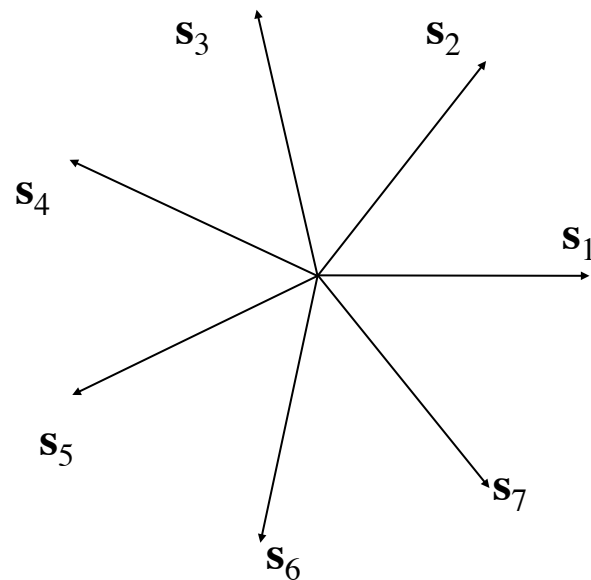
**Achievable error decay  $\varepsilon = O(r^{p+1})$**

- Benedetto J. J., Powell A. M., and Yilmaz O., "Sigma-Delta quantization and finite frames," *IEEE Trans. Info. Theory*, vol. 52, no. 5, pp. 1990–2005, May 2006.

# Example: Simulation Results

## Histogram of the Error Magnitude

Frame: 7th roots of unity



- Random points on the plane, uniform inside the unit circle.
- Quantization points:  $(-7/8, -5/8, -3/8, -1/8, 1/8, 3/8, 5/8, 7/8)$
- Optimal ordering (one of many) is:  $(s_1, s_4, s_7, s_3, s_6, s_2, s_5)$

# Further Reading

---

- Thao N. and Vetterli M., “Reduction of the MSE in R-times oversampled A/D conversion  $O(1/R)$  to  $O(1/R^2)$ ,” *IEEE Trans. Signal Processing*, vol. 42, no. 1, pp. 200–203, Jan 1994.
- Thao N. T. and Vetterli M., “Lower bound on the mean-squared error in oversampled quantization of periodic signals using vector quantization analysis,” *IEEE Trans. Info. Theory*, vol. 42, no. 2, pp. 469–479, Mar. 1996.
- Thao N. T., “Vector quantization analysis of  $\Sigma\Delta$  modulation,” *IEEE Trans. Signal Processing*, vol. 44, no. 4, pp. 808–817, Apr. 1996.
- Goyal V. K., Vetterli M., and Thao N. T., “Quantized overcomplete expansions in  $R^N$  : Analysis, synthesis, and algorithms,” *IEEE Trans. Info. Theory*, vol. 44, no. 1, pp. 16–31, Jan. 1998.
- Boufounos P. and Oppenheim A.V., “Quantization noise shaping on arbitrary frame expansions,” *EURASIP Journal on Applied Signal Processing, Special issue on Frames and Overcomplete Representations in Signal Processing, Communications, and Information Theory*, vol. 2006, pp. Article ID 53 807, 12 pages, DOI:10.1155/ASP/2006/53 807, 2006.
- Boufounos P. T., “Quantization and erasures in frame representations,” *MIT D.Sc. Thesis*, Cambridge, MA, January 2006.
- Benedetto J. J., Powell A. M., and Yilmaz O., “Sigma-Delta quantization and finite frames,” *IEEE Trans. Info. Theory*, vol. 52, no. 5, pp. 1990–2005, May 2006.
- Deift, P., Krahmer, F. and Güntürk, C. S. (2011), “An optimal family of exponentially accurate one-bit Sigma-Delta quantization schemes.” *Comm. Pure Appl. Math.*, 64: 883–919. doi: 10.1002/cpa.20367



# Today's Topics

---

1. Modern Scalar Quantization
2. Compressive Sensing Overview
3. Compressive Sensing and Quantization
4. 1-bit Compressive Sensing
5. Locality Sensitive Hashing and Universal Quantization

# Today's Topics

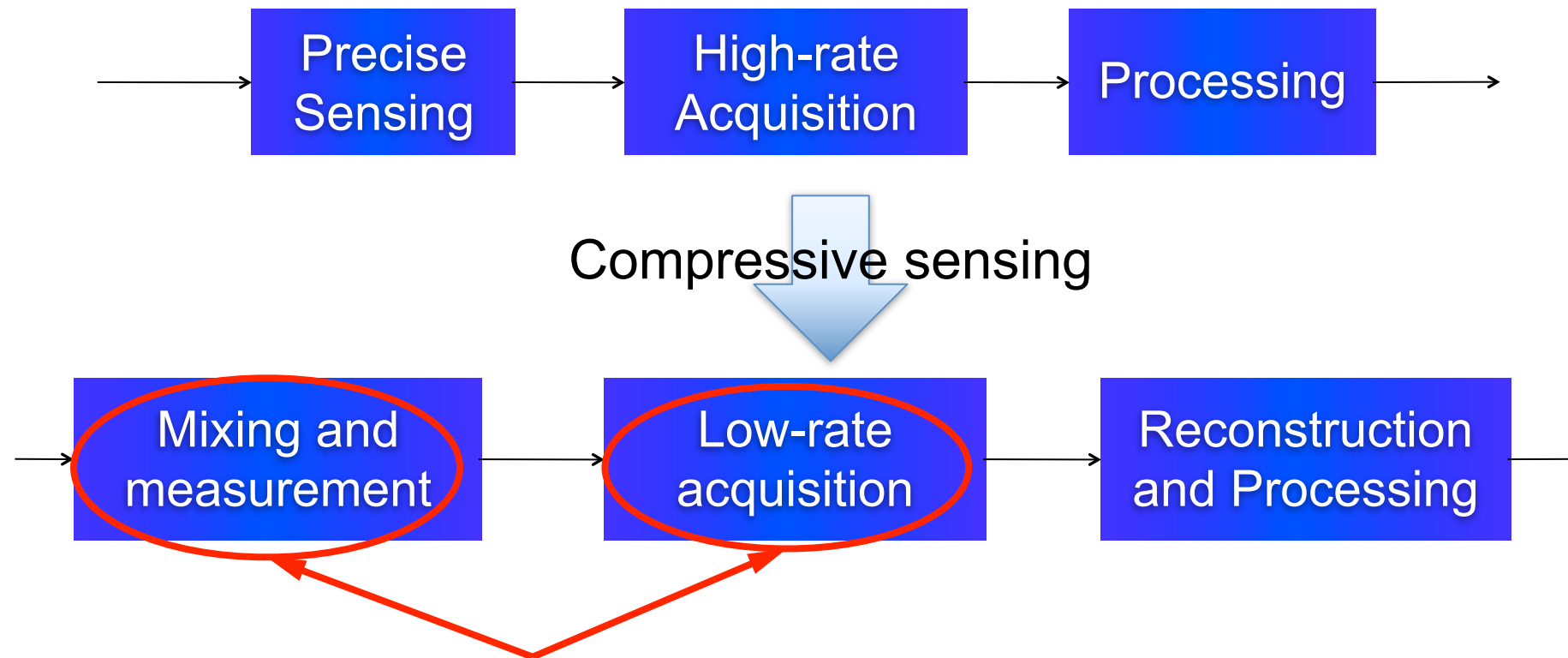
---

1. Modern Scalar Quantization
- 2. Compressive Sensing Overview**
3. Compressive Sensing and Quantization
4. 1-bit Compressive Sensing
5. Locality Sensitive Hashing and Universal Quantization

- Candès, E., Romberg, J., and Tao, T., “Stable signal recovery from incomplete and inaccurate measurements,” *Comm. Pure and Appl. Math.*, vol. 59, no. 8, pp. 1207–1223, 2006.
- Donoho D. , “Compressed sensing,” *IEEE Trans. Info. Theory*, vol. 52, no. 4, pp. 1289–1306, Sept. 2006.

# Sensing Pipeline Paradigm Change

---



**Goal:** exploit mixing to **simplify sensor or improve sensor** specifications (e.g., sensor speed, A/D conversion rate, measured bandwidth/resolution)

- Compressive sensing has significantly improved our sensing capability
- Two fundamental Compressive Sensing research aspects
  - Hardware modifications for efficient acquisition
  - Signal/scene models and processing algorithms

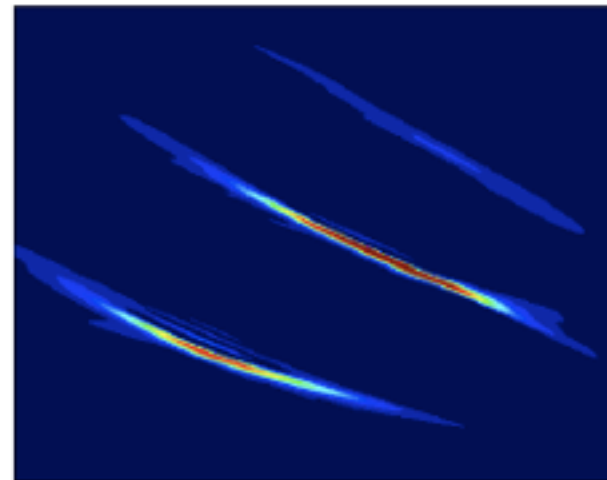
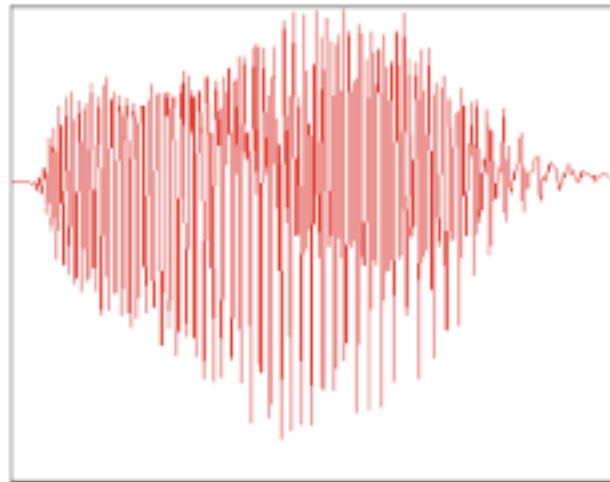
# Signal Structure: Sparsity

$N$   
pixels



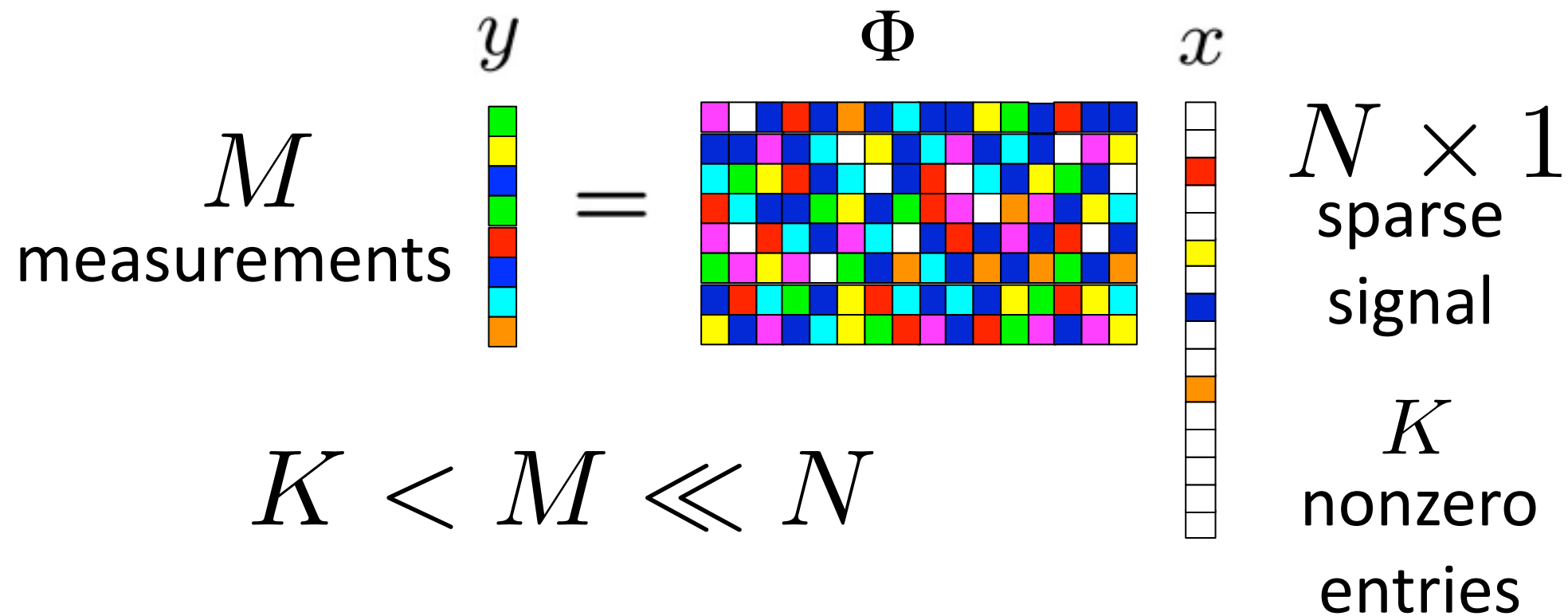
$K \ll N$   
large  
wavelet  
coefficients

$N$   
wideband  
signal  
samples



$K \ll N$   
large  
Gabor  
coefficients

# Measurement Model: Incoherence [Candes et al]



- $x$  is  $K$ -sparse or  $K$ -compressible
- $\Phi$  random, satisfies a *restricted isometry property (RIP)*

$\Phi$  has RIP of order  $2K$  with constant  $\delta$

If there exists  $\delta$  s.t. for all  $2K$ -sparse  $x$ :

$$(1 - \delta)\|x\|_2^2 \leq \|\Phi x\|_2^2 \leq (1 + \delta)\|x\|_2^2$$

- $M = O(K \log N / K)$
- $\Phi$  also has small *coherence*  $\mu \triangleq \max_{i \neq j} |\langle \phi_i, \phi_j \rangle|$

# Measurement Model: Incoherence [Candes et al]

The diagram shows the equation  $y = \Phi x$ . On the left, a vertical vector  $y$  of size 8 is shown with 8 colored squares (green, yellow, blue, green, red, blue, cyan, orange). In the middle is an equals sign. To the right of the equals sign is a matrix  $\Phi$  of size 8x12, represented as a grid of colored squares. Below the matrix is a horizontal vector  $x$  of size 12, represented as a row of colored squares (white, white, red, white, yellow, white, blue, white, orange, white, white, white). The colors in the matrix and vector correspond to the colors in the vector  $y$ .

- $x$  is  $K$ -sparse or  $K$ -compressible
- $\Phi$  random, satisfies a *restricted isometry property (RIP)*

$\Phi$  has RIP of order  $2K$  with constant  $\delta$

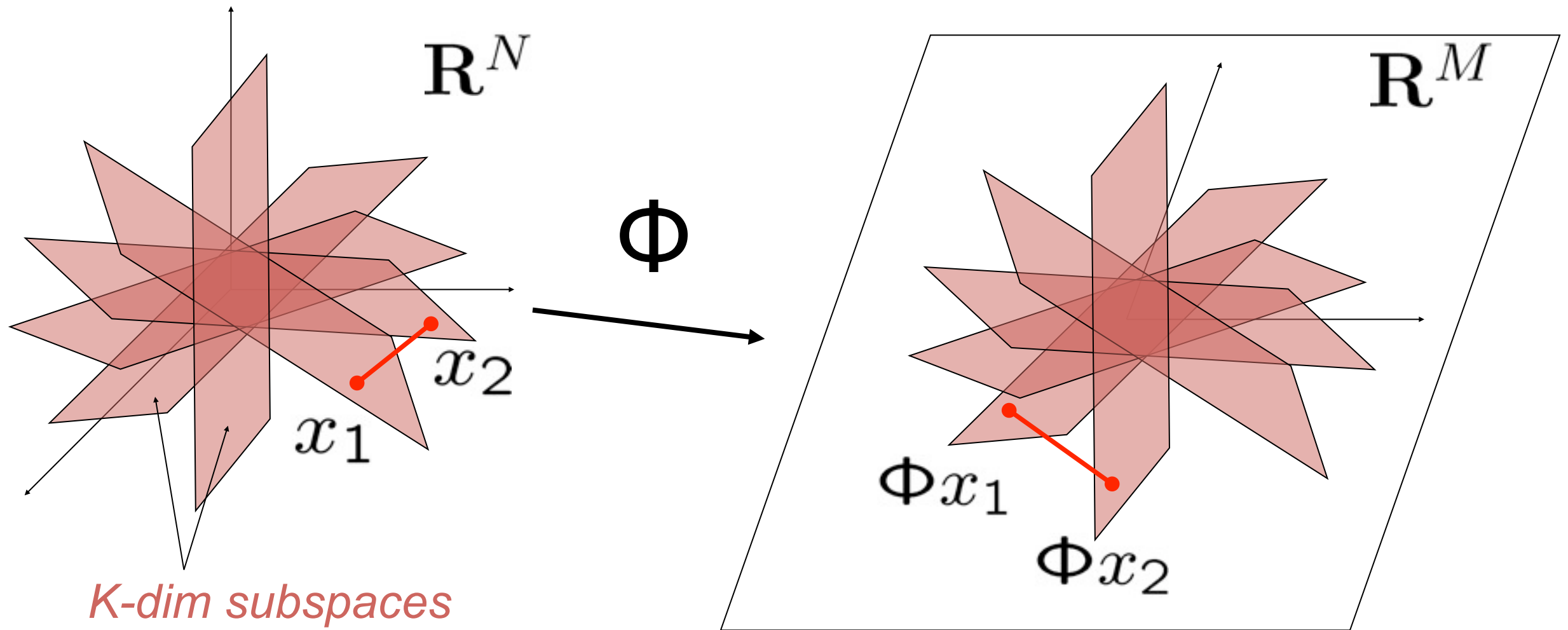
If there exists  $\delta$  s.t. for all  $2K$ -sparse  $x$ :

$$(1 - \delta)\|x\|_2^2 \leq \|\Phi x\|_2^2 \leq (1 + \delta)\|x\|_2^2$$

- $M = O(K \log N / K)$
- $\Phi$  also has small *coherence*  $\mu \triangleq \max_{i \neq j} |\langle \phi_i, \phi_j \rangle|$

# RIP/Stable Embedding

- An information preserving projection  $\Phi$  preserves the **geometry** of the set of sparse signals



Restricted Isometry Property

$$(1 - \delta)\|x\|_2^2 \leq \|\Phi x\|_2^2 \leq (1 + \delta)\|x\|_2^2$$

# Reconstruction: Non-linear, Enforcing Structure

- Reconstruction using **sparse approximation**:

- Find sparsest  $\mathbf{x}$  such that  $\mathbf{y} \approx \mathbf{A}\mathbf{x}$

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_0 \text{ s.t. } \mathbf{y} \approx \Phi \mathbf{x}$$

- **Convex optimization** approach:

- Minimize  $l_1$  norm: e.g.,

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_1 \text{ s.t. } \mathbf{y} \approx \Phi \mathbf{x}$$

- **Greedy algorithms** approach:

- Minimize  $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$  such that  $\mathbf{x}$  is sparse

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 \text{ s.t. } \|\mathbf{x}\|_0 \leq K$$

- MP, OMP, ROMP, StOMP, CoSaMP, SP, ALPS, PYAMP (Pick Your Acronym Matching Pursuit)

- **More general cost functions,**

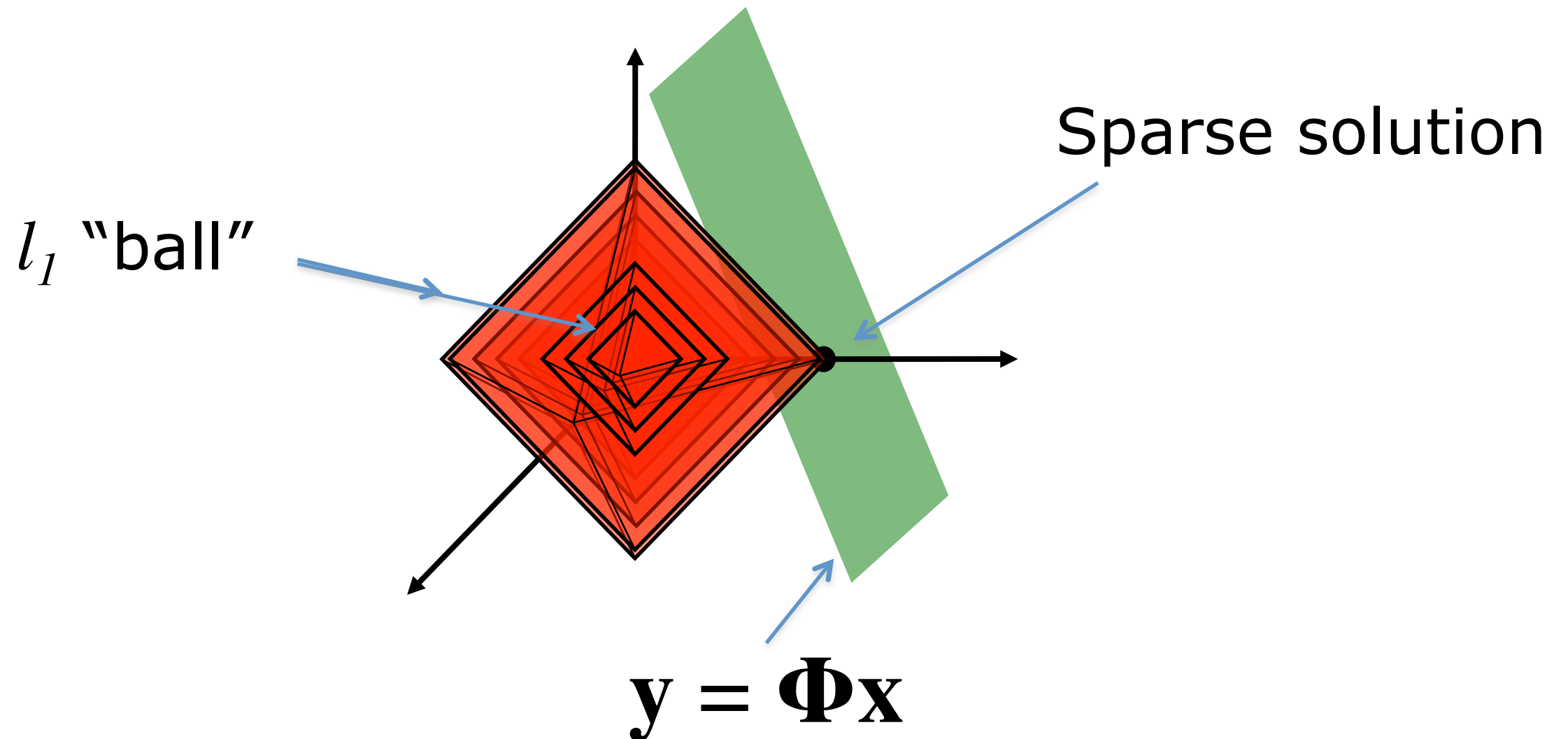
- GraSP, generalization of CoSaMP

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} f(\mathbf{x}) \text{ s.t. } \|\mathbf{x}\|_0 \leq K$$



# Why $l_1$ relaxation works

$$\min \| \mathbf{x} \|_1 \text{ s.t. } \mathbf{y} \approx \Phi \mathbf{x}$$

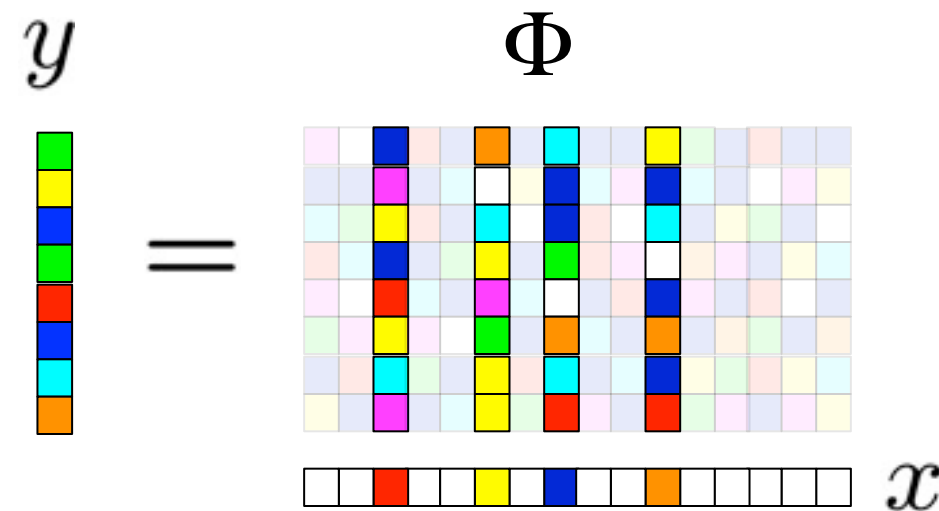


**K-term approximation error**

If  $\Phi$  satisfies the RIP:  $\| \hat{\mathbf{x}} - \mathbf{x} \|_2 \leq c_1 \frac{\| \mathbf{x} - \mathbf{x}_K \|_1}{\sqrt{K}} + c_2 \epsilon$

**Measurement error**

# Greedy Pursuits Core Idea



- $y$  highly correlated with  $\Phi$  at locations where  $x$  is high
- $\Phi^T y$  provides a good idea of these locations
  - This is why low coherence is important

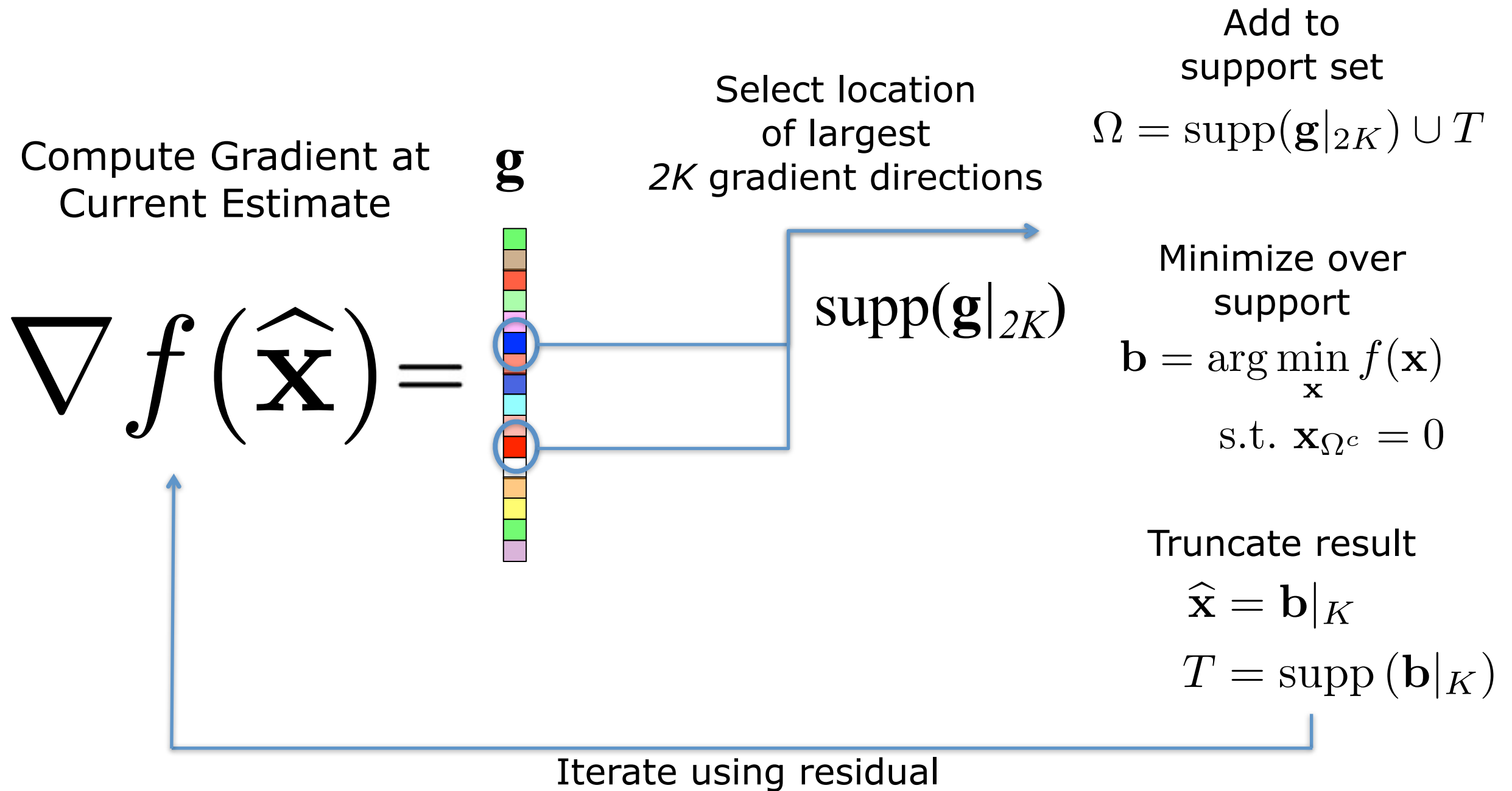
$$\mu \triangleq \max_{i \neq j} |\langle \phi_i, \phi_j \rangle|$$

- $\Phi^T y$  referred to as *proxy* for  $x$
- General Strategy:
  - Identify locations
  - Invert the system only on those locations

# GraSP (Gradient Subspace Pursuit)

**State Variables:** Signal estimate,  $\hat{\mathbf{x}}$  support estimate:  $T$

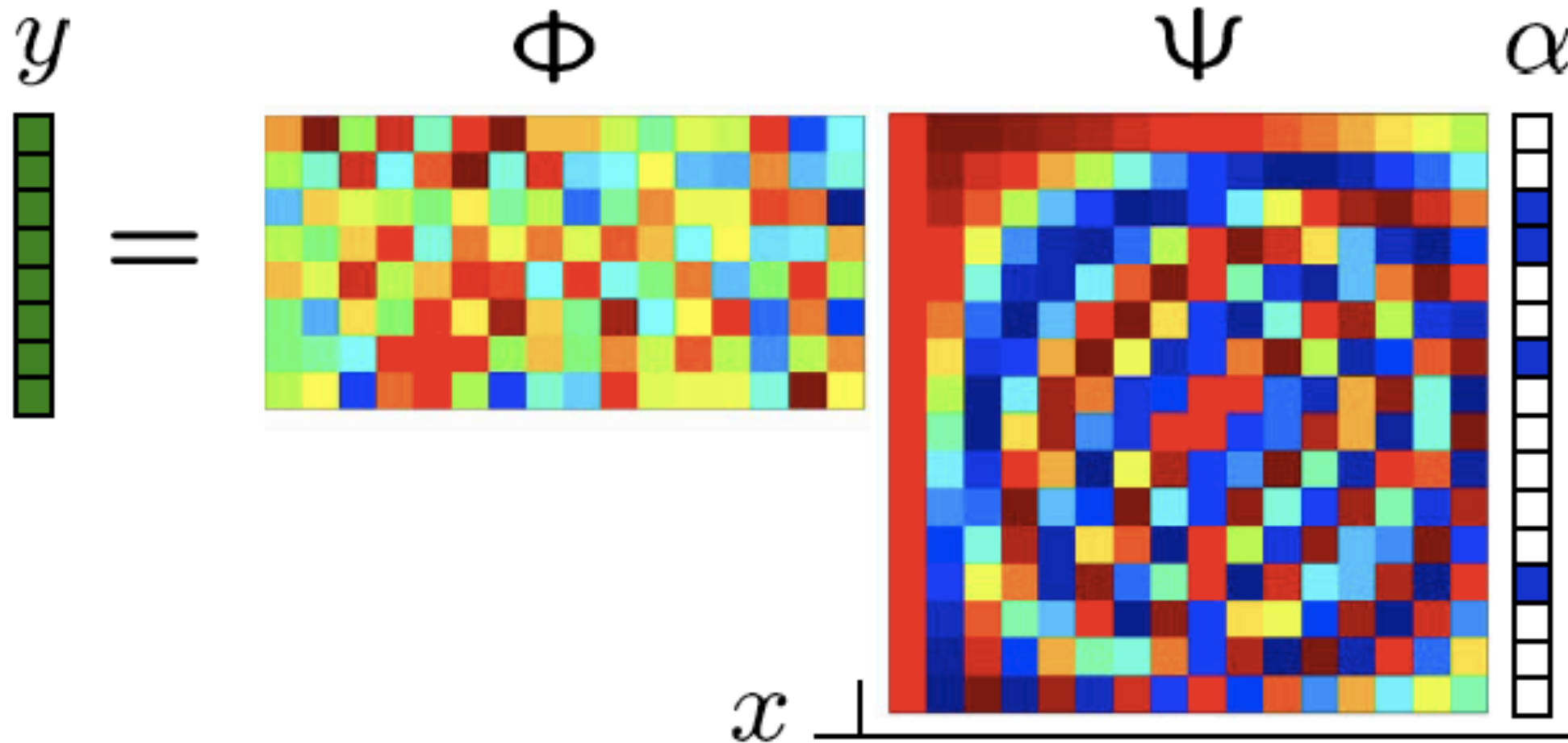
Initialize estimate and support:  $\hat{\mathbf{x}}=0$ ,  $T=\text{supp}(\hat{\mathbf{x}})$



- S. Bahmani, B. Raj, and P. T. Boufounos, "Greedy Sparsity-Constrained Optimization," *Journal of Machine Learning Research*, v. 14, pp. 807-841, March, 2013.

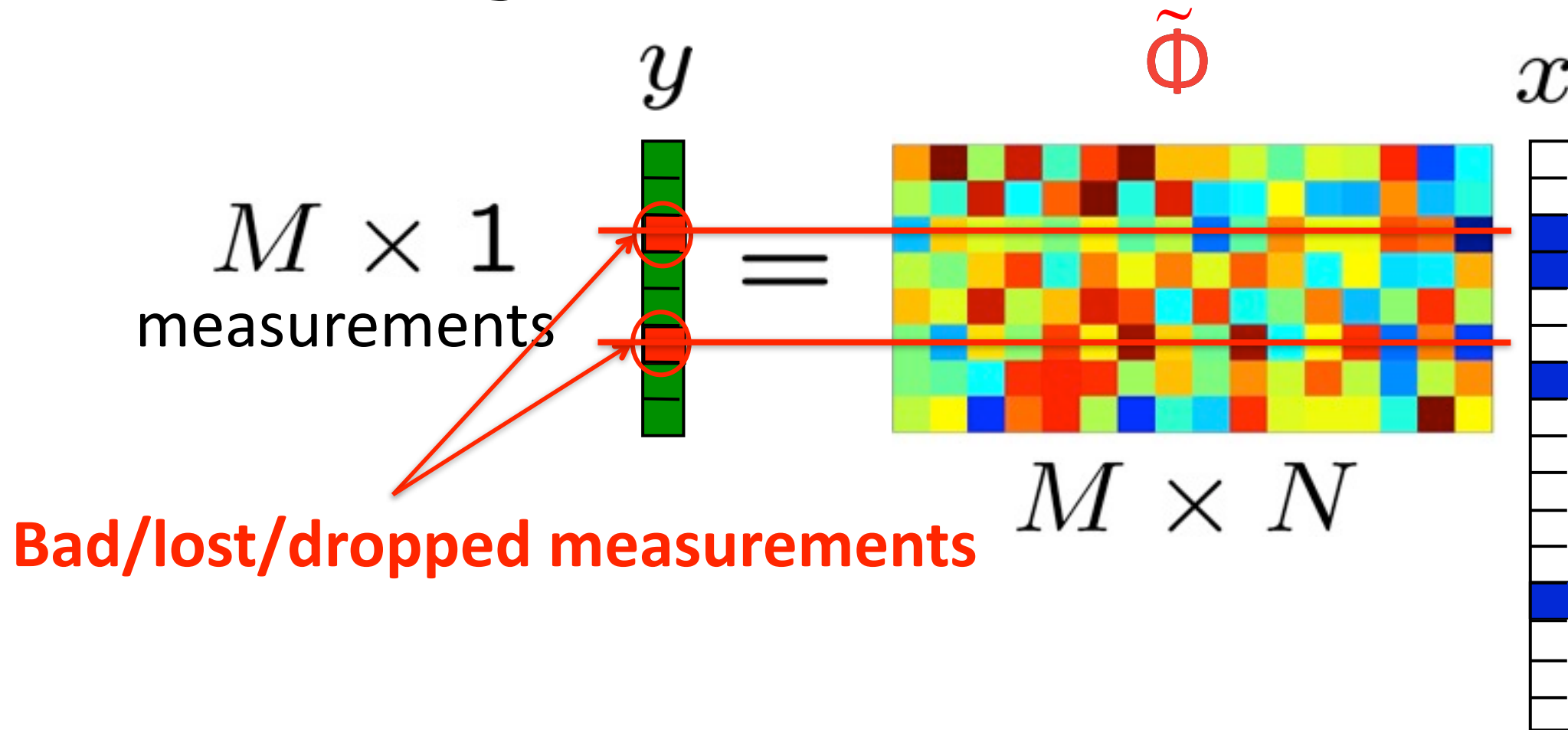
# Universality

- Gaussian white noise basis is incoherent with *any* fixed orthonormal basis (with high probability)
- Signal sparse in frequency domain:  $\psi = \text{idct}$



- Product  $\Phi \Psi$  remains Gaussian white noise

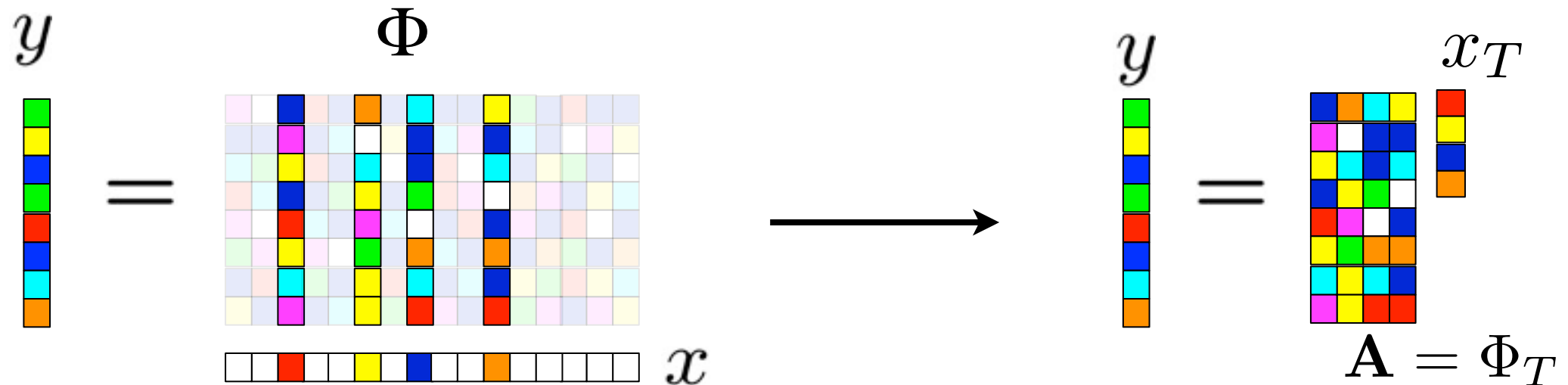
# Democracy



- Measurements are **democratic** [Davenport, Laska, Boufounos, Baraniuk]
  - They are all equally important
  - We can **lose** some **arbitrarily** (i.e. an adversary can choose which ones)
- The  $\tilde{\Phi}$  still satisfies RIP (as long as we don't drop too many)

# Compressive Sensing and Oversampling

Given support of signal  $T$



Resulting system is oversampled:  $\mathbf{A} \in \mathbb{R}^{M \times K}$

$$M = O(K \log N)$$

$$\Rightarrow \textcircled{r} = O(\log N)$$

**Oversampling  
Rate**

**Oversampling** provides **robustness**, but introduces **inefficiency**

- Boufounos P., Baraniuk R. G., "Quantization of Sparse Representations." *Rice University ECE Department Technical Report 0701*. Summary appears in *Proc. of the Data Compression Conference (DCC '07)*, March 27-29 2007, Snowbird, UT.

# Further Reading

---

- Candès, E., Romberg, J., and Tao, T., “Stable signal recovery from incomplete and inaccurate measurements,” *Comm. Pure and Appl. Math.*, vol. 59, no. 8, pp. 1207–1223, 2006.
- Donoho D. , “Compressed sensing,” *IEEE Trans. Info. Theory*, vol. 52, no. 4, pp. 1289–1306, Sept. 2006.
- Candès, Emmanuel J. “Compressive sampling.” *Proceedings oh the International Congress of Mathematicians: invited lectures*, August 22-30, 2006, Madrid, Spain.
- Candes, E.J.; Wakin, M.B., "An Introduction To Compressive Sampling," *IEEE Signal Processing Magazine*, vol.25, no. 2, pp.21,30, March 2008
- M. A. Davenport, J. N. Laska, P. T. Boufounos, and R. G. Baraniuk, “A simple proof that random matrices are democratic,” *Rice University ECE Department Technical Report TREE-0906*, Houston, TX, November, 2009.
- Needell D. and Tropp J.A., “CoSaMP: Iterative signal recovery from incomplete and inaccurate samples,” *Applied and Computational Harmonic Analysis*, vol. 26, no. 3, pp. 301-321, May 2009.
- S. Bahmani, B. Raj, and P. T. Boufounos, “Greedy Sparsity-Constrained Optimization,” *Journal of Machine Learning Research*, v. 14, pp. 807-841, March, 2013.

# Today's Topics

---

1. Modern Scalar Quantization
2. Compressive Sensing Overview
3. Compressive Sensing and Quantization
4. 1-bit Compressive Sensing
5. Locality Sensitive Hashing and Universal Quantization



# Today's Topics

---

1. Modern Scalar Quantization
2. Compressive Sensing Overview
- 3. Compressive Sensing and Quantization**
4. 1-bit Compressive Sensing
5. Locality Sensitive Hashing and Universal Quantization

- Candès, E., Romberg, J., and Tao, T., “Stable signal recovery from incomplete and inaccurate measurements,” *Comm. Pure and Appl. Math.*, vol. 59, no. 8, pp. 1207–1223, 2006.
- Donoho D. , “Compressed sensing,” *IEEE Trans. Info. Theory*, vol. 52, no. 4, pp. 1289–1306, Sept. 2006.

Part III:  
When quantization meets  
compressed sensing

Laurent Jacques, UCL, Belgium  
Petros Boufounos, MERL, USA

# Outline:

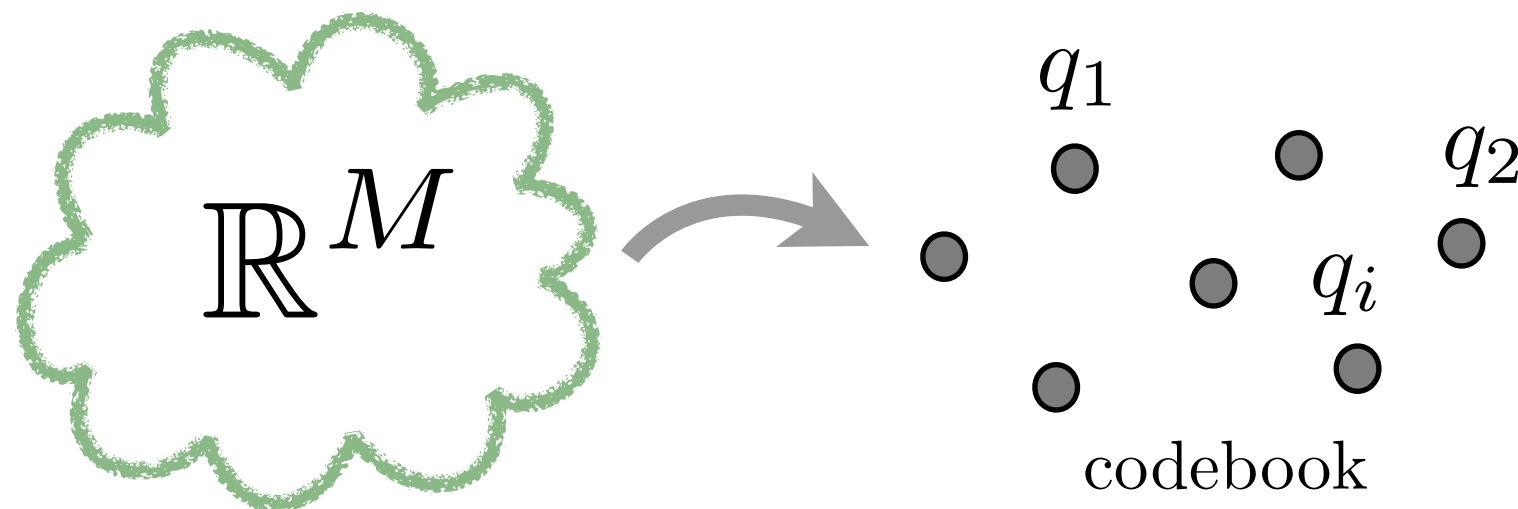
1. Context
2. Former QCS methods and performance limits
3. Consistent Reconstructions
4. Sigma-Delta quantization in CS
5. To saturate or not? And how much?

# 1. Context

# What is quantization?

- Generality:

Intuitively: “Quantization maps a continuous domain to a set of finite elements (or codebook)”



$$Q[x] \in \{q_1, q_2, \dots\}$$

- Oldest example: rounding off  $\lfloor x \rfloor, \lceil x \rceil, \dots \quad \mathbb{R} \rightarrow \mathbb{Z}$

What is quantization? ...

# Example 1: scalar quantization

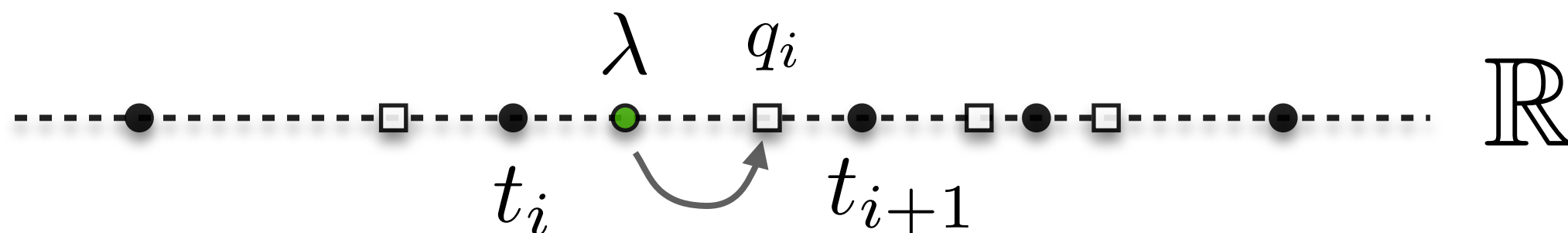
- In  $\mathbb{R}^M$ , on each component of  $M$ -dimensional vectors:

$$\Omega = \{q_i \in \mathbb{R} : 1 \leq i \leq 2^B\}, \quad (\text{levels}) \quad \square$$

$$\mathcal{T} = \{t_i \in \overline{\mathbb{R}} : 1 \leq i \leq 2^B + 1, t_i \leq t_{i+1}\} \quad (\text{thresholds}) \quad \bullet$$

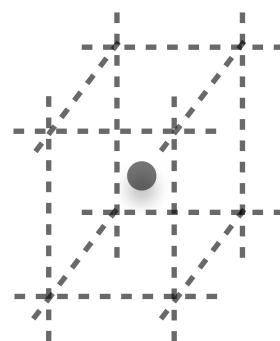
$$\forall \lambda \in \mathbb{R}, \quad \mathcal{Q}[\lambda] = q_i \Leftrightarrow \lambda \in \mathcal{R}_i \triangleq [t_i, t_{i+1}), \quad \text{1-D quantization cell}$$

$$\forall u \in \mathbb{R}^M, \quad (\mathcal{Q}[u])_j = \mathcal{Q}[u_j]$$



- Globally:

$$\mathcal{Q}[z] = \mathbf{q} \in \Omega^M \Leftrightarrow z \in$$



$$\begin{aligned} &M\text{-D quantization cell} \\ &\mathcal{R}_{i_1} \times \mathcal{R}_{i_2} \times \cdots \times \mathcal{R}_{i_M} \\ &:= \mathcal{Q}^{-1}[\mathbf{q}] \end{aligned}$$

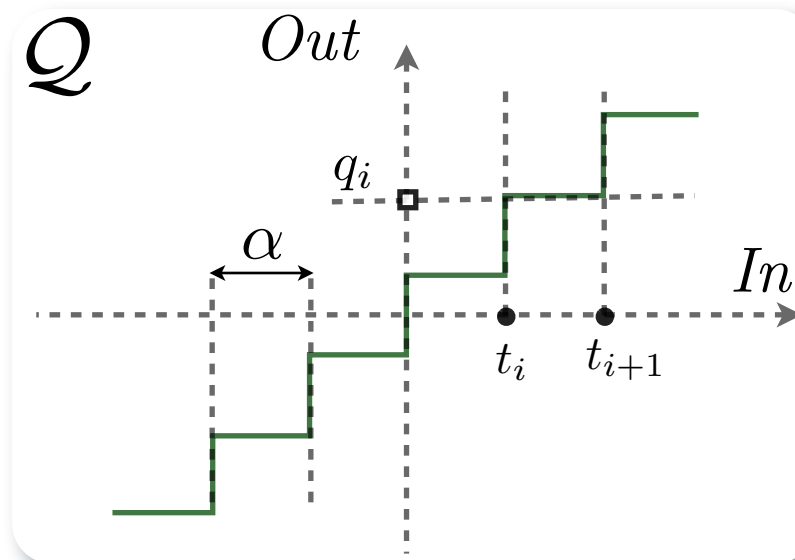
What is quantization? ...

# Example 1: scalar quantization

- Regular uniform

$$q_k = (k + 1/2)\alpha$$

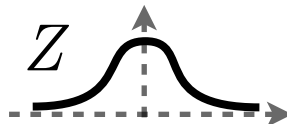
$$t_k = k\alpha$$

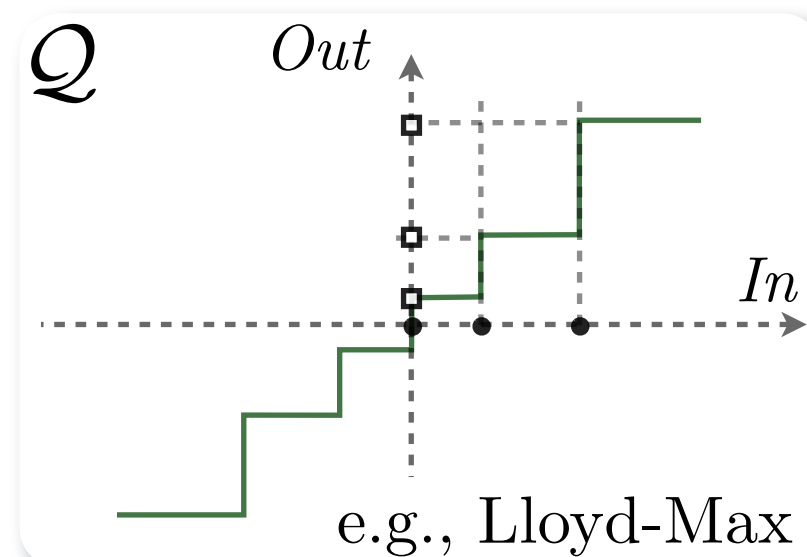


- Regular non-uniform

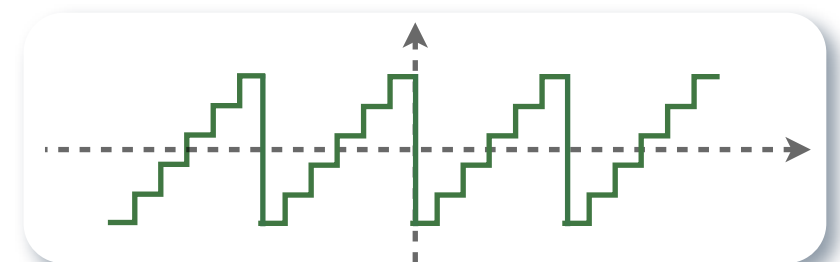
$\Omega$  and  $\mathcal{T}$  optimized

e.g., wrt an input distribution  $Z$   
find minimum distortion, *i.e.*,


$$\argmin_{\mathcal{T}, \Omega} \mathbb{E}_Z \|Z - Q[Z]\|^2$$



- Non-regular  $\rightarrow$  Petros, Part V

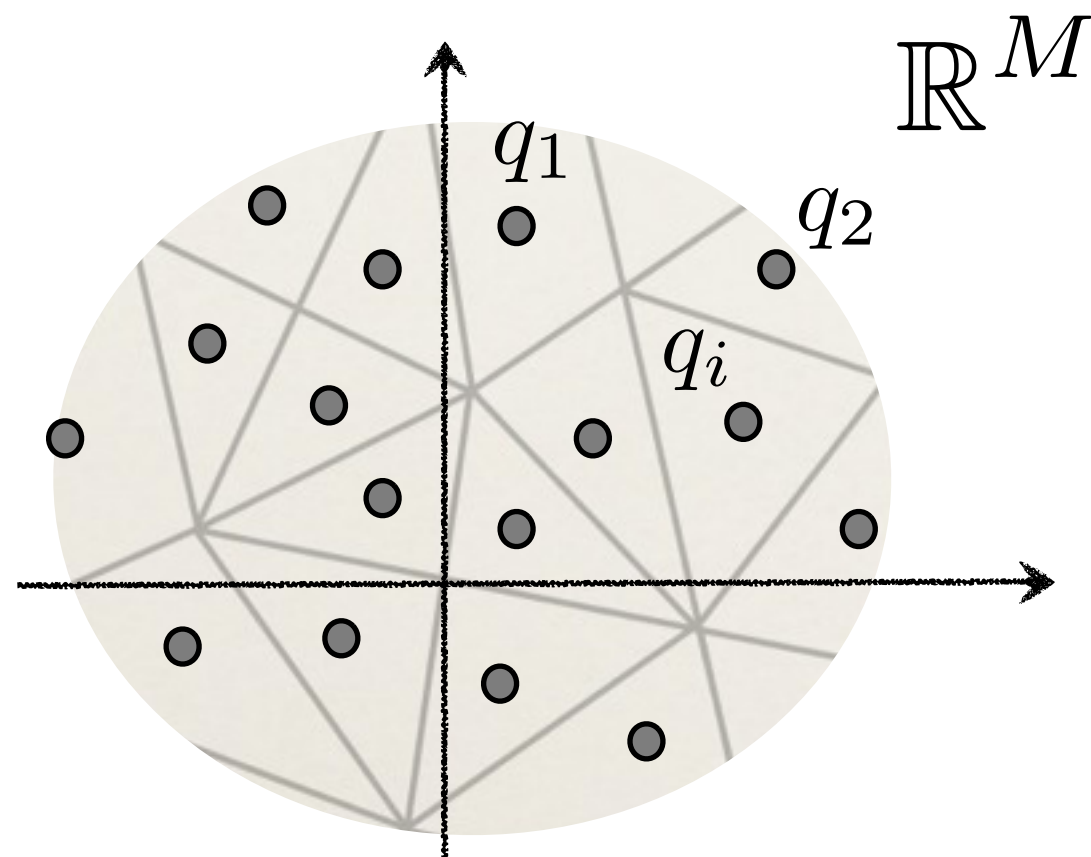


What is quantization? ...

## Example 2: vector quantization

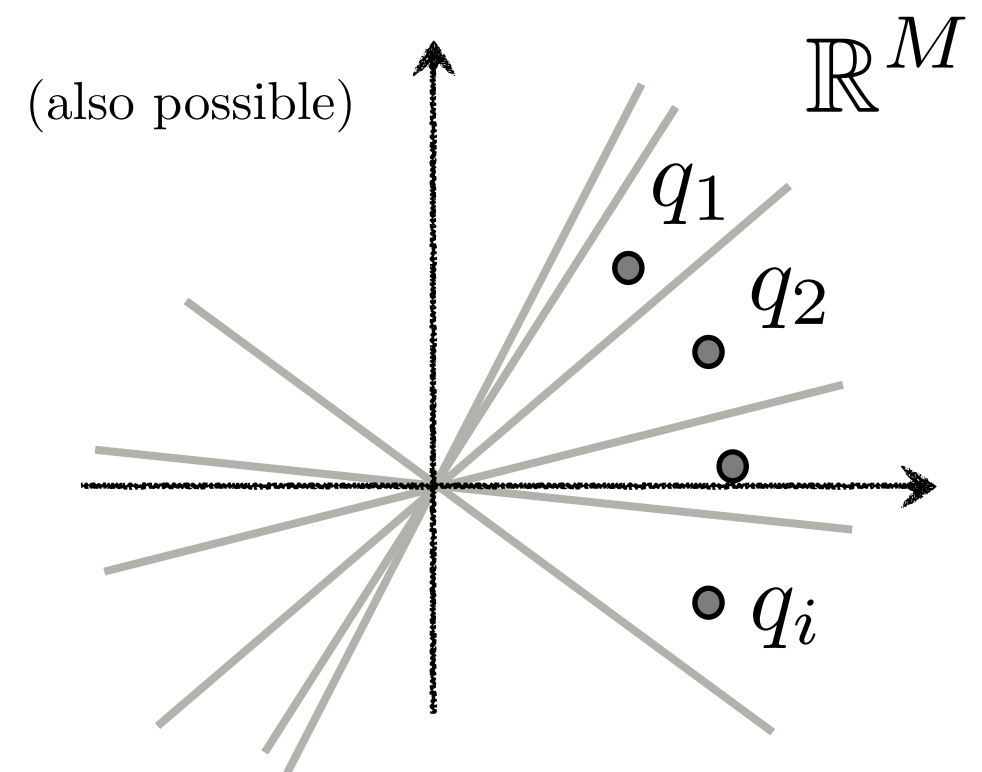
(caveat: not really covered in this tutorial, ... except  $\Sigma\Delta$ , see later)

Quantization = codebook  $\Omega$  + quantization cells  $\mathcal{R} = \{\mathcal{R}_i \subset \mathbb{R}^M\}$



(non-separable quantization)

e.g.,  $\operatorname{argmin}_{\Omega, \mathcal{R}} \mathbb{E}_Z \|Z - Q[Z]\|^2$

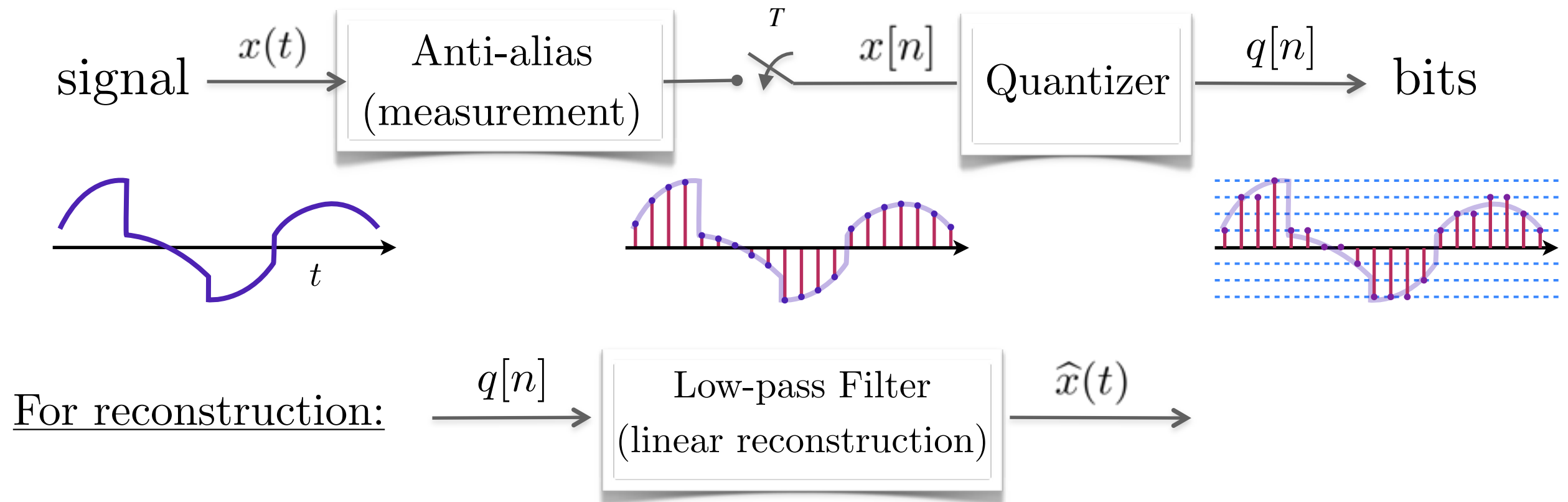


e.g., encoding components ordering + sign  
(*permutation frame quantization*)

(Nguyen et al, Goyal, ...)



# Classical Sampling and Quantization



**Sampling:** discretization in time

Lossless at the Nyquist rate

**Quantization:** discretization in amplitude

Always lossy

Need both for digital data acquisition

# Compressive Sampling and Quantization

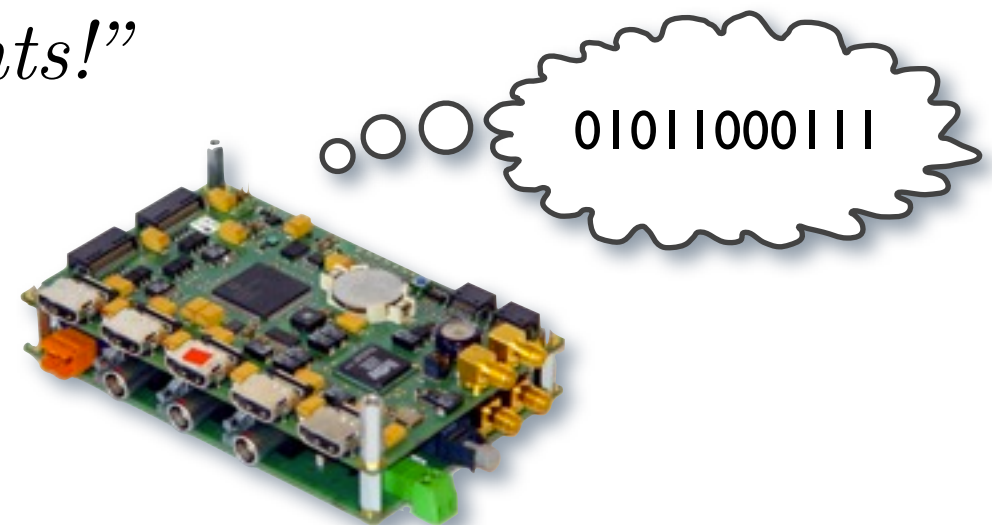
Compressed sensing theory says:

*“Linearly sample a signal  
at a rate function of  
its intrinsic dimensionality”*



Information theory and sensor designer say:

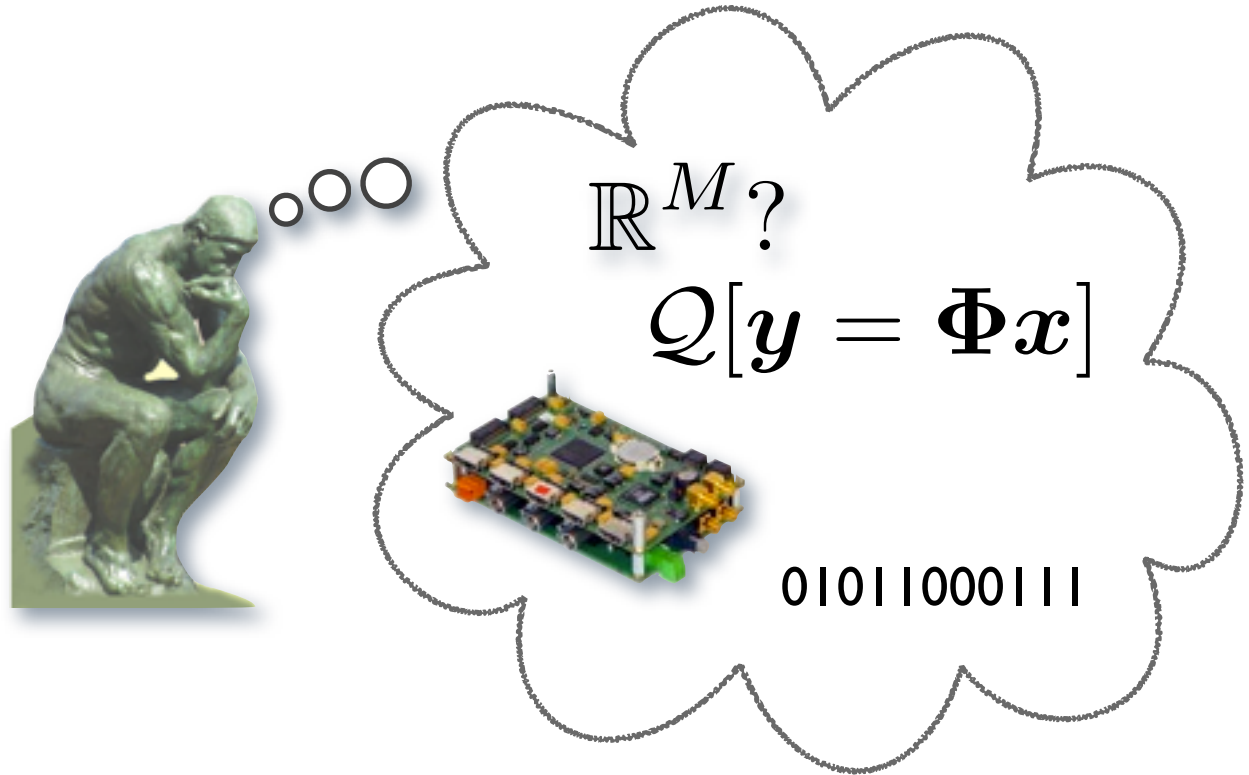
*“Okay, but I need to  
quantize/digitize my measurements!”  
(e.g., in ADC)*



# The Quantized CS Problem (QCS)

## Natural questions:

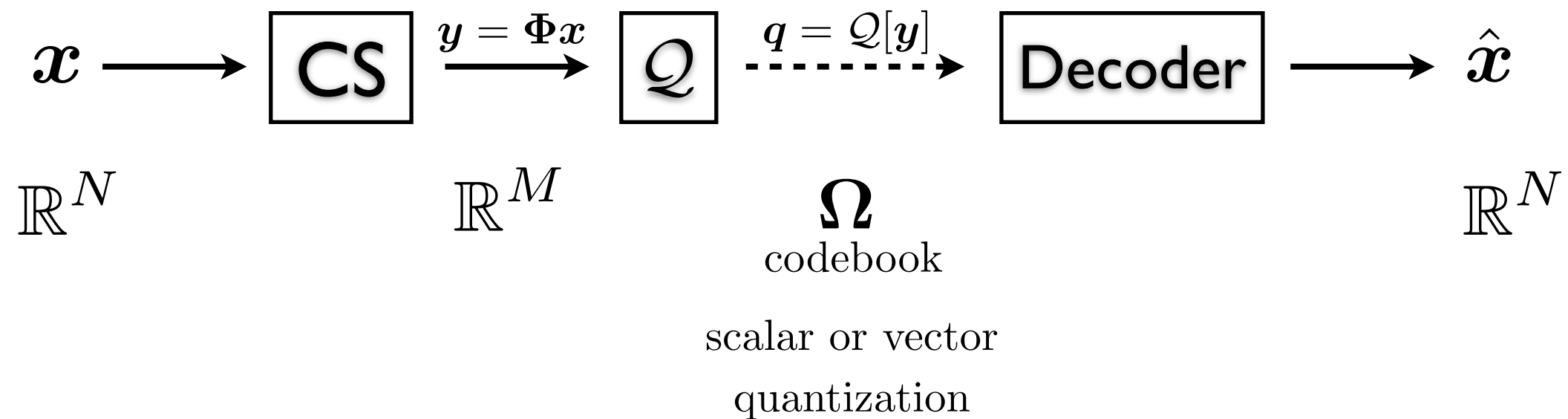
- ▶ How to integrate quantization in CS?
- ▶ What do we loose?
- ▶ Are there some theoretical limitations?  
(related to information theory? geometry?)
- ▶ How to minimize quantization effects in the reconstruction?



# QCS: a system view

With **no additional noise**:

e.g., basis pursuit,  
greedy methods, ...



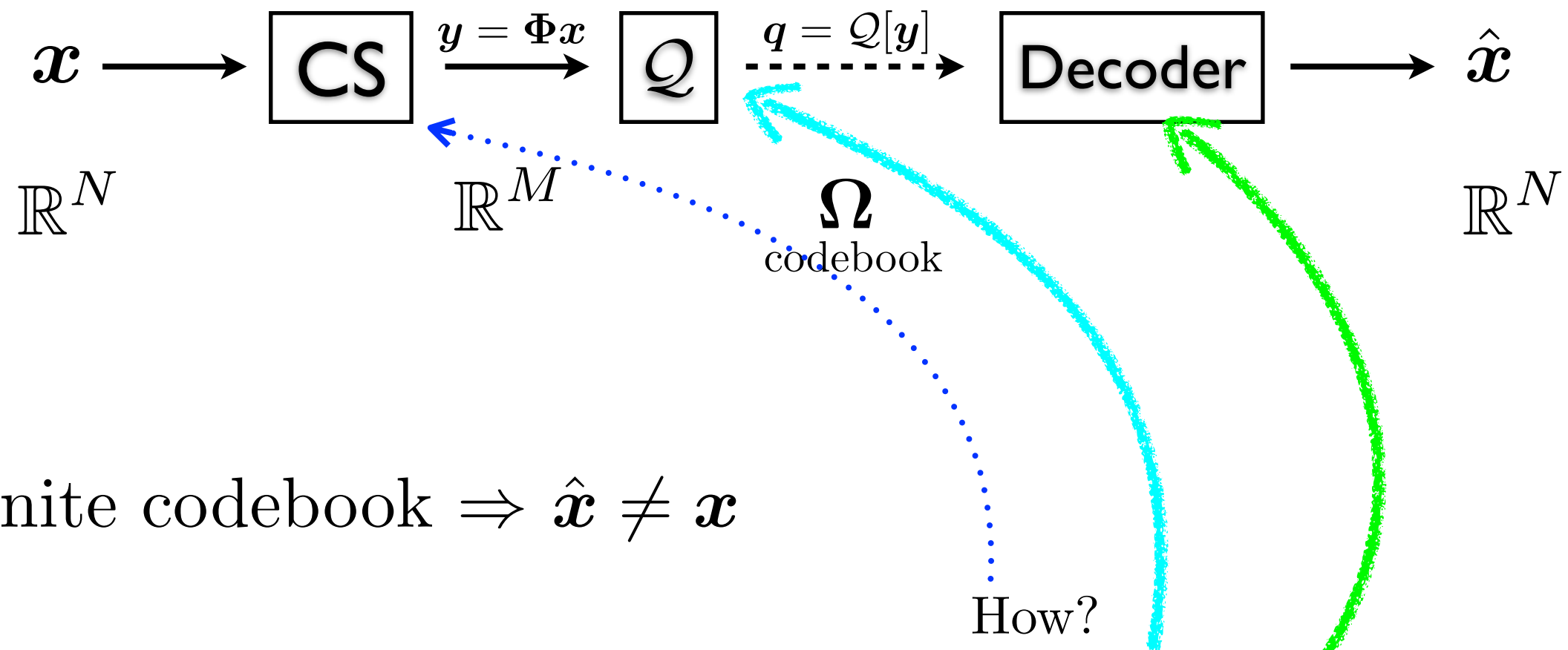
Finite codebook  $\Rightarrow \hat{x} \neq x$

(i.e., impossibility to encode continuous domain in a finite number of elements)

# QCS: a system view

With **no additional noise**:

e.g., basis pursuit,  
greedy methods, ...



Finite codebook  $\Rightarrow \hat{x} \neq x$

Objective: Minimize  $\|\hat{x} - x\|$   
given a certain number of:  
bits, measurements, or bits/meas.

How?

Where to act?

Change CS, Q or decoder?

Some of them? all?

## 2. Former QCS methods and performance limits

# Scalar quantization in CS

Turning measurements into bits  $\rightarrow$  scalar quantization

$$q_i = \mathcal{Q}[(\Phi \mathbf{x})_i] = \mathcal{Q}[\langle \phi_i, \mathbf{x} \rangle] \in \Omega \subset \mathbb{R}$$
$$\mathbf{q} = \mathcal{Q}[\Phi \mathbf{x}] \in \Omega = \Omega^M,$$

Important points:

- Definition of  $\Phi$  independent of  $M$  (e.g.,  $\Phi_{ij} \sim_{\text{iid}} \mathcal{N}(0, 1)$ )  
 $\rightarrow$  preserves measurement dynamic!
- $B$  bits per measurement
- Total bit budget:  $R = BM$
- No further encoding (e.g., entropic)

# Former solution (Candès, Tao, ...)

- Quantization is like a noise

$$q = \mathcal{Q}[\Phi x] = \Phi x + n$$

quantization  
distortion



# Former solution (Candès, Tao, ...)

- Quantization is like a noise

$$\mathbf{q} = \mathcal{Q}[\Phi \mathbf{x}] = \Phi \mathbf{x} + \mathbf{n}$$

and CS is robust (e.g., with *basis pursuit denoise*)

$$\hat{\mathbf{x}} = \underset{\mathbf{u} \in \mathbb{R}^N}{\operatorname{argmin}} \|\mathbf{u}\|_1 \text{ s.t. } \|\Phi \mathbf{u} - \mathbf{q}\| \leq \epsilon \quad (\text{BPDN})$$

$\ell_2 - \ell_1$  instance optimality:

If  $\|\mathbf{n}\| \leq \epsilon$  and  $\frac{1}{\sqrt{M}}\Phi$  is RIP( $\delta, 2K$ ) with  $\delta \leq \sqrt{2} - 1$ , then

How to find it?

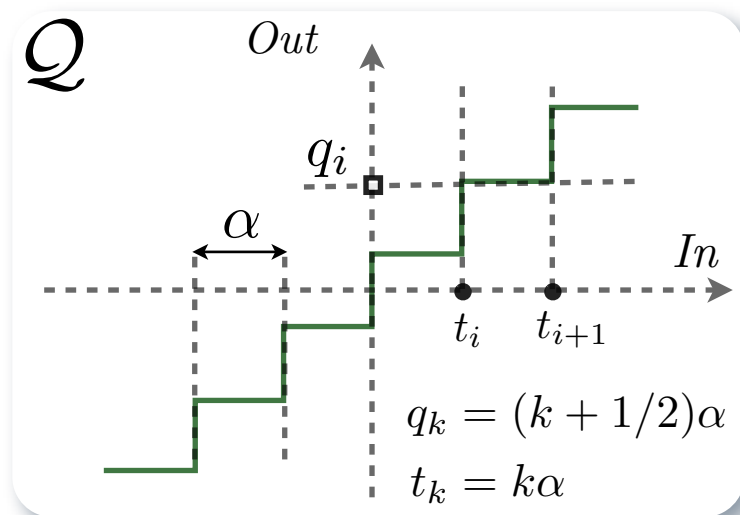
$$\|\hat{\mathbf{x}} - \mathbf{x}\| \leq C \frac{\epsilon}{\sqrt{M}} + D e_0(K),$$

for some  $C, D > 0$  and  $e_0(K) = \|\mathbf{x} - \mathbf{x}_K\|_1 / \sqrt{K}$ .

# Former solution (Candès, Tao, ...)

1. For uniform quantization, by construction:

€?



$$\begin{aligned} n_i &= Q[(\Phi \mathbf{x})_i] - (\Phi \mathbf{x})_i \\ &\in q_{k_i} - \mathcal{R}_{k_i} = [-\alpha/2, \alpha/2] \\ &\Rightarrow \|\mathbf{n}\|_{\infty} \leq \alpha/2 \end{aligned}$$

$$\Rightarrow \|\mathbf{n}\|^2 \leq M \|\mathbf{n}\|_{\infty}^2 \leq M \alpha^2 / 4$$

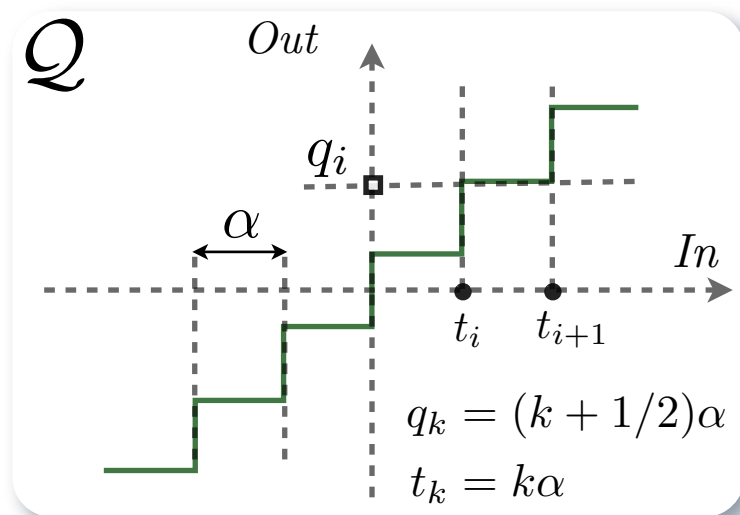
and plug this upper bound in BPDN

can be improved!

# Former solution (Candès, Tao, ...)

2. For uniform quantization, **uniform model!**

€?

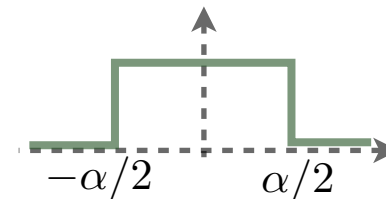


$$n_i = Q[(\Phi x)_i] - (\Phi x)_i$$

$$\in q_{k_i} - \mathcal{R}_{k_i} = [-\alpha/2, \alpha/2]$$

$$\sim_{\text{iid}} \text{Uniform}([- \alpha/2, \alpha/2])$$

(HRA - high resolution assumption)



$$\Rightarrow \mathbb{E}|n_i|^2 = \alpha^2/12$$

$$\Rightarrow \|\mathbf{n}\|^2 \leq \mathbb{E}\|\mathbf{n}\|^2 + \kappa \sqrt{\text{Var}\|\mathbf{n}\|^2} \quad (\text{Chernoff-Hoeffding, bounded RVs})$$

$$\leq M \frac{\alpha^2}{12} + \kappa \sqrt{M} \frac{\alpha^2}{6\sqrt{5}} = \epsilon_2^2 \simeq M \frac{\alpha^2}{12}$$

$$\text{with } \Pr > 1 - e^{-2\kappa^2}$$

and plug this upper bound in BPDN

# Former solution (Candès, Tao, ...)

- Therefore, from BPDN  $\ell_2 - \ell_1$  instance optimality:

$$\|\hat{\mathbf{x}} - \mathbf{x}\| \lesssim C \alpha + D e_0(K), \quad \text{for } C, D > 0$$

(for BPDN with  $\epsilon_2$ , under prev. cond.)

- Assuming :

- bounded dynamics:  $\|\Phi \mathbf{x}\|_\infty = \max_j |(\Phi \mathbf{x})_i| \leq \rho$  (e.g., by discarding saturation)  
(see later)
- $B$  bits per measurements  $\Rightarrow \alpha \simeq \rho 2^{1-B}$

$$\Rightarrow \text{BPDN RMSE} \lesssim C' 2^{-B} + D e_0(K) \quad \text{for } C', D > 0$$

as soon as RIP holds:  $M = O(K \log N/K)$

- Equivalently: BPDN RMSE  $\simeq O(2^{-R/M}) + e_0(K)$   
for a rate  $R = BM$  bits (total "bid budget" for all meas.)

# RMSE Lower bound?

- ▶ Let a fixed  $K$ -sparse  $\mathbf{x} \in \mathbb{R}^N$
- ▶ Oracle: you know  $T = \text{supp } \mathbf{x}$
- ▶ Noisy measurements (random noise):

Given  $\Phi \in \mathbb{R}^{M \times N}$  with  $\Phi_{ij} \sim_{\text{iid}} N(0, 1)$

$$\mathbf{y} = \Phi_T \mathbf{x} + \mathbf{n}, \quad \text{with } \mathbb{E} \mathbf{n} \mathbf{n}^T = \sigma^2 \mathbf{Id}_{M \times M}$$

- ▶ Assume:  $\frac{1}{\sqrt{M}} \Phi$  is RIP( $K, \delta_K$ ) and RIP( $1, \delta_1$ )
- ▶ Compute LS solution:  $\hat{\mathbf{x}}_T = \Phi_T^\dagger \mathbf{y} = (\Phi_T^* \Phi_T)^{-1} \Phi_T^* \mathbf{y}$   
 $\hat{\mathbf{x}}_{T^c} = 0$   
pseudo-inverse

- ▶ Then:  $\text{MSE} = \mathbb{E}_{\mathbf{n}} \|\mathbf{x} - \hat{\mathbf{x}}\|^2 \geq r^{-1} \sigma^2 \left( \frac{1 - \delta_1}{1 + \delta_K} \right)$   
for oversampling factor  $r = M/K$

(as for BPDN)  
 $\& \text{MSE} \leq \frac{1}{1 - \delta_K} \sigma^2$   
from [Needell, Tropp, 08]

- ▶ for QCS:  $\Rightarrow \text{RMSE} = \Omega(r^{-1/2} 2^{-B})$   $\& \text{RMSE} = O(2^{-B})$



# 3. Consistent Reconstructions

# Consistent reconstructions in CS?

- ▶ **Problem in previous case:** if  $\hat{\mathbf{x}}$  solution of BPDN,
- ▶ no **Quantization Consistency** (QC):  $\mathcal{Q}[\Phi\hat{\mathbf{x}}] \neq \mathcal{Q}[\Phi\mathbf{x}]$

$$\|\Phi\hat{\mathbf{x}} - \mathcal{Q}[\Phi\mathbf{x}]\| \leq \epsilon_2 \not\Rightarrow \mathcal{Q}[\Phi\hat{\mathbf{x}}] = \mathcal{Q}[\Phi\mathbf{x}]$$

(from BPDN constraint)

$\Rightarrow$  sensing information is fully not exploited!

- ▶  $\ell_2$  constraint  $\approx$  Gaussian distribution (MAP - cond. log. lik.)
- ▶ But why looking for consistency ?

**Proposition (Goyal, Vetterli, Thao, 98)** *If  $T$  is known (with  $|T| = K$ ), the best decoder  $\text{Dec}()$  provides a  $\hat{\mathbf{x}} = \text{Dec}(\mathbf{y}, \Phi)$  such that:*

$$\text{RMSE} = (\mathbb{E}\|\mathbf{x} - \hat{\mathbf{x}}\|^2)^{1/2} \gtrsim r^{-1}\alpha,$$

where  $\mathbb{E}$  is wrt a probability measure on  $\mathbf{x}_T$  in a bounded set  $\mathcal{S} \subset \mathbb{R}^K$ .

*This bound is achieved, at least, for  $\Phi_T = \text{DFT} \in \mathbb{R}^{M \times K}$ , when  $\text{Dec}()$  is **consistent**.*





# In quest of consistency...

$$\ell_2 \rightarrow \ell_\infty$$

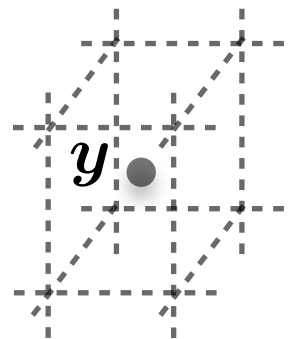
- Modify BPDN [W. Dai, O. Milenkovic, 09]

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^N} \|\mathbf{u}\|_1 \text{ s.t. } \mathcal{Q}[\Phi \mathbf{u}] = \mathbf{q}$$

+ modified greedy algo:  
“*subspace pursuit*”

$$\Leftrightarrow \Phi \mathbf{u} \in \mathcal{Q}^{-1}[\mathbf{q}]$$

convex set in  $\mathbb{R}^M$



$$\Leftrightarrow \|\Phi \mathbf{u} - \mathbf{q}\|_\infty \leq \alpha/2$$

(if uniform quant.)

$\exists$  numerical methods



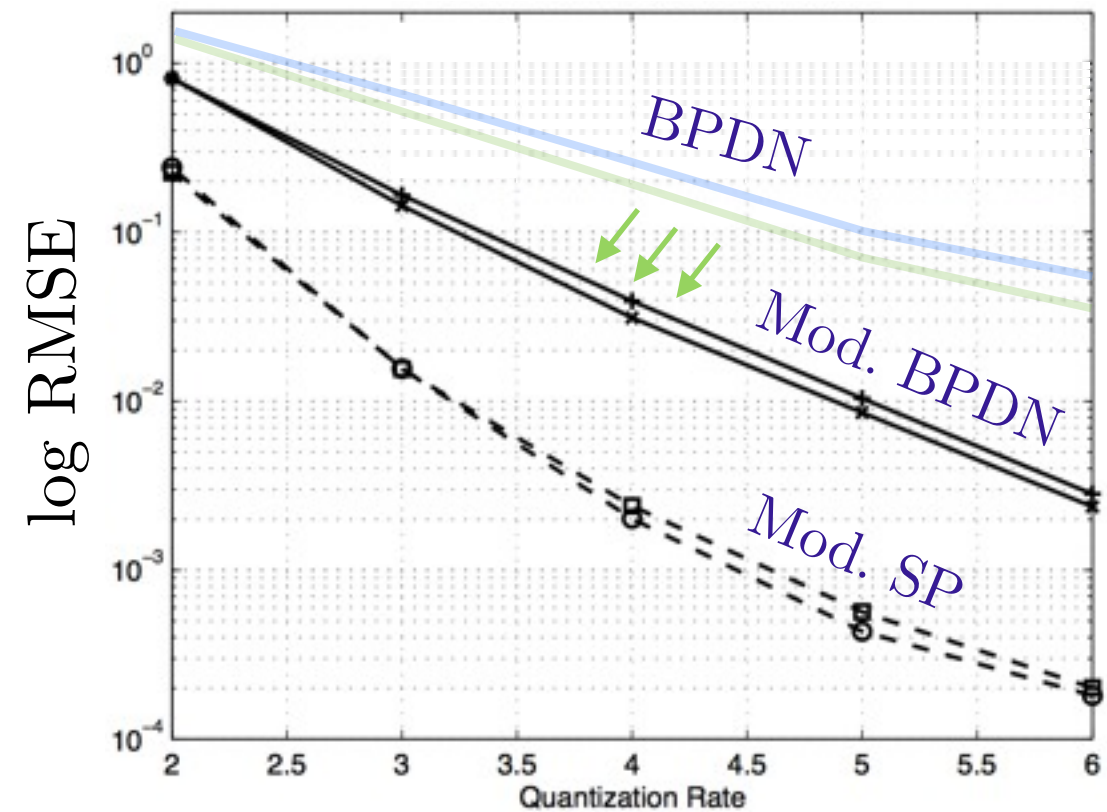
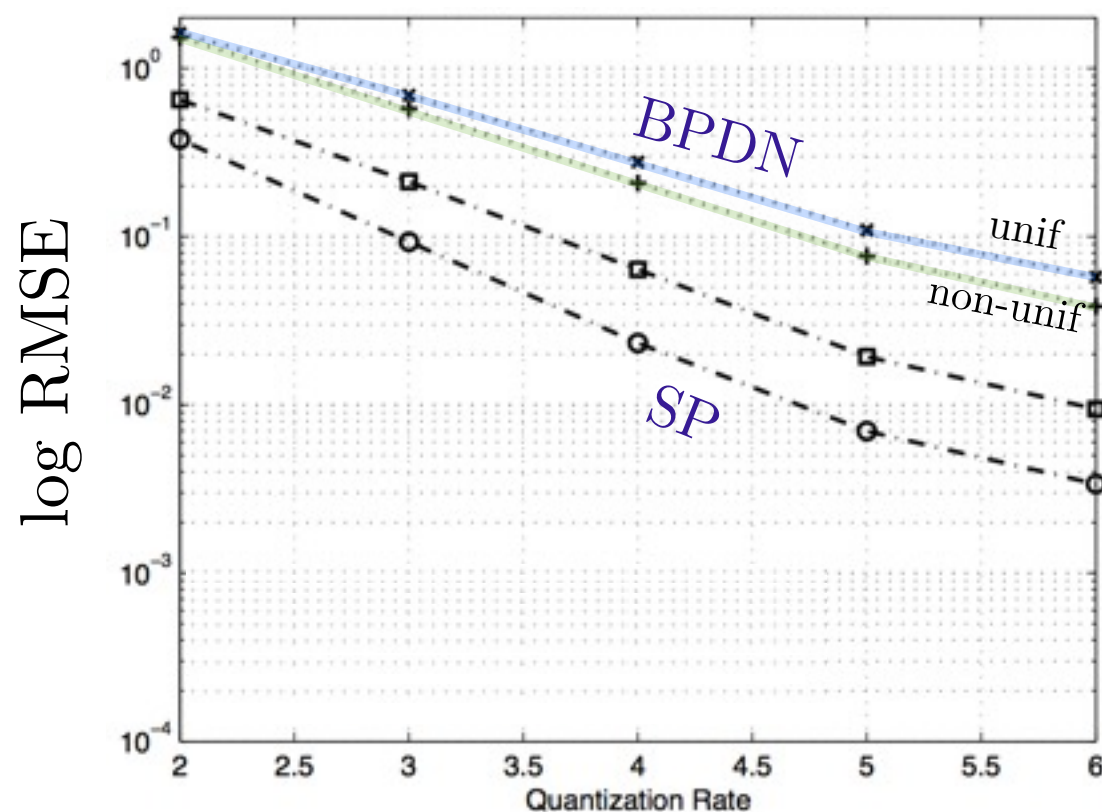
# In quest of consistency...

$$\ell_2 \rightarrow \ell_\infty$$

- Modify BPDN [W. Dai, O. Milenkovic, 09]

$$\hat{\mathbf{x}} = \underset{\mathbf{u} \in \mathbb{R}^N}{\operatorname{argmin}} \|\mathbf{u}\|_1 \text{ s.t. } \mathcal{Q}[\Phi \mathbf{u}] = \mathbf{q}$$

Simulations:  $M = 128, N = 256, K = 6, 1000$  trials  $\Rightarrow \lambda \simeq 20$



W. Dai, H. V. Pham, and O. Milenkovic, "Quantized Compressive Sensing", preprint, 2009

# Dequantizing CS?

[LJ, Hammond, Fadili, 2009, 2011]

- Distortion model:

$$\mathbf{q} = \mathcal{Q}[\Phi \mathbf{x}] = \Phi \mathbf{x} + \mathbf{n}, \quad n_i \sim U(-\frac{\alpha}{2}, \frac{\alpha}{2})$$

- Observation:  $\|\Phi \mathbf{x} - \mathbf{q}\|_\infty \leq \alpha/2$

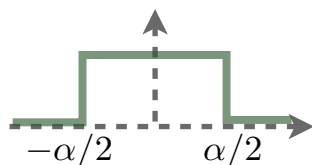
- Reconstruction: Generalizing BPDN with BPDQ

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{u} \in \mathbb{R}^N} \|\mathbf{u}\|_1 \text{ s.t. } \|\mathbf{q} - \Phi \mathbf{u}\|_p \leq \epsilon_p$$

Towards  $p = \infty$   
Related to GGD MAP

How to find it? again, uniform model:

$$\begin{aligned} n_i &= \mathcal{Q}[(\Phi \mathbf{x})_i] - (\Phi \mathbf{x})_i \\ &\in q_{k_i} - \mathcal{R}_{k_i} = [-\alpha/2, \alpha/2] \\ &\sim_{\text{iid}} \text{Uniform}([- \alpha/2, \alpha/2]) \end{aligned}$$



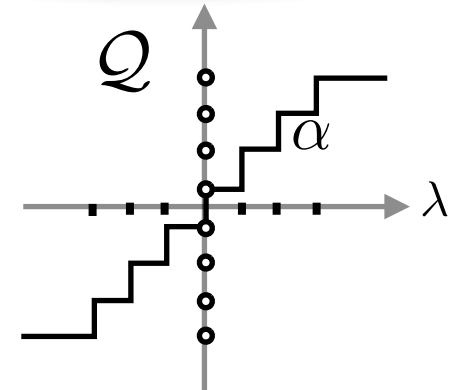
Estimating  $p^{\text{th}}$  moment:

$$\epsilon_p(\alpha) = \frac{\alpha}{2(p+1)^{1/p}} \left( M + \kappa(p+1)\sqrt{M} \right)^{1/p}$$

works with  $\Pr \geq 1 - e^{-2\kappa^2}$

Note:  $\epsilon_p(\alpha) \xrightarrow{p \rightarrow \infty} \frac{\alpha}{2} = \text{QC!}$

$$\ell_2 \rightarrow \ell_p \ (p \geq 2)$$



# Dequantizing CS?

[LJ, Hammond, Fadili, 2009, 2011]

- Distortion model:

$$\mathbf{q} = \mathcal{Q}[\Phi \mathbf{x}] = \Phi \mathbf{x} + \mathbf{n}, \quad n_i \sim U(-\frac{\alpha}{2}, \frac{\alpha}{2})$$

- Observation:  $\|\Phi \mathbf{x} - \mathbf{q}\|_\infty \leq \alpha/2$

- Reconstruction: Generalizing BPDN with BPDQ

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{u} \in \mathbb{R}^N} \|\mathbf{u}\|_1 \text{ s.t. } \|\mathbf{q} - \Phi \mathbf{u}\|_p \leq \epsilon_p$$

Towards  $p = \infty$   
Related to GGD MAP

BPDQ Stability ?

If  $\Phi$  is  $\text{RIP}_p$  of order  $K$ , *i.e.*,

$$\exists \mu_p > 0, \delta \in (0, 1),$$

$$\sqrt{1 - \delta} \|\mathbf{v}\|_2 \leq \frac{1}{\mu_p} \|\Phi \mathbf{v}\|_p \leq \sqrt{1 + \delta} \|\mathbf{v}\|_2,$$

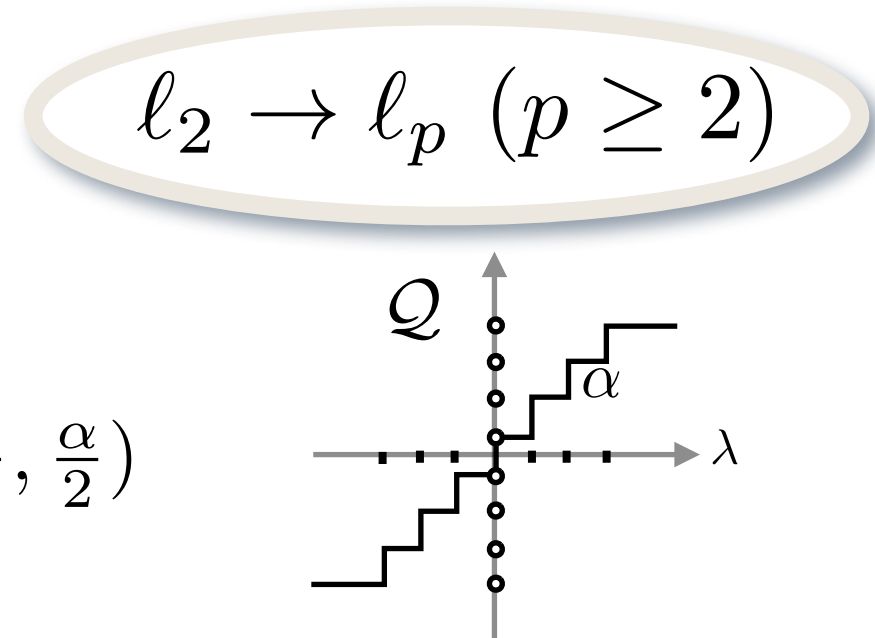
for all  $K$  sparse signals  $\mathbf{v}$ .

Gain over BPDN (for tight  $\epsilon_p(\alpha, M)$ )  
 $\Rightarrow \|\mathbf{x} - \hat{\mathbf{x}}\| = O(\epsilon_p / \mu_p)$

$$\Rightarrow \|\mathbf{x} - \hat{\mathbf{x}}\| = O(\alpha / \sqrt{p + 1})$$

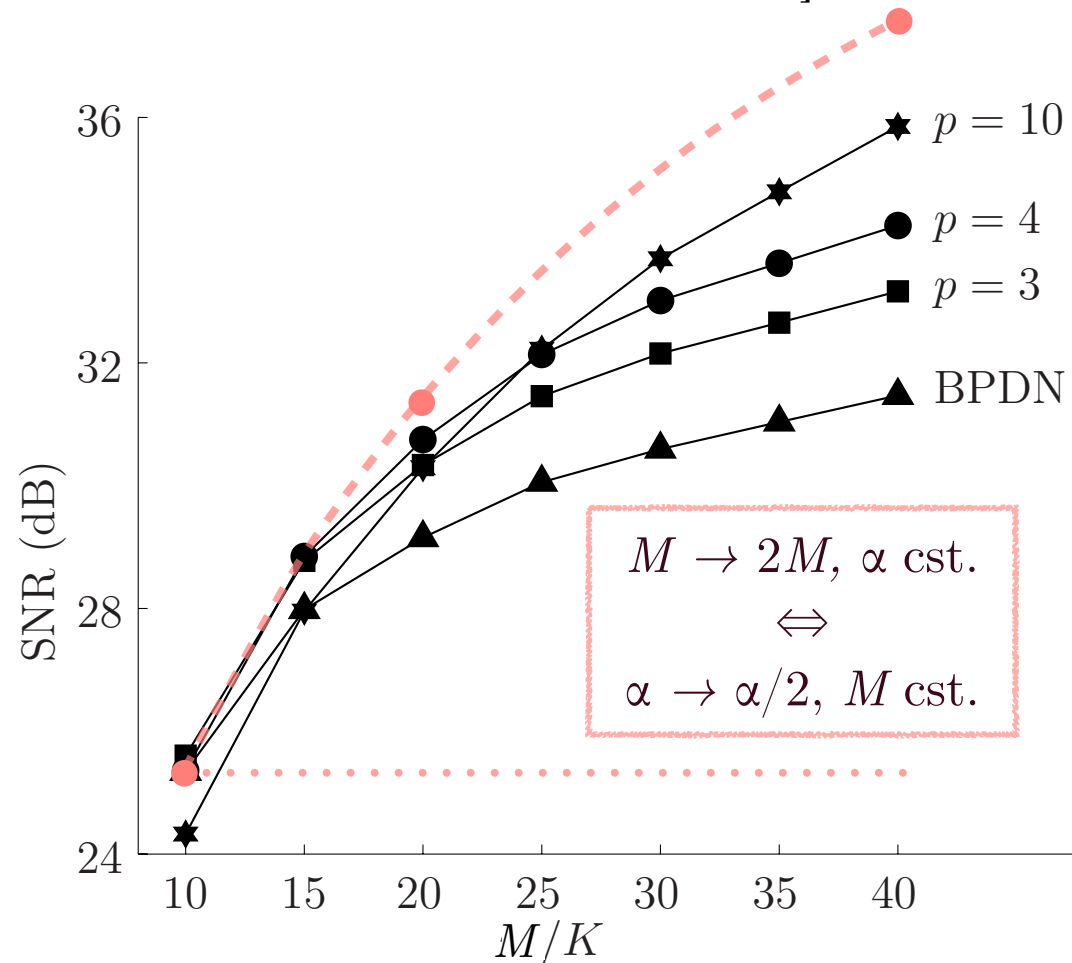
But no free lunch: for  $\Phi$  Gaussian  
 $M = O((K \log N / K)^{p/2})$

$\Rightarrow$  Another reading: limited range of valid  $p$  for a given  $M$  (and  $K$ )!



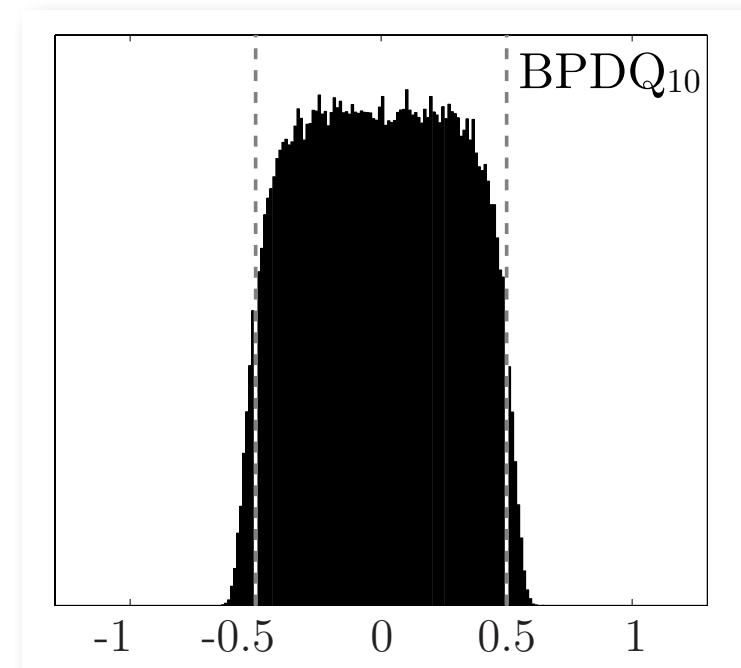
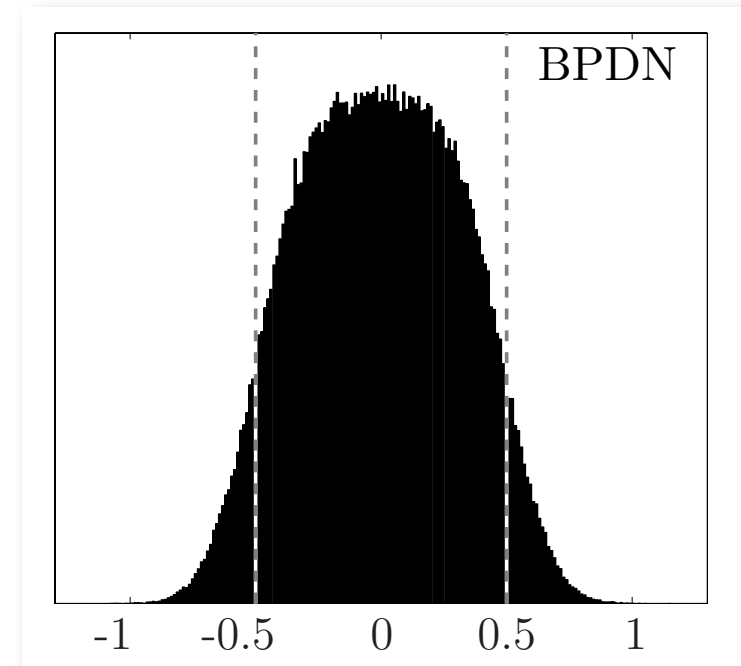
# Dequantizing CS?

[LJ, Hammond, Fadili, 2009, 2011]



- \*  $N=1024$ ,  $K=16$ , Gaussian  $\Phi$
- \* 500  $K$ -sparse (canonical basis)
- \* Non-zero components follow  $\mathcal{N}(0, 1)$
- \* Quantiz. bin width  $\alpha = \|\Phi \mathbf{x}\|_{\infty}/40$

Histograms of  
 $\alpha^{-1}(\mathbf{q} - \Phi \hat{\mathbf{x}})_i$

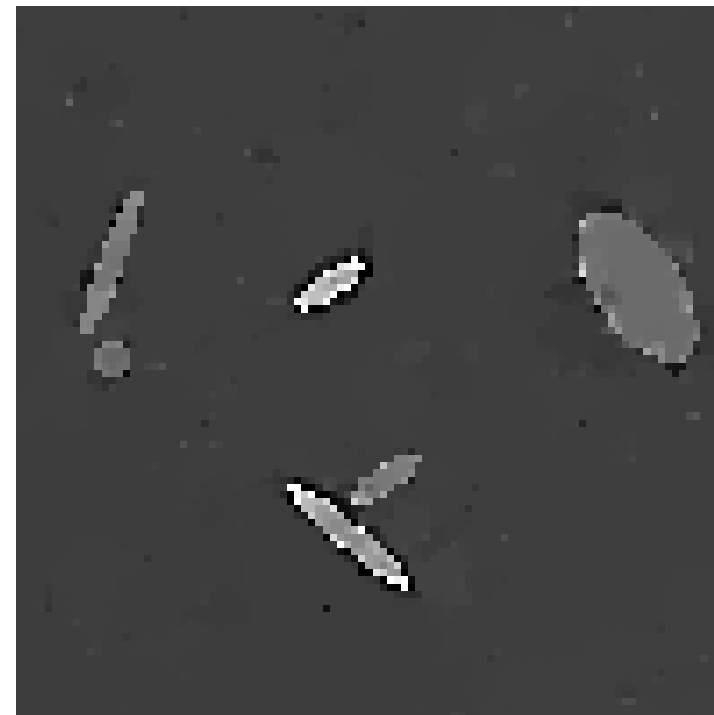
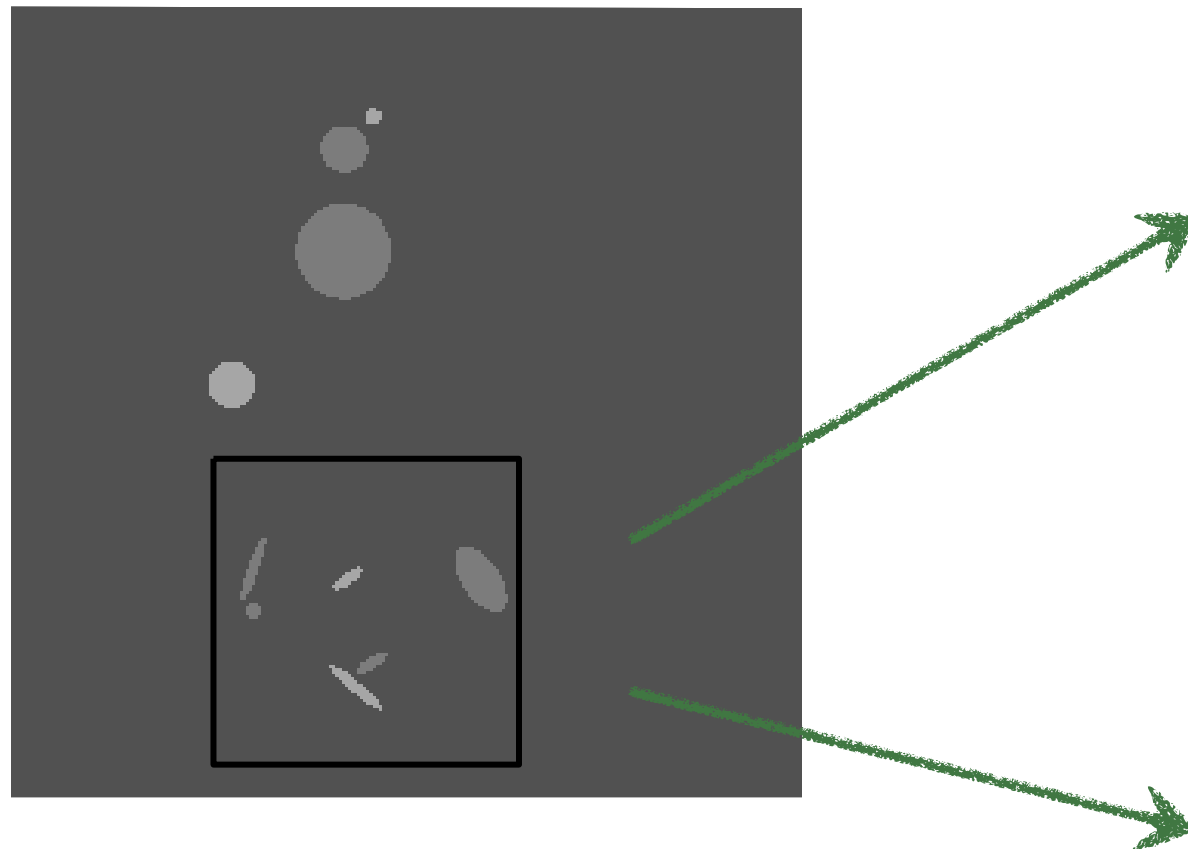


LJ, D. Hammond, J. Fadili "Dequantizing compressed sensing: When oversampling and non-gaussian constraints combine." Information Theory, IEEE Transactions on, 57(1), 559-571.

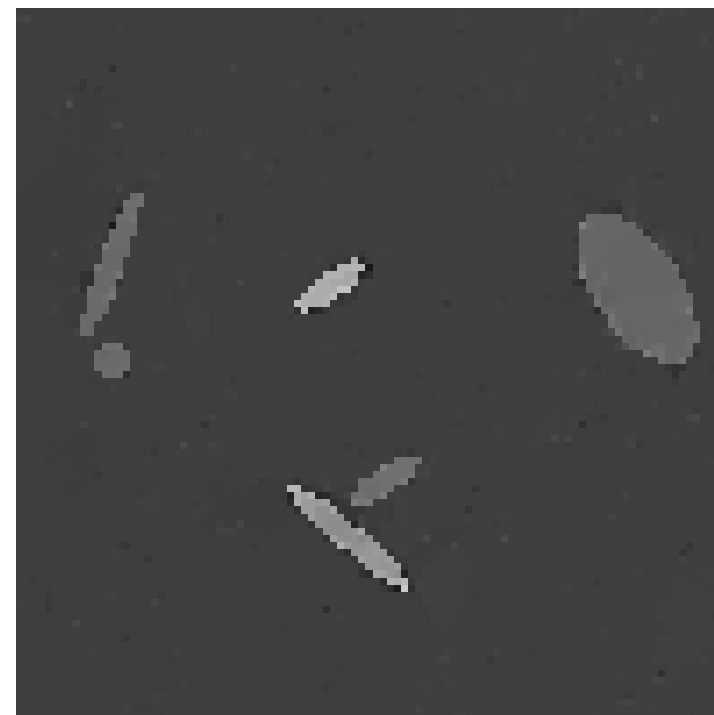
# Dequantizing CS?

[LJ, Hammond, Fadili, 2009, 2011]

**A bit outside the theory...**



BPDN-TV  
SNR: 8.96 dB



BPDQ<sub>10</sub>-TV  
SNR: 12.03 dB

- \* Synthetic Angiogram [Michael Lustig 07, SPARCO],
- \*  $\Phi$ : **Random Fourier Ensemble**
- \*  $N/M = 8$
- \* Decoder:  $\Delta_{TV,p}(y, \epsilon_p)$
- \* Quantiz. bin width = 50 (i.e. 12 bins)

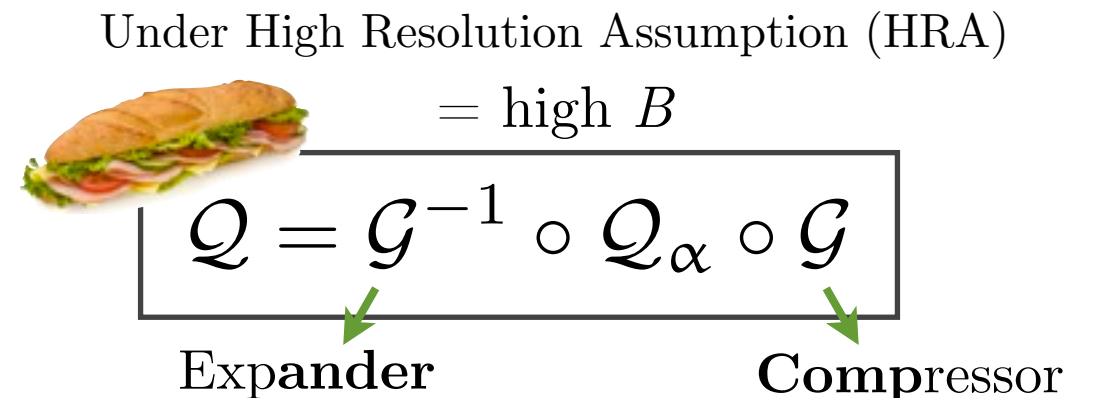
LJ, D. Hammond, J. Fadili “Dequantizing compressed sensing: When oversampling and non-gaussian constraints combine.” Information Theory, IEEE Transactions on, 57(1), 559-571.



# Non-uniform dequantization?

Possible!

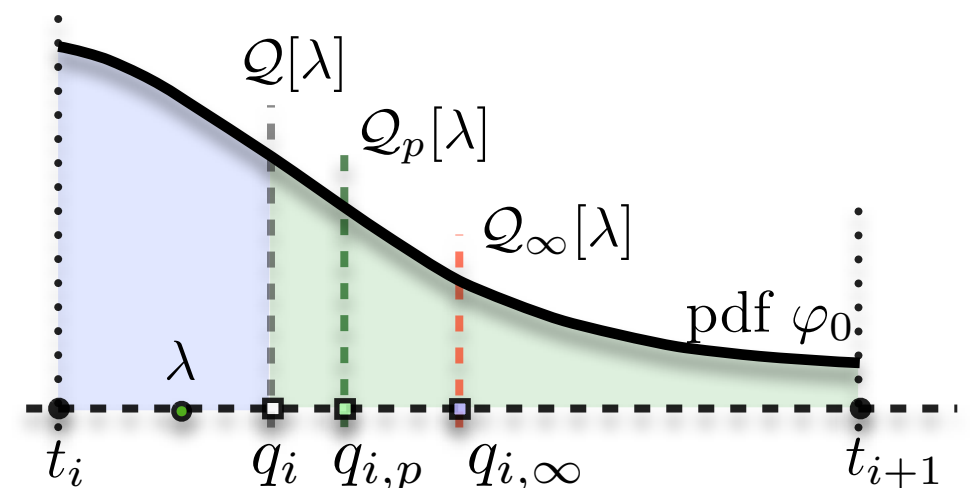
1. Use compander formalism:



2. Redefine  $\mathbf{q}$  : (post-sensing)

$$\mathbf{q} \rightarrow \mathbf{q}_p = \mathcal{Q}_p[\mathbf{q}] = \mathcal{Q}_p[\Phi \mathbf{x}]$$

with  $(q_p)_i$  minimizing  $p^{\text{th}}$  moment in each bin.



3. Reweight the bins:  $\|\cdot\|_p \rightarrow \|\text{diag}(\mathbf{w}) \cdot \cdot\|_p =: \|\cdot\|_{p,\mathbf{w}}$

$$\text{with: } w_i(p) := \mathcal{G}'((q_p)_i)^{(p-2)/p}$$

→ kind of noise stabilization operation (“equi- $p$ -distortion”)

4. Solve:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{u} \in \mathbb{R}^N} \|\mathbf{u}\|_1 \text{ s.t. } \|\mathbf{q} - \Phi \mathbf{u}\|_{p,\mathbf{w}} \leq \epsilon_p$$

# Non-uniform dequantization?

- Stability? Well ... need a more general RIP 😊

$$\text{RIP}(\ell_{p,\mathbf{w}}, \ell_2 | K, \delta, \mu) \quad \begin{array}{l} \exists \mu > 0, \delta \in (0, 1) \\ \sqrt{1 - \delta} \|\mathbf{v}\|_2 \leq \frac{1}{\mu} \|\Phi \mathbf{v}\|_{p,\mathbf{w}} \leq \sqrt{1 + \delta} \|\mathbf{v}\|_2 \\ \text{for all } K \text{ sparse signals } \mathbf{v}. \end{array}$$

$$\Rightarrow M = O((\theta(\mathbf{w}) K \log N/K)^{p/2})$$

$$\text{with: } \theta(\mathbf{w}) \simeq M^{2/p} \|\mathbf{w}\|_\infty^2 / \|\mathbf{w}\|_p^2 \quad (= 1 \text{ if } w_i = \text{cst})$$

Then,

Given  $\mathcal{Q}_p[\cdot]$ ,  $\mathbf{w}(p) \in \mathbb{R}_+^M$  and  $\epsilon_p$  as before, GBPDN robustness provides:

$$\|\mathbf{x}^* - \mathbf{x}\| \underset{B, M}{\lesssim} 4 c' \frac{2^{-B}}{\sqrt{p+1}} + 2 e_0(K),$$

$$\text{with } c' = (9/8)(e\pi/3)^{1/2} < 1.8981.$$

# 4. Sigma-Delta quantization in CS



# Context:

- ▶ **Former attempts:** (see prev. slides)

CS + uniform scalar quantization (or pulse code modulation - PCM)

For  $K$ -sparse signals:  $\|\mathcal{Q}_\alpha[\Phi\mathbf{x}] - \Phi\mathbf{x}\|_2 \leq c\sqrt{M}\alpha \Rightarrow \|\mathbf{x}^* - \mathbf{x}\| \leq C\alpha$  (with RIP)

and for high  $\lambda$ ,  $\|\mathcal{Q}_\alpha[\Phi\mathbf{x}] - \Phi\mathbf{x}\|_p \leq cM^{1/p}\alpha \Rightarrow \|\mathbf{x}^* - \mathbf{x}\| \leq C\alpha/\sqrt{p+1}$  (with  $\text{RIP}_p$ )

- ▶ **No improvement if  $M$  increases!**

- ▶ **Can we do better?**

Can we have  $\|\mathbf{x}^* - \mathbf{x}\| \leq O(r^{-s}\alpha)$  for some  $s > 0$  ?

- ▶ **Staying with PCM,**  $s \leq 1$  (Goyal-Vetterli-Thao lower bound)

- ▶ **Solution: replacing PCM by  $\Sigma\Delta$  quantization!**

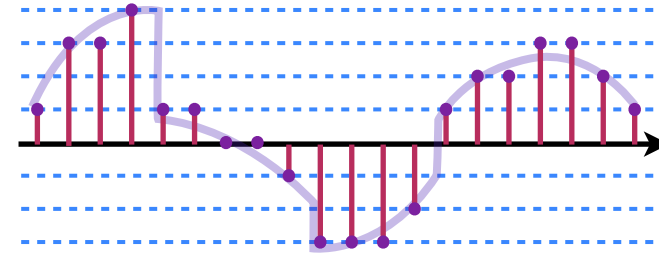
[S. Güntürk, A. Powell, R. Saab, Ö. Yılmaz]

# $\Sigma\Delta$ quantization (reminder)

- PCM: Signal sensing + unif. quantization (step  $\alpha$ )

$$\mathbf{x} \in \mathbb{R}^K \rightarrow \mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^M$$

$$\mathbf{q} = \mathcal{Q}_{\text{PCM}}[\mathbf{y}] \text{ with}$$



$$q_k = \mathcal{Q}_{\text{PCM}}[y_k] := \underset{u \in \alpha\mathbb{Z}}{\operatorname{argmin}} |y_k - u|, \quad 1 \leq k \leq M$$

Let  $\mathbf{A}^\#$ , a left inverse of  $\mathbf{A}$ , *i.e.*,  $\mathbf{A}^\# \mathbf{A} = \mathbf{Id}$ .

$$\hat{\mathbf{x}} := \mathbf{A}^\# \mathbf{q} \Rightarrow \|\mathbf{x} - \hat{\mathbf{x}}\| = \underbrace{\|\mathbf{A}^\# (\mathbf{y} - \mathbf{q})\|}_{\text{quant. noise}}$$

Taking (Moore-Penrose) pseudo-inverse:  $\mathbf{A}^\# = \mathbf{A}^\dagger = (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^*$   
(or canonical dual of the frame  $\mathbf{A}$ )

minimize  $\|\mathbf{A}^\# (\mathbf{y} - \mathbf{q})\|$ ! (least square solution)

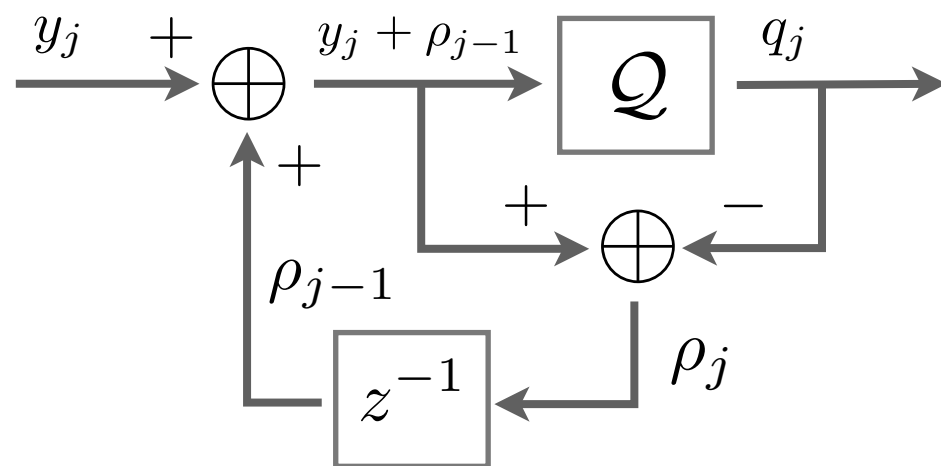
- In CS, this could be used if signal support was known (see before)

# $\Sigma\Delta$ quantization (reminder)

- ▶  $\Sigma\Delta \equiv$  noise shaping! Enjoy of:
  - ▶ **freedom** to pick  $\mathbf{q} \in \alpha\mathbb{Z}^M$
  - ▶ **freedom** to take another left inverse  $\mathbf{A}^\#$
- ▶ 1<sup>st</sup> order  $\Sigma\Delta$ : (in 1-D) Quantizing the sequence  $\{y_j : j \geq 0\}$

Use of state variables  $\{\rho_j\}$  (1-step memory):

$$\begin{aligned} \text{find } q_j: \quad q_j &= \mathcal{Q}_{\Sigma\Delta}^{(1)}[y_j] := \operatorname{argmin}_{u \in \alpha\mathbb{Z}} |\rho_{j-1} + y_j - u| = \mathcal{Q}_{\text{PCM}}[\rho_{j-1} + y_j] \\ \text{find } \rho_j: \quad (\Delta\rho)_j &= \rho_j - \rho_{j-1} = y_j - q_j \quad (\text{difference eq.}) \end{aligned}$$



with:  $|\rho_j| \leq \alpha$   
 $|y_j - q_j| \leq 2\alpha$

bigger than  $\alpha$  but still  $O(\alpha)$

# $\Sigma\Delta$ quantization (reminder)

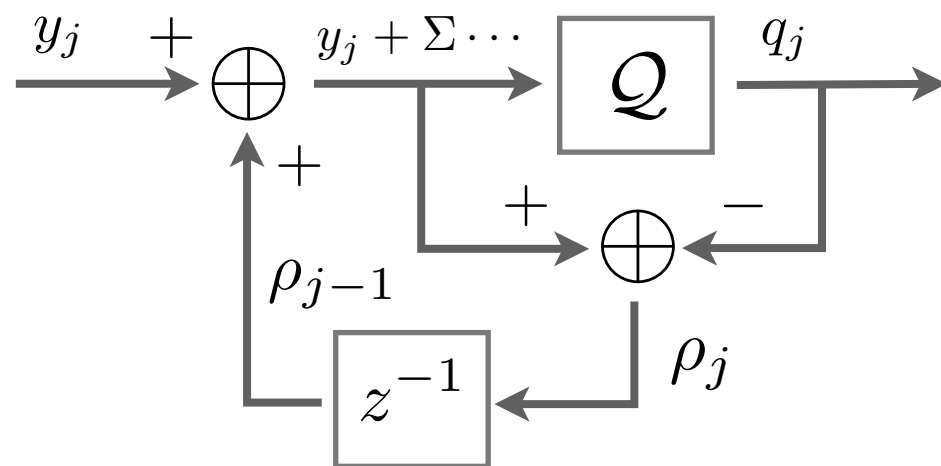
- ▶  $\Sigma\Delta \equiv$  noise shaping! Enjoy of:
- ▶ **freedom** to pick  $\mathbf{q} \in \alpha\mathbb{Z}^M$
- ▶ **freedom** to take another left inverse  $\mathbf{A}^\#$
- ▶  $s^{\text{th}}$  order  $\Sigma\Delta$ : (in 1-D) Quantizing the sequence  $\{y_j : j \geq 0\}$   
Use of state variables  $\{\rho_j\}$  (s-step memory):

Remark:

PCM is  
0<sup>th</sup> order  $\Sigma\Delta$

$$\text{find } q_j: \quad q_j = \mathcal{Q}_{\Sigma\Delta}^{(s)}[y_j] := \operatorname{argmin}_{u \in \alpha\mathbb{Z}} \left| \sum_{i=1}^s (-1)^{i-1} \binom{s}{i} \rho_{j-i} + y_j - u \right|$$

$$\text{find } \rho_j: \quad (\Delta^s \rho)_j = y_j - q_j \quad (s^{\text{th}} \text{ order difference eq.})$$



with:  $|\rho_j| \leq \alpha$

$$|y_j - q_j| \leq 2^{s-1} \alpha$$

bigger than  $\alpha$  but still  $O(\alpha)$

# $\Sigma\Delta$ quantization (reminder)

- ▶  $\Sigma\Delta \equiv$  noise shaping! Enjoy of:
  - ▶ freedom to pick  $\mathbf{q} \in \alpha\mathbb{Z}^M$
  - ▶ freedom to take another left inverse  $\mathbf{A}^\#$
- ▶  $s^{\text{th}}$  order  $\Sigma\Delta$ :

Most important fact:  $(\Delta^s \rho)_j = y_j - q_j \Leftrightarrow \mathbf{D}^s \boldsymbol{\rho} = \mathbf{y} - \mathbf{q}$

$$\hat{\mathbf{x}} := \mathbf{A}^\# \mathbf{q} \Rightarrow \|\mathbf{x} - \hat{\mathbf{x}}\| = \|\mathbf{A}^\# \mathbf{D}^s (\mathbf{y} - \mathbf{q})\|$$

$$\text{minimize } \|\mathbf{A}^\# \mathbf{D}^s (\mathbf{y} - \mathbf{q})\|!$$

~~Pseudo-inverse~~

$$\mathbf{A}^\dagger = (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^*$$

Sobolev duals

$$\mathbf{A}_{\text{sob},s} = (\mathbf{D}^{-s} \mathbf{A})^\dagger \mathbf{D}^{-s}$$

# $\Sigma\Delta$ quantization (reminder)

- ▶  $\Sigma\Delta \equiv$  noise shaping! Enjoy of:
  - ▶ **freedom** to pick  $\mathbf{q} \in \alpha\mathbb{Z}^M$
  - ▶ **freedom** to take another left inverse  $\mathbf{A}^\#$
- ▶  $s^{\text{th}}$  order  $\Sigma\Delta$ :

Most important fact:  $(\Delta^s \rho)_j = y_j - q_j \Leftrightarrow \mathbf{D}^s \boldsymbol{\rho} = \mathbf{y} - \mathbf{q}$

$$\hat{\mathbf{x}} = \mathbf{A}_{\text{sob},s} \mathbf{q}$$

$$\mathbf{A}_{\text{sob},s} = (\mathbf{D}^{-s} \mathbf{A})^\dagger \mathbf{D}^{-s}$$

**Proposition** Let  $\mathbf{A} \in \mathbb{R}^{M \times K}$  with  $A_{ij} \sim_{\text{iid}} \mathcal{N}(0, 1)$ .

For any  $\kappa \in (0, 1)$ , if  $r := M/K \geq c(\log M)^{1/(1-\kappa)}$ , then with  $Pr > 1 - e^{-c' M/r^\kappa}$ ,

$$\|\hat{\mathbf{x}} - \mathbf{x}\| \leq C_s r^{-\kappa(s-\frac{1}{2})} \alpha,$$

for some  $c, c', C_s > 0$ .

proof: show that

$$\sigma_{\min}(\mathbf{D}^{-s} \mathbf{A}) > C'_s r^{\kappa(s-\frac{1}{2})} \sqrt{M}$$

# $\Sigma\Delta$ quantization in CS

$$\mathbf{x} \in \Sigma_K \subset \mathbb{R}^N \rightarrow \mathbf{y} = \Phi \mathbf{x} \in \mathbb{R}^M \rightarrow \mathbf{q} = \mathcal{Q}_{\Sigma\Delta}^{(s)}[\mathbf{y}]$$

$$\|\mathbf{y} - \mathbf{q}\| \leq 2^{s-1} \alpha \sqrt{M}$$

Two-steps procedure:

1. find the support  $T$  of  $\mathbf{x}$  : coarse approx. with BPDN
2. compute  $\hat{\mathbf{x}} := (\Phi_T)_{\text{sob},s} \mathbf{q} = (\mathbf{D}^{-s} \Phi_T)^\dagger \mathbf{D}^{-s} \mathbf{q}$

**Proposition** Let  $\Phi \in \mathbb{R}^{M \times K}$  with  $\Phi_{ij} \sim_{\text{iid}} \mathcal{N}(0, 1)$ . Suppose  $\kappa \in (0, 1)$  and  $r := M/K \geq c(\log M)^{1/(1-\kappa)}$  for  $c > 0$ . Then,  $\exists c', C, C_s > 0$  such that, with  $Pr > 1 - e^{-c' M/r^\kappa}$ , for all  $\mathbf{x} \in \Sigma_K$  s.t.  $\min_{i \in \text{supp } \mathbf{x}} |x_i| \geq C\alpha$ ,

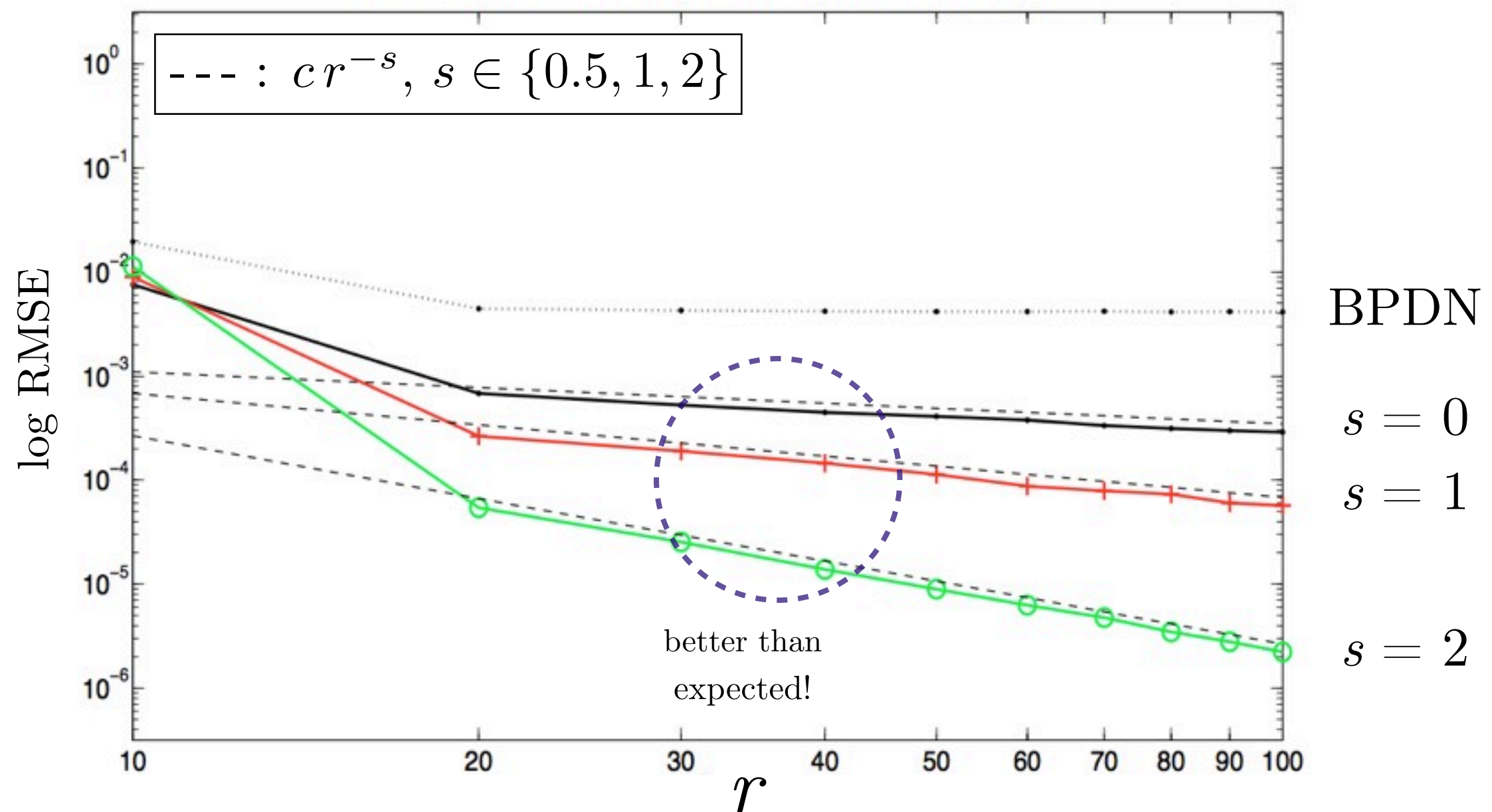
$$\|\hat{\mathbf{x}} - \mathbf{x}\| \leq C_s r^{-\kappa(s-\frac{1}{2})} \alpha.$$

proof: Union bound on any  $K$ -column subset of  $\Phi$   
+ proba having good support.



# $\Sigma\Delta$ quantization in CS (Simulations)

$M \in \{100, 200, \dots, 1000\}$ ,  $K = 10$  and 1000 trials ( $x_i \in \{0, \pm 1/\sqrt{K}\}$ ,  $\|\mathbf{x}\| \simeq 1$ ,  $\alpha = 10^{-2}$ )



Güntürk, C. S., Lammers, M., Powell, A. M., Saab, R., & Yilmaz, Ö. (2013). **Sobolev duals for random frames and  $\Sigma\Delta$  quantization of compressed sensing measurements.** Foundations of Computational Mathematics, 13(1), 1-36.



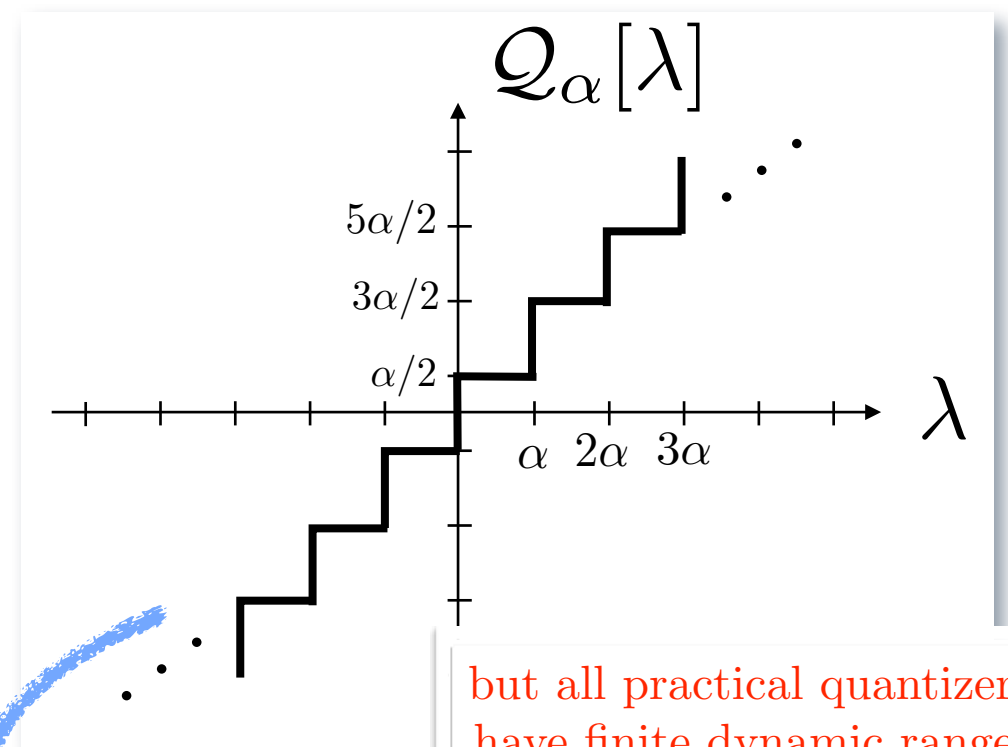
5. To saturate or not?  
And how much?

# Saturation phenomenon:

Uniform quantization:

- $\alpha$  quantization interval
- error per measurement bounded:

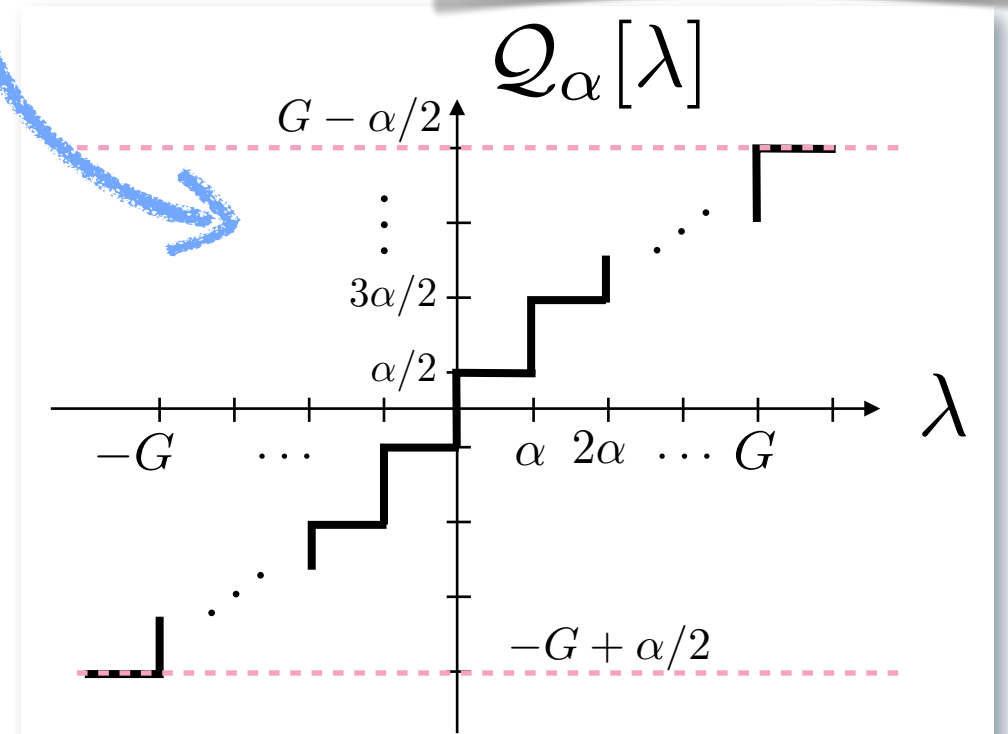
$$|\lambda - Q_\alpha[\lambda]| \leq \alpha/2$$



Finite Dynamic Range Quantization:

- $G$  “saturation level”
- $B$  bit rate (bits per measurement)
- quantization interval is  $\alpha = 2^{-B+1}G$
- measurements above  $G$  saturate
- saturation error is *unbounded*

*CS guarantees are for bounded errors only!*

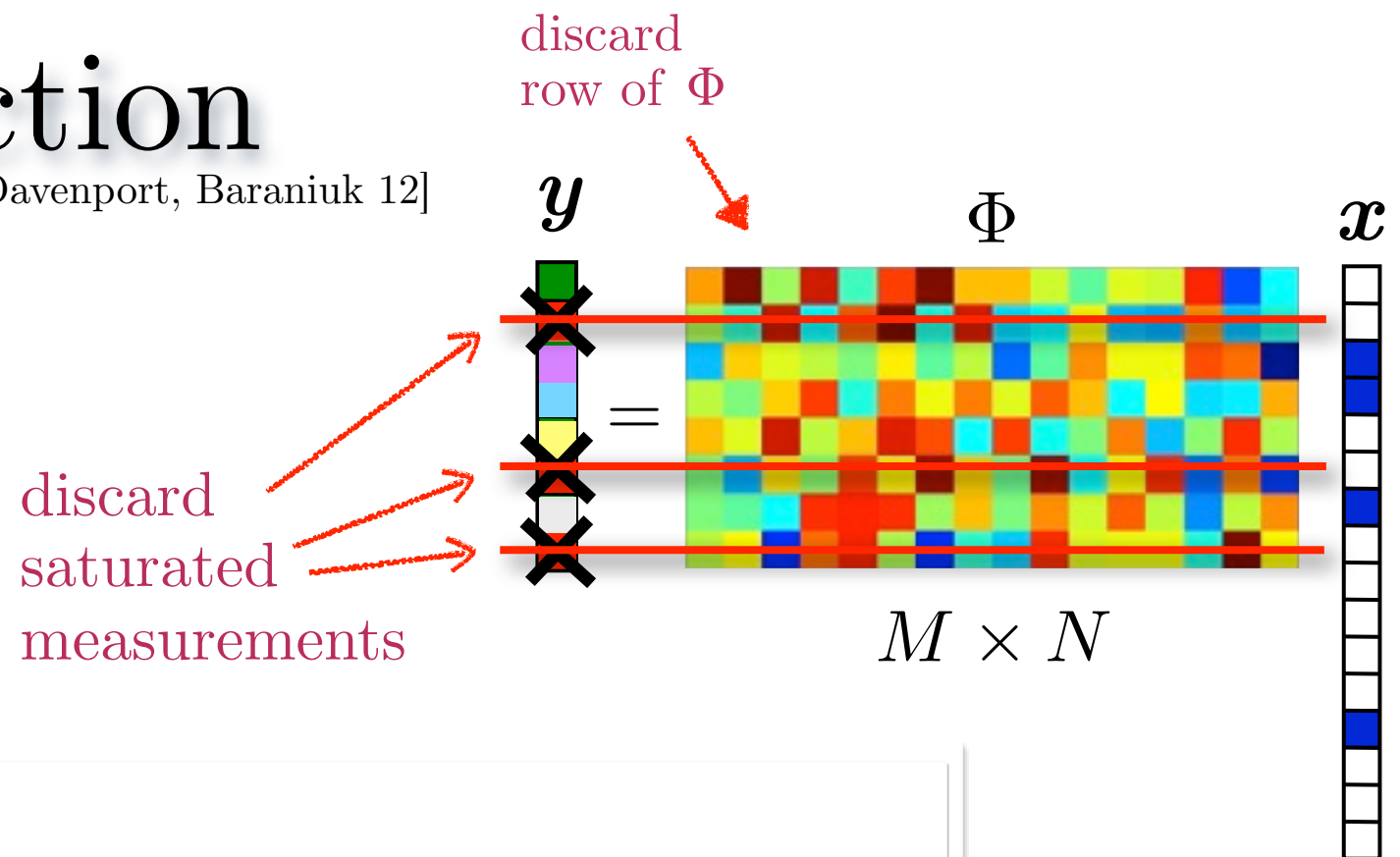


# Democracy in Action

[Laska, Boufounos, Davenport, Baraniuk 12]

## (i) Saturation Rejection:

Simply discard saturated measurements and corresponding rows of  $\Phi$



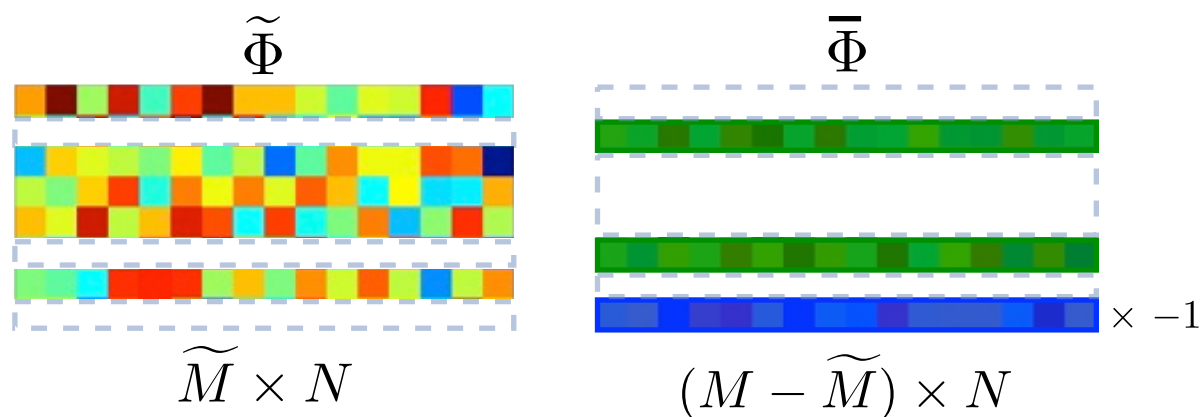
*“democratic measurements”*

each measurement has roughly same amount of information

RIP holds on row subsets of  $\Phi$

## (ii) Saturation Consistency:

Include saturated measurements as inequality constraint



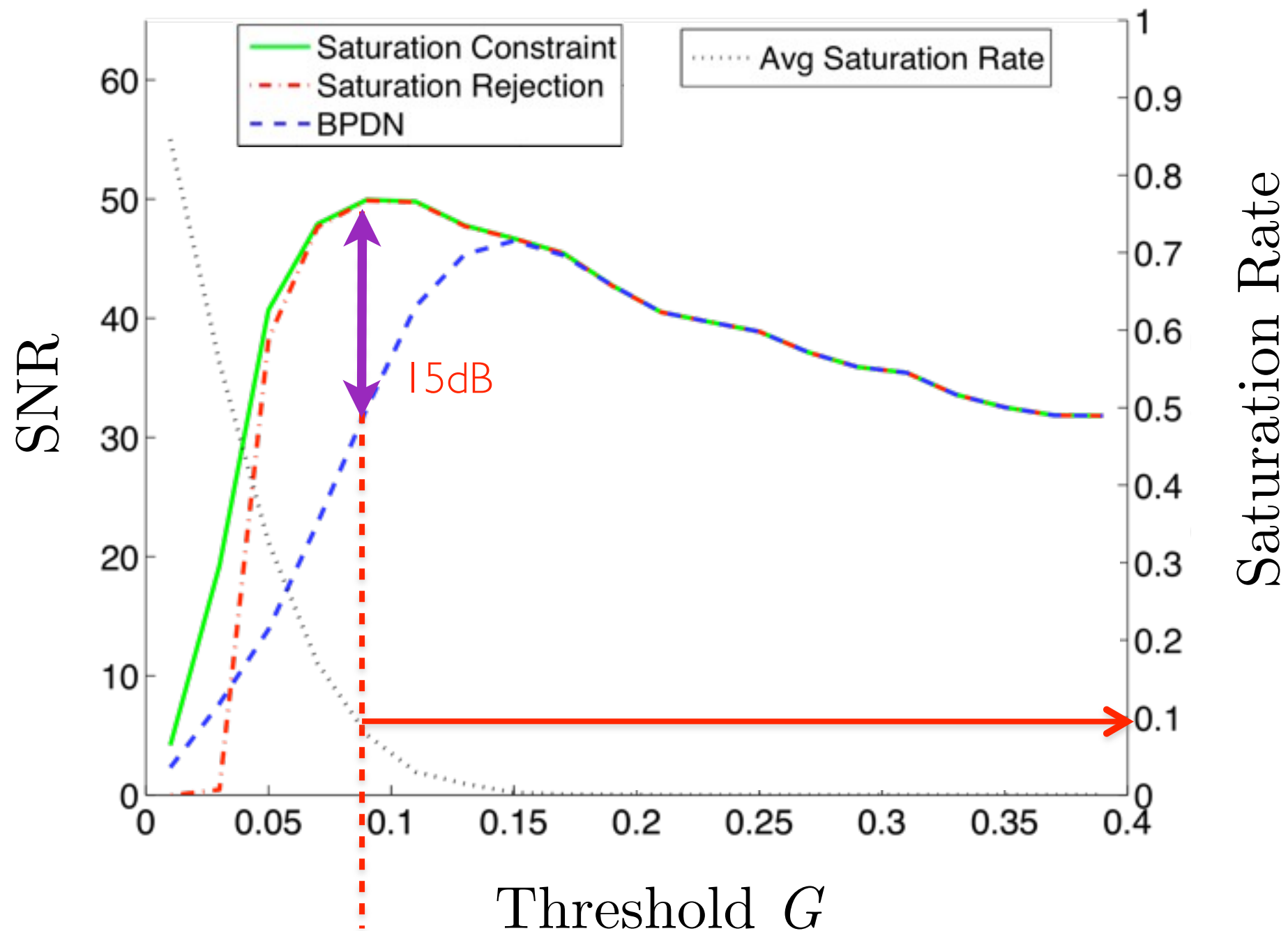
$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \|\tilde{\Phi}\mathbf{x} - \tilde{\mathbf{y}}\|_2 < \epsilon$$

Measurement error  
term (quantization)

Saturation consistency  
constraint

and  $\bar{\Phi}\mathbf{x} \geq G \cdot \mathbf{1}$

# Experimental Results

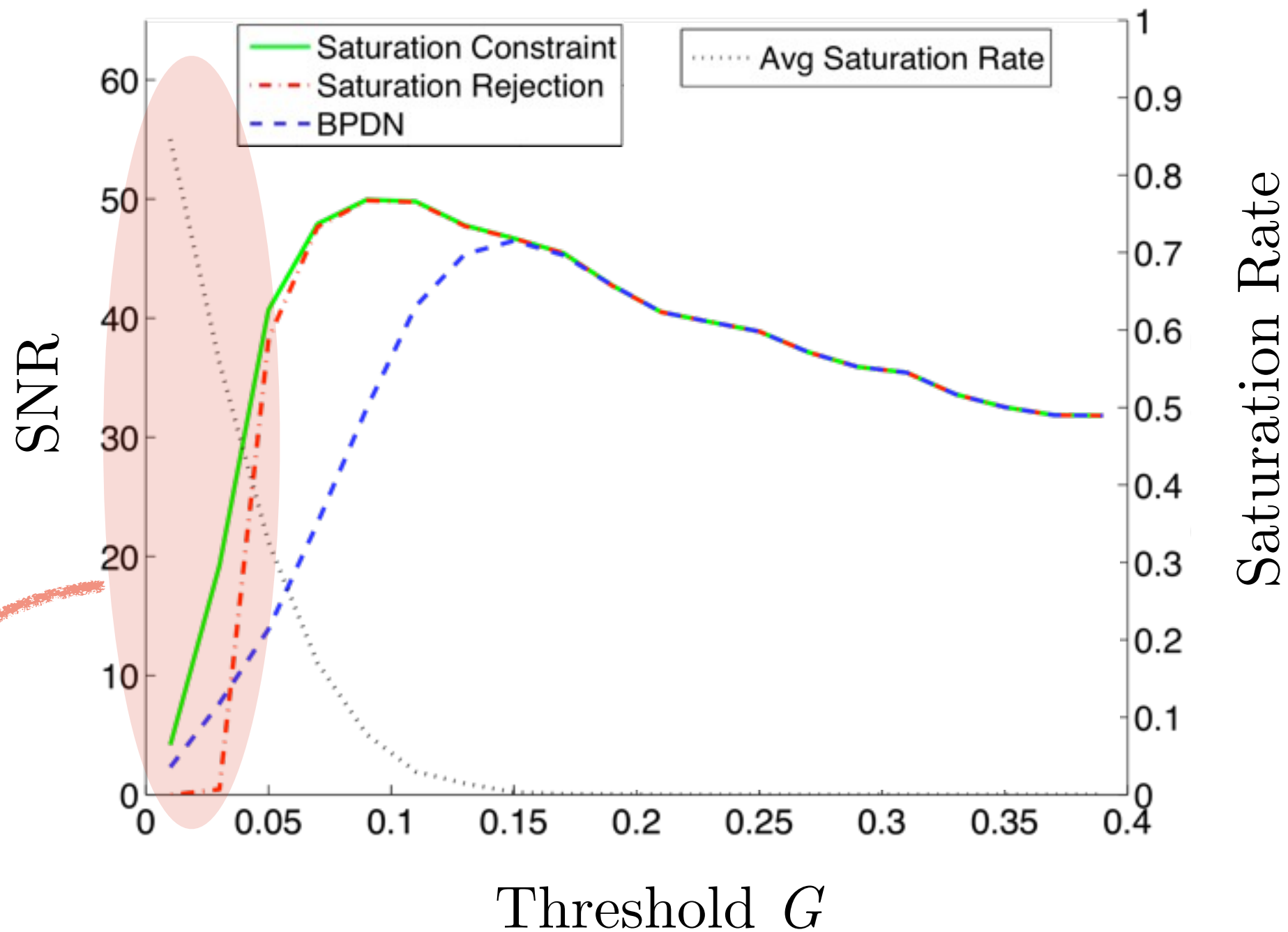


Note: optimal performance **requires** 10% saturation

J.N. Laska, P.T. Boufounos, M.A. Davenport, R.G. Baraniuk, "Democracy in action: Quantization, saturation, and compressive sensing". *Applied and Computational Harmonic Analysis*, 31(3), 429-443. (2011)

# Experimental Results

The “saturation gap”

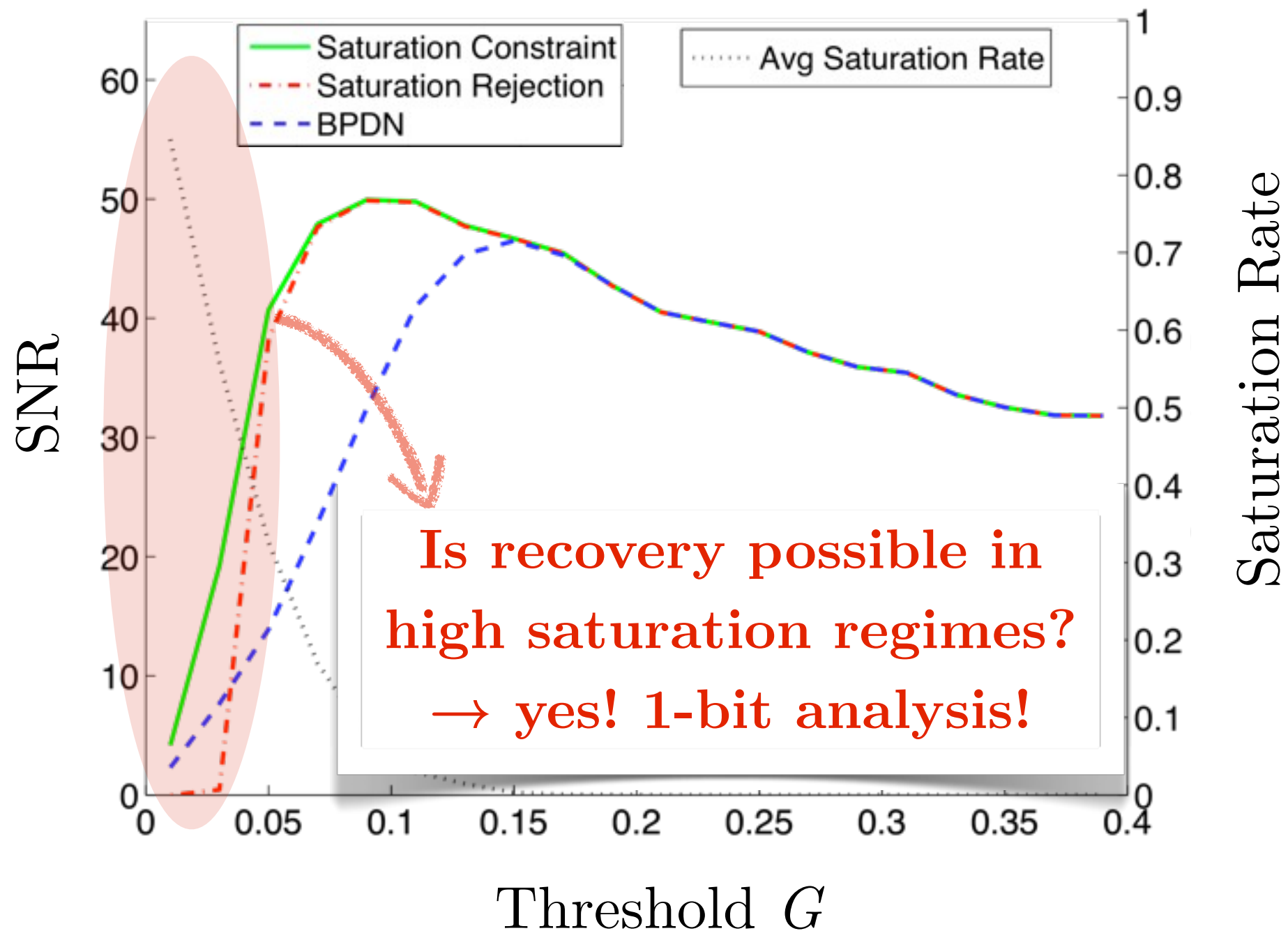


- Majority of measurements saturate
- Recovery fails

J.N. Laska, P.T. Boufounos, M.A. Davenport, R.G. Baraniuk, “Democracy in action: Quantization, saturation, and compressive sensing”. *Applied and Computational Harmonic Analysis*, 31(3), 429-443. (2011)

# Experimental Results

The “saturation gap”



- Majority of measurements saturate
- Recovery fails

J.N. Laska, P.T. Boufounos, M.A. Davenport, R.G. Baraniuk, “Democracy in action: Quantization, saturation, and compressive sensing”. *Applied and Computational Harmonic Analysis*, 31(3), 429-443. (2011)



# Further Reading

- ▶ V. K Goyal, M. Vetterli, N. T. Thao, “Quantized Overcomplete Expansions in RN: Analysis, Synthesis, and Algorithms”, *IEEE Trans. Info. Theory*, 44(1), 1998
- ▶ P. T. Boufounos and R. G. Baraniuk, “Quantization of sparse representations,” *Rice University ECE Department Technical Report 0701*. Summary appears in *Proc. Data Compression Conference (DCC)*, Snowbird, UT, March 27-29, 2007
- ▶ W. Dai, H. V. Pham, and O. Milenkovic, “Quantized Compressive Sensing”, preprint, 2009
- ▶ L. Jacques, D. Hammond, J. Fadili “Dequantizing compressed sensing: When oversampling and non-gaussian constraints combine.” *IEEE Transactions on Information Theory*, 57(1), 559-571, 2011
- ▶ J.N. Laska, P.T. Boufounos, M.A. Davenport, R.G.Baraniuk, “Democracy in action: Quantization, saturation, and compressive sensing”. *Applied and Computational Harmonic Analysis*, 31(3), 429-443, 2011
- ▶ L. Jacques, D. Hammond, J. Fadili, “Stabilizing Nonuniformly Quantized Compressed Sensing with Scalar Companders”, arXiv:1206.6003, 2012
- ▶ Güntürk, C. S., Lammers, M., Powell, A. M., Saab, R., & Yilmaz, Ö. “Sobolev duals for random frames and  $\Sigma\Delta$  quantization of compressed sensing measurements”. *Foundations of Computational Mathematics*, 13(1), 1-36, 2013

Part IV:  
Extreme quantization:  
1-bit compressed sensing

Laurent Jacques, UCL, Belgium  
Petros Boufounos, MERL, USA

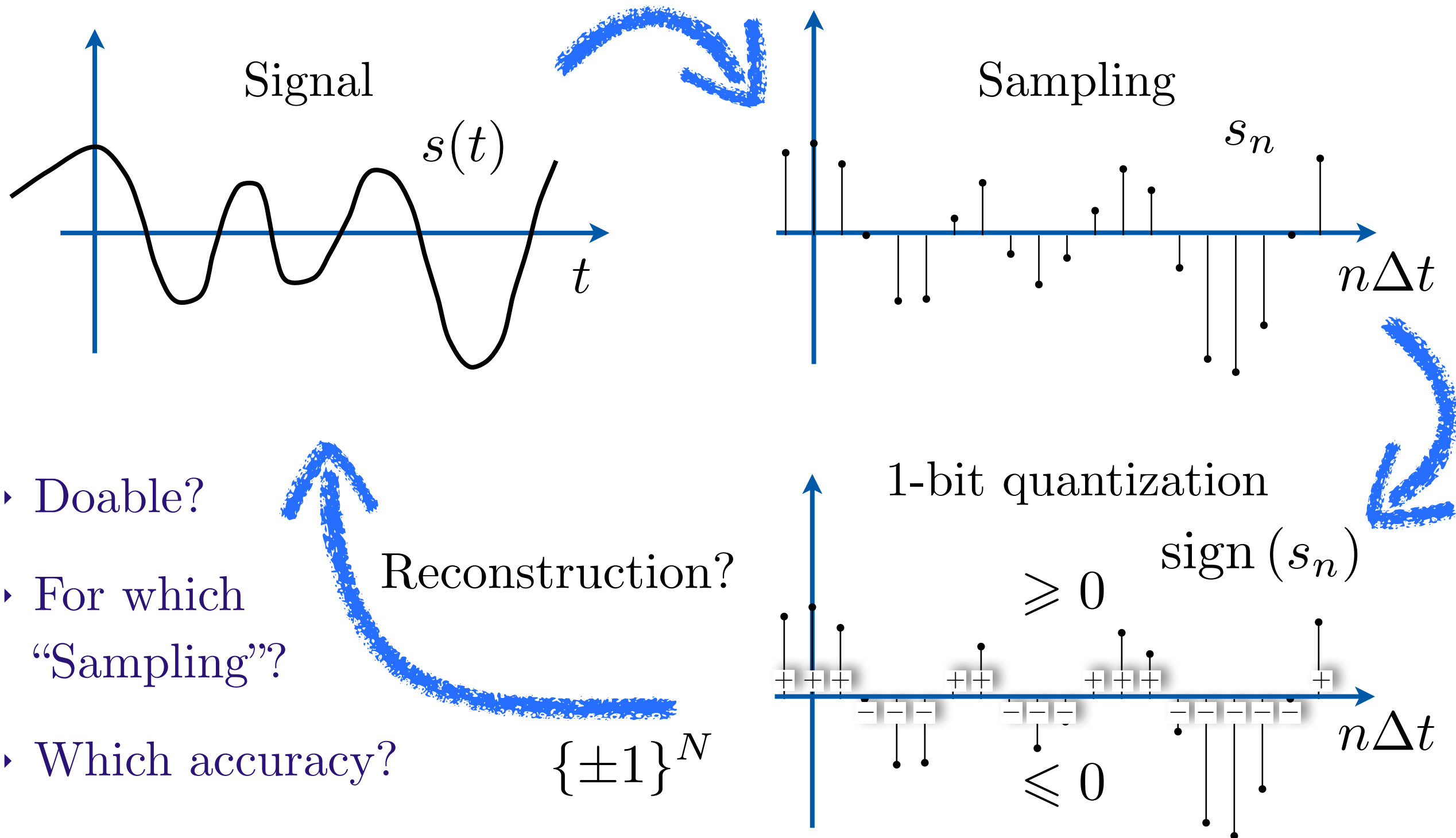


# Outline:

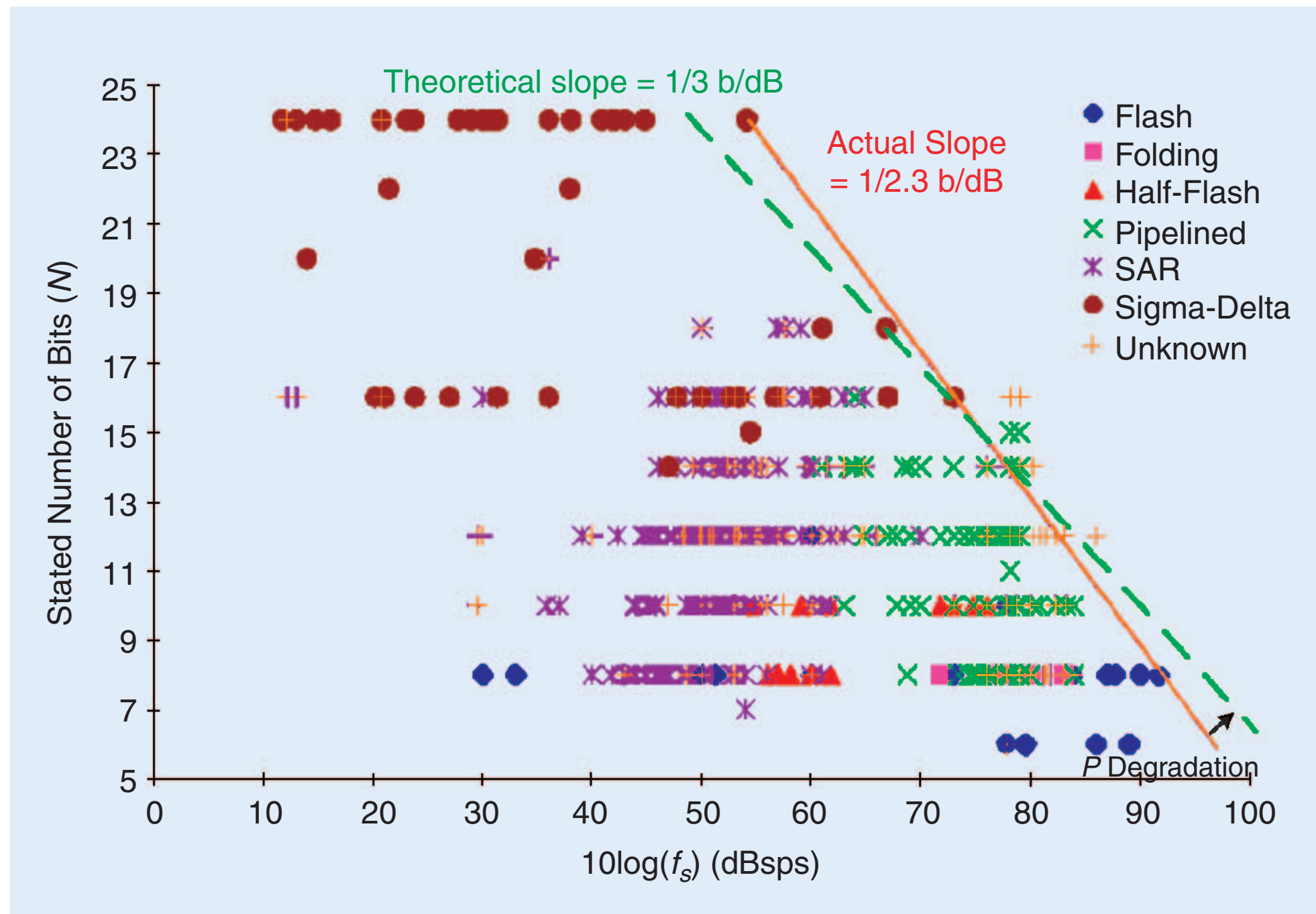
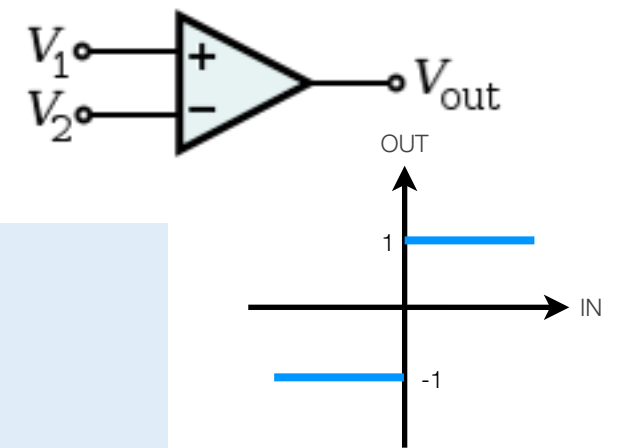
1. Context
2. Theoretical performance limits
3. Stable embeddings: angles are preserved
4. Generalized Embeddings
5. 1-bit CS Reconstructions?
6. Playing with thresholds in 1-bit CS

# 1. Context

# Central question: 1-bit sampling?



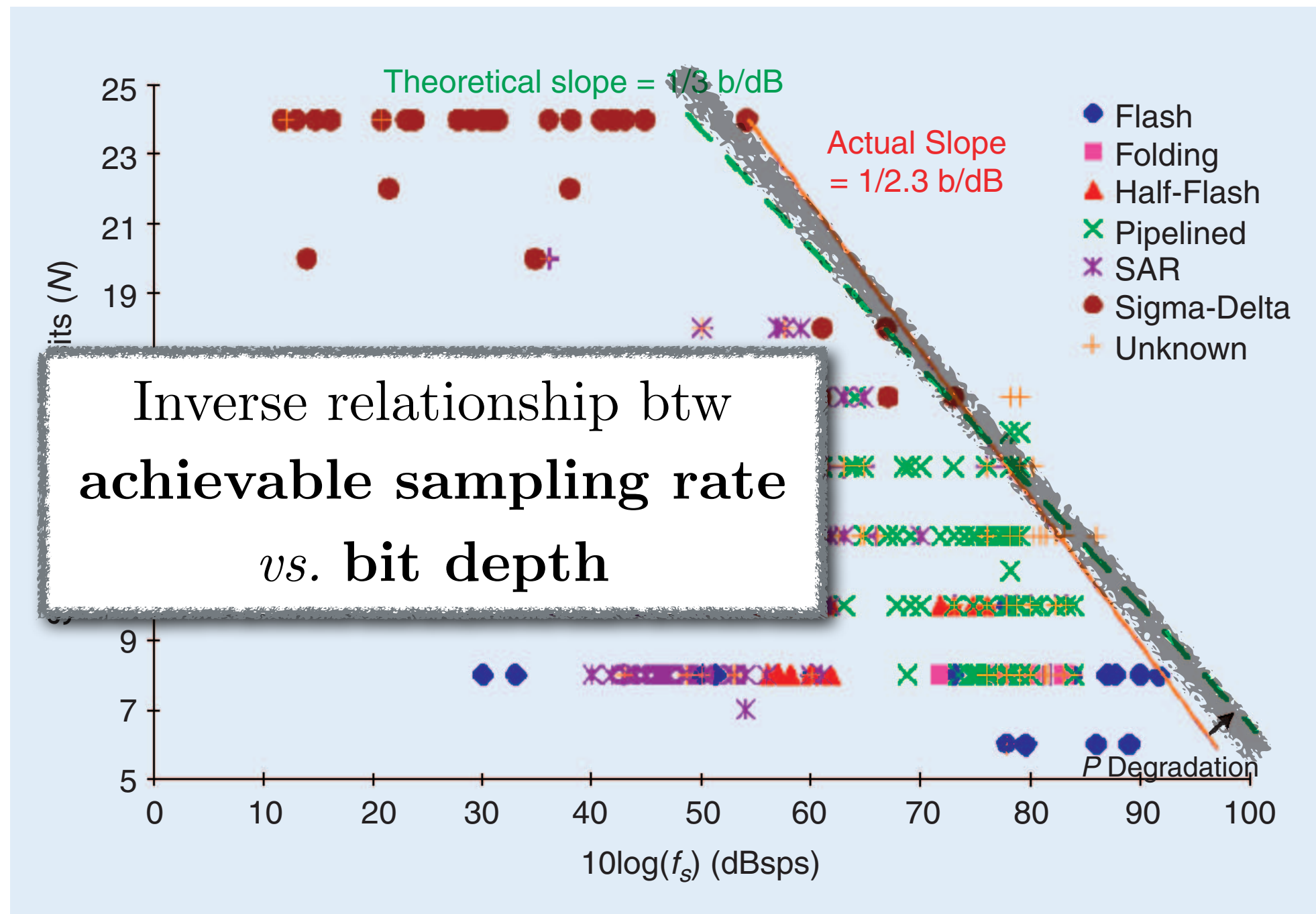
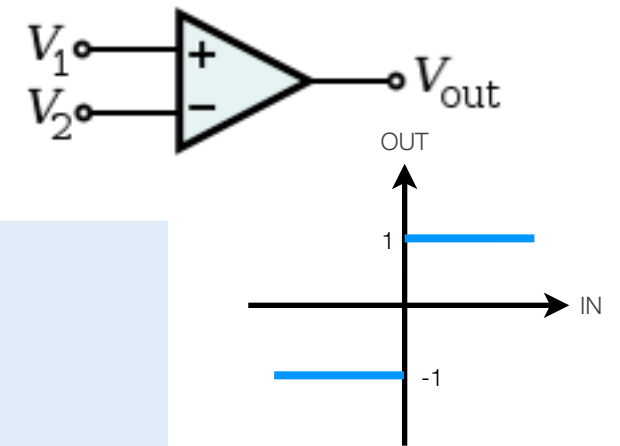
# Why 1-bit? Very Fast Quantizers!



**[FIG1] Stated number of bits versus sampling rate.**

[From "Analog-to-digital converters" B. Le, T.W. Rondeau, J.H. Reed, and C.W.Bostian, IEEE Sig. Proc. Magazine, Nov 2005]

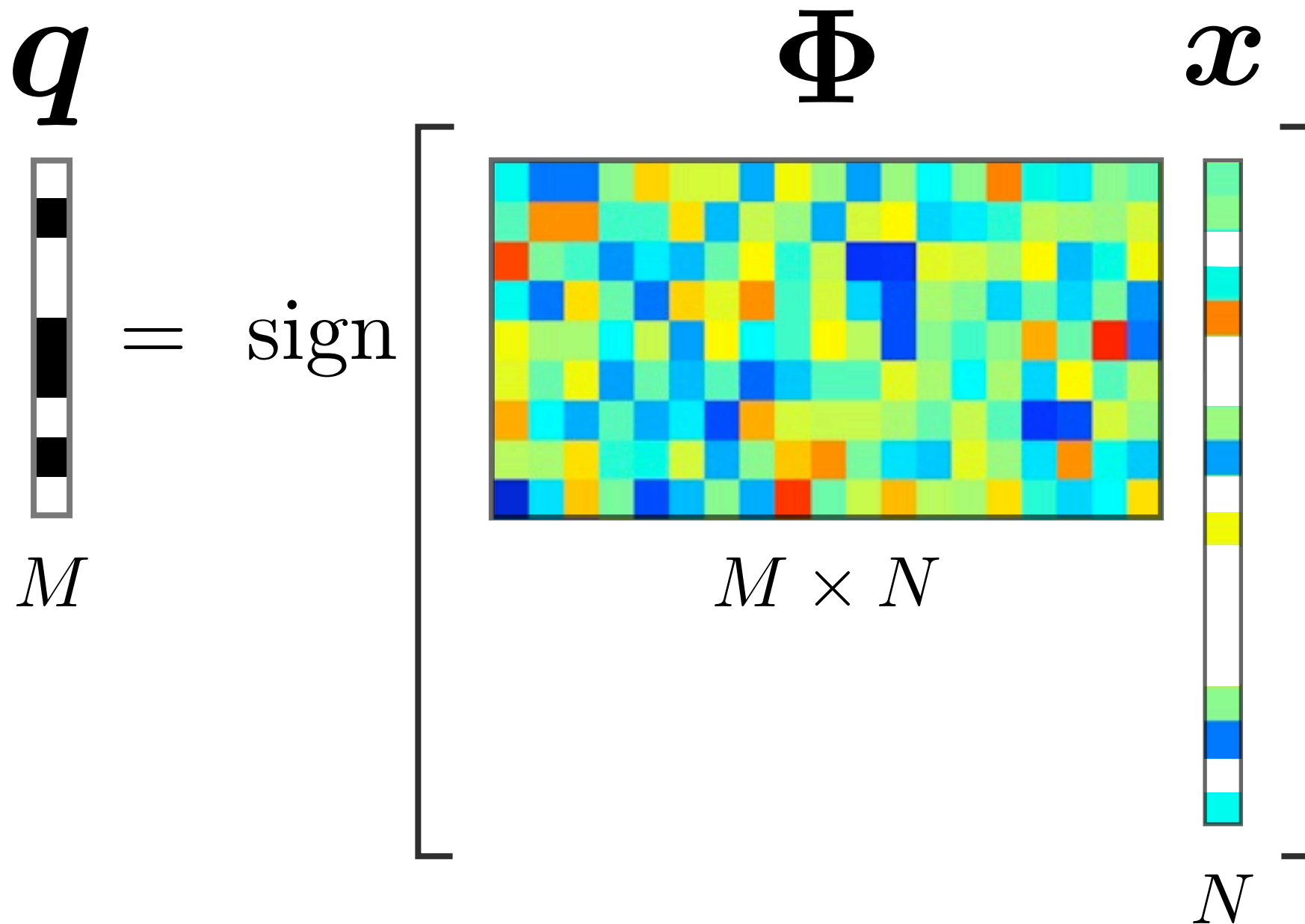
# Why 1-bit? Very Fast Quantizers!



**[FIG1] Stated number of bits versus sampling rate.**

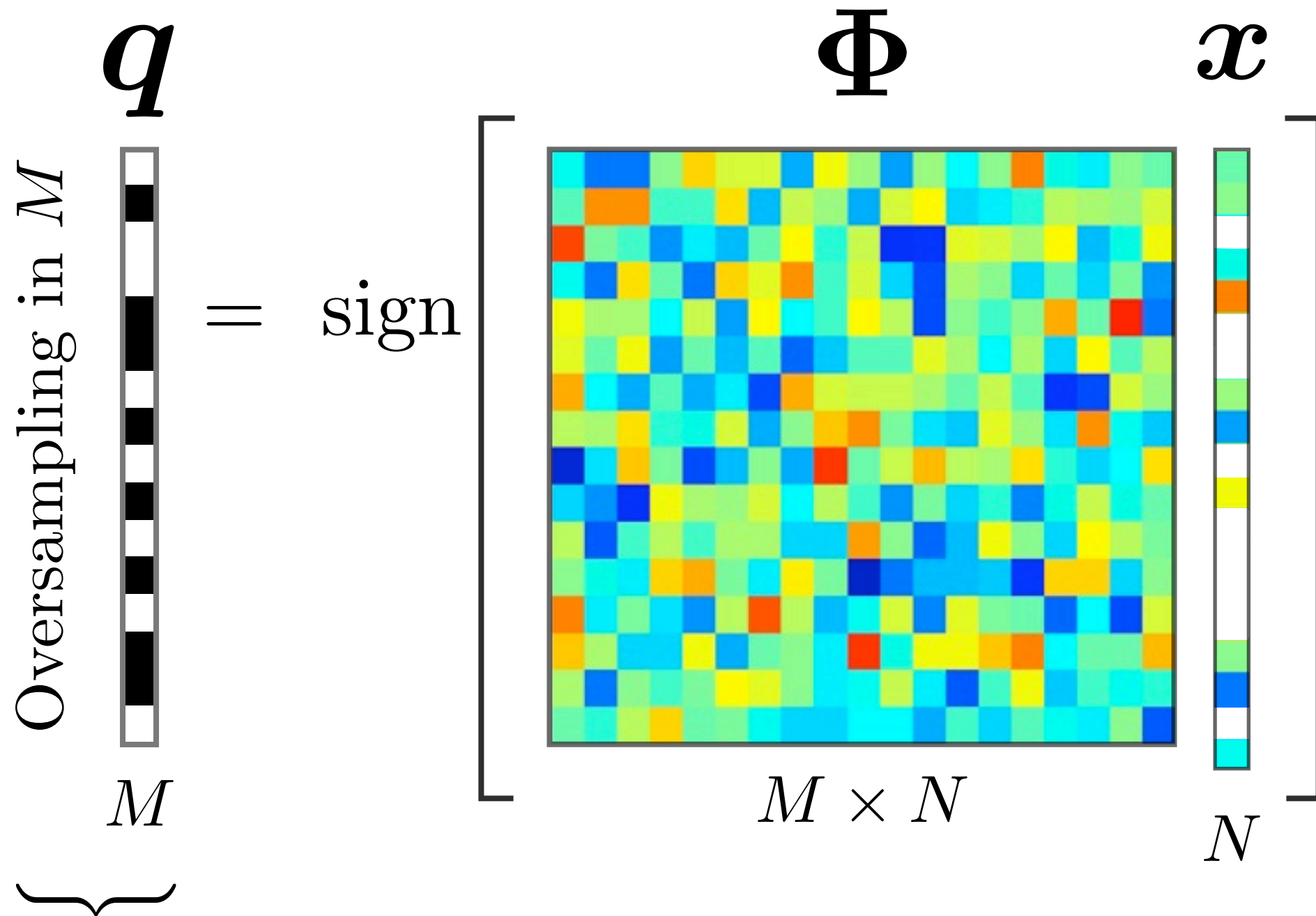
[From "Analog-to-digital converters" B. Le, T.W. Rondeau, J.H. Reed, and C.W. Bostian, IEEE Sig. Proc. Magazine, Nov 2005]

# 1-bit Compressed Sensing



with:  $\text{sign } t = \begin{cases} 1 & \text{if } t > 0 \\ -1 & \text{if } t \leq 0 \end{cases}$  component-wise

# 1-bit Compressed Sensing

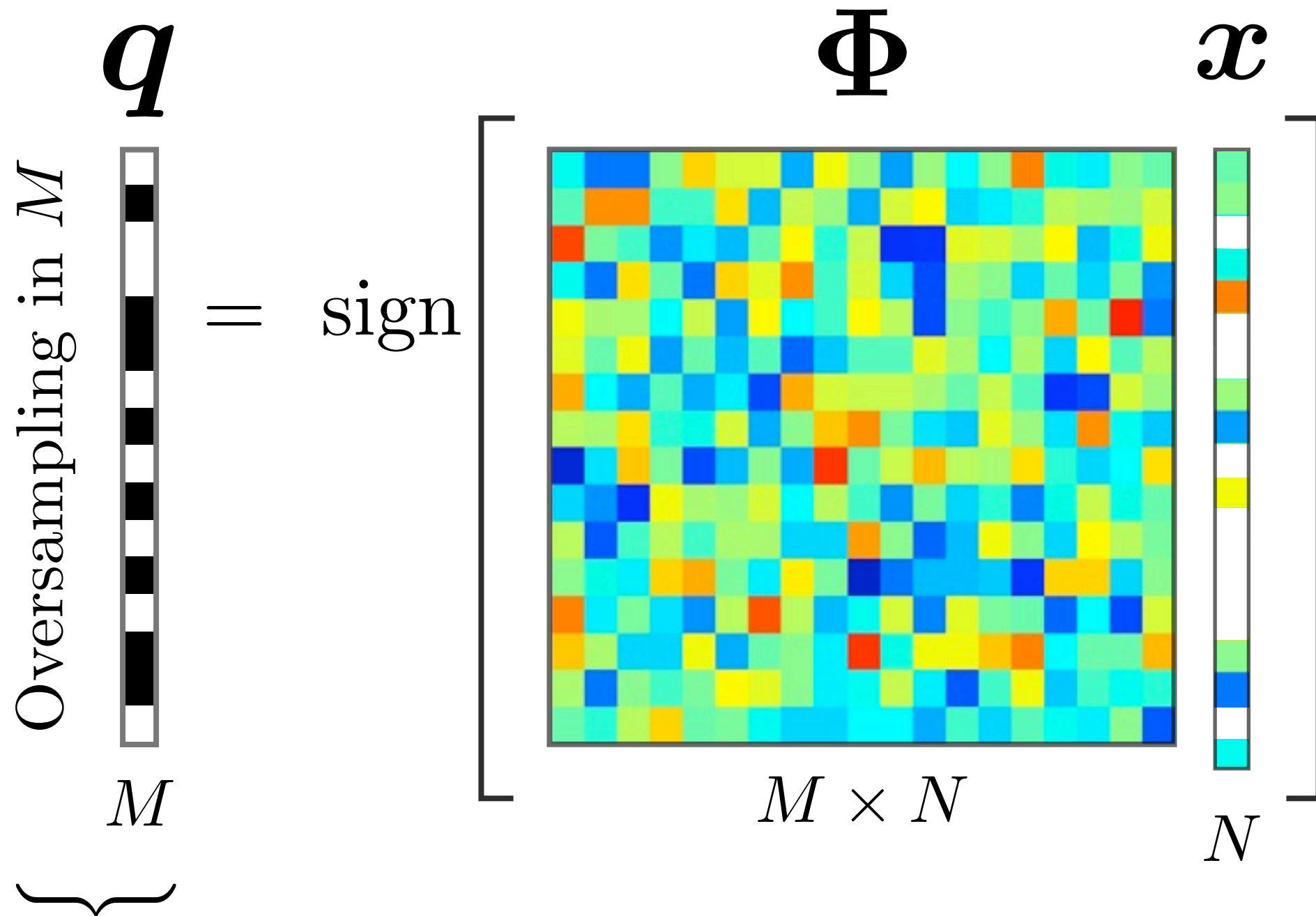


$M$ -bits! But, which information inside  $q$  ?



# 1-bit Computational ~~pressed~~ Sensing

bits matter!

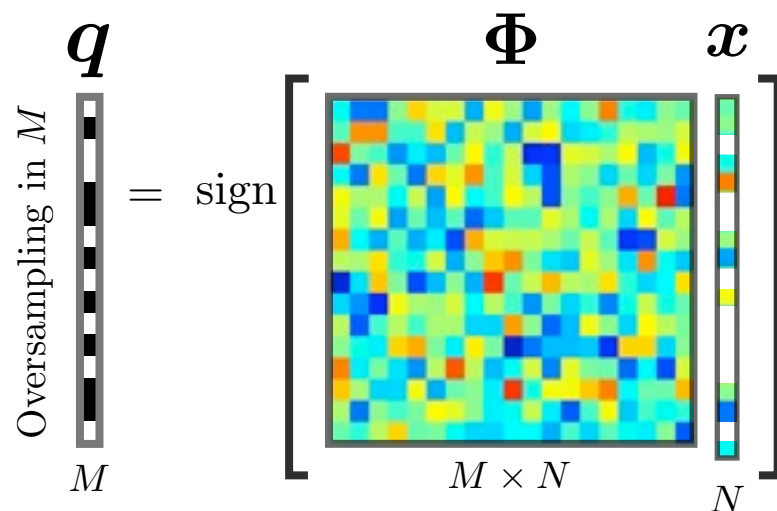


$M$ -bits! But, which information inside  $q$ ?



# 1-bit Computational Compressed Sensing

bits matter!



Warning 1: signal amplitude is lost!

$$q = \text{sign}(\Phi(\lambda x)) = \text{sign}(\Phi x), \quad \forall \lambda > 0$$

$\Rightarrow$  Amplitude is arbitrarily fixed

Examples :  $\|x\| = 1$  or  $\|\Phi x\|_1 = 1$

# 1-bit Computational Compressed Sensing

bits matter!

$$\underset{M}{\overset{\text{Oversampling in } M}{q}} = \text{sign} \left[ \underset{M \times N}{\Phi} \underset{N}{x} \right]$$

[Plan, Vershynin, 11]

## Warning 2: $\exists$ forbidden sensing!

Let  $\mathbf{x}_\lambda := (1, \lambda, 0, \dots, 0)^T \in \mathbb{R}^N$   
and  $\Phi \in \{\pm 1\}^{M \times N}$  (e.g., Bernoulli).

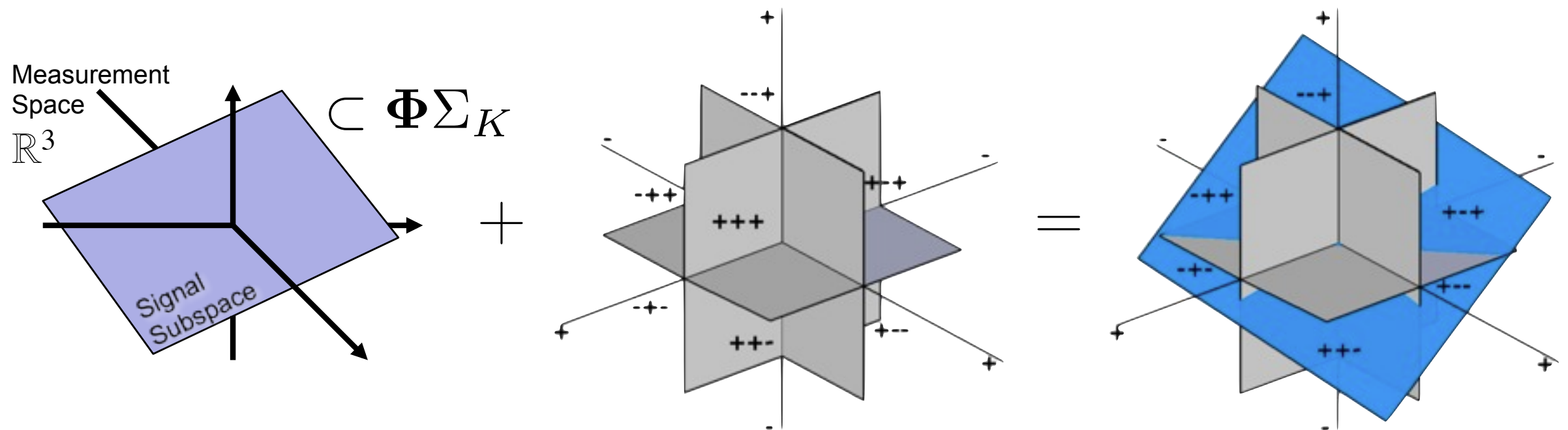
We have  $\|\mathbf{x}_0 - \mathbf{x}_\lambda\| = \lambda$

but  $\mathbf{q} = \text{sign}(\Phi \mathbf{x}_0) = \text{sign}(\Phi \mathbf{x}_\lambda), \forall |\lambda| < 1$

$\Rightarrow$  No hope to distinguish them by increasing  $M$ !

## 2. Theoretical performance limits

# Lower bound: cell intersection viewpoint

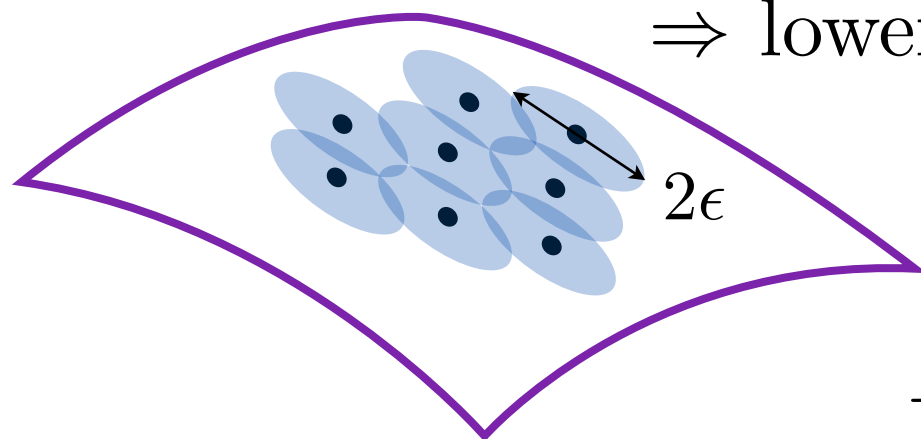


Not all quantization cells intersected!

no more than  $C = 2^K \binom{N}{K} \binom{M}{K}$

Most efficient  $\epsilon$ -covering of  $S^{N-1} \cap \Sigma_K$  with  $\epsilon$ -caps

$\Rightarrow$  lower bound on  $C$  by “ $\text{vol}(S^{N-1} \cap \Sigma_K) / \text{vol}(\epsilon\text{-cap})$ ”



$$\Rightarrow \epsilon = \Omega(K/M)$$

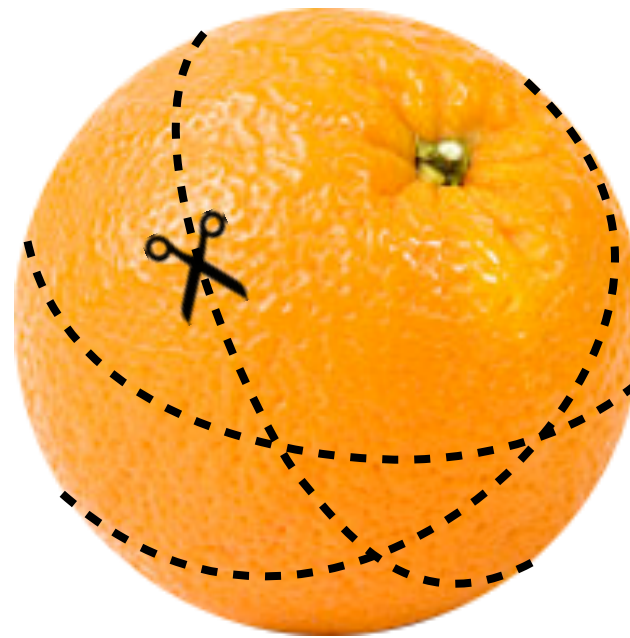
$\rightarrow$  Lower bound on any 1-bit reconstruction error

# Reaching this bound ?



Carl Friedrich Gauss:  
“1-bit CS? I solved it at  
breakfast by randomly  
slicing my orange!”

<http://www.gaussfacts.com>



# Reaching this bound ?

$x$  on  $S^2$

$M$  vectors:

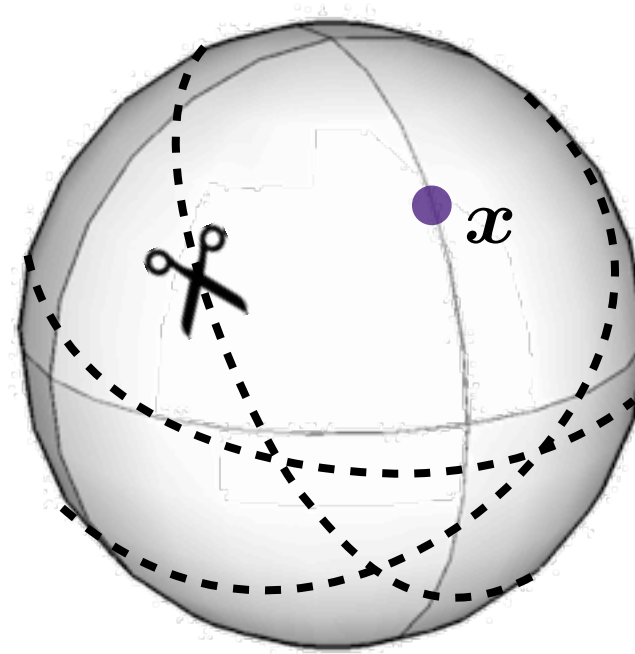
$$\{\varphi_i : 1 \leq i \leq M\}$$

iid Gaussian



Carl Friedrich Gauss:  
“1-bit CS? I solved it at  
breakfast by randomly  
slicing my orange!”

<http://www.gaussfacts.com>





# Reaching this bound ?

$x$  on  $S^2$

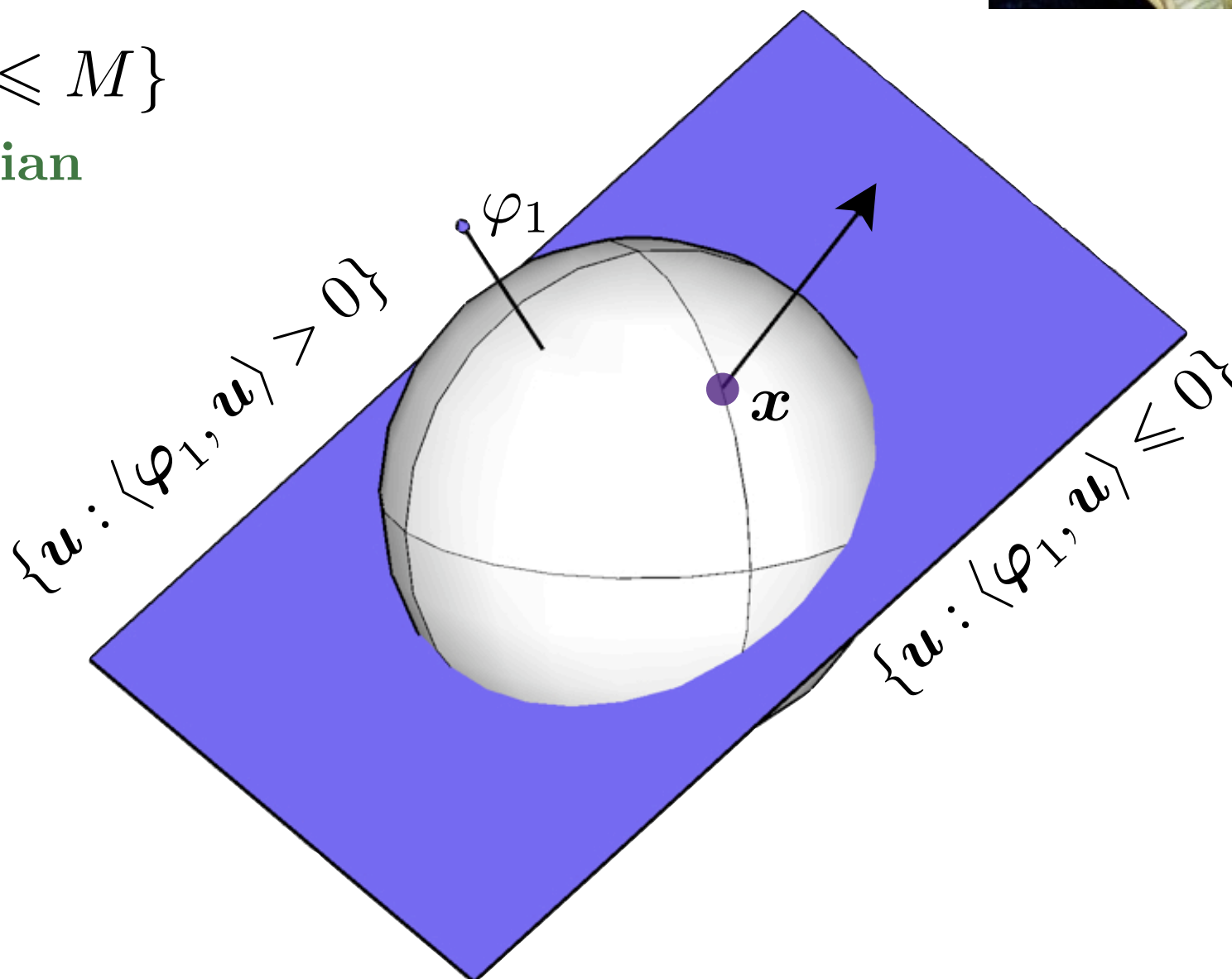
$M$  vectors:

$$\{\varphi_i : 1 \leq i \leq M\}$$

iid Gaussian

1-bit Measurements

$$\langle \varphi_1, x \rangle > 0$$



Carl Friedrich Gauss:

“1-bit CS? I solved it at breakfast by randomly slicing my orange!”

<http://www.gaussfacts.com>

# Reaching this bound ?

$\mathbf{x}$  on  $S^2$

$M$  vectors:

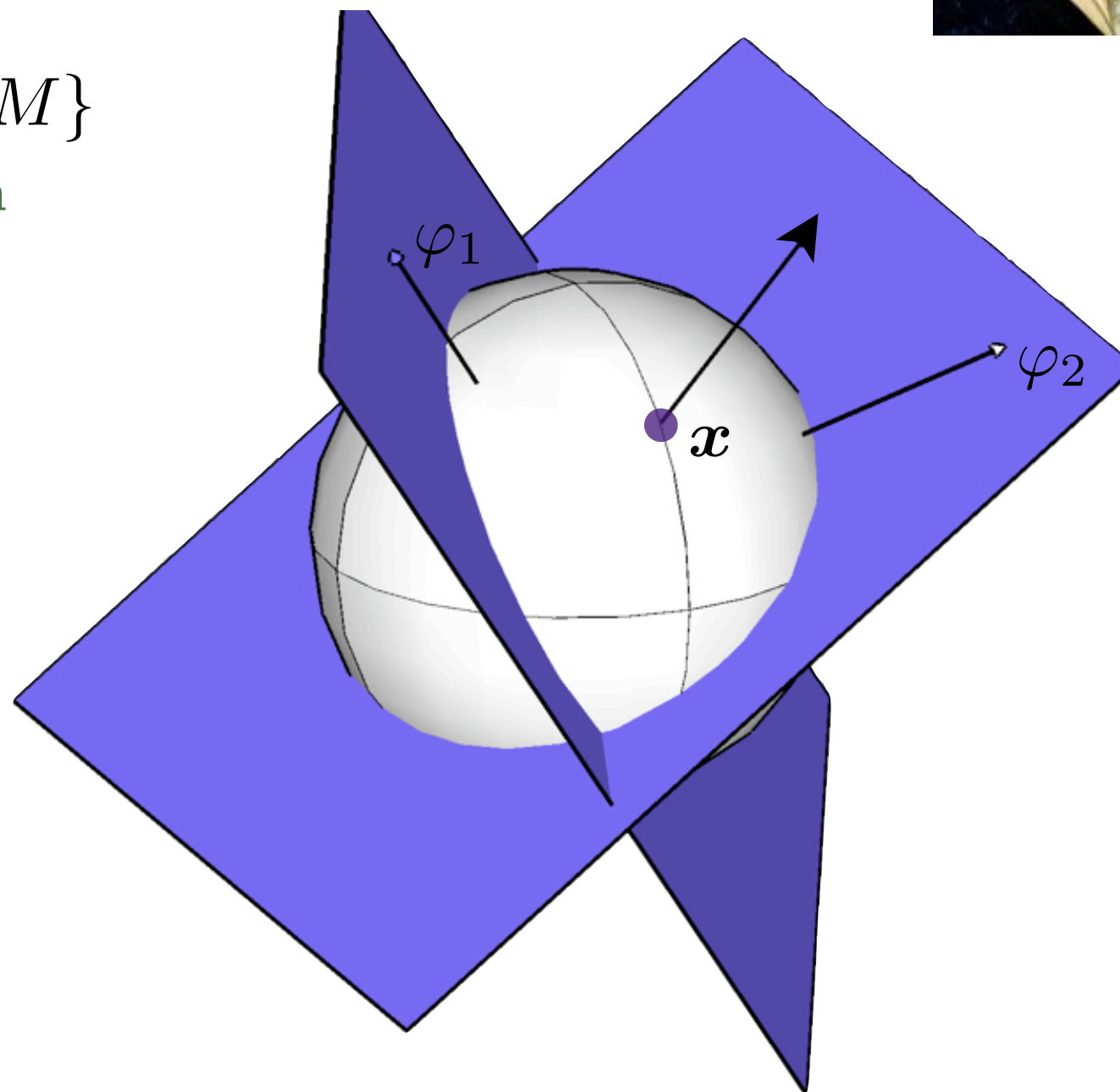
$$\{\varphi_i : 1 \leq i \leq M\}$$

iid Gaussian

1-bit Measurements

$$\langle \varphi_1, \mathbf{x} \rangle > 0$$

$$\langle \varphi_2, \mathbf{x} \rangle > 0$$



Carl Friedrich Gauss:

“1-bit CS? I solved it at breakfast by randomly slicing my orange!”

<http://www.gaussfacts.com>



# Reaching this bound ?

$x$  on  $S^2$

$M$  vectors:

$$\{\varphi_i : 1 \leq i \leq M\}$$

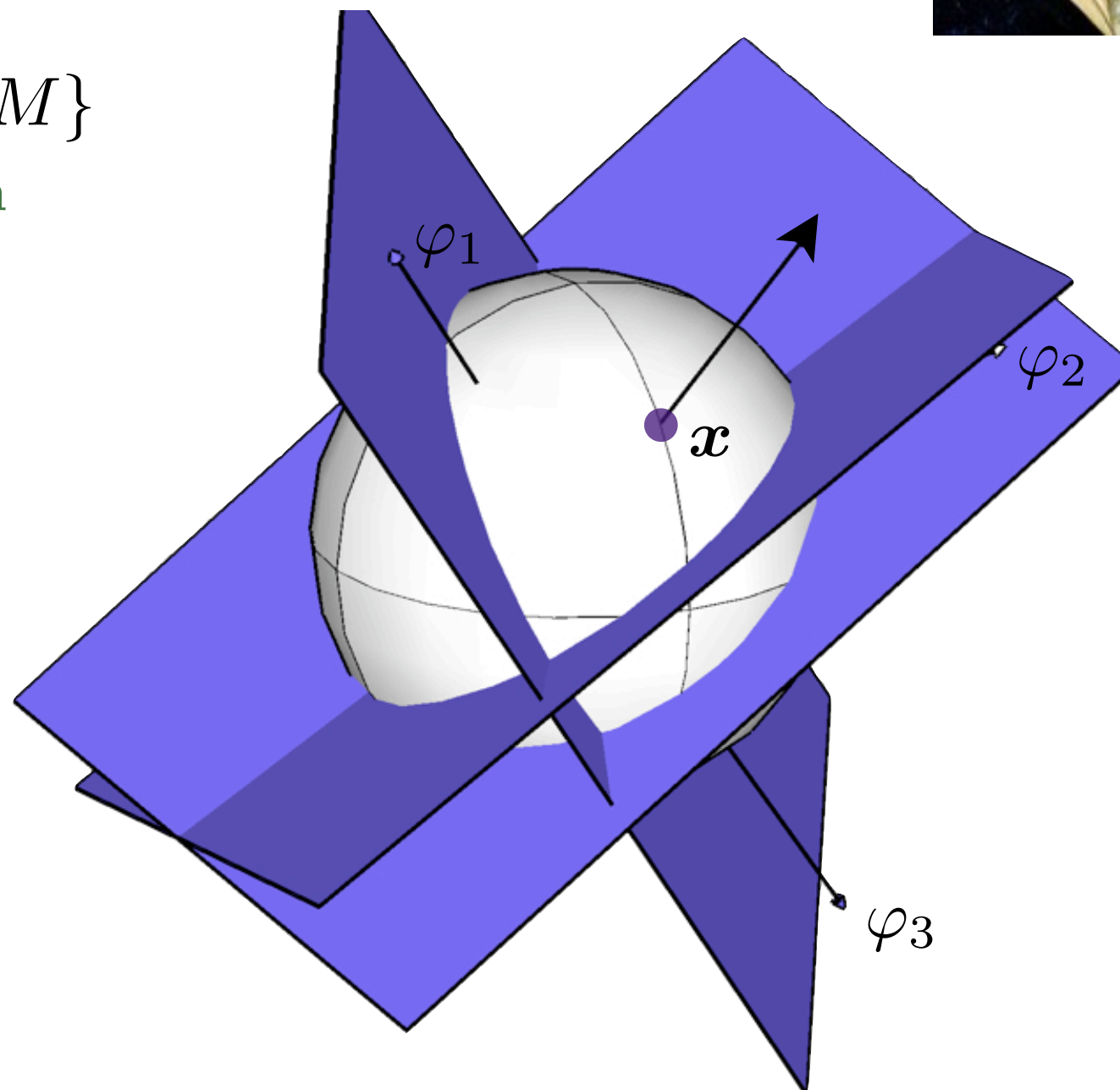
iid Gaussian

1-bit Measurements

$$\langle \varphi_1, x \rangle > 0$$

$$\langle \varphi_2, x \rangle > 0$$

$$\langle \varphi_3, x \rangle \leq 0$$



Carl Friedrich Gauss:

“1-bit CS? I solved it at breakfast by randomly slicing my orange!”

<http://www.gaussfacts.com>

# Reaching this bound ?

$\mathbf{x}$  on  $S^2$

$M$  vectors:

$$\{\varphi_i : 1 \leq i \leq M\}$$

iid Gaussian

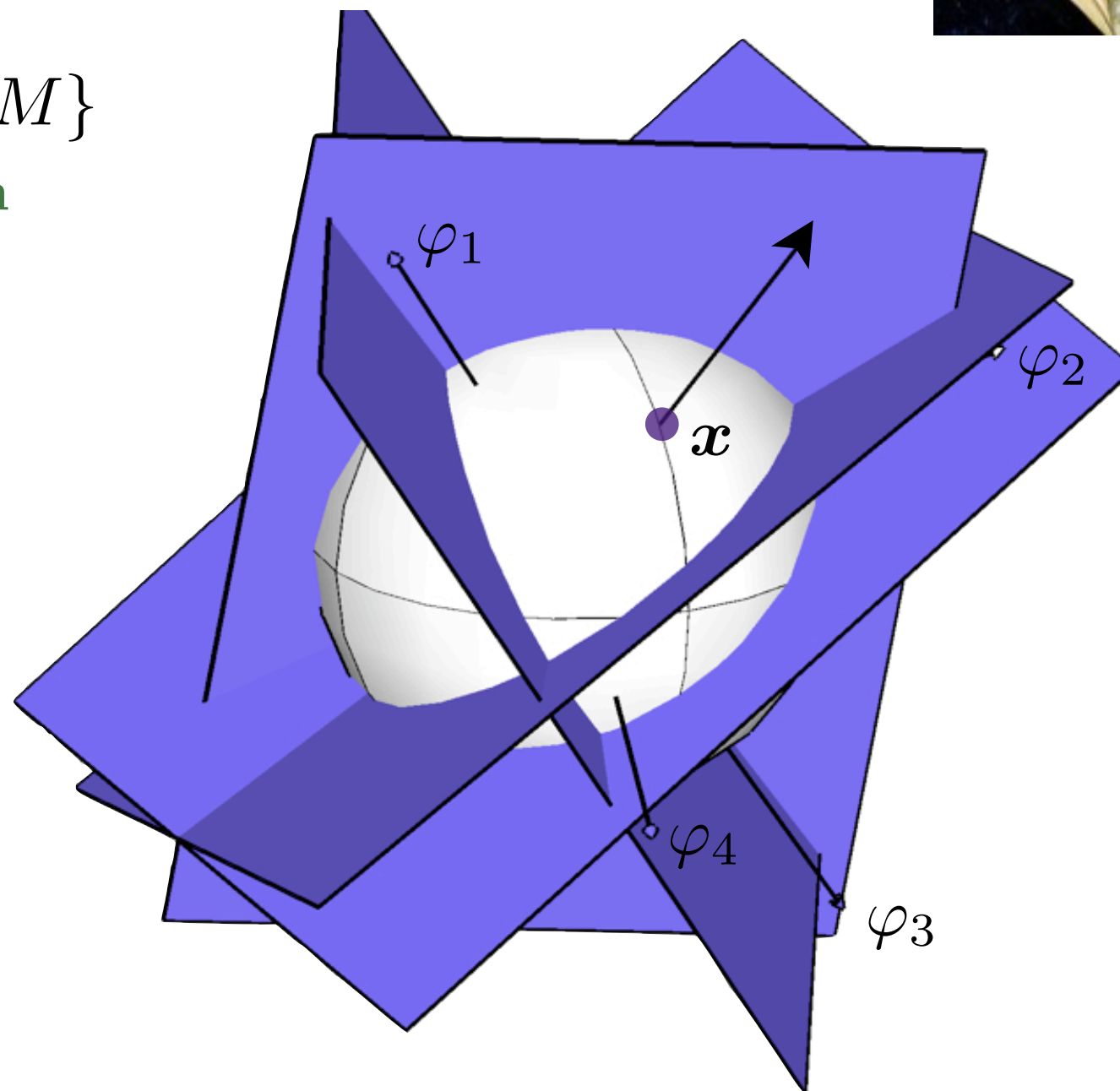
1-bit Measurements

$$\langle \varphi_1, \mathbf{x} \rangle > 0$$

$$\langle \varphi_2, \mathbf{x} \rangle > 0$$

$$\langle \varphi_3, \mathbf{x} \rangle \leq 0$$

$$\langle \varphi_4, \mathbf{x} \rangle > 0$$



Carl Friedrich Gauss:

“1-bit CS? I solved it at breakfast by randomly slicing my orange!”

<http://www.gaussfacts.com>

# Reaching this bound ?

$x$  on  $S^2$

$M$  vectors:

$$\{\varphi_i : 1 \leq i \leq M\}$$

iid Gaussian

1-bit Measurements

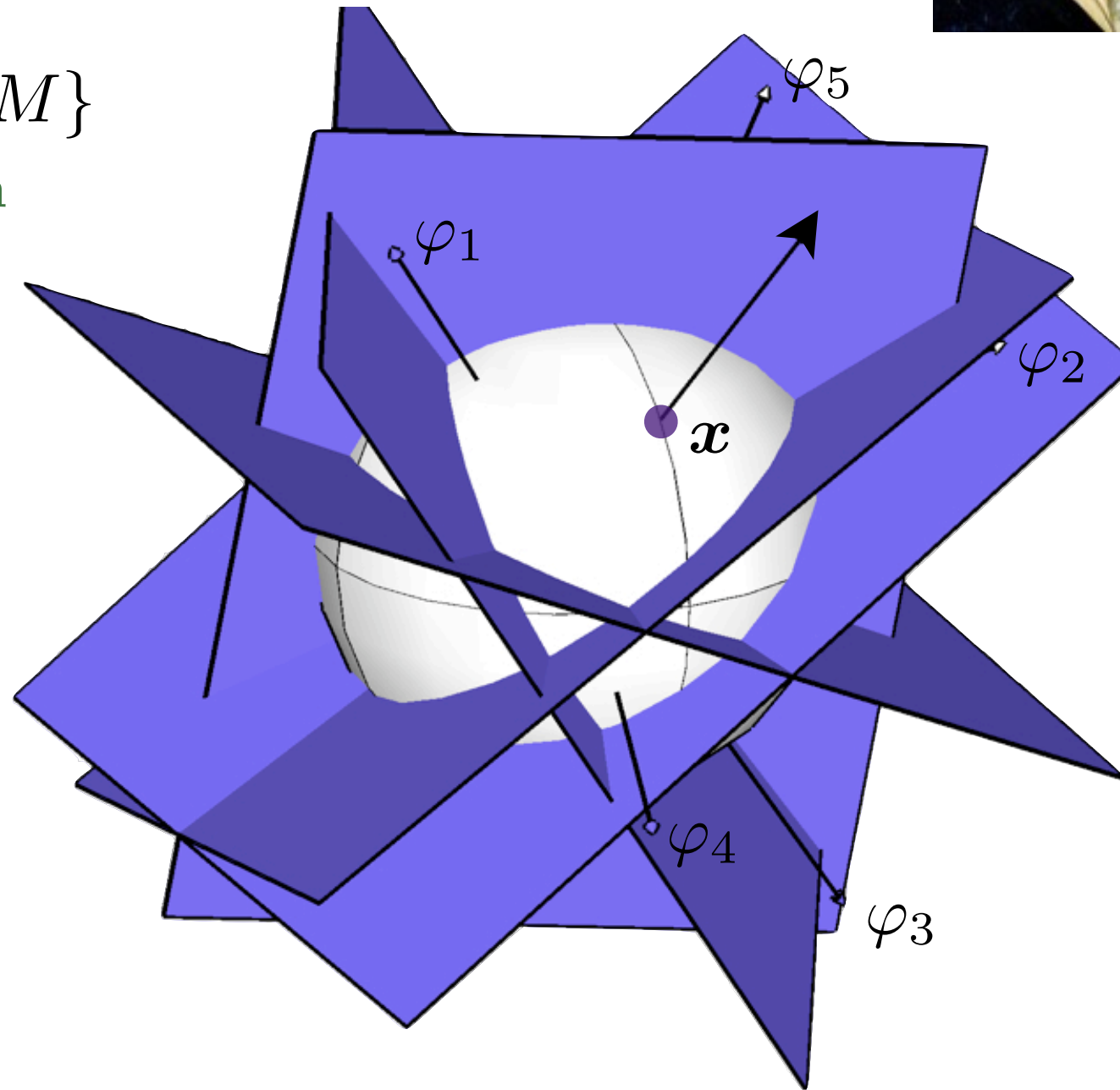
$$\langle \varphi_1, x \rangle > 0$$

$$\langle \varphi_2, x \rangle > 0$$

$$\langle \varphi_3, x \rangle \leq 0$$

$$\langle \varphi_4, x \rangle > 0$$

$$\langle \varphi_5, x \rangle > 0$$



Carl Friedrich Gauss:

“1-bit CS? I solved it at breakfast by randomly slicing my orange!”

<http://www.gaussfacts.com>

# Reaching this bound ?

$x$  on  $S^2$

$M$  vectors:

$$\{\varphi_i : 1 \leq i \leq M\}$$

iid Gaussian

1-bit Measurements

$$\langle \varphi_1, x \rangle > 0$$

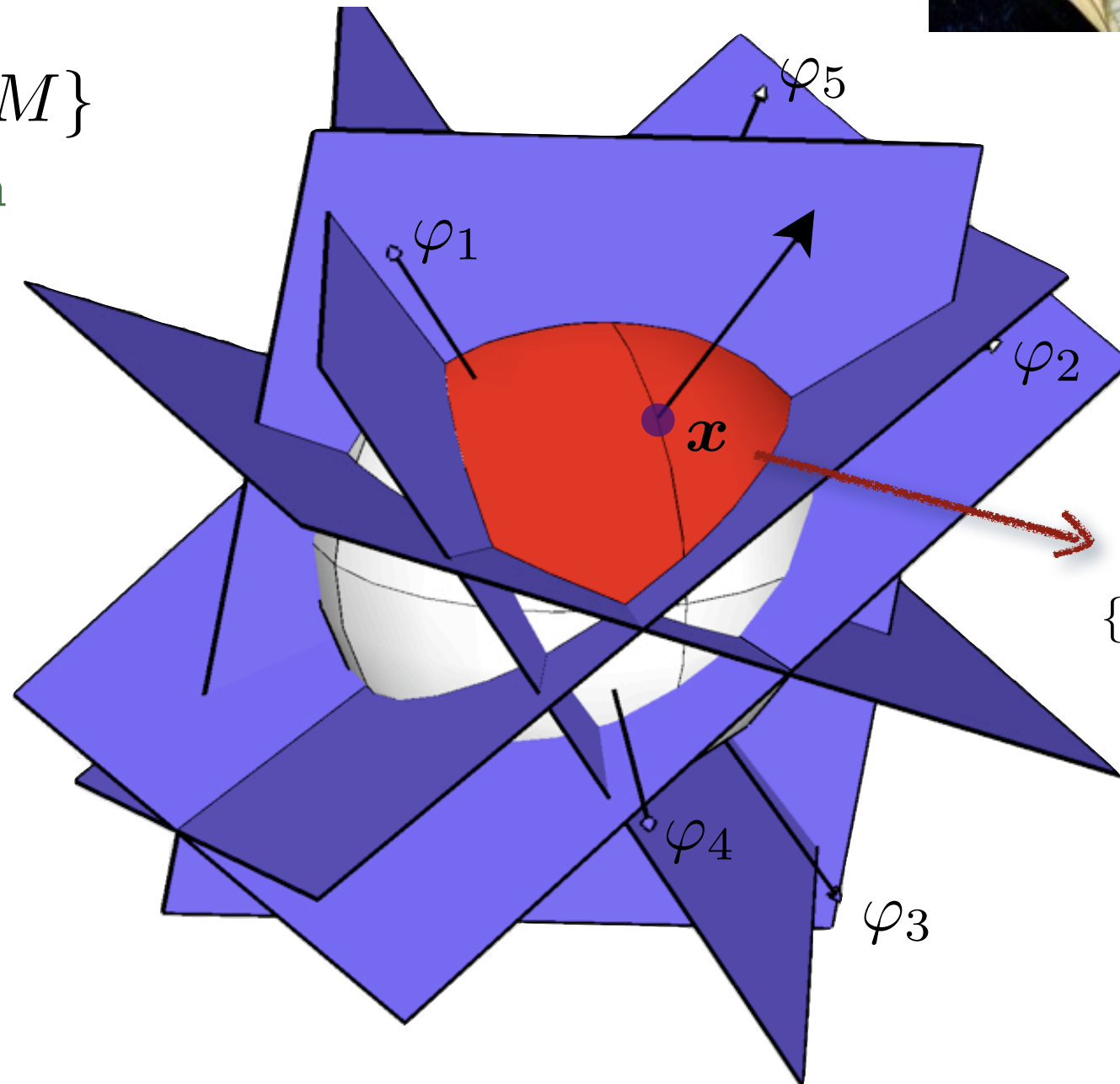
$$\langle \varphi_2, x \rangle > 0$$

$$\langle \varphi_3, x \rangle \leq 0$$

$$\langle \varphi_4, x \rangle > 0$$

$$\langle \varphi_5, x \rangle > 0$$

$\vdots$



Smaller and smaller  
when  $M$  increases  
 $\{u : \text{sign}(\Phi u) = \text{sign}(\Phi x)\}$



Carl Friedrich Gauss:  
“1-bit CS? I solved it at  
breakfast by randomly  
slicing my orange!”

<http://www.gaussfacts.com>



# Reaching this bound ?

$\mathbf{x}$  on  $S^2$

$M$  vectors:

$$\{\varphi_i : 1 \leq i \leq M\}$$

iid Gaussian

1-bit Measurements

$$\langle \varphi_1, \mathbf{x} \rangle > 0$$

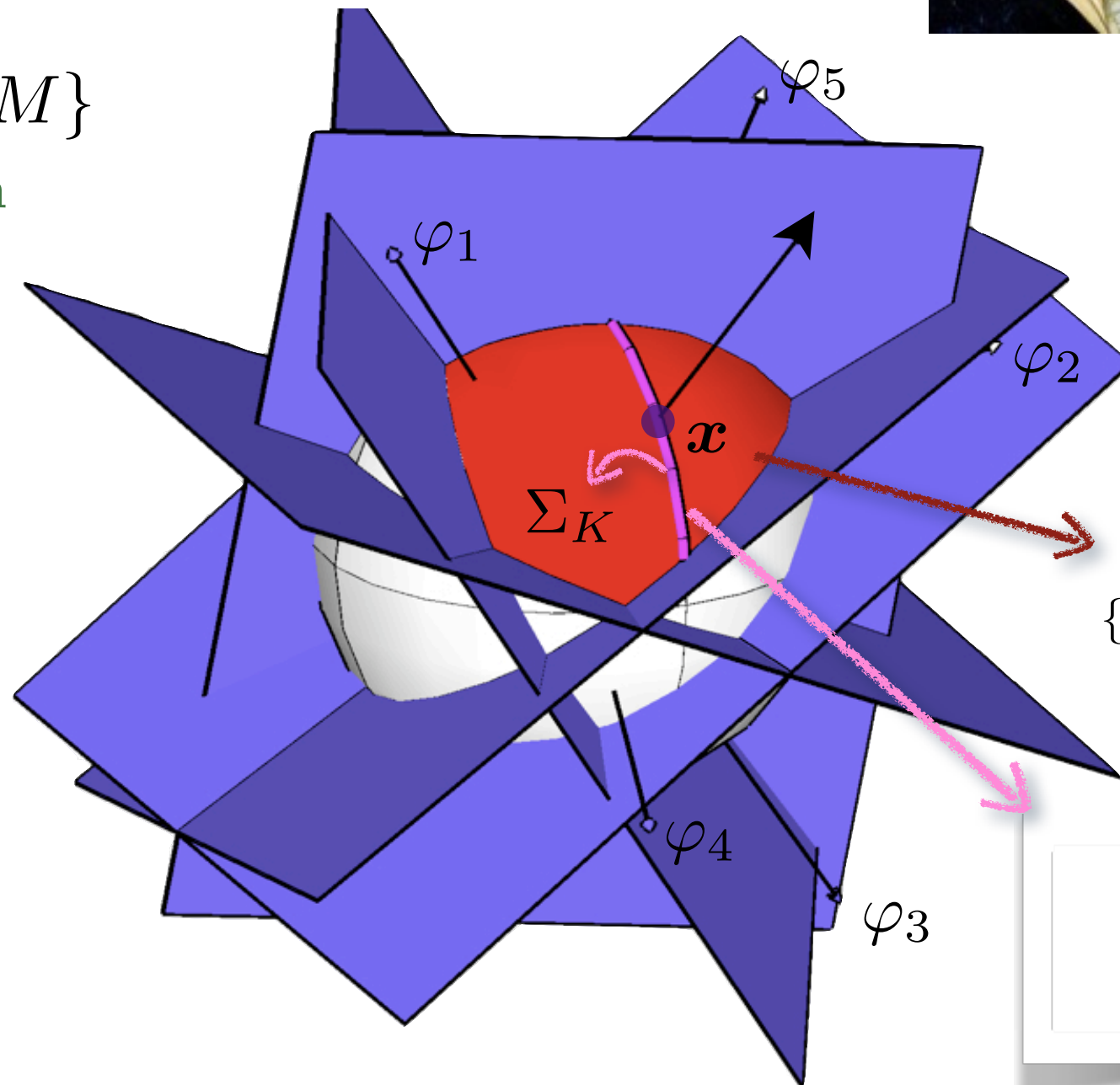
$$\langle \varphi_2, \mathbf{x} \rangle > 0$$

$$\langle \varphi_3, \mathbf{x} \rangle \leq 0$$

$$\langle \varphi_4, \mathbf{x} \rangle > 0$$

$$\langle \varphi_5, \mathbf{x} \rangle > 0$$

$\vdots$



Smaller and smaller  
when  $M$  increases  
 $\{u : \text{sign}(\Phi u) = \text{sign}(\Phi x)\}$

Lower bound on  
this width?



Carl Friedrich Gauss:

“1-bit CS? I solved it at  
breakfast by randomly  
slicing my orange!”

<http://www.gaussfacts.com>

# Reaching this bound ?



Carl Friedrich Gauss:  
“1-bit CS? I solved it at  
breakfast by randomly  
slicing my orange!”  
<http://www.gaussfacts.com>

Let  $A(\cdot) := \text{sign}(\Phi \cdot)$  with  $\Phi \sim \mathcal{N}^{M \times N}(0, 1)$ .

If  $M = O(\epsilon^{-1} K \log N)$ , then, w.h.p,

for any two unit  $K$ -sparse vectors  $\mathbf{x}$  and  $\mathbf{s}$ ,

$$A(\mathbf{x}) = A(\mathbf{s}) \quad \Rightarrow \quad \|\mathbf{x} - \mathbf{s}\| \leq \epsilon$$

$$\Leftrightarrow \epsilon = O\left(\frac{K}{M} \log \frac{MN}{K}\right)$$

almost optimal

Note: You can even afford a small error, *i.e.*,  
if only  $b$  bits are different  
between  $A(\mathbf{x})$  and  $A(\mathbf{s})$   $\Rightarrow \|\mathbf{x} - \mathbf{s}\| \leq \frac{K+b}{K} \epsilon$

### 3. Stable embeddings: angles are preserved

# Starting point: Hamming/Angle Concentration

- Metrics of interest:

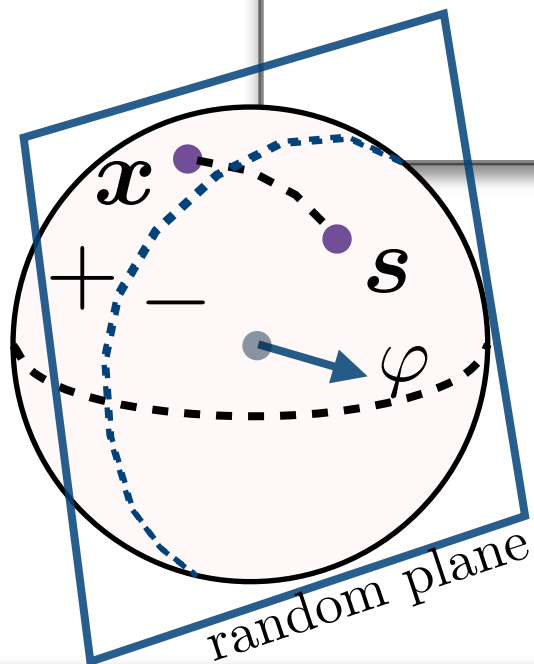
$$d_H(\mathbf{u}, \mathbf{v}) = \frac{1}{M} \sum_i (u_i \oplus v_i) \quad (\text{norm. Hamming})$$

$$d_{\text{ang}}(\mathbf{x}, \mathbf{s}) = \frac{1}{\pi} \arccos(\langle \mathbf{x}, \mathbf{s} \rangle) \quad (\text{norm. angle})$$

- Known fact: if  $\Phi \sim \mathcal{N}^{M \times N}(0, 1)$  [e.g., Goemans, Williamson 1995]

Let  $\Phi \sim \mathcal{N}^{M \times N}(0, 1)$ ,  $A(\cdot) = \text{sign}(\Phi \cdot) \in \{-1, 1\}^M$  and  $\epsilon > 0$ .  
For any  $\mathbf{x}, \mathbf{s} \in S^{N-1}$ , we have

$$\mathbb{P}_{\Phi} \left[ \left| d_H(A(\mathbf{x}), A(\mathbf{s})) - d_{\text{ang}}(\mathbf{x}, \mathbf{s}) \right| \leq \epsilon \right] \geq 1 - 2e^{-2\epsilon^2 M}.$$



Thanks to  $A(\cdot)$ , Hamming distance concentrates around vector angles!



# Binary $\epsilon$ Stable Embedding ( $B_{\epsilon}SE$ )

A mapping  $A : \mathbb{R}^N \rightarrow \{\pm 1\}^M$  is a **binary  $\epsilon$ -stable embedding** ( $B_{\epsilon}SE$ ) of order  $K$  for sparse vectors if

$$|d_{\text{ang}}(\mathbf{x}, \mathbf{s}) - \epsilon| \leq d_H(A(\mathbf{x}), A(\mathbf{s})) \leq d_{\text{ang}}(\mathbf{x}, \mathbf{s}) + \epsilon$$

for all  $\mathbf{x}, \mathbf{s} \in S^{N-1}$  with  $\mathbf{x} \pm \mathbf{s}$   $K$ -sparse.

kind of “binary restricted (quasi) isometry”

- ▶ *Corollary*: for any algorithm with output  $\mathbf{x}^*$  jointly  $K$ -sparse and consistent (*i.e.*,  $A(\mathbf{x}^*) = A(\mathbf{x})$ ),

$$d_{\text{ang}}(\mathbf{x}, \mathbf{x}^*) \leq 2\epsilon!$$

- ▶ If limited binary noise,  $d_{\text{ang}}$  still bounded
- ▶ If not exactly sparse signals (but almost),  $d_{\text{ang}}$  still bounded

# B $\epsilon$ SE existence? Yes!

Let  $\Phi \sim \mathcal{N}^{M \times N}(0, 1)$ , fix  $0 \leq \eta \leq 1$  and  $\epsilon > 0$ . If

$$M \geq \frac{4}{\epsilon^2} \left( K \log(N) + 2K \log\left(\frac{50}{\epsilon}\right) + \log\left(\frac{2}{\eta}\right) \right),$$

then  $\Phi$  is a B $\epsilon$ SE with  $\Pr > 1 - \eta$ .

$$M = O(\epsilon^{-2} K \log N)$$

Proof sketch:

1) Generalize

$$\mathbb{P}_{\Phi} \left[ \left| d_H(A(\mathbf{x}), A(\mathbf{s})) - d_{\text{ang}}(\mathbf{x}, \mathbf{s}) \right| \leq \epsilon \right] \geq 1 - 2e^{-2\epsilon^2 M}.$$

to

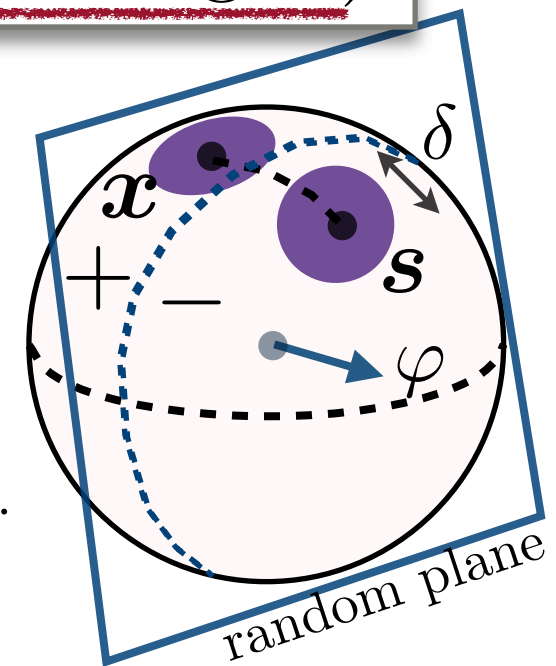
$$\mathbb{P}_{\Phi} \left[ \left| d_H(A(\mathbf{u}), A(\mathbf{v})) - d_{\text{ang}}(\mathbf{x}, \mathbf{s}) \right| \leq \epsilon + \left(\frac{\pi}{2} D\right)^{1/2} \delta \right] \geq 1 - 2e^{-2\epsilon^2 M}.$$

for  $\mathbf{u}, \mathbf{v}$  in a  $D$ -dimensional neighborhood of width  $\delta$  around  $\mathbf{x}$  and  $\mathbf{s}$  resp.

2) Covers the space of " $K$ -sparse signal pairs" in  $\mathbb{R}^N$  by

$$O\left(\binom{N}{K} \delta^{-2K}\right) = O\left(\left(\frac{eN}{K\delta^2}\right)^K\right) \text{ neighborhoods.}$$

3) Apply Point 1 with union bound, and "stir until the proof thickens"



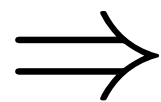
# B $\epsilon$ SE existence? Yes!

Let  $\Phi \sim \mathcal{N}^{M \times N}(0, 1)$ , fix  $0 \leq \eta \leq 1$  and  $\epsilon > 0$ . If

$$M \geq \frac{4}{\epsilon^2} \left( K \log(N) + 2K \log\left(\frac{50}{\epsilon}\right) + \log\left(\frac{2}{\eta}\right) \right),$$

then  $\Phi$  is a B $\epsilon$ SE with  $\Pr > 1 - \eta$ .

$$M = O(\epsilon^{-2} K \log N)$$



B $\epsilon$ SE consistency “width”:

$$\epsilon = O\left(\left(\frac{K}{M} \log \frac{MN}{K}\right)^{1/2}\right)$$

not as optimal but  
stronger result!

$$d_H \leftrightarrow d_{\text{ang}}$$

# 4. Generalized Embeddings

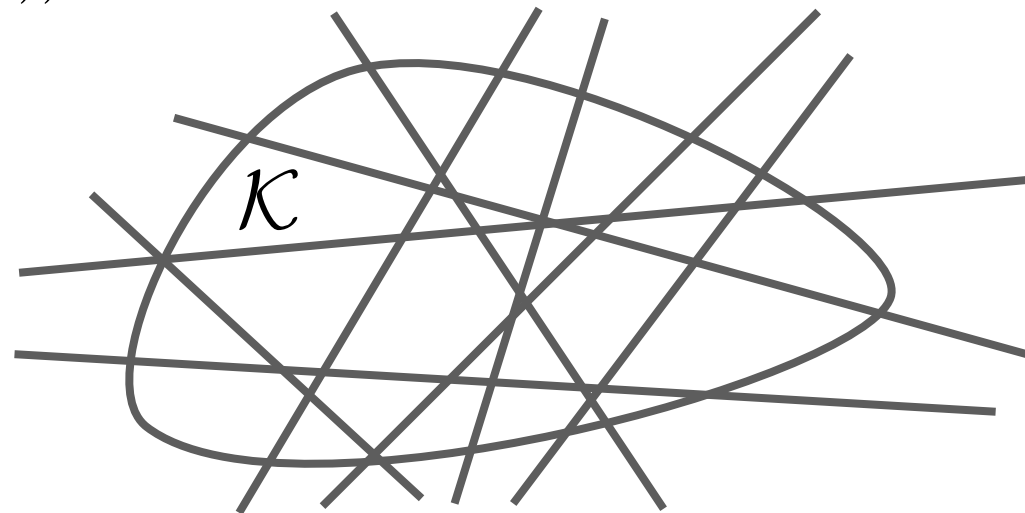
# Beyond strict sparsity ...

Let  $\mathcal{K} \subset S^{N-1}$  (e.g., compressible signals s.t.  $\|\mathbf{x}\|_2/\|\mathbf{x}\|_1 \leq \sqrt{K}$ )  
 $\neq \Sigma_K$

What can we say on  $d_H(A(\mathbf{x}), A(\mathbf{s}))$  for  $\mathbf{x}, \mathbf{s} \in \mathcal{K}$ ?

*Uniform tessellation:* [Plan, Vershynin, 11]

$P(\frac{\# \text{ random hyperplanes btw } \mathbf{x} \text{ and } \mathbf{s}}{d_H(A(\mathbf{x}), A(\mathbf{s}))} \propto d_{\text{ang}}(\mathbf{x}, \mathbf{s})) ?$



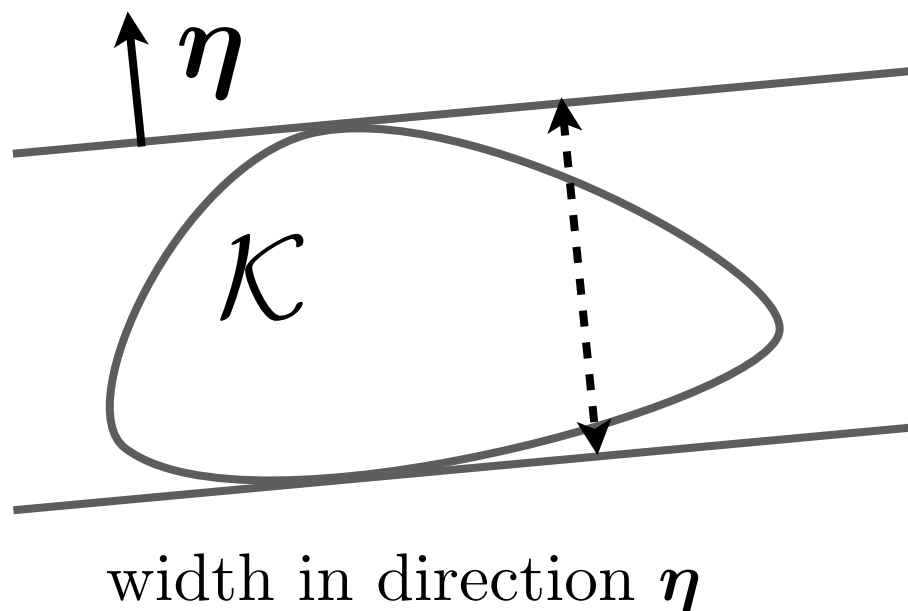
Y. Plan, R. Vershynin, "Dimension reduction by random hyperplane tessellations", 2011, arXiv:1111.4452

Y. Plan, R. Vershynin, "Robust 1-bit compressed sensing and sparse logistic regression: a convex programming approach", IEEE TIT 2012, arXiv:1202.1212.

# Beyond strict sparsity ...

Measuring the “dimension” of  $\mathcal{K} \rightarrow$  Gaussian mean width:

$$w(\mathcal{K}) := \mathbb{E} \sup_{\mathbf{u} \in \mathcal{K}} \langle \mathbf{g}, \mathbf{u} \rangle, \text{ with } g_k \sim_{\text{iid}} \mathcal{N}(0, 1)$$



Examples:

$$w^2(\mathcal{S}^{N-1}) \leq 4N$$

$$w^2(\mathcal{K}) \leq C \log |\mathcal{K}| \quad (\text{for finite sets})$$

$$w^2(\mathcal{K}) \leq L \quad \text{if subspace with } \dim \mathcal{K} = L$$

$$w^2(\Sigma_K) \simeq K \log(2N/K)$$

⋮

Y. Plan, R. Vershynin, “Dimension reduction by random hyperplane tessellations”, 2011, arXiv:1111.4452

Y. Plan, R. Vershynin, “Robust 1-bit compressed sensing and sparse logistic regression: a convex programming approach”, IEEE TIT 2012, arXiv:1202.1212.

# Beyond strict sparsity ...

**Proposition** Let  $\Phi \sim \mathcal{N}^{M \times N}(0, 1)$  and  $\mathcal{K} \subset \mathbb{R}^N$ . Then, for some  $C, c > 0$ , if

$$M \geq C \epsilon^{-6} w^2(\mathcal{K}),$$

not as optimal but  
stronger result!

then, with  $Pr \geq 1 - e^{-c\epsilon^2 M}$ , we have

$$d_{\text{ang}}(\mathbf{x}, \mathbf{s}) - \epsilon \leq d_H(A(\mathbf{x}), A(\mathbf{s})) \leq d_{\text{ang}}(\mathbf{x}, \mathbf{s}) + \epsilon, \quad \forall \mathbf{x}, \mathbf{s} \in \mathcal{K}.$$

Generalize B $\epsilon$ SE to more general sets.

In particular, to

$$\mathcal{C}_K = \{\mathbf{u} \in \mathbb{R}^N : \|\mathbf{u}\|_2 / \|\mathbf{u}\|_1 \leq \sqrt{K}\} \supset \Sigma_K$$

$$\text{with } w^2(\mathcal{C}_K) \leq cK \log N/K.$$

$\Rightarrow$  Extension to “1-bit Matrix Completion” possible!

$$\text{i.e., } w^2(r\text{-rank } N_1 \times N_2 \text{ matrix}) \leq cr(N_1 + N_2)!$$

Y. Plan, R. Vershynin, “Dimension reduction by random hyperplane tessellations”, 2011, arXiv:1111.4452

Y. Plan, R. Vershynin, “Robust 1-bit compressed sensing and sparse logistic regression: a convex programming approach”, IEEE TIT 2012, arXiv:1202.1212.

## 5. 1-bit CS Reconstructions?



# Dumbest 1-bit reconstruction

Fact:

If  $M = O(\epsilon^{-2} K \log N/K)$  (for  $\mathbf{x} \in \Sigma_K$  fixed,  $\forall \mathbf{s} \in \Sigma_K$ )

or, if  $M = O(\epsilon^{-6} K \log N/K)$  ( $\forall \mathbf{x}, \mathbf{s} \in \Sigma_K$ ), then, w.h.p,

$$|\frac{\sqrt{\pi}/2}{M} \langle \text{sign}(\Phi \mathbf{x}), \Phi \mathbf{s} \rangle - \langle \mathbf{x}, \mathbf{s} \rangle| \leq \epsilon \quad [\text{Plan, Vershynin, 12}]$$

► Implication? [LJ, Degraux, De Vleeschouwer, 13]

Let  $\mathbf{x} \in \Sigma_K \cap S^{N-1}$  and  $\mathbf{q} = \text{sign}(\Phi \mathbf{x})$ .  
Compute

$$\hat{\mathbf{x}} = \frac{\pi}{2M} \mathcal{H}_K(\Phi^* \mathbf{q})$$

Then, if previous property holds,

$$\|\mathbf{x} - \hat{\mathbf{x}}\| \leq 2\epsilon.$$

Non-uniform case ( $\mathbf{x}$  given):

$$\Rightarrow \epsilon = O\left(\left(\frac{K}{M} \log \frac{MN}{K}\right)^{1/2}\right)$$

Uniform case:

$$\Rightarrow \epsilon = O\left(\left(\frac{K}{M} \log \frac{MN}{K}\right)^{1/6}\right)$$

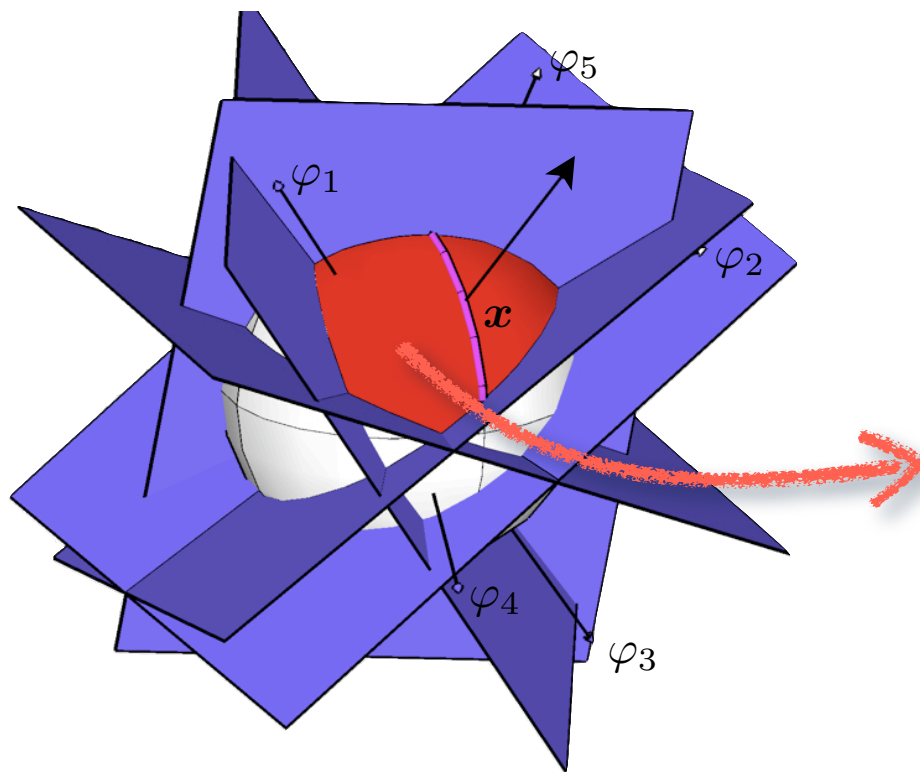
Y. Plan, R. Vershynin, "Robust 1-bit compressed sensing and sparse logistic regression: a convex programming approach", IEEE TIT 2012, arXiv:1202.1212.

LJ, K. Degraux, C. De Vleeschouwer, "Quantized Iterative Hard Thresholding: Bridging 1-bit and High-Resolution Quantized Compressed Sensing", SAMPTA2013

# Initial approach

- ▶ Let  $\mathbf{q} = \text{sign}(\Phi \mathbf{x}) =: A(\mathbf{x})$
- ▶ Initially: [Boufounos, Baraniuk 2008]

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{u}} \|\mathbf{u}\|_1 \quad \text{s.t.} \quad \text{diag}(\mathbf{q}) \Phi \mathbf{u} > 0 \quad \text{and} \quad \|\mathbf{u}\|_2 = 1$$



**Non-convex!** 2 numerical choices :

1. relax + projection on  $S^{N-1}$
2. “trust region methods”  
→ *Restricted-Step Shrinkage (RSS)*

**Consistency constraint:**

$$\begin{aligned} & \{\mathbf{u} \in \mathbb{R}^N \cap S^{N-1} : \mathbf{q} = A(\mathbf{u})\} \\ & \Leftrightarrow \{\mathbf{u} \in \mathbb{R}^N \cap S^{N-1} : \text{diag}(\mathbf{q}) \Phi \mathbf{u} > 0\} \\ & \ni \mathbf{x} \end{aligned}$$

# Initial approach

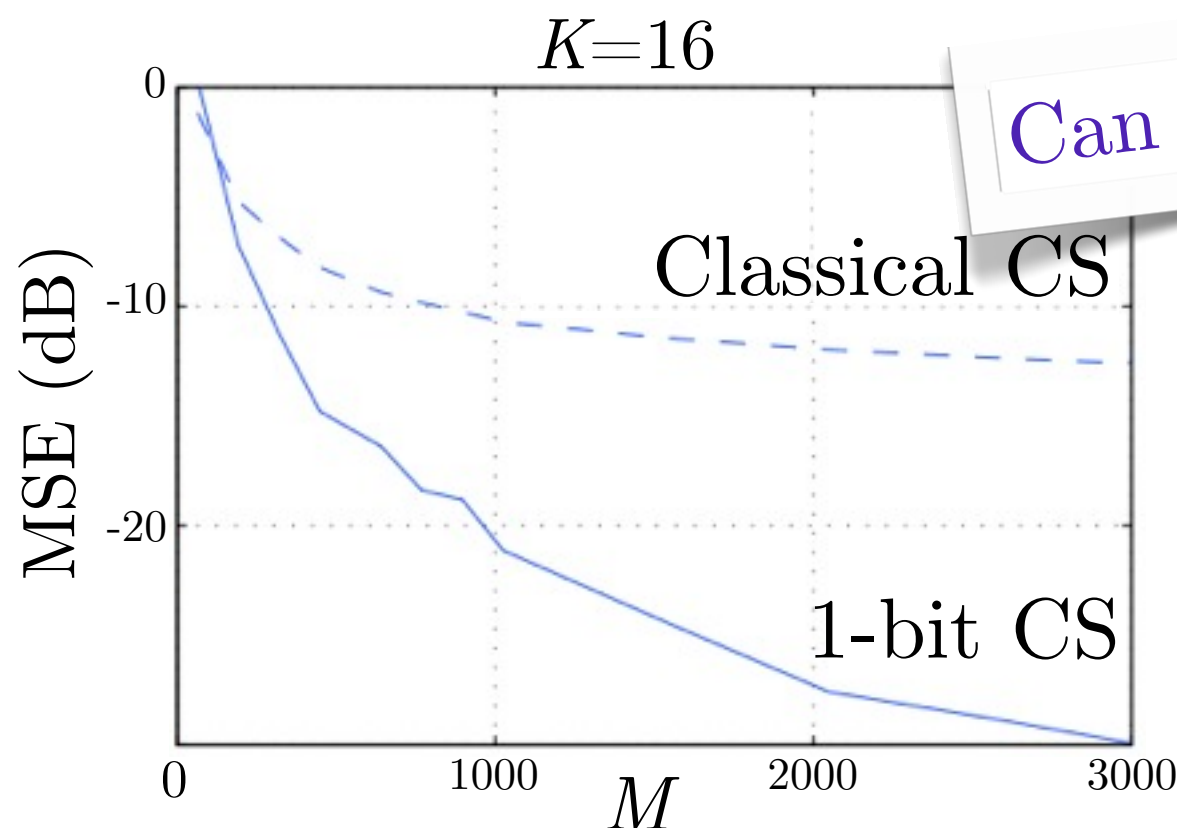
- ▶ Let  $\mathbf{q} = \text{sign}(\Phi \mathbf{x}) =: A(\mathbf{x})$
- ▶ Initially: [Boufounos, Baraniuk 2008]

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{u}} \|\mathbf{u}\|_1 \quad \text{s.t.} \quad \text{diag}(\mathbf{q}) \Phi \mathbf{u} > 0 \quad \text{and} \quad \|\mathbf{u}\|_2 = 1$$

(e.g., take  
the 1<sup>st</sup> choice)

(relaxed) 
$$\hat{\mathbf{x}} = \arg \min_{\mathbf{u}} \|\mathbf{u}\|_1 + \lambda \|(\text{diag}(\mathbf{q}) \Phi \mathbf{u})_-\|^2 \quad \text{s.t.} \quad \|\mathbf{u}\|_2 = 1$$

→ Solved by projected gradient descent



Can we do better?

gain brought  
by (almost)  
consistency

# Other methods:

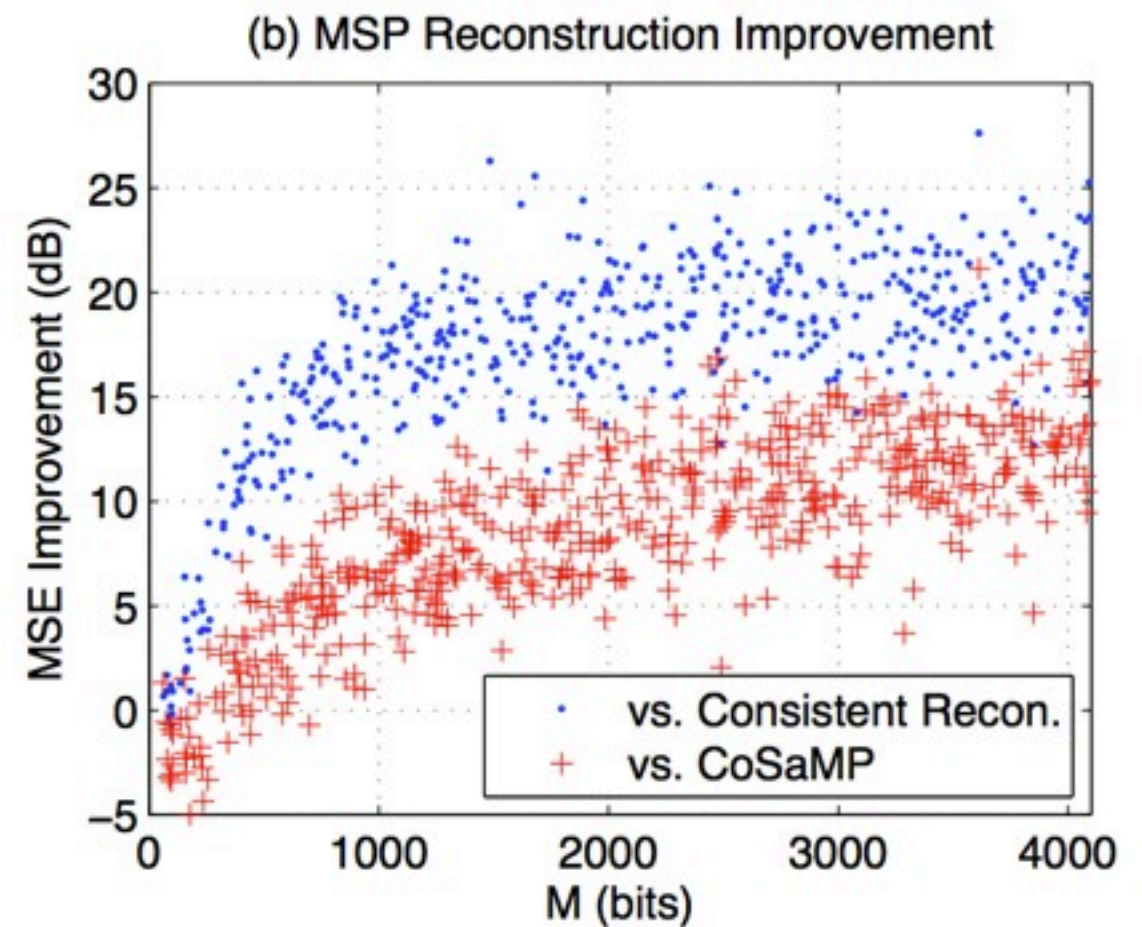
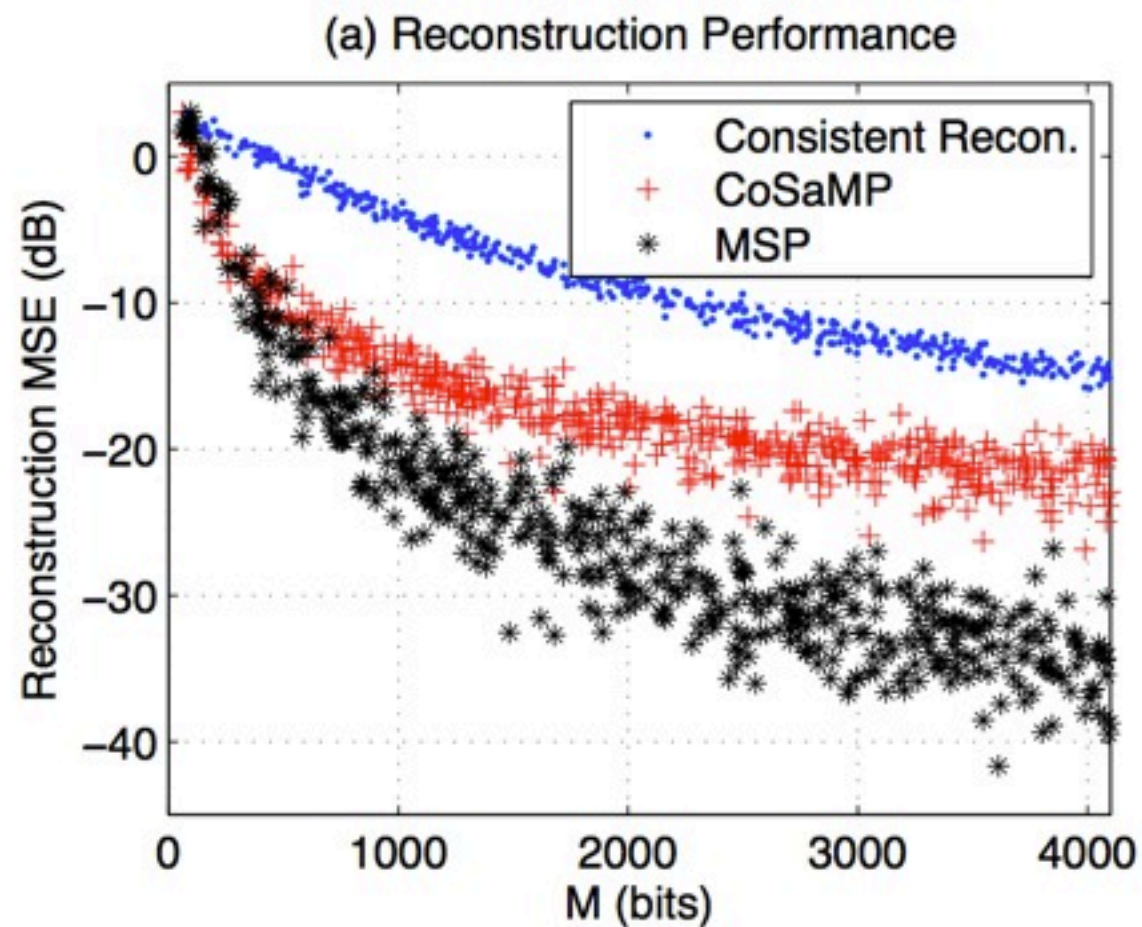
- ✓ ▶ Matching Sign Pursuit [Boufounos]
- ▶ Restricted-Step Shrinkage (RSS) [Laska, We, Yin, Baraniuk]
- ✓ ▶ Binary Iterative Hard Thresholding [Jacques, Laska, Boufounos, Baraniuk]
- ✓ ▶ Convex Optimization [Plan, Vershynin]
- ▶ ...

# Matching Sign Pursuit (MSP)

- ▶ Iterative greedy algorithm, similar to CoSaMP [Needell, Tropp, 08]
- ▶ Maintains running signal estimate and its support  $T$ .
- ▶ MSP iteration:
  - ▶ Identify **sign violations**  $\rightarrow \mathbf{r} = (\text{diag}(\mathbf{y}) \Phi \hat{\mathbf{x}})_-$
  - ▶ Compute **proxy**  $\rightarrow \mathbf{p} = \Phi^T \mathbf{r}$
  - ▶ Identify **support**  $\rightarrow \Omega = \text{supp } \mathbf{p}|_{2K} \cup T$
  - ▶ **Consistent Reconstruction** over support estimate:
$$\mathbf{b}|_{\Omega} = \arg \min_{\mathbf{u} \in \mathbb{R}^N} \|(\text{diag}(\mathbf{y}) \Phi \mathbf{u})_-\|_2^2 \text{ s.t. } \|\mathbf{u}\|_2 = 1 \text{ and } \mathbf{u}|_{T^c} = 0$$
  - ▶ Truncate, normalize, and **update** estimate:  $\hat{\mathbf{x}} \leftarrow \mathbf{b}|_K / \|\mathbf{b}|_K\|_2$



# Matching Sign Pursuit (MSP)



Boufounos, P. T. (2009, November). "Greedy sparse signal reconstruction from sign measurements".

In Signals, Systems and Computers, 2009 Conference Record of the Forty-Third Asilomar Conference on (pp. 1305-1309). IEEE.

# Binary Iterative Hard Thresholding

Given  $\mathbf{q} = A(\mathbf{x})$  and  $K$ , set  $l = 0$ ,  $\mathbf{x}^0 = 0$ :

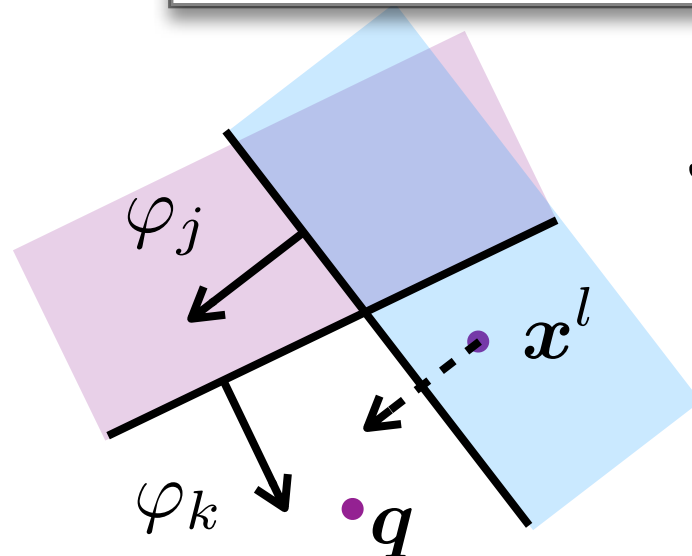
$$\begin{aligned} \mathbf{a}^{l+1} &= \mathbf{x}^l + \frac{\tau}{2} \Phi^T (\mathbf{q} - A(\mathbf{x}^l)), \\ \mathbf{x}^{l+1} &= \mathcal{H}_K(\mathbf{a}^{l+1}), \quad l \leftarrow l + 1 \end{aligned}$$

(“gradient” towards consistency)  
( $\tau > 0$  controls gradient descent)  
(proj.  $K$ -sparse signal set)

with  $\mathcal{H}_K(\mathbf{u}) = K$ -term hard thresholding

Stop when  $d_H(\mathbf{q}, A(\mathbf{x}^{l+1})) = 0$  or  $l = \text{max. iter.}$

minimizes  $\mathcal{J}(\mathbf{x}') = \|\text{diag}(\mathbf{q})(\Phi \mathbf{x}')\|_1$  with  $(\lambda)_- = (\lambda - |\lambda|)/2$



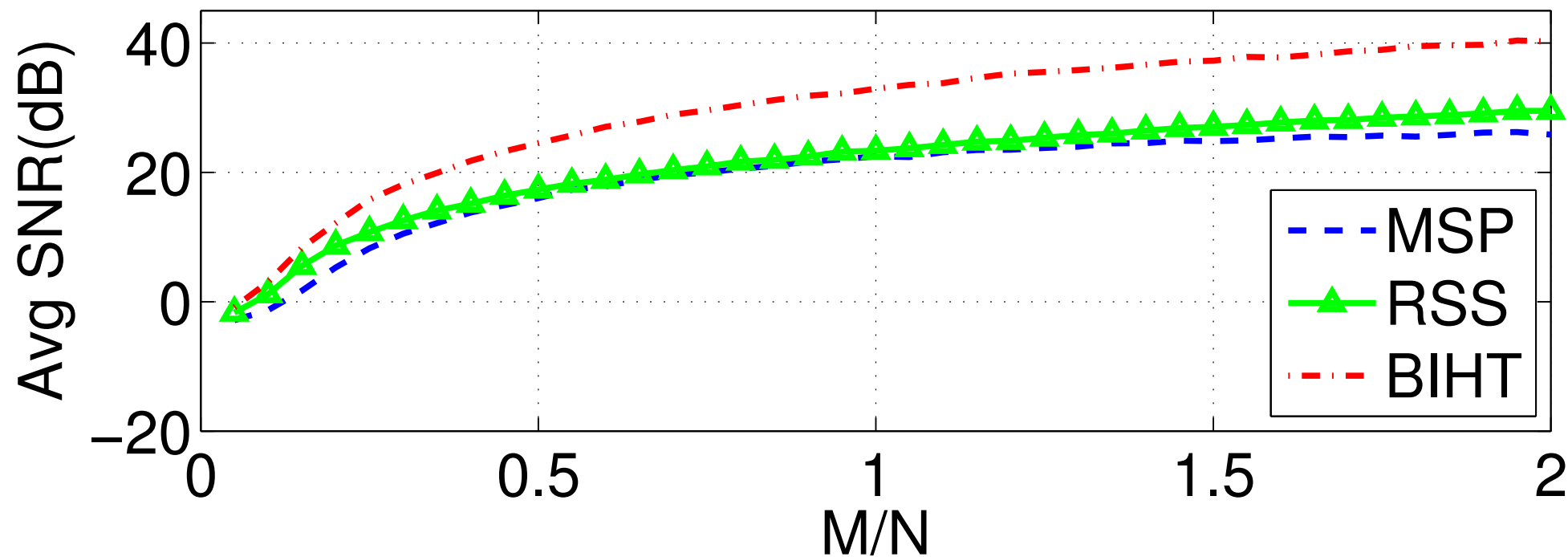
$$\mathcal{J}(\mathbf{x}') = \sum_{j=1}^M \left| \overbrace{(\text{sign}(\langle \varphi_j, \mathbf{x} \rangle) \langle \varphi_j, \mathbf{x}' \rangle)}^{q_j} \right|_-$$

$$q_k - A(\mathbf{x}^l)_k = 0$$

$$q_j - A(\mathbf{x}^l)_j > 0$$

(connections with ML hinge loss, 1-bit classification)

# Binary Iterative Hard Thresholding



$N = 1000, K = 10$

Bernoulli-Gaussian model

normalized signals

1000 trials

Matching Sign pursuit (MSP)

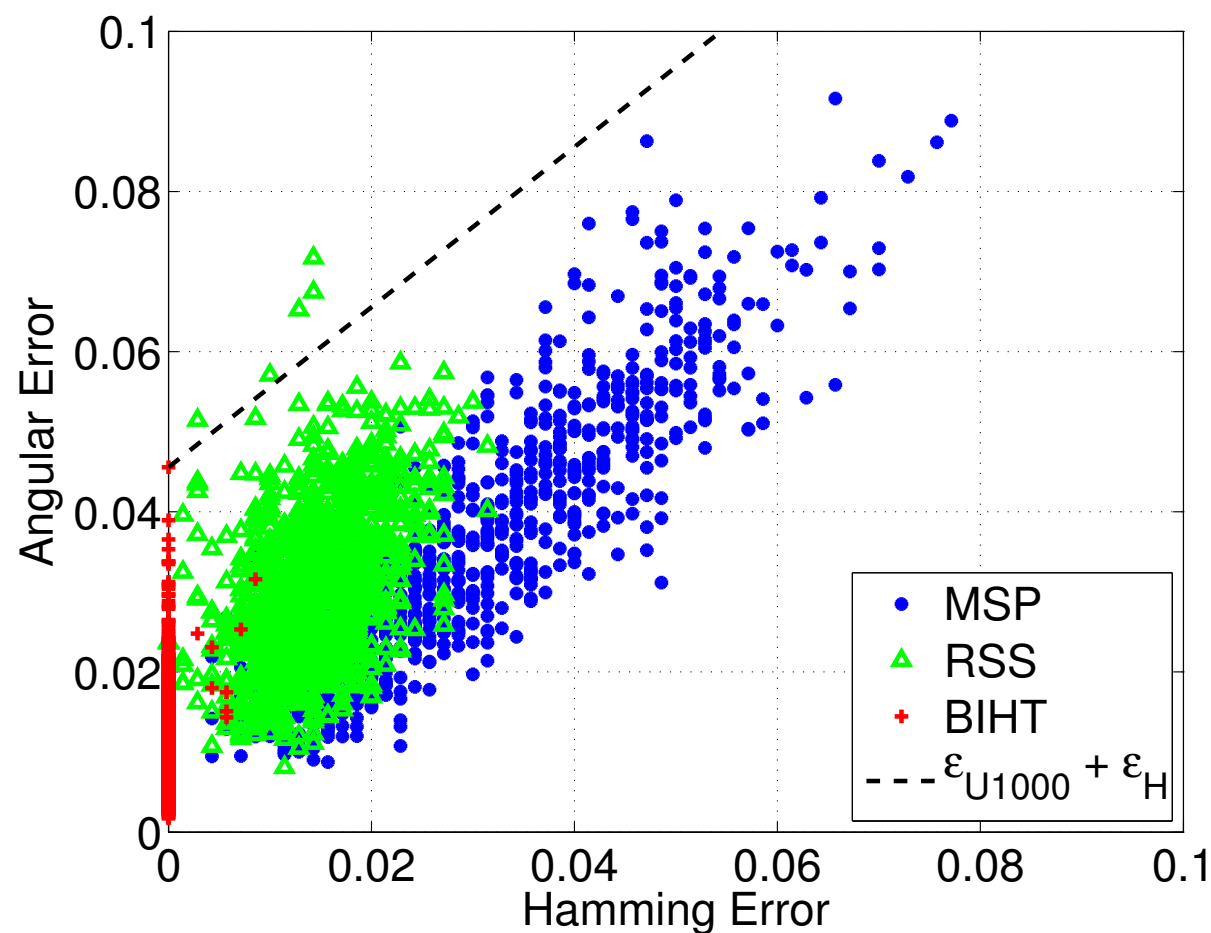
Restricted-Step Shrinkage (RSS)

Binary Iterative Hard Thresholding (BIHT)

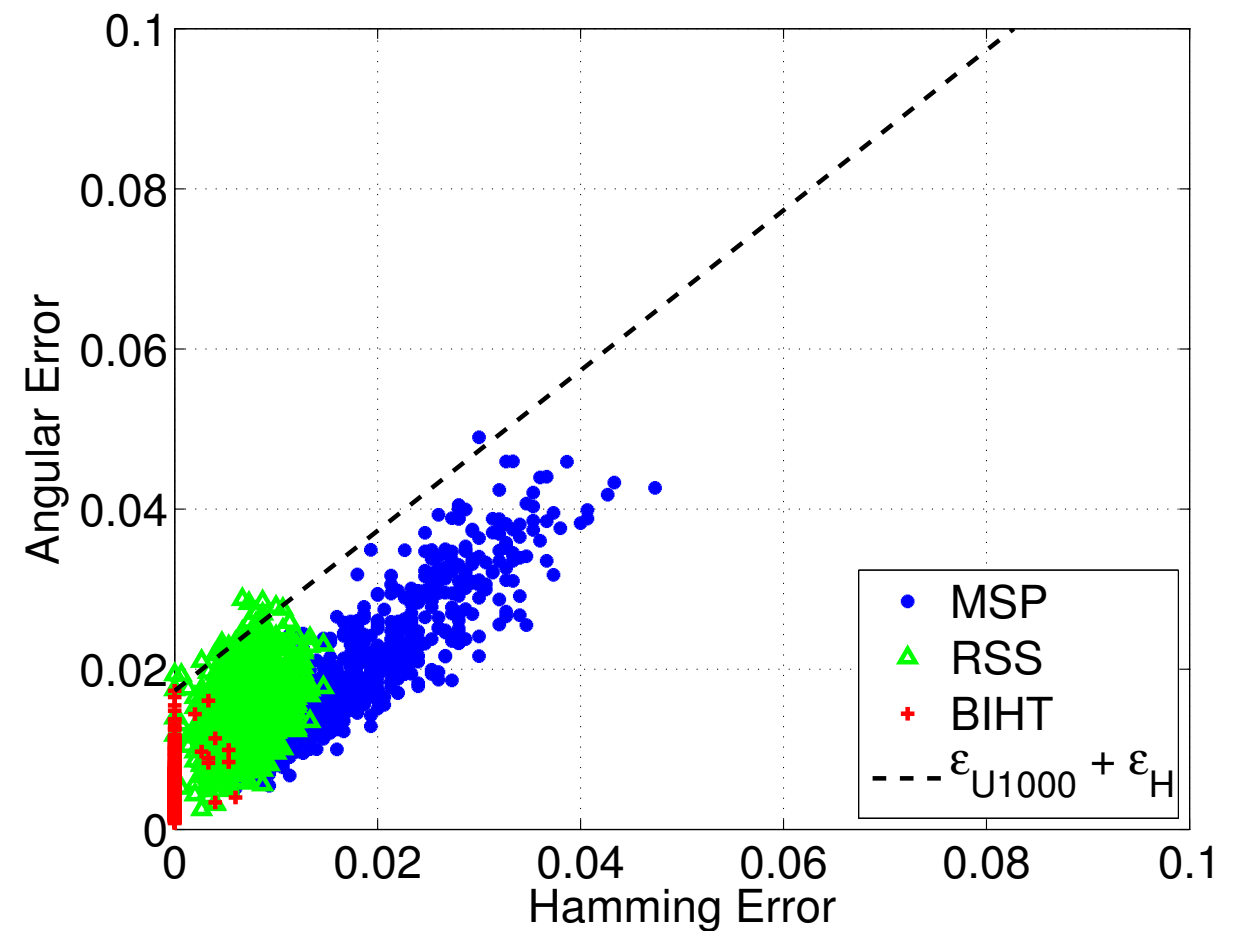


# Binary Iterative Hard Thresholding

- Testing B $\epsilon$ SE:  $d_{\text{ang}}(\mathbf{x}, \mathbf{x}^*) \leq d_H(A(\mathbf{x}), A(\mathbf{x}^*)) + \epsilon(M)$

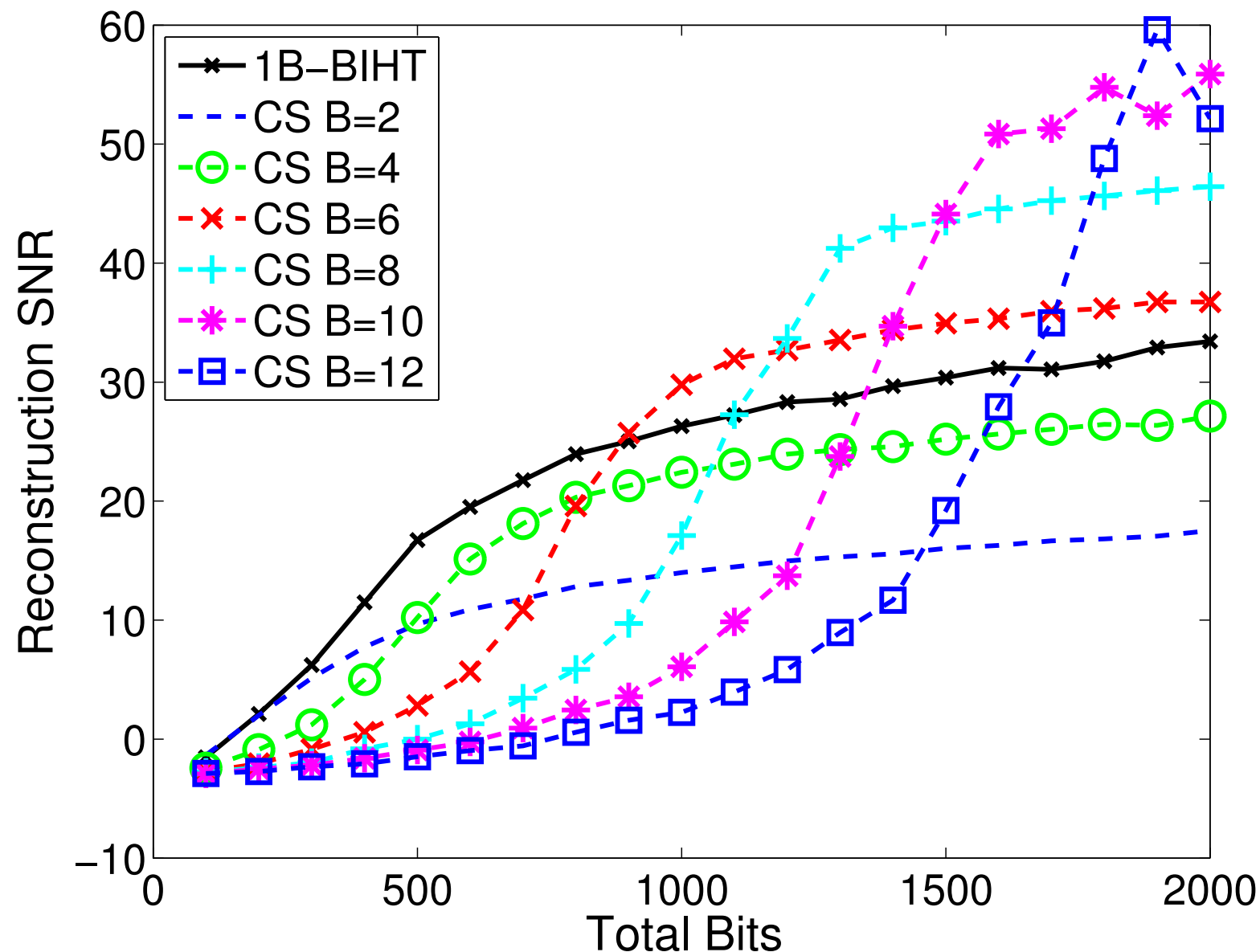


$$M/N = 0.7$$



$$M/N = 1.5$$

# Remark: CS vs bits/meas.



$$N = 2000, K = 20$$

Bernoulli-Gaussian model  
normalized signals

$B$  bits/measurement

$$B = 1, \dots, 12$$

$$M = \text{Total Bits} / B$$

1000 trials

# Convex Optimization [Plan, Vershynin, 12]

Let  $\mathbf{q} = \text{sign}(\Phi \mathbf{x})$  for some signal  $\mathbf{x} \in \mathcal{K} \subset B_2^N$  e.g., sparse,  
compressible,  
low-rank matrix

Compute  $\hat{\mathbf{x}} = \arg \max_{\mathbf{u} \in \mathbb{R}^N} \mathbf{q}^T \Phi \mathbf{u} \quad \text{s.t.} \quad \mathbf{u} \in \mathcal{K}$

 maximize consistency

Convex problem if  $\mathcal{K}$  convex!

No ambiguous amplitude definition

( $\mathbf{u} = 0$  avoided)

Remark: (PV-L0 problem) [Bahmani, Boufounos, Raj, 13]

$$\hat{\mathbf{x}} = \frac{1}{\|\mathcal{H}_K(\Phi^* \mathbf{q})\|} \mathcal{H}_K(\Phi^* \mathbf{q}) \text{ if } \mathcal{K} = \Sigma_K !!$$

# Convex Optimization

[Plan, Vershynin, 12]

Let  $\mathbf{q} = \text{sign}(\Phi \mathbf{x})$  for some signal  $\mathbf{x} \in \mathcal{K} \subset B_2^N$

e.g., sparse,  
compressible,  
low-rank matrix

Compute  $\hat{\mathbf{x}} = \arg \max_{\mathbf{u} \in \mathbb{R}^N} \mathbf{q}^T \Phi \mathbf{u} \quad \text{s.t.} \quad \mathbf{u} \in \mathcal{K}$

maximize  
consistency

**Proposition** (assuming  $\|\mathbf{x}\| = 1$ ) For some  $C, c > 0$ , if  $M \geq C\epsilon^{-6}w^2(\mathcal{K})$ , then, with  $Pr \geq 1 - e^{-c\epsilon^2 M}$ , we have  $\|\hat{\mathbf{x}} - \mathbf{x}\|^2 \leq \sqrt{\frac{\pi}{2}} \epsilon$ .

-2 if  $\mathbf{x}$  is fixed

# Convex Optimization

[Plan, Vershynin, 12]

Let  $\mathbf{q} = \text{sign}(\Phi \mathbf{x})$  for some signal  $\mathbf{x} \in \mathcal{K} \subset B_2^N$

Compute  $\hat{\mathbf{x}} = \arg \max_{\mathbf{u} \in \mathbb{R}^N} \mathbf{q}^T \Phi \mathbf{u} \quad \text{s.t.} \quad \mathbf{u} \in \mathcal{K}$

**Proposition** (assuming  $\|\mathbf{x}\| = 1$ ) For some  $C, c > 0$ , if  $M \geq C\epsilon^{-6}w^2(\mathcal{K})$ , then, with  $Pr \geq 1 - e^{-c\epsilon^2 M}$ , we have  $\|\hat{\mathbf{x}} - \mathbf{x}\|^2 \leq \sqrt{\frac{\pi}{2}} \epsilon$ .

+ Robust to noise: noise (bit flip)

Let  $\mathbf{q}_n = \text{diag}(\boldsymbol{\eta}) \mathbf{q}$  with  $\eta_i \in \{\pm 1\}^M$ , and assume  $d_H(\mathbf{q}, \mathbf{q}_n) \leq p$  noise power

(under the same conditions)

$$\|\hat{\mathbf{x}} - \mathbf{x}\|^2 \leq \epsilon \sqrt{\log e / \epsilon} + 11 p \sqrt{\log e / p}$$

Note: if  $M = O(\epsilon^{-2}(p - 1/2)^{-2} K \log N / K)$   
this term disappear if  $\eta_i = \pm 1$  are iid RVs (with  $P(\eta_i = 1) = p$ )

# 5. Playing with thresholds in 1-bit CS

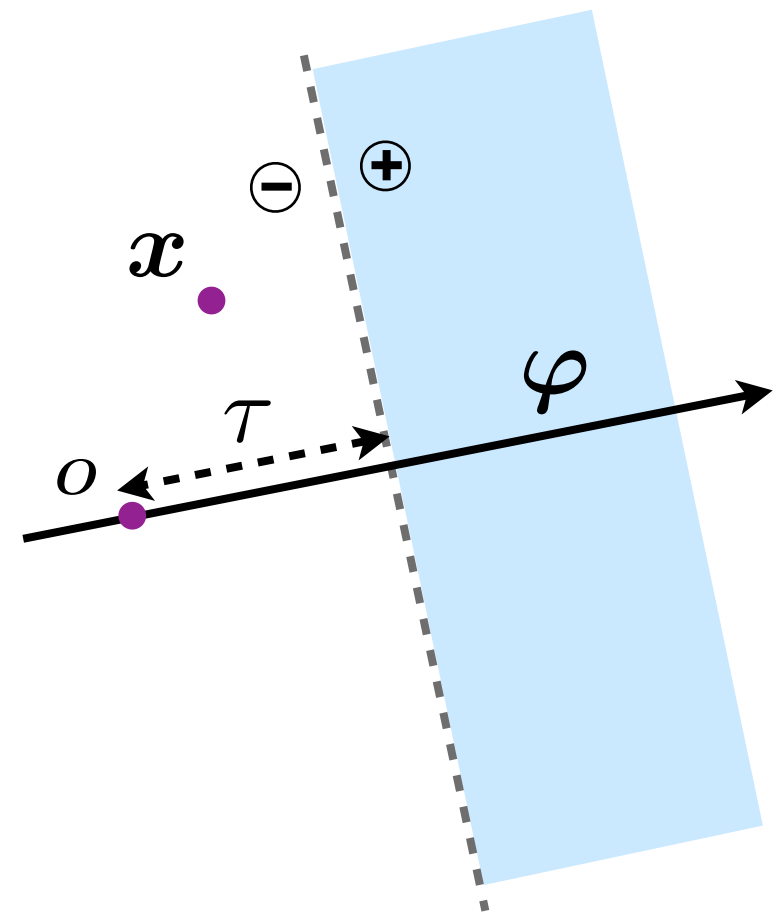
# Thresholds?

- ▶ Given  $\mathbf{x} \in \mathbb{R}^N$  (e.g., sparse)  
Is there an interest in sensing

$$\text{sign}(\langle \boldsymbol{\varphi}, \mathbf{x} \rangle - \tau)$$

for some (random)  $\boldsymbol{\varphi}$  and  $\tau \in \mathbb{R}$ ?

- ▶ Two recent applications:
  - ▶ adaptive thresholds [Kamilov, Bourquard, Amini, Unser, 12]
  - ▶ bridging 1-bit and  $B$ -bits QCS [LJ, Degraux, De Vleeschouwer, 13]

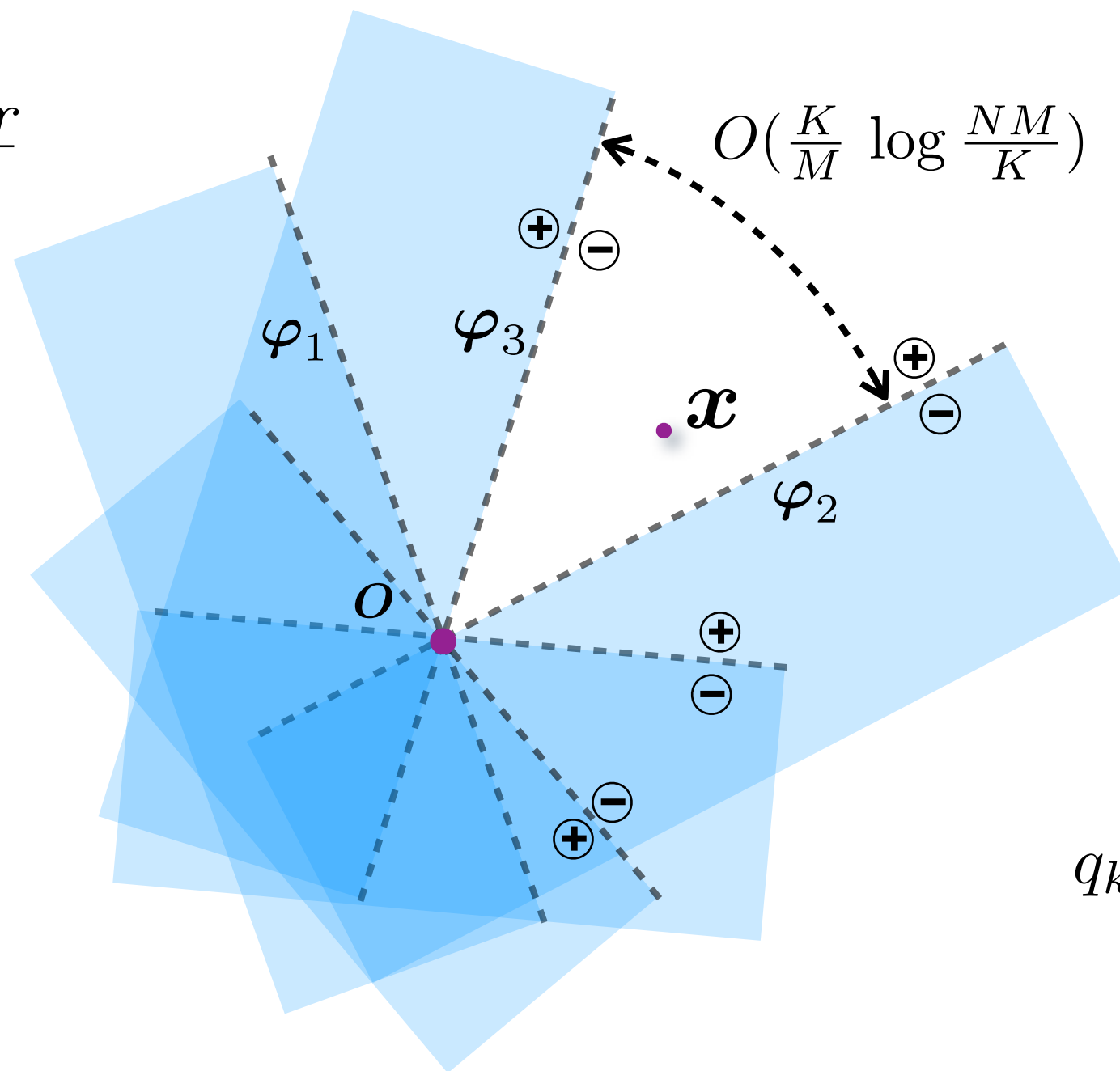




# 1-bit CS with adaptive thresholds

Non-adaptive 1-bit CS ( $\tau = 0$ )

Reminder



$$q_k = \text{sign}(\langle \varphi_k, x \rangle)$$

# 1-bit CS with adaptive thresholds

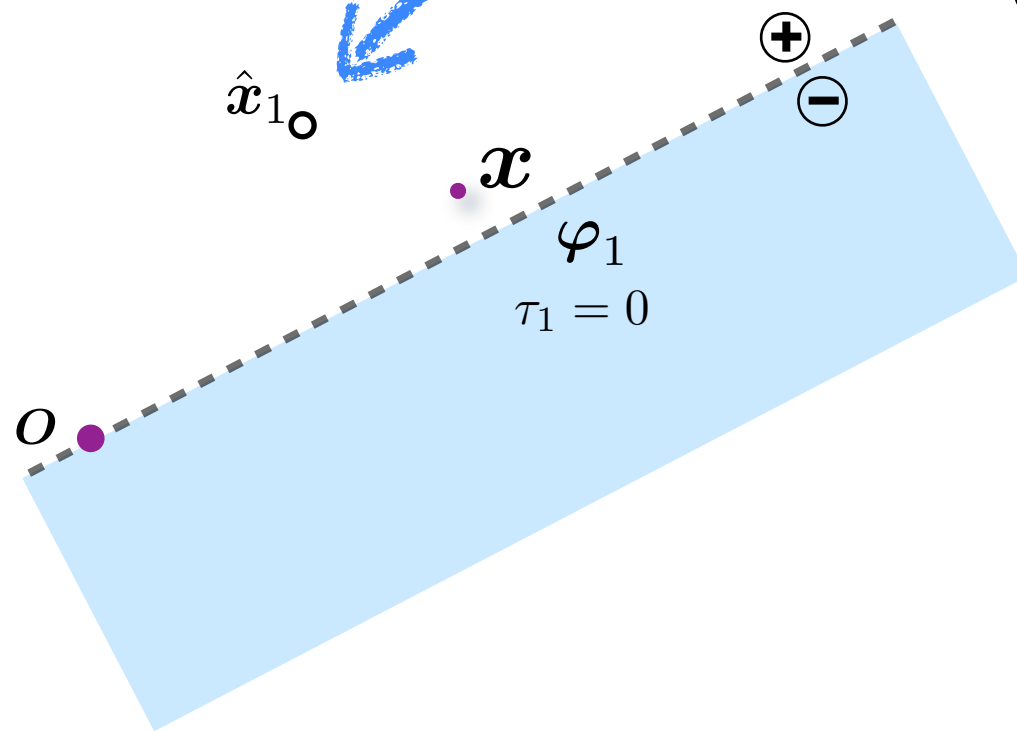
## Adaptive 1-bit CS [Kamilov, Bourquard, Amini, Unser, 12]

Given a decoder  $\text{Rec}()$

adapted from prev. meas.

$$q_k = \text{sign} \left( \langle \varphi_k, \mathbf{x} \rangle - \tau_k \right)$$

$$\begin{cases} \hat{\mathbf{x}}_k := \text{Rec}(y_1, \dots, y_k, \varphi_1, \dots, \varphi_k, \tau_1, \dots, \tau_k) \\ \tau_{k+1} \text{ s.t. } \langle \varphi_{k+1}, \hat{\mathbf{x}}_k \rangle - \tau_{k+1} = 0 \end{cases}$$



U.S. Kamilov, A. Bourquard, A. Amini, M. Unser,

“One-bit measurements with adaptive thresholds”. *Signal Processing Letters, IEEE*, 19(10), 607-610.

# 1-bit CS with adaptive thresholds

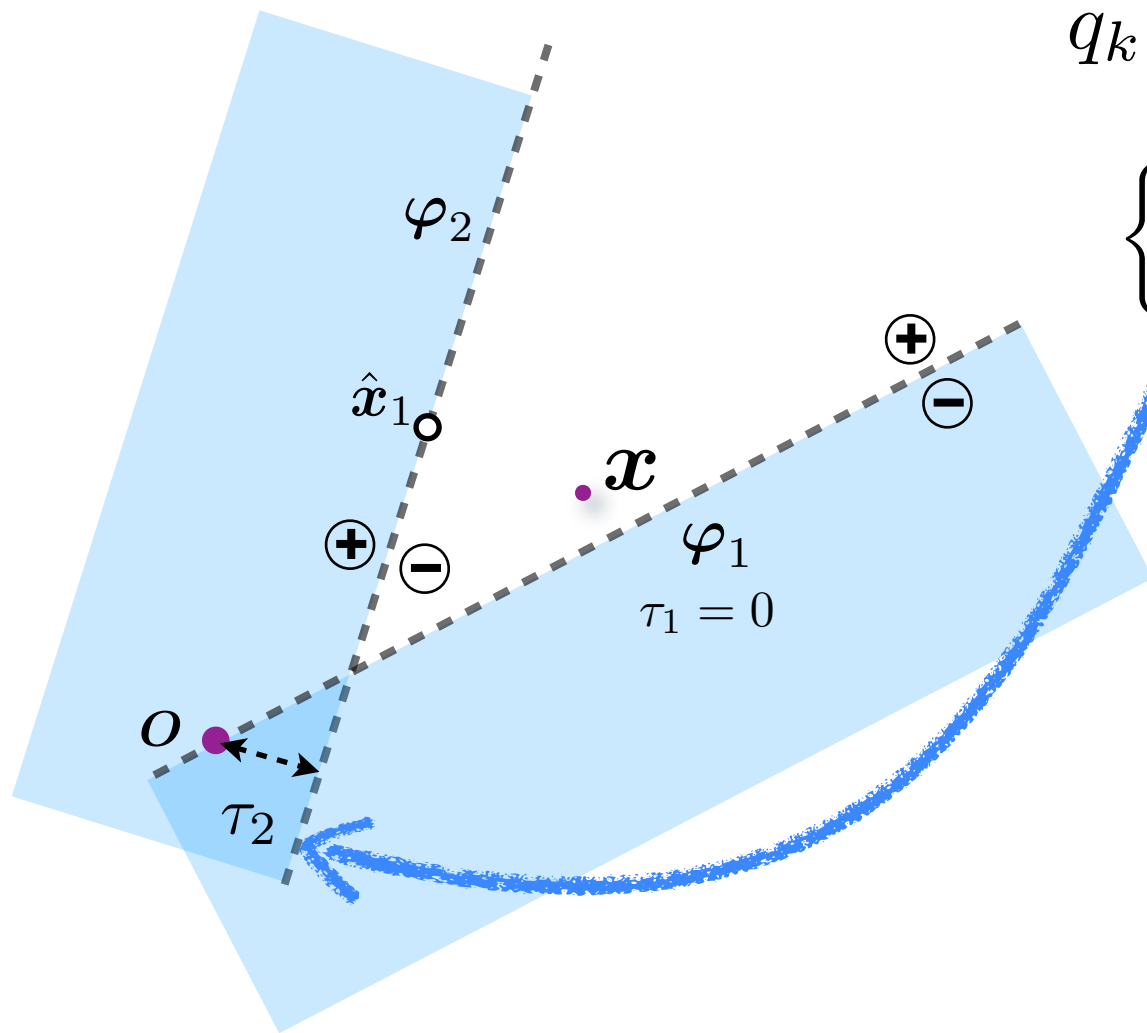
Adaptive 1-bit CS [Kamilov, Bourquard, Amini, Unser, 12]

Given a decoder  $\text{Rec}()$

adapted from prev. meas.

$$q_k = \text{sign}(\langle \varphi_k, \mathbf{x} \rangle - \tau_k)$$

$$\begin{cases} \hat{\mathbf{x}}_k := \text{Rec}(y_1, \dots, y_k, \varphi_1, \dots, \varphi_k, \tau_1, \dots, \tau_k) \\ \tau_{k+1} \text{ s.t. } \langle \varphi_{k+1}, \hat{\mathbf{x}}_k \rangle - \tau_{k+1} = 0 \end{cases}$$



U.S. Kamilov, A. Bourquard, A. Amini, M. Unser,

“One-bit measurements with adaptive thresholds”. Signal Processing Letters, IEEE, 19(10), 607-610.

# 1-bit CS with adaptive thresholds

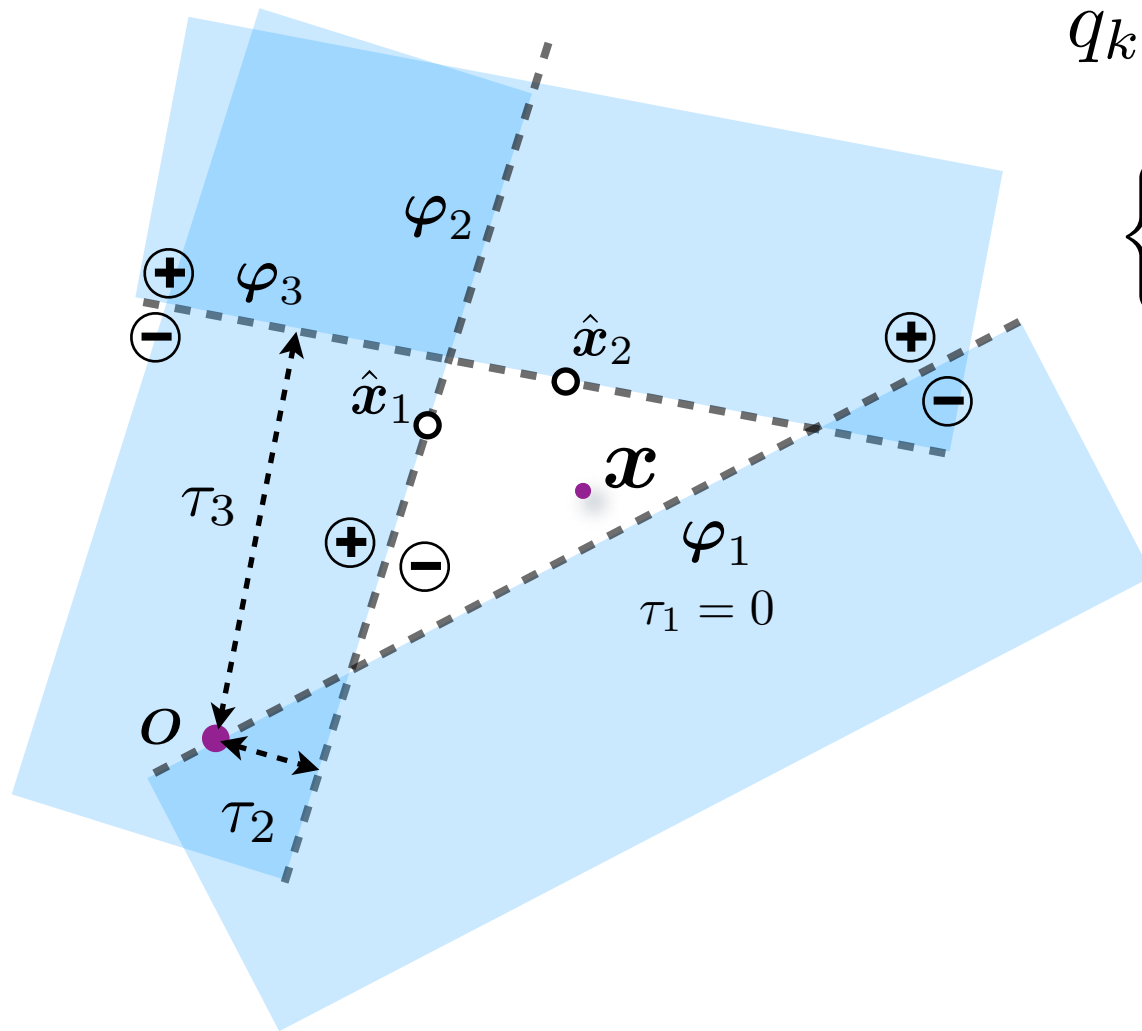
Adaptive 1-bit CS [Kamilov, Bourquard, Amini, Unser, 12]

Given a decoder  $\text{Rec}()$

adapted from prev. meas.

$$q_k = \text{sign}(\langle \varphi_k, \mathbf{x} \rangle - \tau_k)$$

$$\begin{cases} \hat{\mathbf{x}}_k := \text{Rec}(y_1, \dots, y_k, \varphi_1, \dots, \varphi_k, \tau_1, \dots, \tau_k) \\ \tau_{k+1} \text{ s.t. } \langle \varphi_{k+1}, \hat{\mathbf{x}}_k \rangle - \tau_{k+1} = 0 \end{cases}$$



U.S. Kamilov, A. Bourquard, A. Amini, M. Unser,

“One-bit measurements with adaptive thresholds”. Signal Processing Letters, IEEE, 19(10), 607-610.

# 1-bit CS with adaptive thresholds

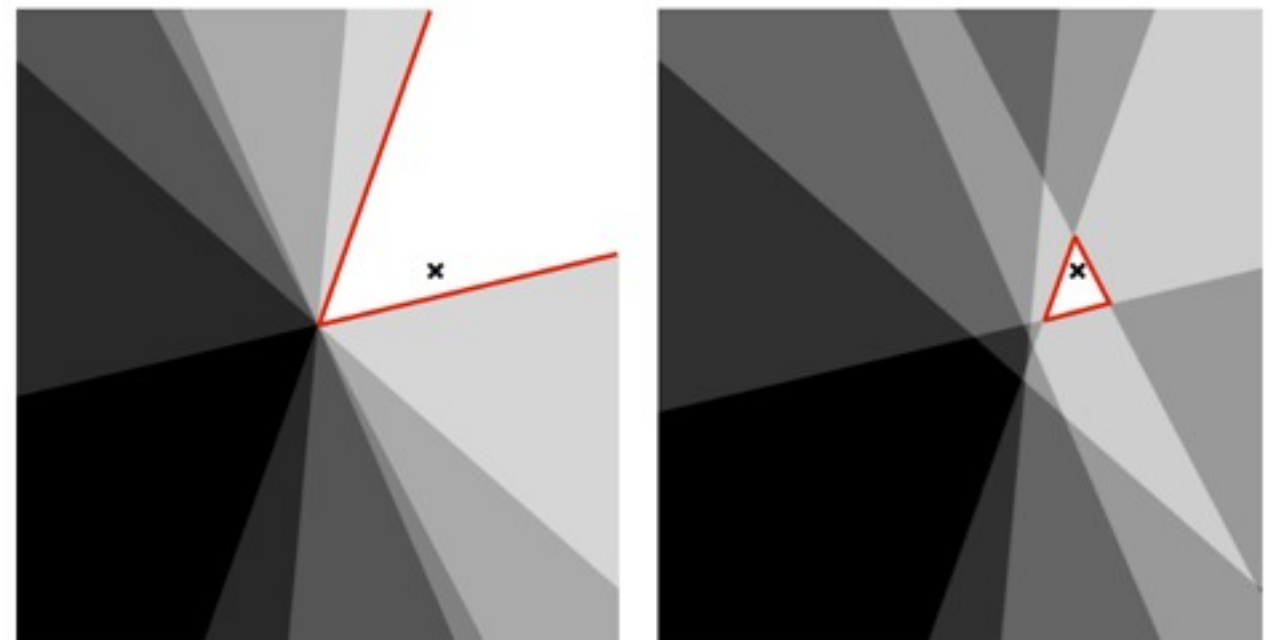
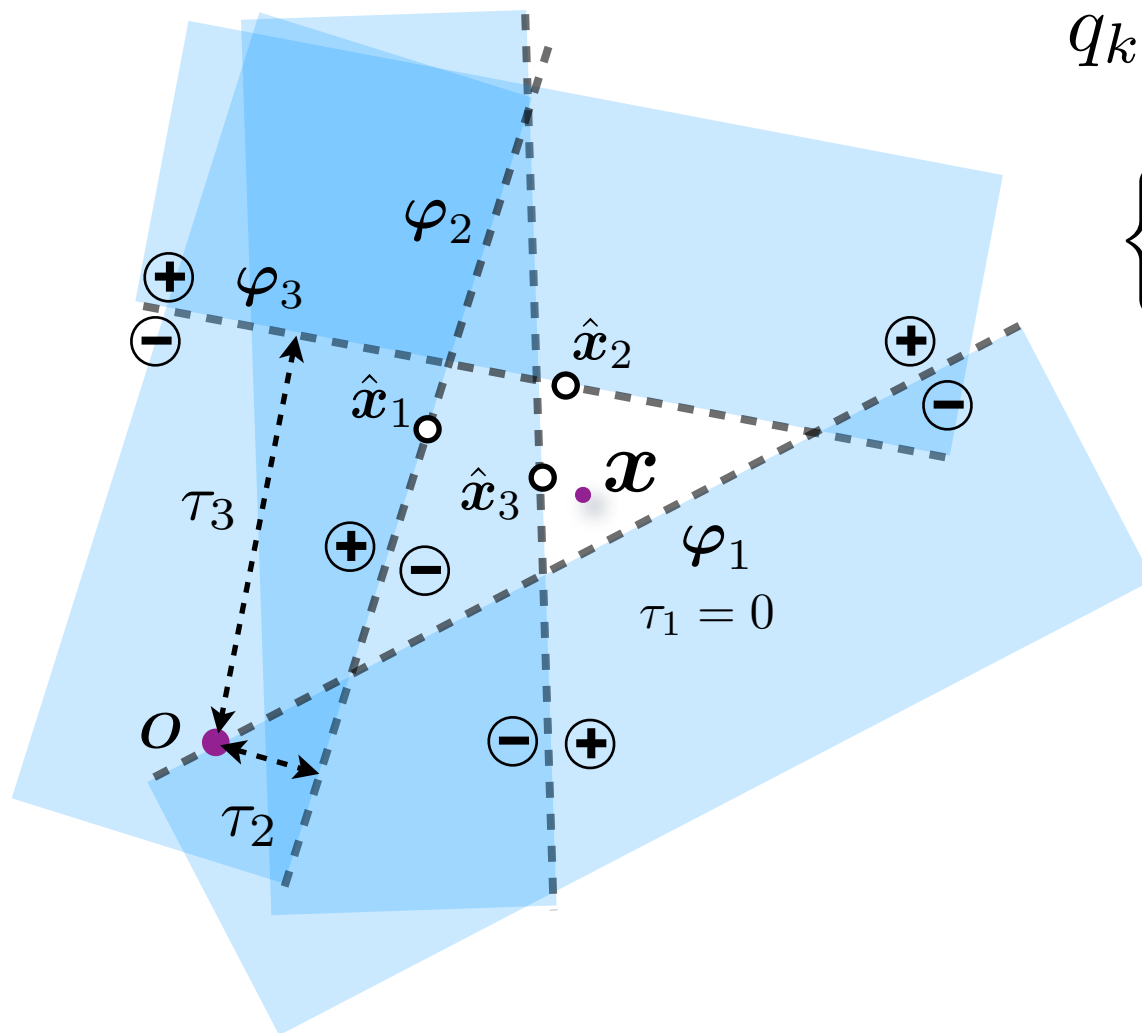
Adaptive 1-bit CS [Kamilov, Bourquard, Amini, Unser, 12]

Given a decoder  $\text{Rec}()$

adapted from prev. meas.

$$q_k = \text{sign}(\langle \varphi_k, \mathbf{x} \rangle - \tau_k)$$

$$\begin{cases} \hat{\mathbf{x}}_k := \text{Rec}(y_1, \dots, y_k, \varphi_1, \dots, \varphi_k, \tau_1, \dots, \tau_k) \\ \tau_{k+1} \text{ s.t. } \langle \varphi_{k+1}, \hat{\mathbf{x}}_k \rangle - \tau_{k+1} = 0 \end{cases}$$

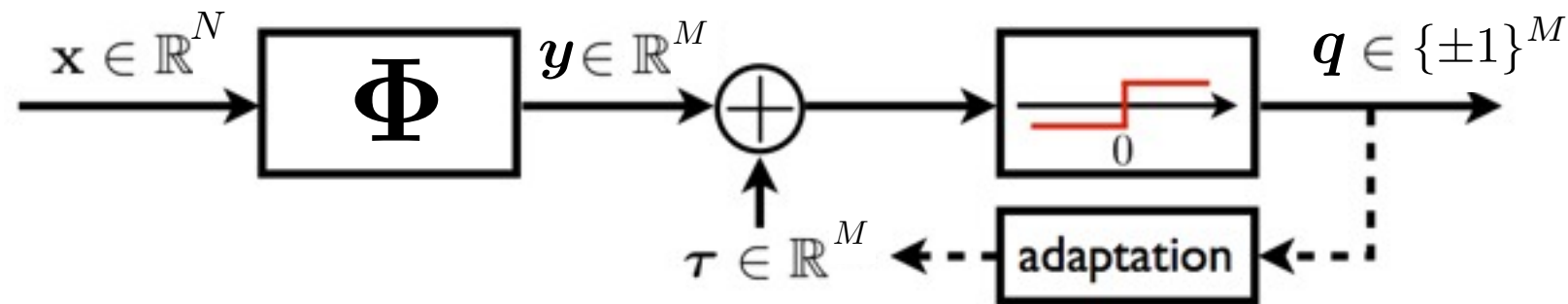


U.S. Kamilov, A. Bourquard, A. Amini, M. Unser,

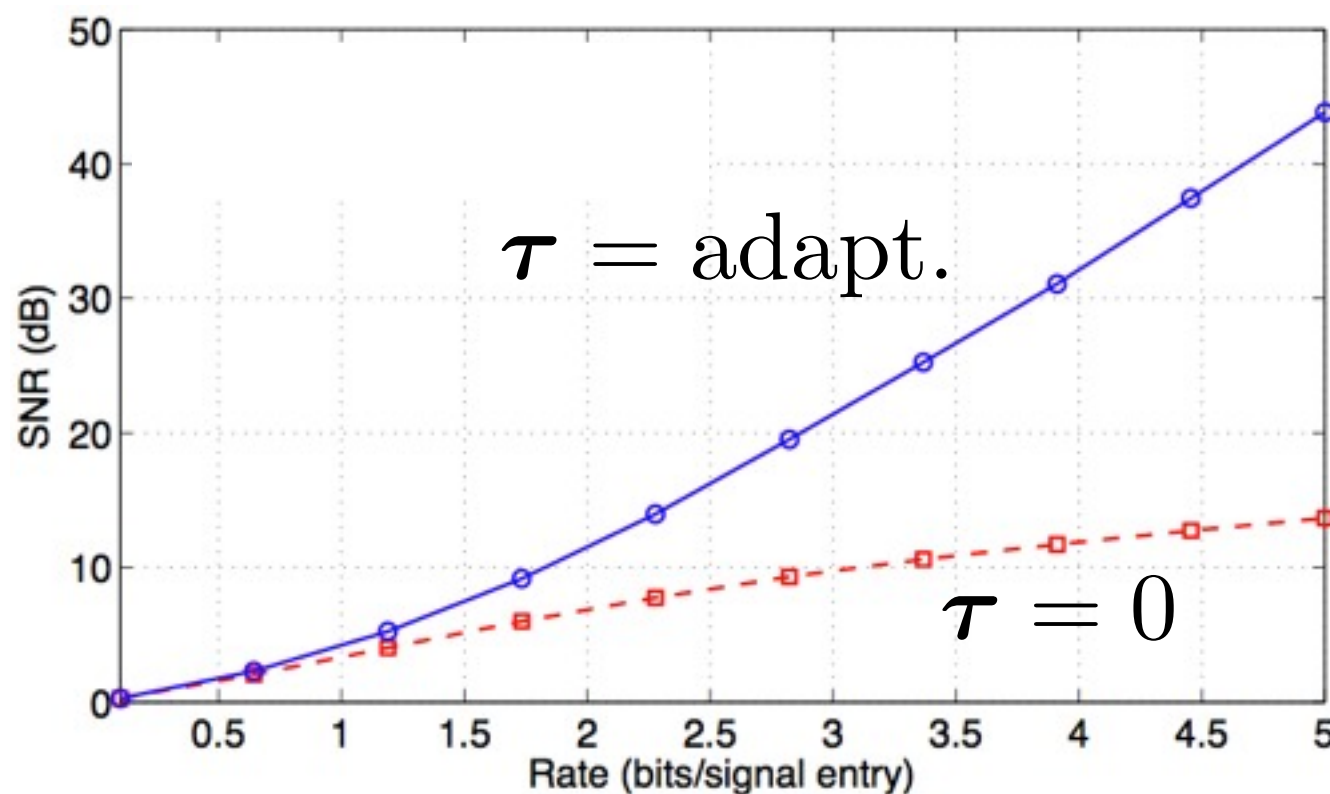
“One-bit measurements with adaptive thresholds”. Signal Processing Letters, IEEE, 19(10), 607-610.

# 1-bit CS with adaptive thresholds

System view:



Kind of  $\Sigma\Delta$  loop



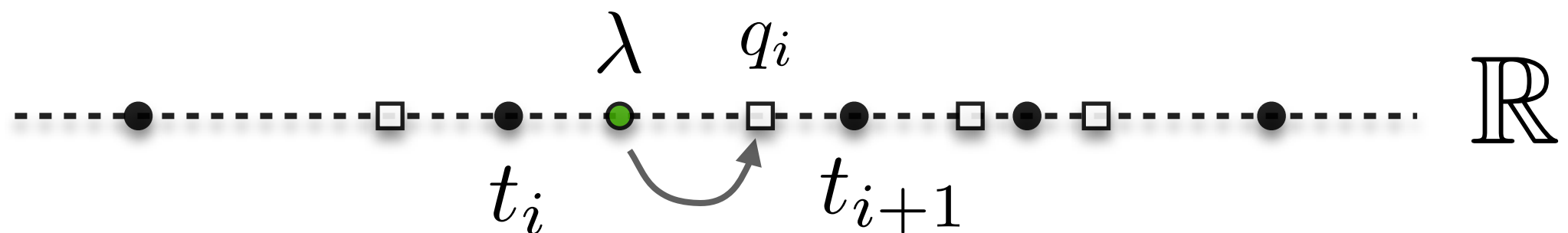
Rec() set to  
Generalized  
Approximate  
Message  
Passing

U.S. Kamilov, A. Bourquard, A. Amini, M. Unser,

“One-bit measurements with adaptive thresholds”. Signal Processing Letters, IEEE, 19(10), 607-610.

# Bridging 1-bit & $B$ -bit CS?

- ▶  $B$ -bit quantizer defined with thresholds:



$$\lambda \in \mathcal{R}_i = [t_i, t_{i+1}) \Leftrightarrow \text{sign}(\lambda - t_i) = +1 \text{ \& \; } \text{sign}(\lambda - t_{i+1}) = -1$$

- ▶ Can we combine multiple thresholds in 1-bit CS?



# Bridging 1-bit & $B$ -bit CS?

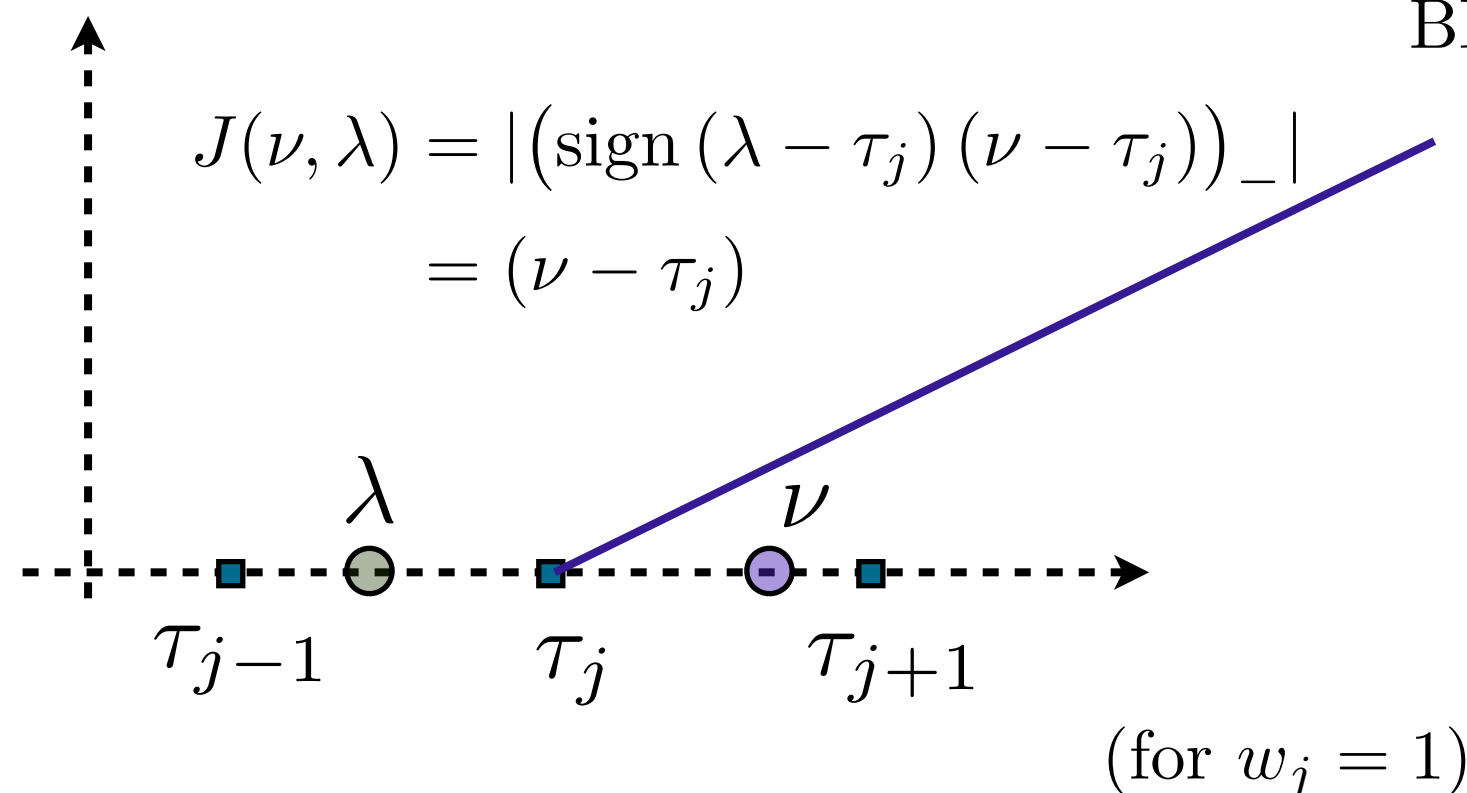
Given  $\mathcal{T} = \{\tau_j\}$  and  $\Omega = \{q_j\}$  ( $|\mathcal{T}| = 2^B + 1 = |\Omega| + 1$ ), let's define

$$J(\nu, \lambda) = \sum_{j=2}^{2^B} w_j \left| \left( \text{sign}(\lambda - \tau_j) (\nu - \tau_j) \right)_- \right|,$$

with  $w_j = q_j - q_{j-1}$ .

Illustration:  $\lambda \in [\tau_{j-1}, \tau_j)$ ,  $\nu \in [\tau_j, \tau_{j+1})$

“delocalized”  
BIHT  $\ell_1$ -sided norm





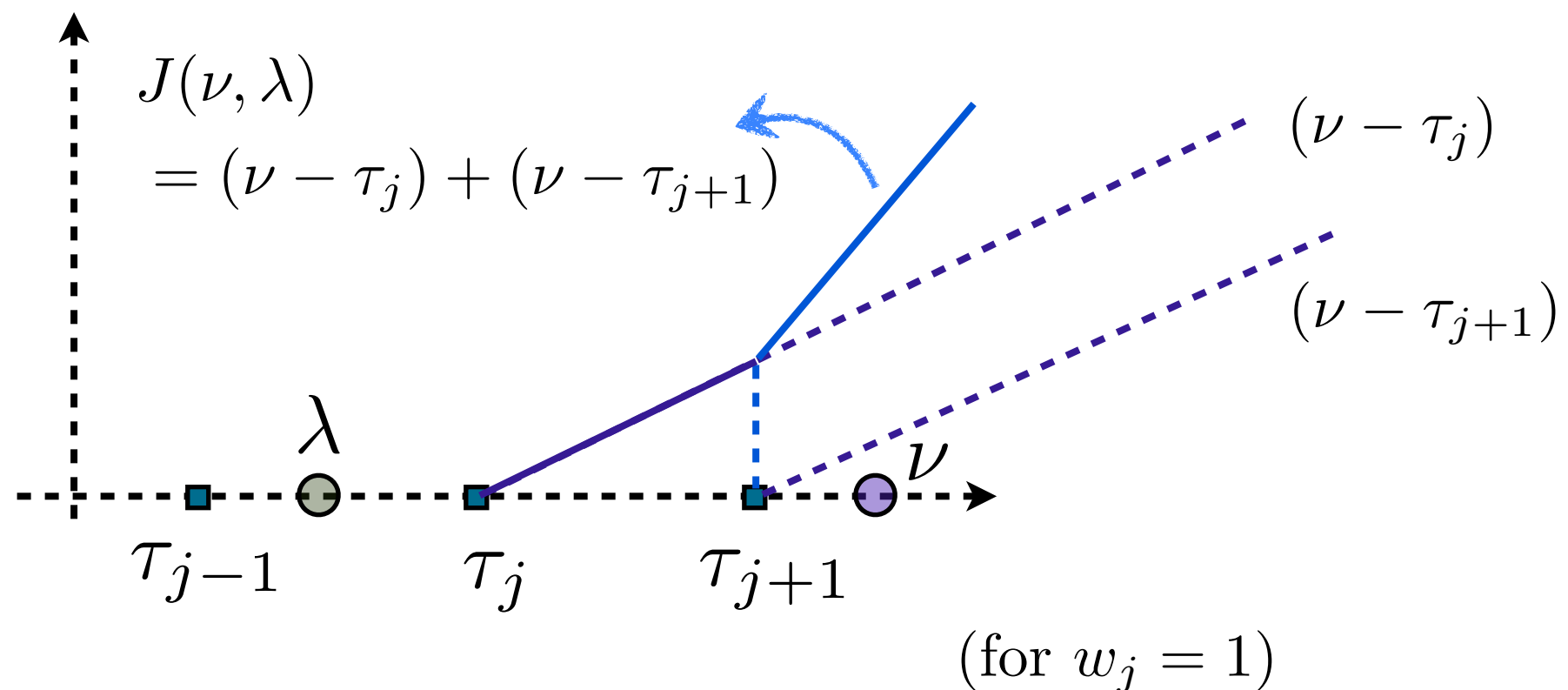
# Bridging 1-bit & $B$ -bit CS?

Given  $\mathcal{T} = \{\tau_j\}$  and  $\Omega = \{q_j\}$  ( $|\mathcal{T}| = 2^B + 1 = |\Omega| + 1$ ), let's define

$$J(\nu, \lambda) = \sum_{j=2}^{2^B} w_j \left| \left( \text{sign}(\lambda - \tau_j) (\nu - \tau_j) \right)_- \right|,$$

with  $w_j = q_j - q_{j-1}$ .

Illustration:  $\lambda \in [\tau_{j-1}, \tau_j)$ ,  $\nu \in [\tau_{j+1}, \tau_{j+2})$



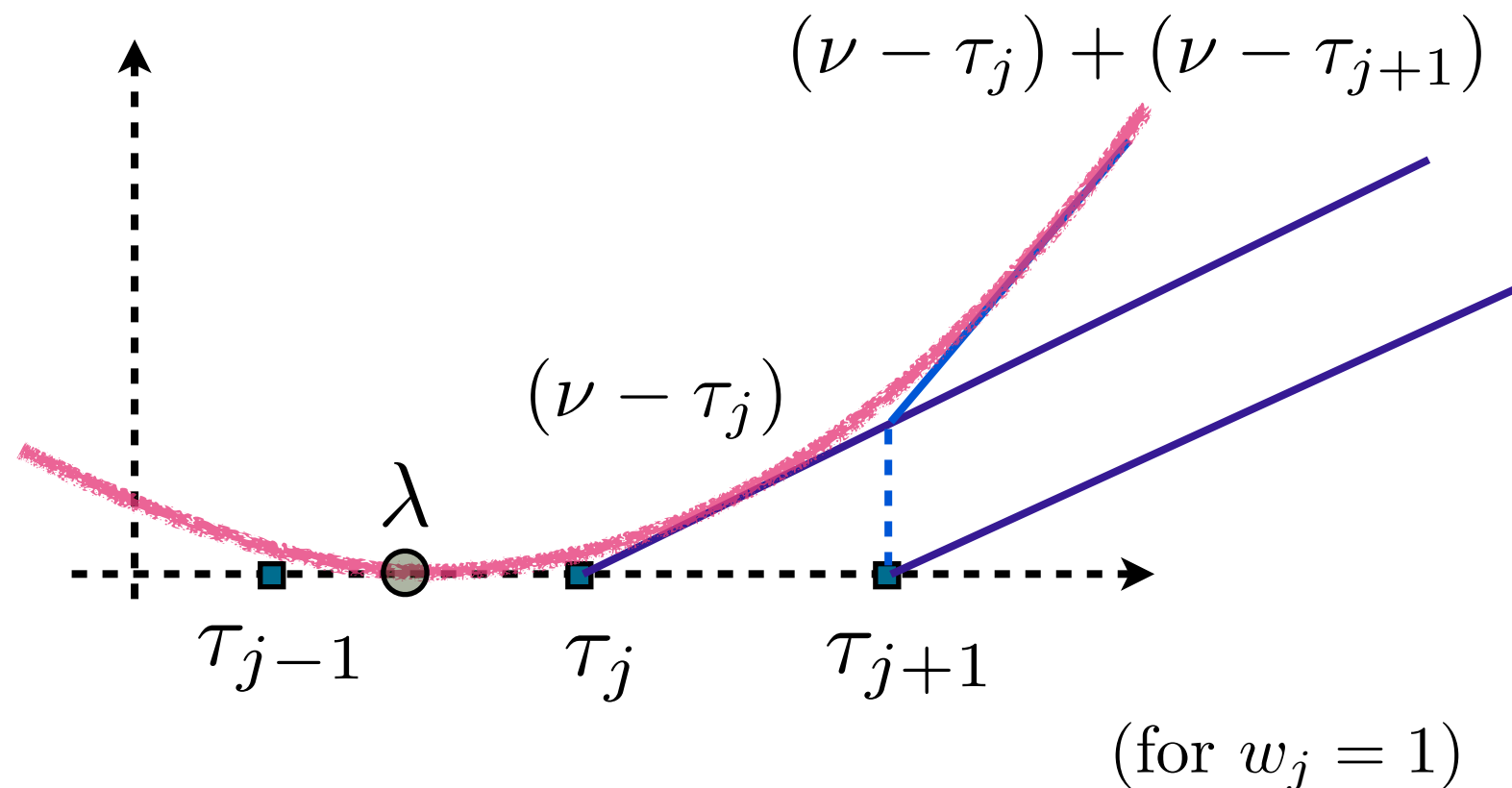
# Bridging 1-bit & $B$ -bit CS?

Given  $\mathcal{T} = \{\tau_j\}$  and  $\Omega = \{q_j\}$  ( $|\mathcal{T}| = 2^B + 1 = |\Omega| + 1$ ), let's define

$$J(\nu, \lambda) = \sum_{j=2}^{2^B} w_j \left| \left( \text{sign}(\lambda - \tau_j) (\nu - \tau_j) \right)_- \right|,$$

with  $w_j = q_j - q_{j-1}$ .

Illustration:



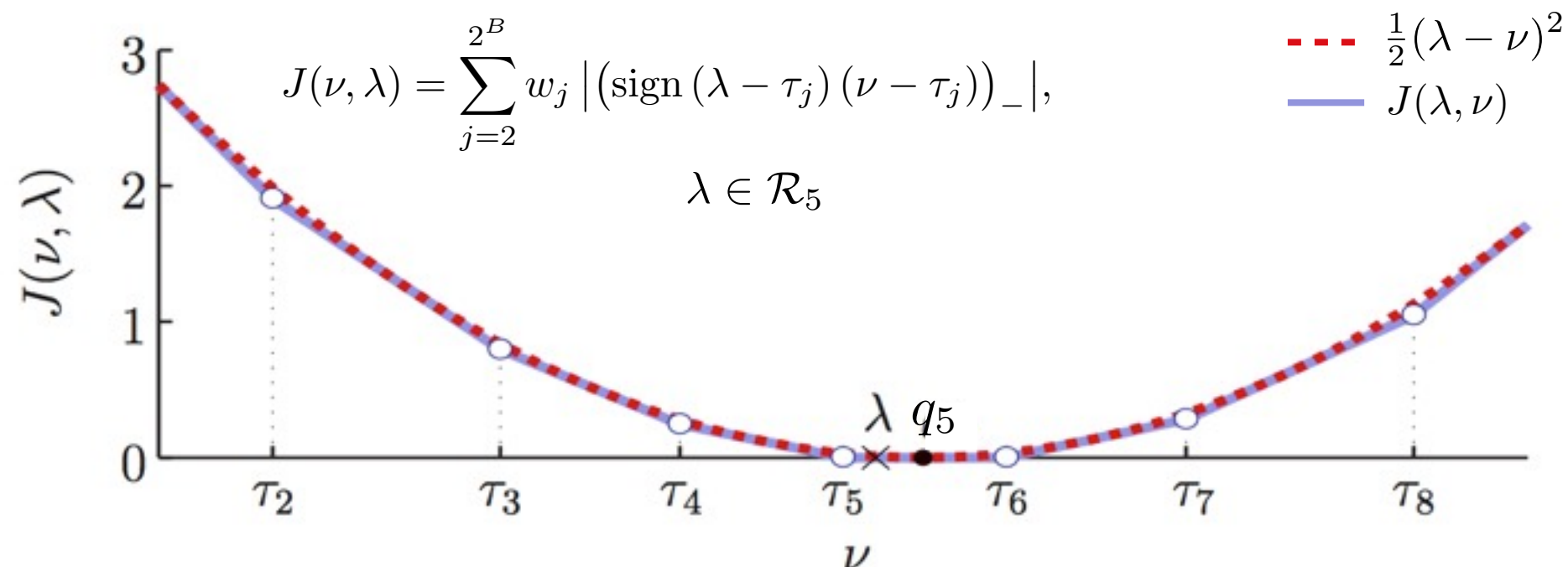
# Bridging 1-bit & $B$ -bit CS?

Given  $\mathcal{T} = \{\tau_j\}$  and  $\Omega = \{q_j\}$  ( $|\mathcal{T}| = 2^B + 1 = |\Omega| + 1$ ), let's define

$$J(\nu, \lambda) = \sum_{j=2}^{2^B} w_j \left| \left( \text{sign}(\lambda - \tau_j) (\nu - \tau_j) \right)_- \right|,$$

with  $w_j = q_j - q_{j-1}$ .

Illustration: more bins



# Bridging 1-bit & $B$ -bit CS?

Given  $\mathcal{T} = \{\tau_j\}$  and  $\Omega = \{q_j\}$  ( $|\mathcal{T}| = 2^B + 1 = |\Omega| + 1$ ), let's define

$$J(\nu, \lambda) = \sum_{j=2}^{2^B} w_j \left| \left( \text{sign}(\lambda - \tau_j) (\nu - \tau_j) \right)_- \right|,$$

with  $w_j = q_j - q_{j-1}$ .

For  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^M$ :  $\mathcal{J}(\mathbf{u}, \mathbf{v}) := \sum_{k=1}^M J(u_k, v_k)$

## Remarks:

- ▶  $J$  is convex in  $\nu$
- ▶ For  $B = 1$  ( $j = 2$  only):  
 $\mathcal{J}(\mathbf{u}, \mathbf{v}) \propto \|(\text{sign}(\mathbf{v}) \odot \mathbf{u})_-\|_1 \rightarrow \ell_1$ -sided 1-bit energy

- ▶ For  $B \gg 1$ :

$$J(\nu, \lambda) \rightarrow \frac{1}{2}(\nu - \lambda)^2 \text{ and } \mathcal{J}(\mathbf{u}, \mathbf{v}) \rightarrow \frac{1}{2}\|\mathbf{u} - \mathbf{v}\|^2 \text{ (quadratic energy)}$$

# Bridging 1-bit & $B$ -bit CS?

- Let's define an *inconsistency* energy:

$$\mathcal{E}_B(\mathbf{u}) := \mathcal{J}(\Phi \mathbf{u}, \mathbf{q}) \text{ with } \mathbf{q} = \mathcal{Q}_B[\Phi \mathbf{x}] \text{ and } \mathcal{E}_B(\mathbf{x}) = 0$$

- Idea: Minimize it in  $\Sigma_K$  (as for Iterative Hard Thresholding)

[Blumensath, Davies, 08]

$$\min_{\mathbf{u} \in \mathbb{R}^N} \mathcal{E}_B(\mathbf{u}) \text{ s.t. } \|\mathbf{u}\|_0 \leq K,$$

- NP Hard but greedy solution (as for IHT):

$$\mathbf{x}^{(n+1)} = \mathcal{H}_K[\mathbf{x}^{(n)} - \underbrace{\mu \partial \mathcal{E}_B(\mathbf{x}^{(n)})}_{\text{(sub) gradient}}] \text{ and } \mathbf{x}^{(0)} = 0.$$

$$\Phi^*(\text{sign}(\Phi \mathbf{u}) - \text{sign}(\Phi \mathbf{x})) \xleftarrow{B=1} \partial \mathcal{E}_B(\mathbf{u}) = \Phi^*(\mathcal{Q}_B(\Phi \mathbf{u}) - \mathbf{q}) \xrightarrow{B \gg 1} \Phi^*(\Phi \mathbf{u} - \mathbf{q})$$

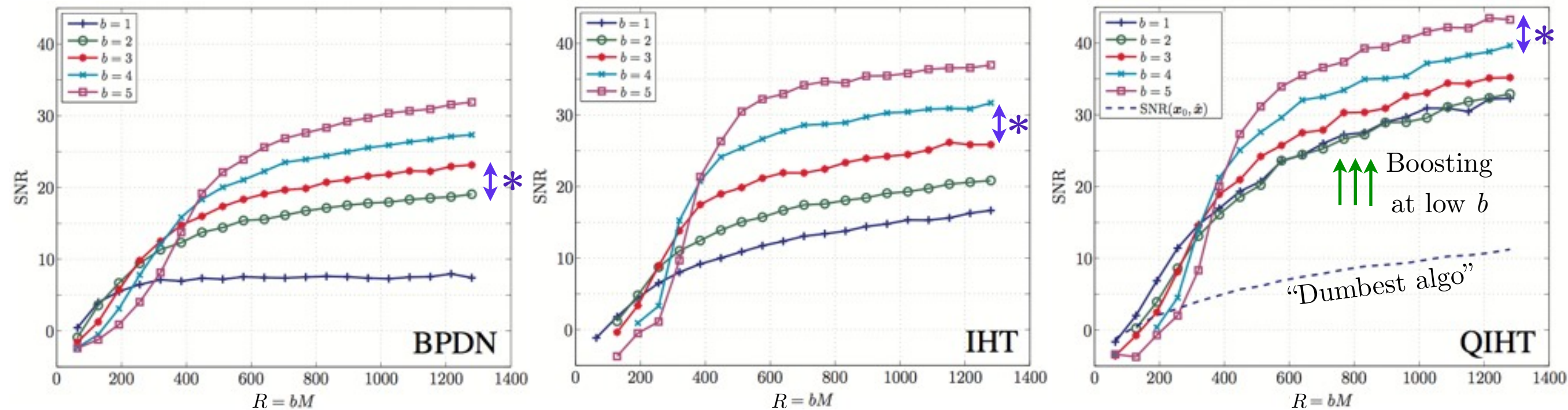
BIHT!
Quantized IHT (QIHT)
IHT!

T. Blumensath, M.E. Davies, "Iterative thresholding for sparse approximations". *Journal of Fourier Analysis and Applications*, 14(5-6), 629-654. (2008).

LJ, K. Degraux, C. De Vleeschouwer, "Quantized Iterative Hard Thresholding: Bridging 1-bit and High-Resolution Quantized Compressed Sensing", SAMPTA2013

# Bridging 1-bit & $B$ -bit CS?

$N = 1024$ ,  $K = 16$ ,  $R = BM \in \{64, 128, \dots, 1280\}$ , 100 trials (+ Lloyd-Max Gauss. Q.)



$R$ : total bit budget ( $BM$ )

\*: almost "6dB per bit" gain

$$\mu = \frac{1}{M}(1 - \sqrt{2K/M})$$

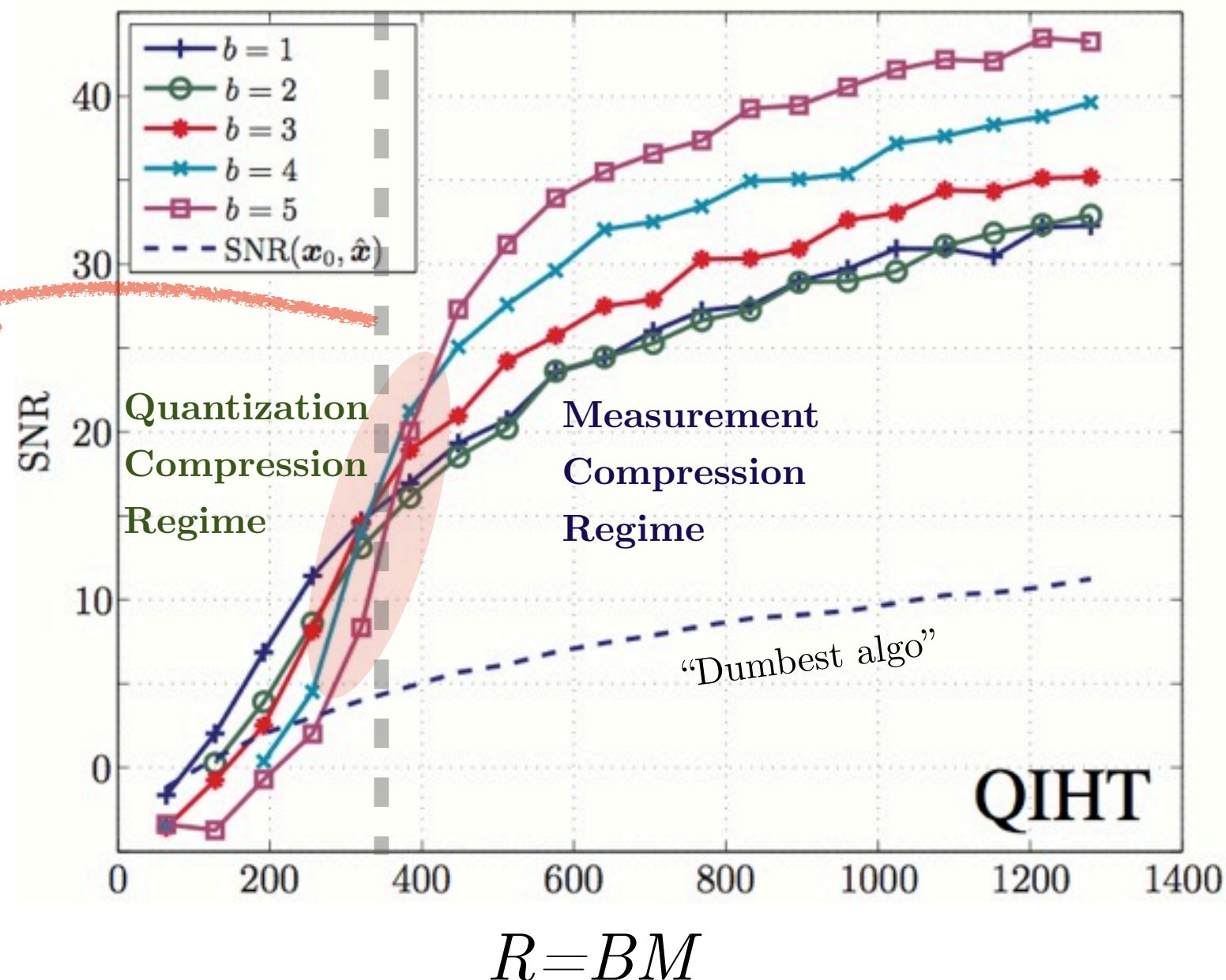
Adjusted by limit case  
analysis: BIHT and IHT

Note: entropy could be computed instead of  $B$  (*e.g.*, for further efficient coding)



# Bridging 1-bit & $B$ -bit CS?

$N = 1024$ ,  $K = 16$ ,  $R = BM \in \{64, 128, \dots, 1280\}$ , 100 trials



Interesting transition at

$R_0 \simeq 375$

“Regime Change?”

[Laska, Baraniuk, 12]

$R_0$  could increase with input noise power.

# Further Reading

- ▶ T. Blumensath, M.E. Davies, “Iterative thresholding for sparse approximations”. *Journal of Fourier Analysis and Applications*, 14(5-6), pp. 629-654, 2008
- ▶ P. T. Boufounos and R. G. Baraniuk, “1-Bit compressive sensing,” *Proc. Conf. Inform. Science and Systems (CISS)*, Princeton, NJ, March 19-21, 2008.
- ▶ Boufounos, P. T. (2009, November). “Greedy sparse signal reconstruction from sign measurements”. In *Conference Record of the Forty-Third Asilomar Conference on Signals, Systems and Computers*, 2009
- ▶ Y. Plan, R. Vershynin, “Dimension reduction by random hyperplane tessellations”, arXiv:1111.4452, 2011.
- ▶ Y. Plan, R. Vershynin, “Robust 1-bit compressed sensing and sparse logistic regression: a convex programming approach”, *IEEE Trans. Info. Theory*, arXiv:1202.1212, 2012.
- ▶ J. N. Laska, R. G. Baraniuk, ‘Regime change: Bit-depth versus measurement-rate in compressive sensing’, *IEEE Trans. Signal Processing*, 60(7), pp. 3496-3505, 2012.
- ▶ U.S. Kamilov, A. Bourquard, A. Amini, M. Unser, “One-bit measurements with adaptive thresholds”. *IEEE Signal Processing Letters*, 19(10), pp. 607-610, 2012
- ▶ L. Jacques, J. N. Laska, P. T. Boufounos, and R. G. Baraniuk, “Robust 1-Bit Compressive Sensing via Binary Stable Embeddings of Sparse Vectors,” *IEEE Trans. Info. Theory*, 59(4), 2013.
- ▶ L. Jacques, K. Degraux, C. De Vleeschouwer, “Quantized Iterative Hard Thresholding: Bridging 1-bit and High-Resolution Quantized Compressed Sensing”, SAMPTA 2013, to appear.



# Today's Topics

---

1. Modern Scalar Quantization
2. Compressive Sensing Overview
3. Compressive Sensing and Quantization
4. 1-bit Compressive Sensing
5. Locality Sensitive Hashing and Universal Quantization

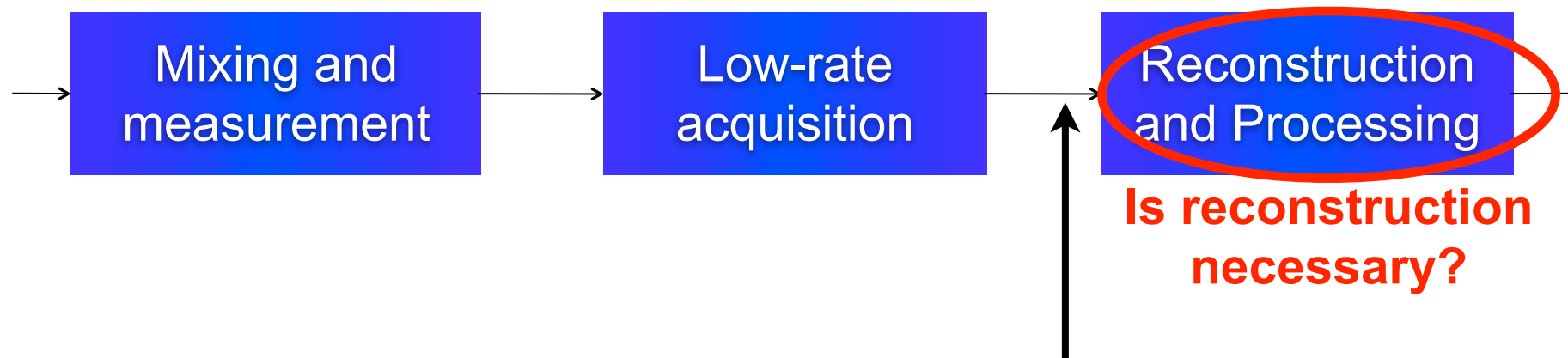
# Today's Topics

---

1. Modern Scalar Quantization
2. Compressive Sensing Overview
3. Compressive Sensing and Quantization
4. 1-bit Compressive Sensing
- 5. Locality Sensitive Hashing and Universal Quantization**

# **INFORMATION EMBEDDING**

# Compressive Domain Processing



**No:** operate in compressive domain

**Moreover:** signal does not have to be sparse  
(as long as it has some structure)

Compressive operations: **detection, estimation, filtering**

Randomized projection **embeds** signal **information**.

**Main benefits:** Computation, Memory

**Questions:**

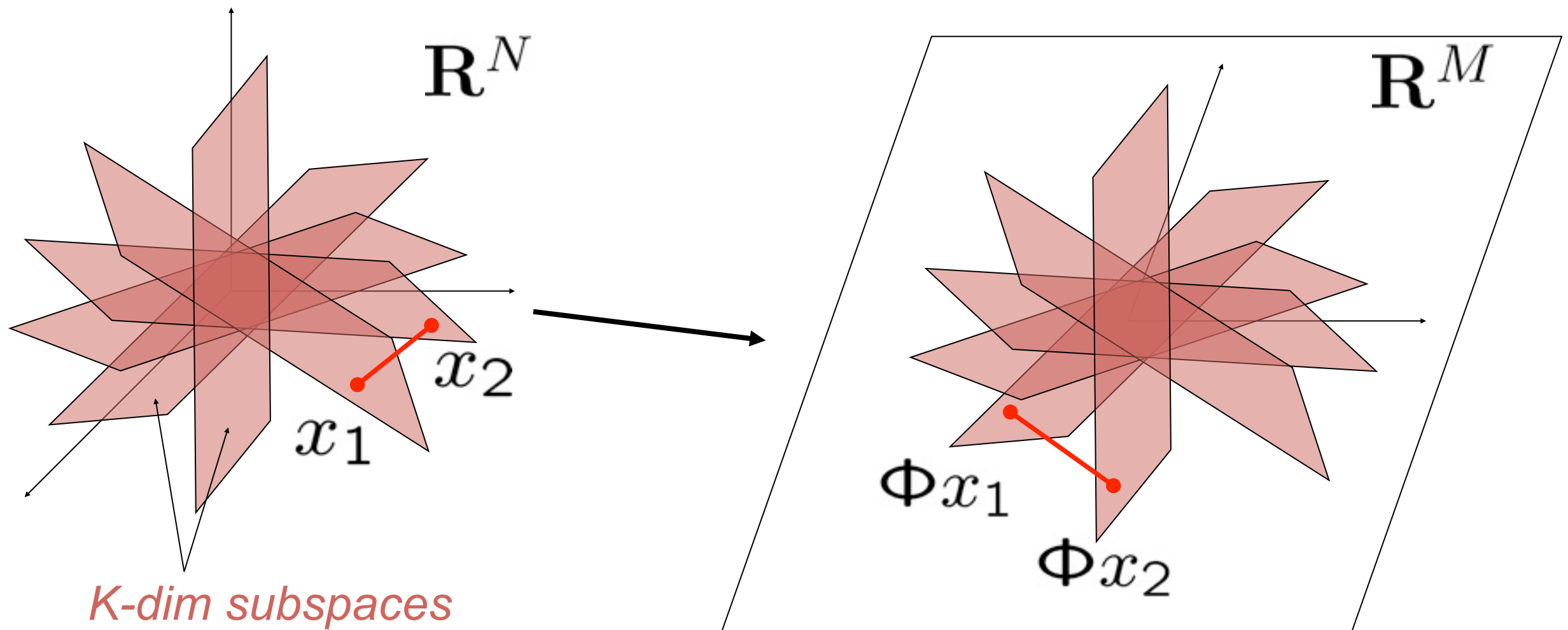
**What information is embedded?**

**How to best embed information?**

- Davenport M. A., Boufounos P. T., Wakin M. B., and Baraniuk R. G., "Signal processing with compressive measurements," *IEEE Journal of Selected Topics in Signal Processing*, v. 4, no. 2, pp. 445-460, April, 2010.

# RIP/Stable Embedding

- An information preserving projection  $\mathbf{A}$  preserves the **geometry** of the set of sparse signals

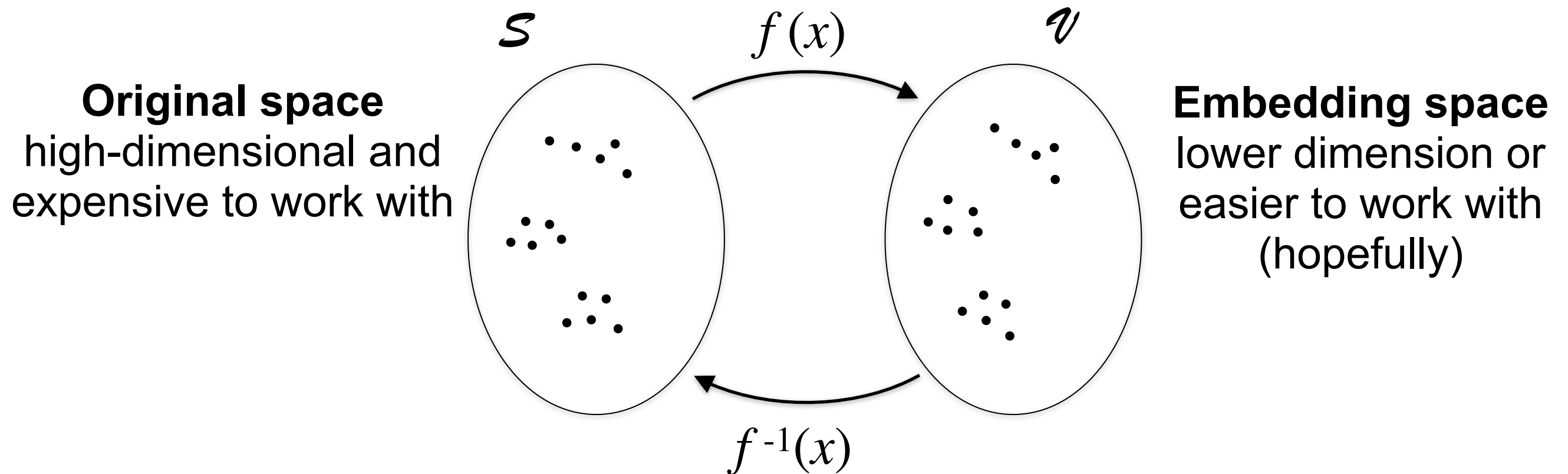


Restricted Isometry Property

$$(1 - \delta)\|x\|_2^2 \leq \|\Phi x\|_2^2 \leq (1 + \delta)\|x\|_2^2$$

# **GEOMETRY-PRESERVING EMBEDDINGS**

# Isometric (approximate) embeddings

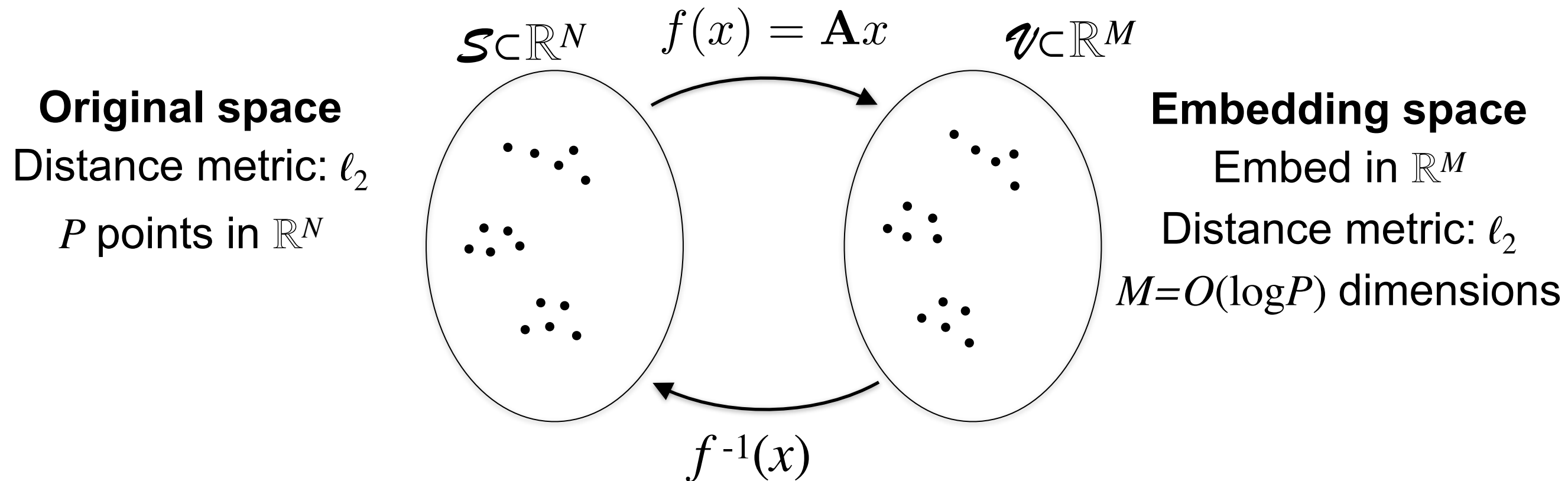


Transformations that preserve distances

$$\text{For all } x, y \text{ in } \mathcal{S} : d_{\mathcal{S}}(x, y) \approx d_{\mathcal{V}}(f(x), f(y))$$



# Johnson-Lindenstrauss embeddings



Transformations that preserve distances

For all  $x, y$  in  $\mathcal{S}$ :

$$(1 - \epsilon) \|x - y\|_2^2 \leq \|f(x) - f(y)\|_2^2 \leq (1 + \epsilon) \|x - y\|_2^2$$

- Johnson W. and Lindenstrauss J., “Extensions of Lipschitz mappings into a Hilbert space,” *Contemporary Mathematics*, vol. 26, pp. 189 – 206, 1984.

# Johnson-Lindenstrauss Lemma

---

Consider  $\mathcal{S} \subset \mathbb{R}^N$  containing  $P$  points.

We can embed  $\mathcal{S}$  in  $\mathbb{R}^M$  such that for all  $x, y$  in  $\mathcal{S}$ :

$$(1 - \epsilon) \|x - y\|_2^2 \leq \|f(x) - f(y)\|_2^2 \leq (1 + \epsilon) \|x - y\|_2^2$$

using only  $M = O\left(\frac{\log P}{\epsilon^2}\right)$  dimensions

## Later results

$f(x)$  can be linear  $f(x) = Ax$ , randomized  $A$  achieves bound  
(e.g., entries Gaussian,  $\pm 1/-1$  Bernoulli, etc.)

Bound (almost) tight:  $M = O\left(\frac{\log P}{\epsilon^2 \log \frac{1}{\epsilon}}\right)$  dimensions necessary

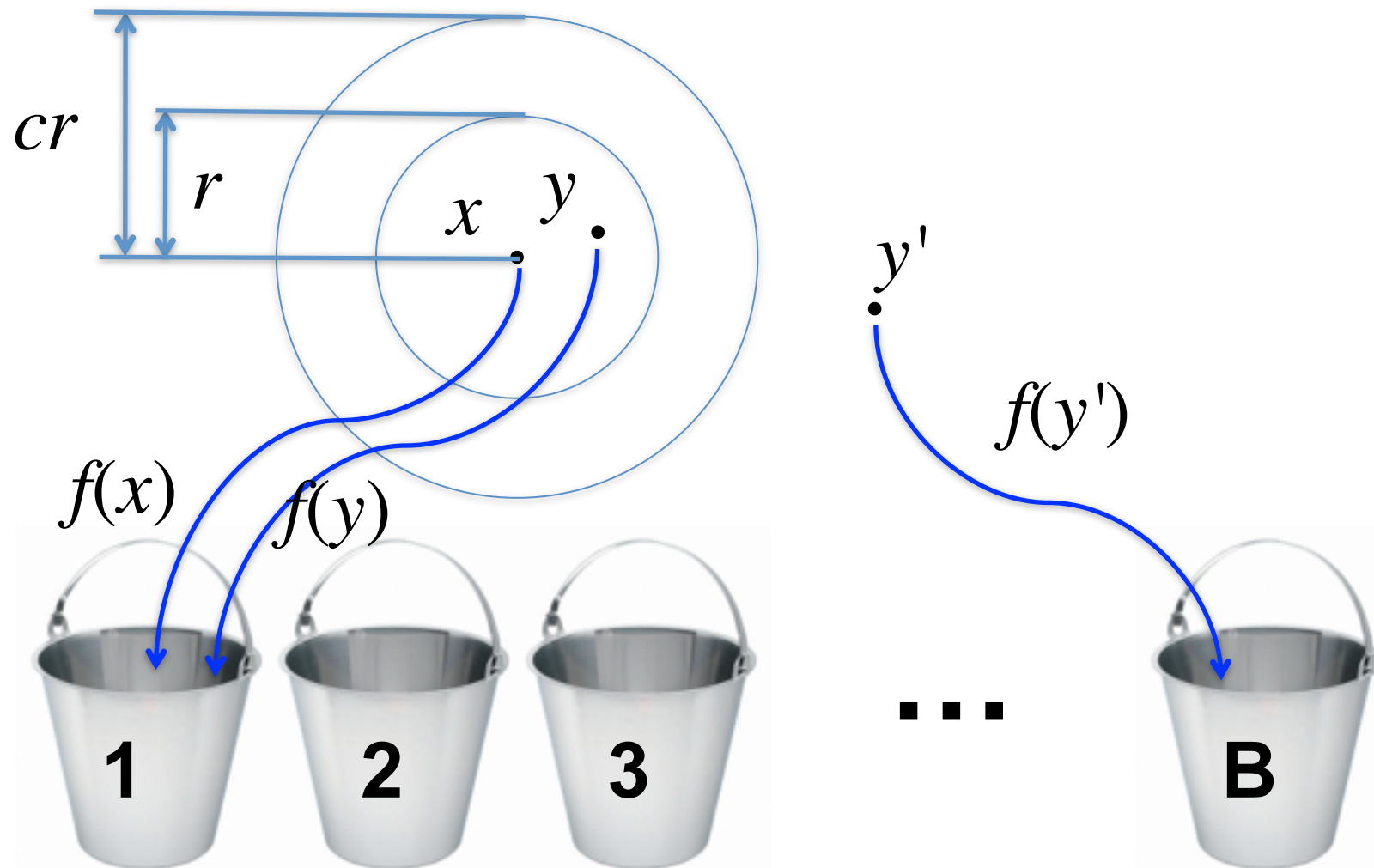
**BUT: Quantization is necessary for transmission!**  
**Are J-L Embeddings still appropriate?**

# Locality Sensitive Hashing

Randomized signal hash  $f: \mathbb{R}^N \rightarrow \mathbb{N}$  such that:

$d(x, y) \leq r \Rightarrow f(x) = f(y)$  with high probability

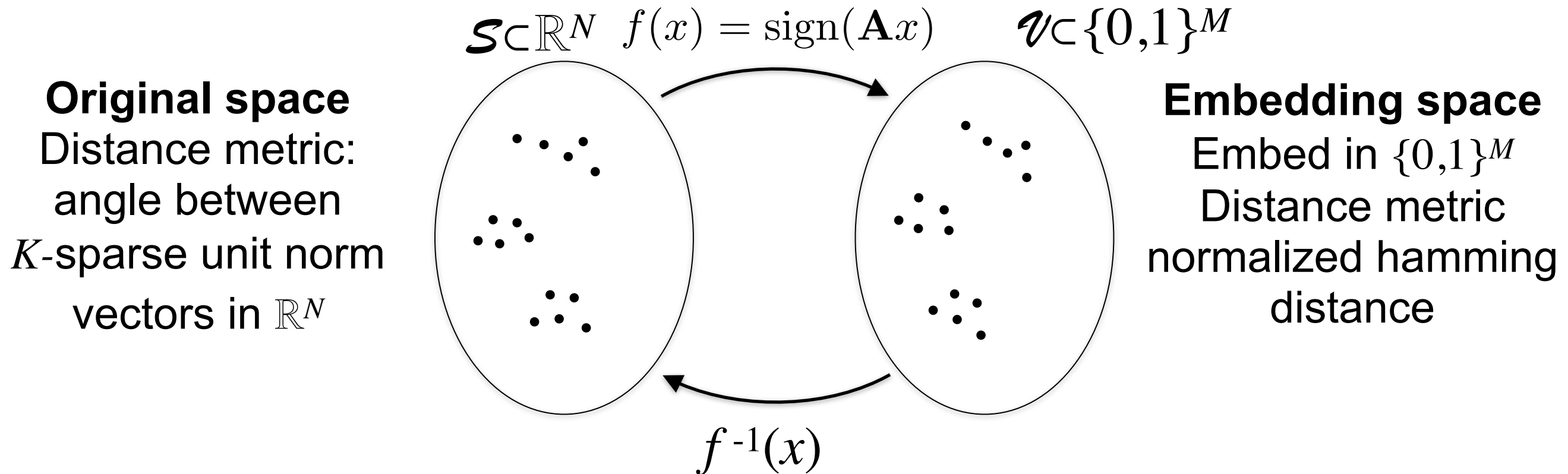
$d(x, y) \geq cr \Rightarrow f(x) \neq f(y)$  with high probability



**(One) LSH approach: random projection and quantization,  
i.e., Quantized Johnson-Lindenstrauss**

- Andoni, A. and Indyk, P., “Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions,” *Commun. ACM*, vol. 51, no. 1, pp. 117–122, 2008.
- Datar M., Immorlica N., Indyk P., and Mirrokni V., “Locality-Sensitive Hashing Scheme Based on p-Stable Distributions,” *Proc. Symposium on Computational Geometry*, 2004

# Binary Stable Embedding



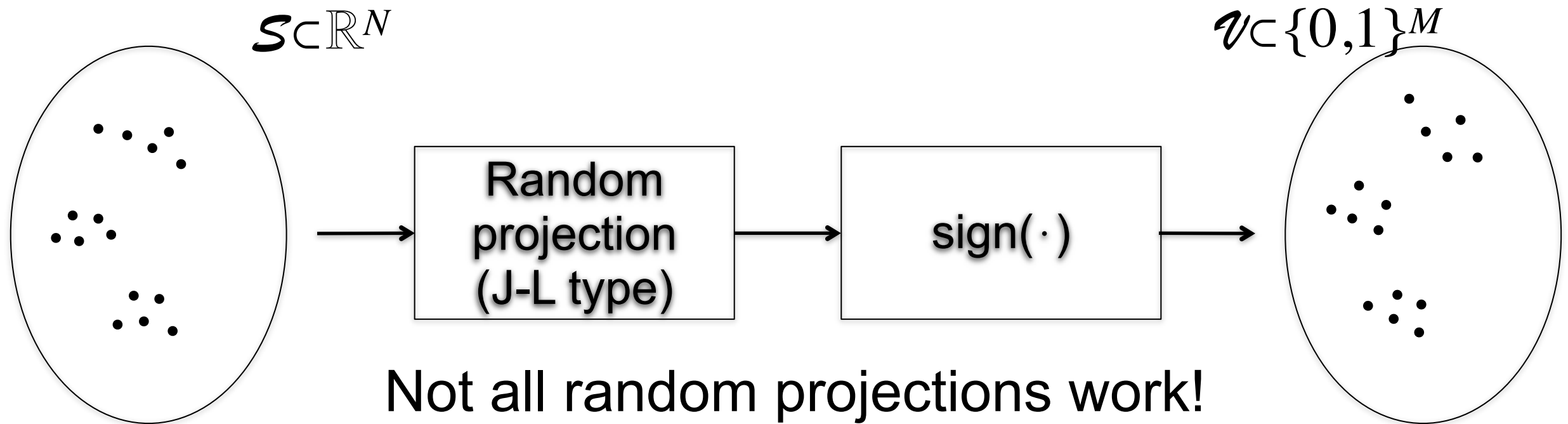
Embedding preserves angles for all  $K$ -sparse  $x, y$  in  $\mathbb{R}^N$ :

$$\arccos \left( \frac{\langle x, y \rangle}{\|x\| \|y\|} \right) - \epsilon \leq d_H(f(x), f(y)) \leq \arccos \left( \frac{\langle x, y \rangle}{\|x\| \|y\|} \right) + \epsilon$$

using only  $M = O \left( \frac{1}{\epsilon^2} \left( K \log N + K \log \frac{1}{\epsilon} \right) \right)$  measurements

- Jacques L., Laska J. N., Boufounos P. T., Baraniuk R. J., "Robust 1-Bit Compressive Sensing via Binary Stable Embeddings of Sparse Vectors," *IEEE Trans. Info. Theory*, v. 59, no. 4, April, 2013.

# Binary Stable Embedding



Not all random projections work!  
Matrices w/ i.i.d. Gaussian entries work  
w/ i.i.d. Bernoulli they don't always

Sufficient information for sparse recovery (previous part)

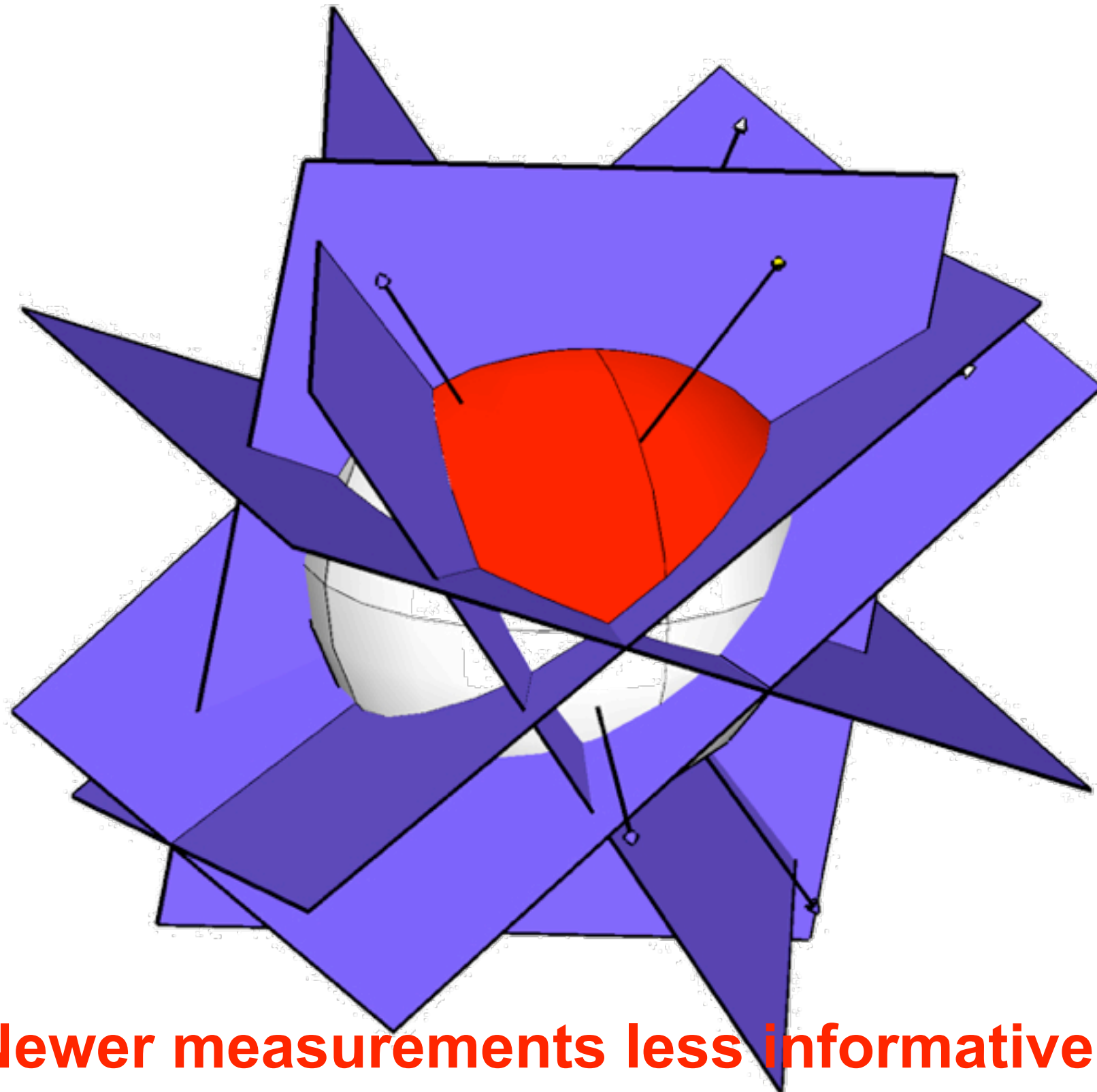
Embedding does not preserve amplitudes

**Is embedding rate-efficient?**

- Plan, Y. and Vershynin, R., “Dimension reduction by random hyperplane tessellations,” preprint, arXiv:1111.4452, 2011.
- Ai, A., Lapanowski, A., Plan, Y., Vershynin, R., “One-bit compressed sensing with non-Gaussian measurements”, *Linear Algebra and Applications*, to appear.

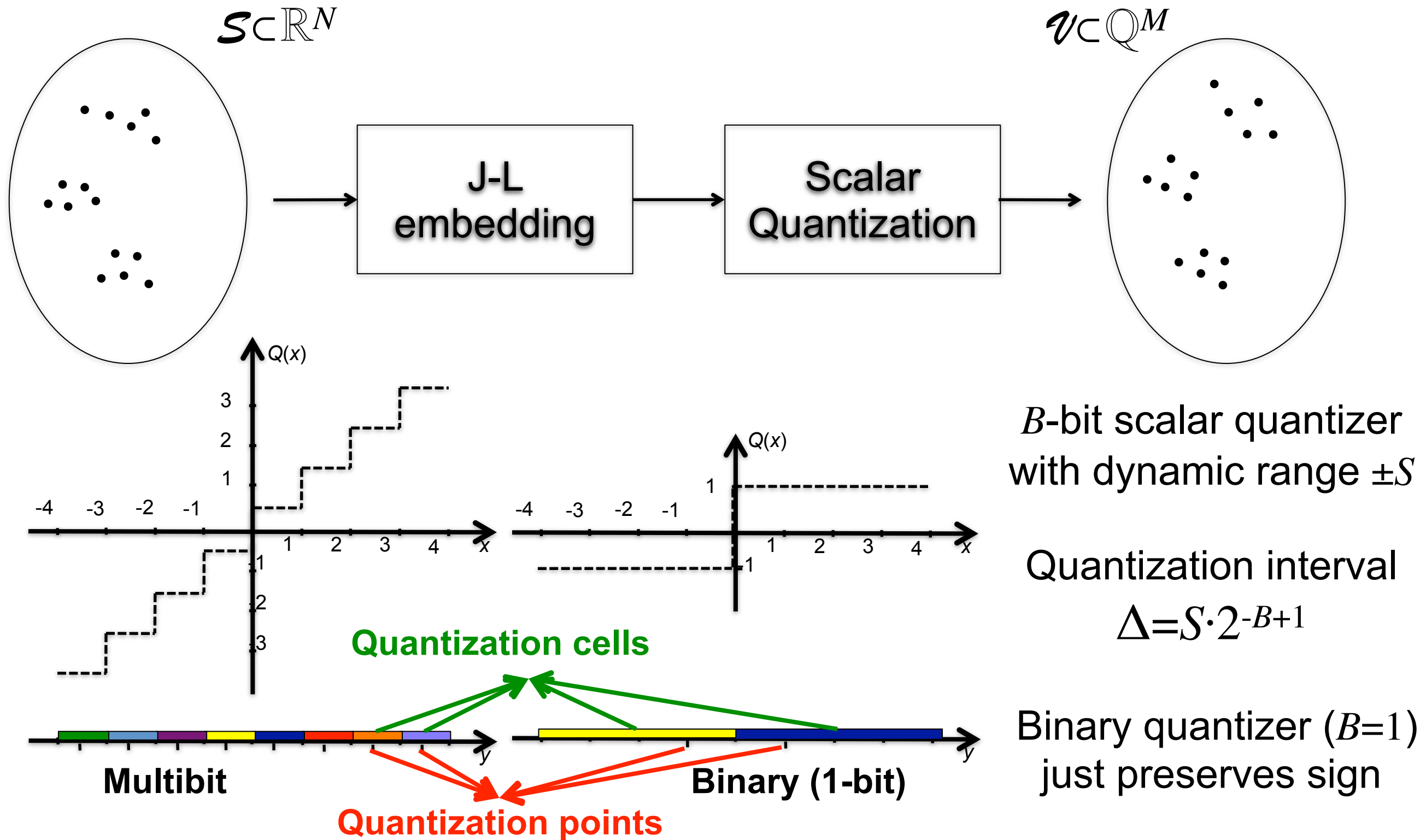
# Information in 1-bit Measurements

---



**Newer measurements less informative.  
Chance of intersection increasingly smaller  
Embedding rate-inefficient!**

# Quantized J-L Embeddings



- Li M., Rane S., and Boufounos P. T., "Quantized embeddings of scale-invariant image features for mobile augmented reality," *IEEE 14th International Workshop on Multimedia Signal Processing (MMSP)*, Banff, Canada, Sept. 17-19, 2012



# Johnson-Lindenstrauss With Quantization [w/ Li, Rane]

---

Consider  $\mathcal{S} \subset \mathbb{R}^N$  containing  $P$  points.

We can embed  $\mathcal{S}$  in  $\mathbb{R}^M$  such that for all  $x, y$  in  $\mathcal{S}$ :

$$\begin{aligned} (1 - \epsilon) \|x - y\|_2 - 2^{-B+1} S &\leq \\ &\|Q(f(x)) - Q(f(y))\|_2 \\ &\leq (1 + \epsilon) \|x - y\|_2 + 2^{-B+1} S \end{aligned}$$

using only  $M = O\left(\frac{\log P}{\epsilon^2}\right)$  dimensions

and  $B$  bits per dimension  
(with appropriate normalizations/saturation levels)

**Total rate:  $R=BM$**



# Quantized J-L at Fixed Rate

Given total rate:  $R=MB$

How to assign  $B$  and  $M$ ? More  $M$  or more  $B$ ?

Larger  $B$ , less quantization distortion

$$(1 - \epsilon) \|x - y\|_2 - 2^{-\frac{R}{M} + 1} S \leq \|Q(f(x)) - Q(f(y))\|_2 \leq (1 + \epsilon) \|x - y\|_2 + 2^{-\frac{R}{M} + 1} S$$

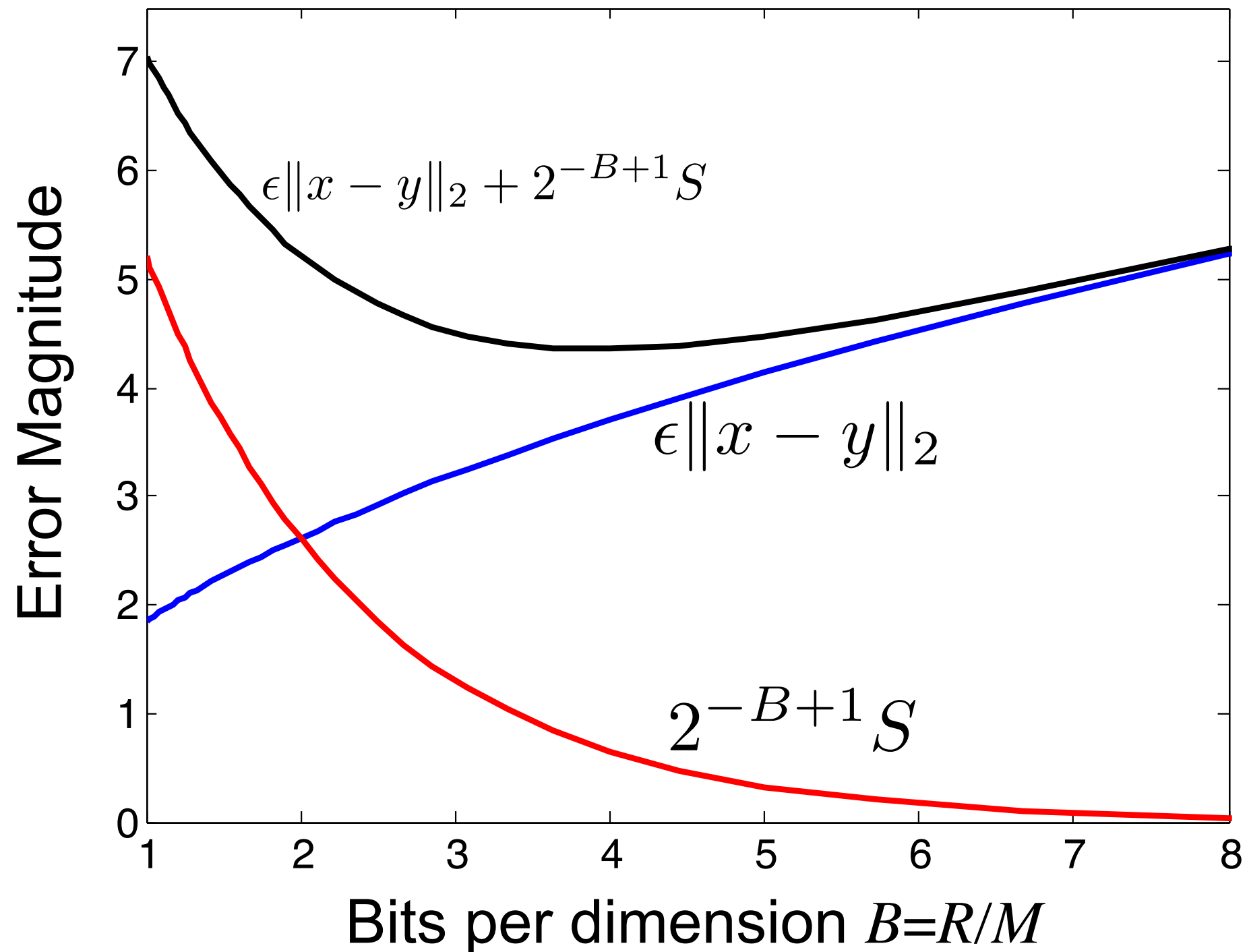
Larger  $M$ , less J-L type distortion  $\epsilon$

$$\epsilon = O(1/\sqrt{M})$$

Design tradeoff:

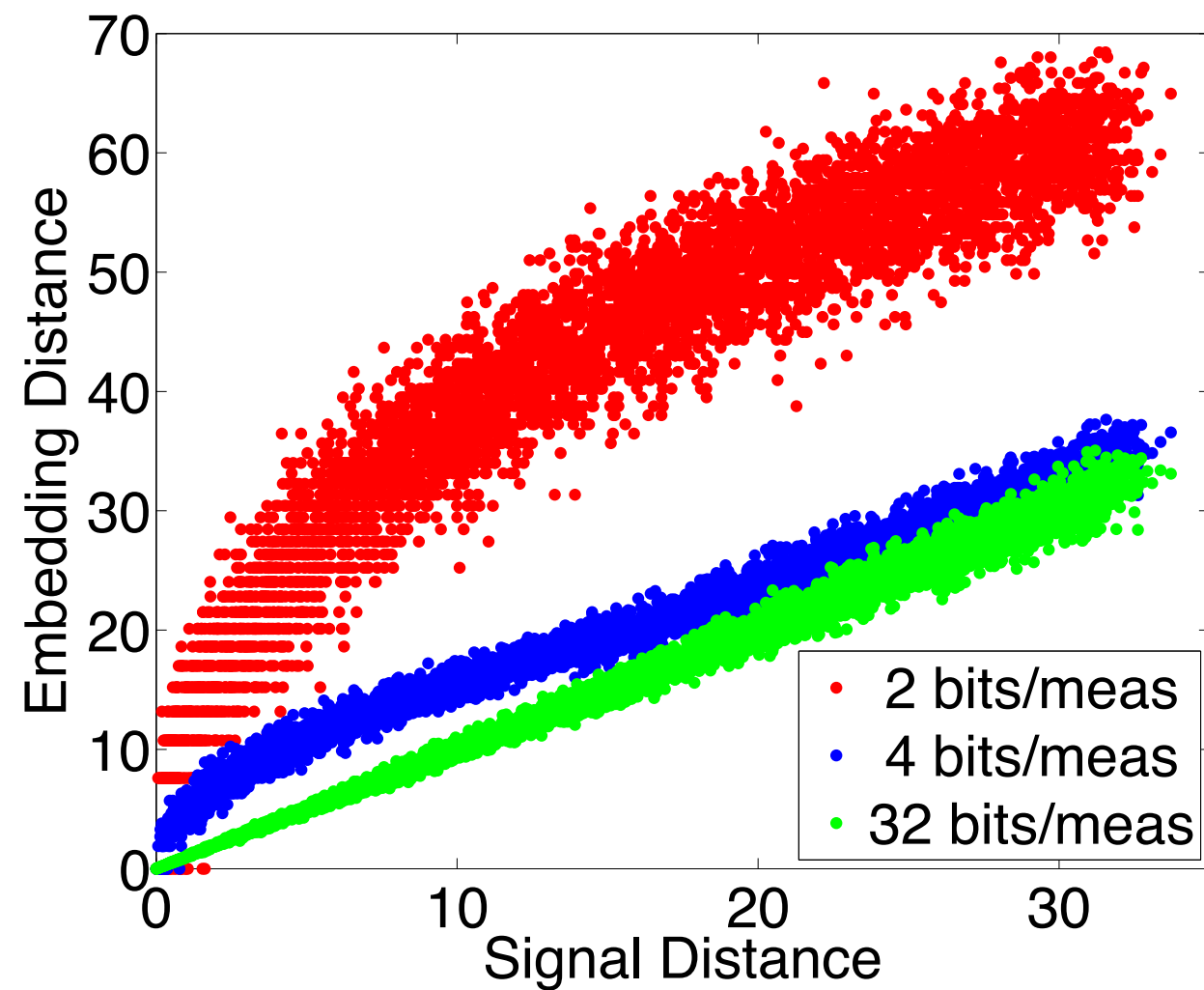
Number of projections vs. bits per projection

# Exploring the Design Trade-off

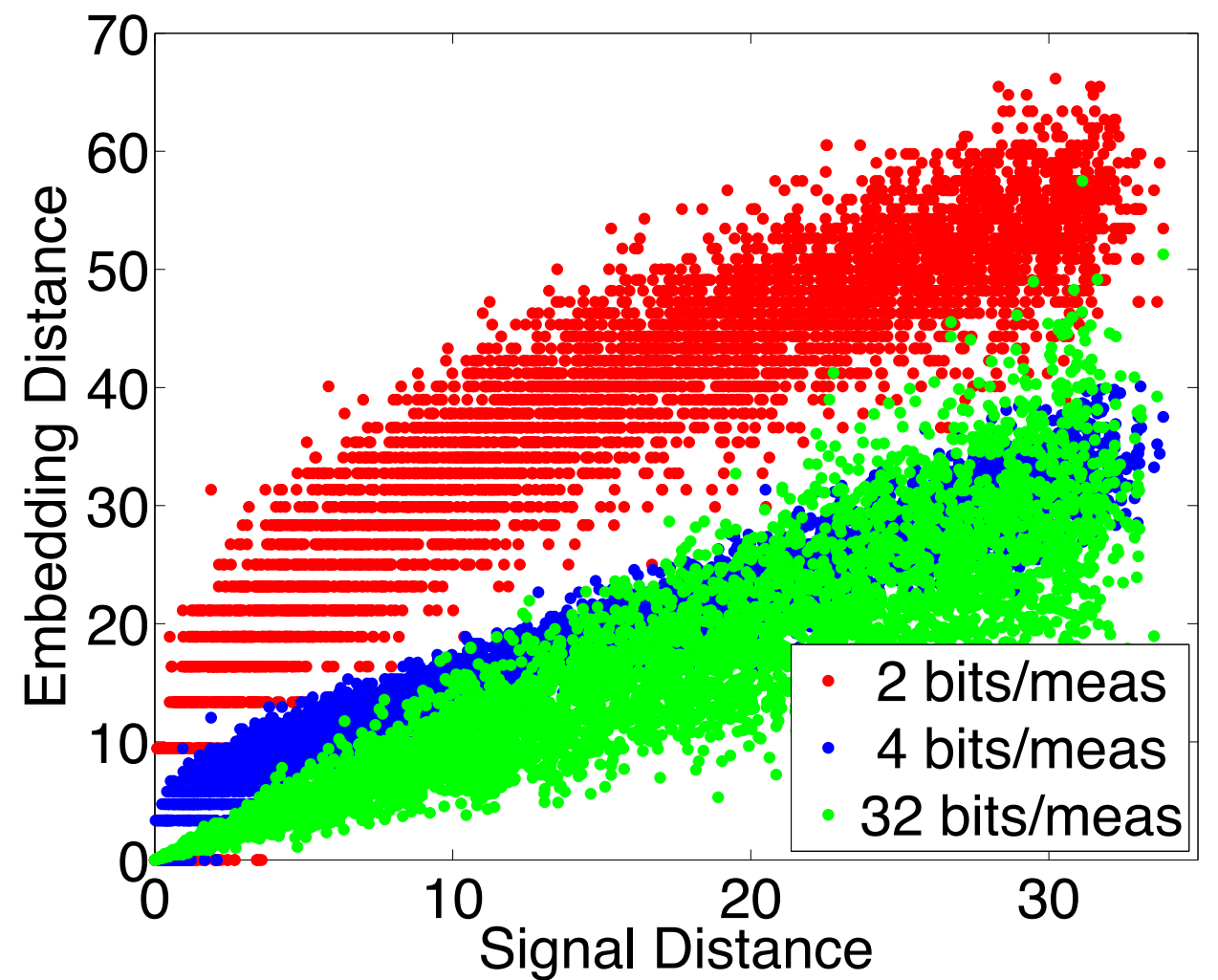


# Exploring the Design Trade-off

Fixed  $M=256$



Fixed  $R=MB=256$



**IN PRACTICE**

# The Augmented Reality Problem



**Server-side processing** increasingly important  
(e.g. cloud computing, augmented reality)

**Compression is necessary**

**Goal: detection;** not image transmission

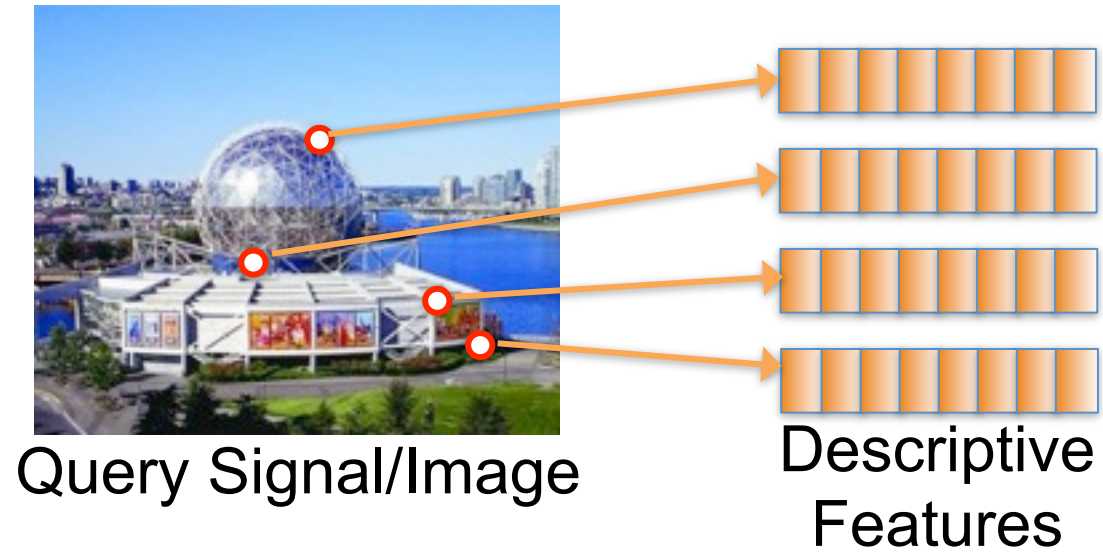
**Q: Should we transmit the signal?**

**Can we reduce the rate?**

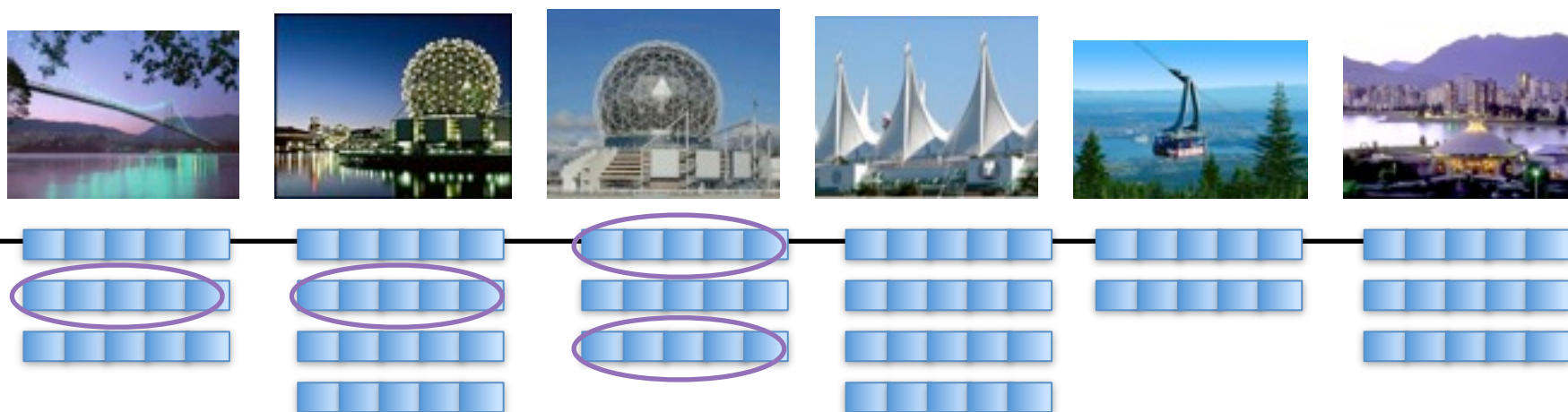


# Signal/Image-based Retrieval

## Feature Extraction

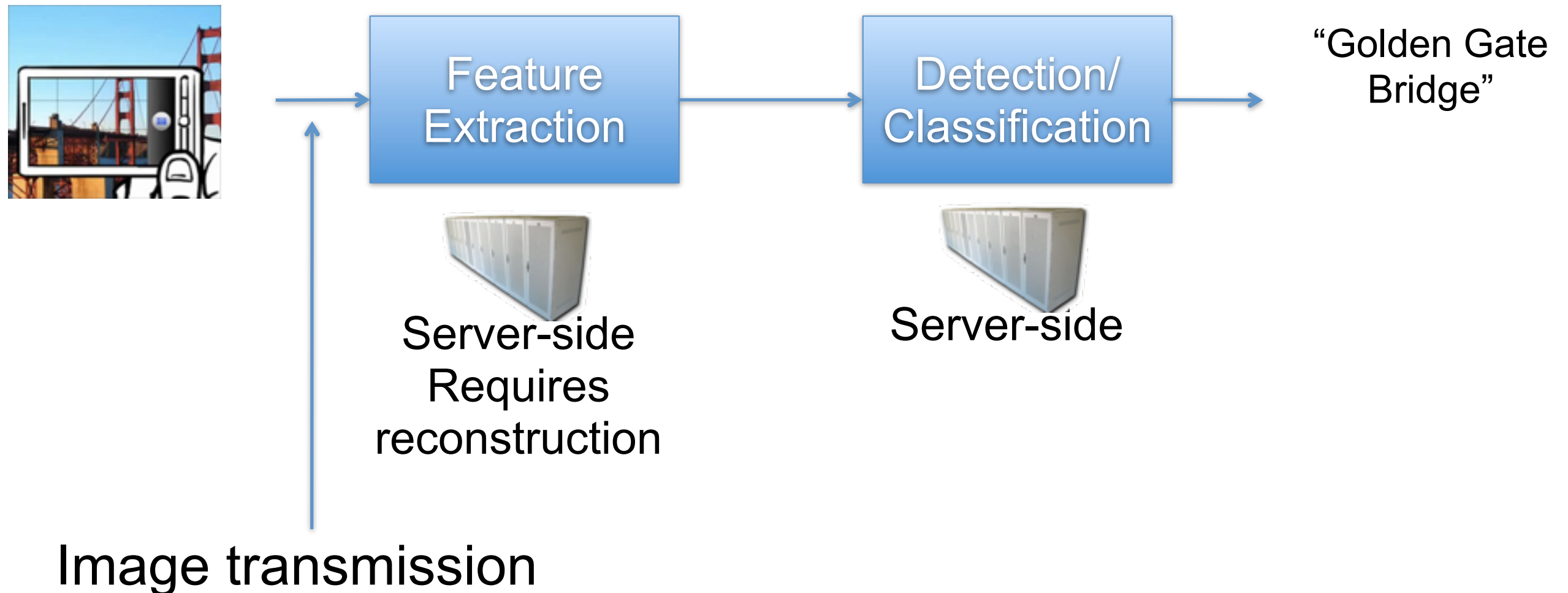


## Signal/Image Database



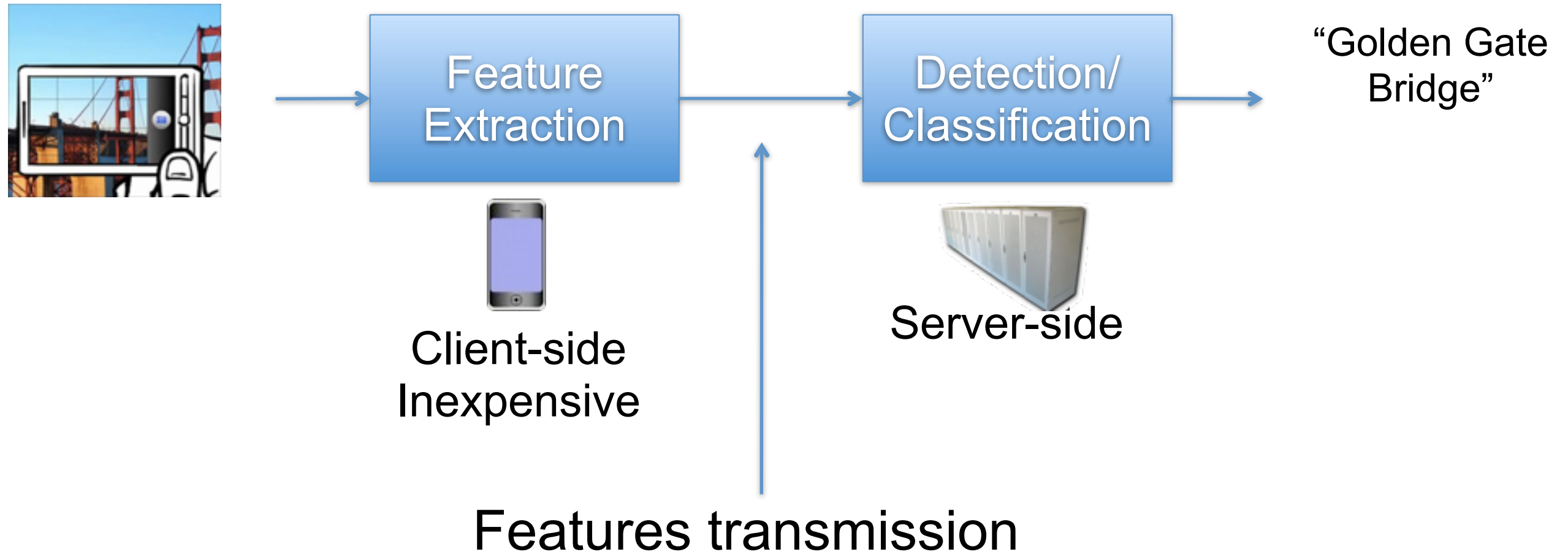


# Detection/Classification Pipeline (typical)



Detection/Classification: Based on distance/inner product

# Detection/Classification Pipeline (efficient)

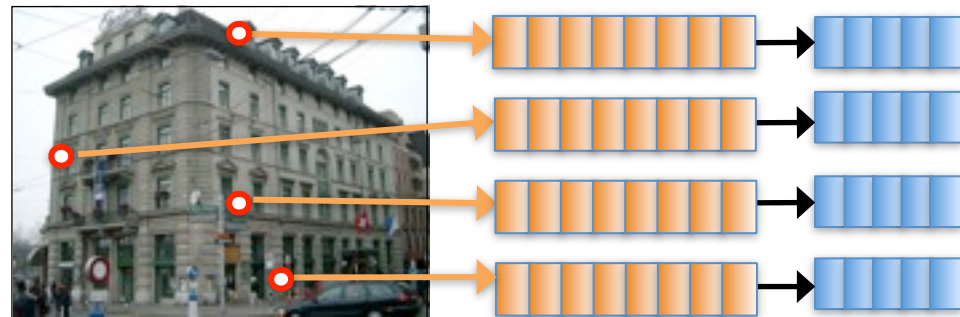


Detection/Classification: Based on distance/inner product

**Goal: rate-efficient distance-preserving transmission**



# ZuBuD: Zurich Buildings Database

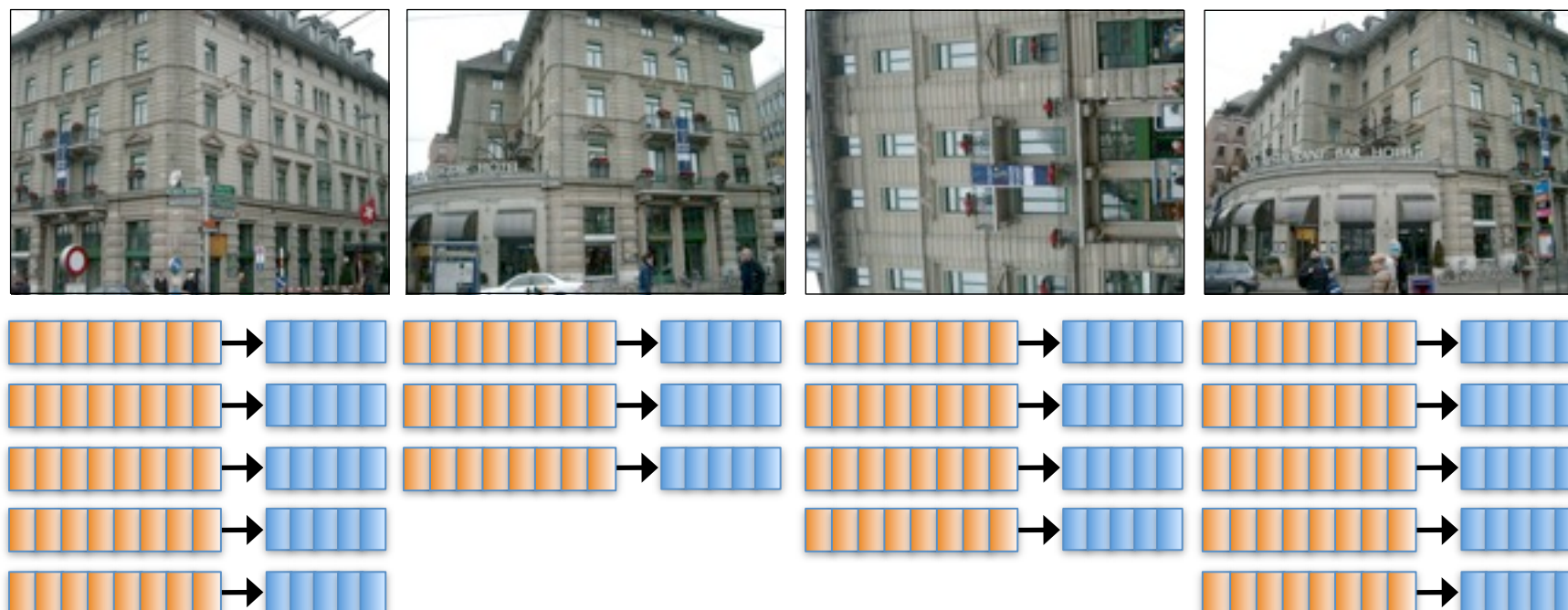


1005 images: 201 buildings from 5 viewpoints each  
804 images (4 viewpoints per building) in server  
201 query images (1 viewpoint per building)

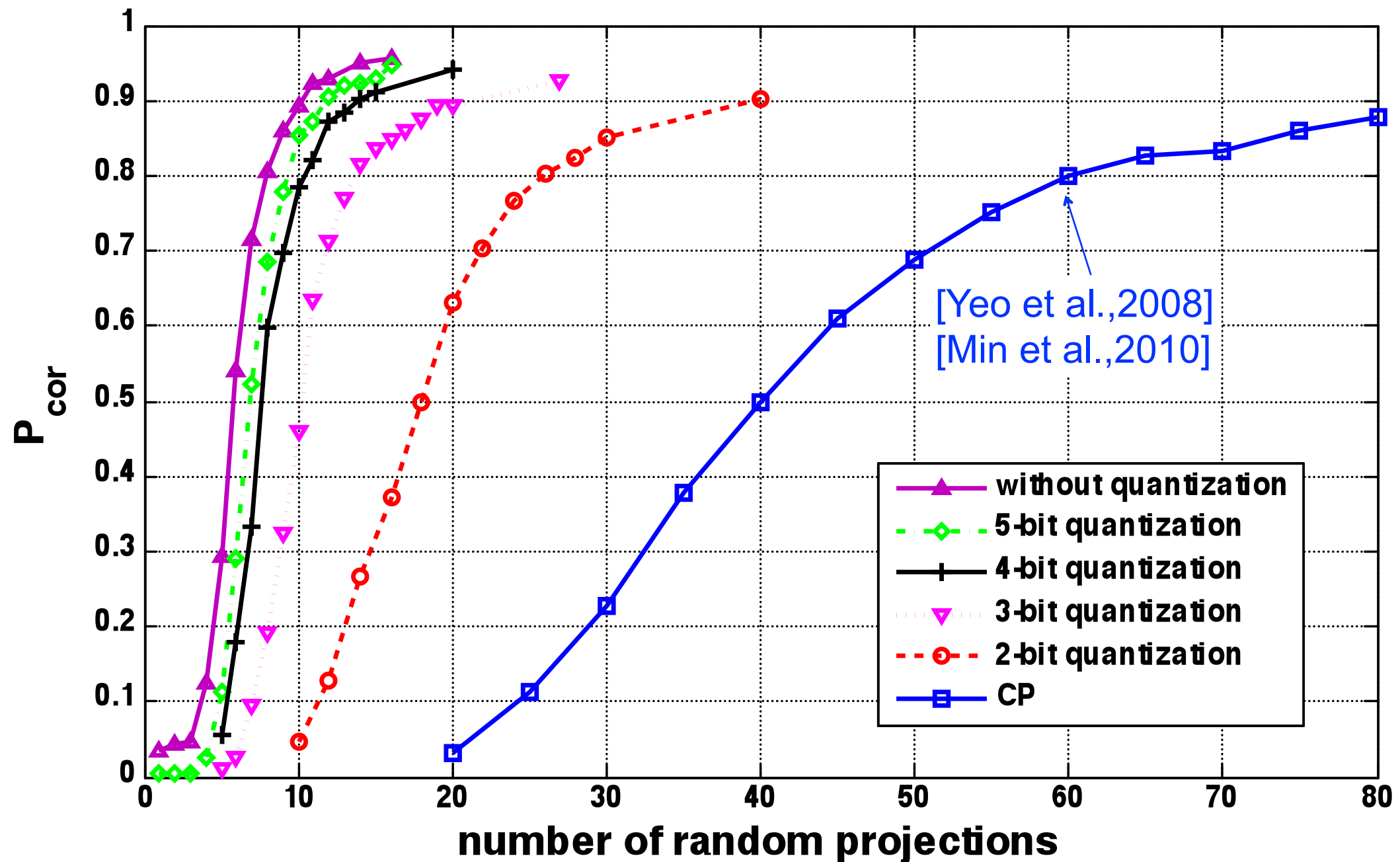


Building ID

## Server Database

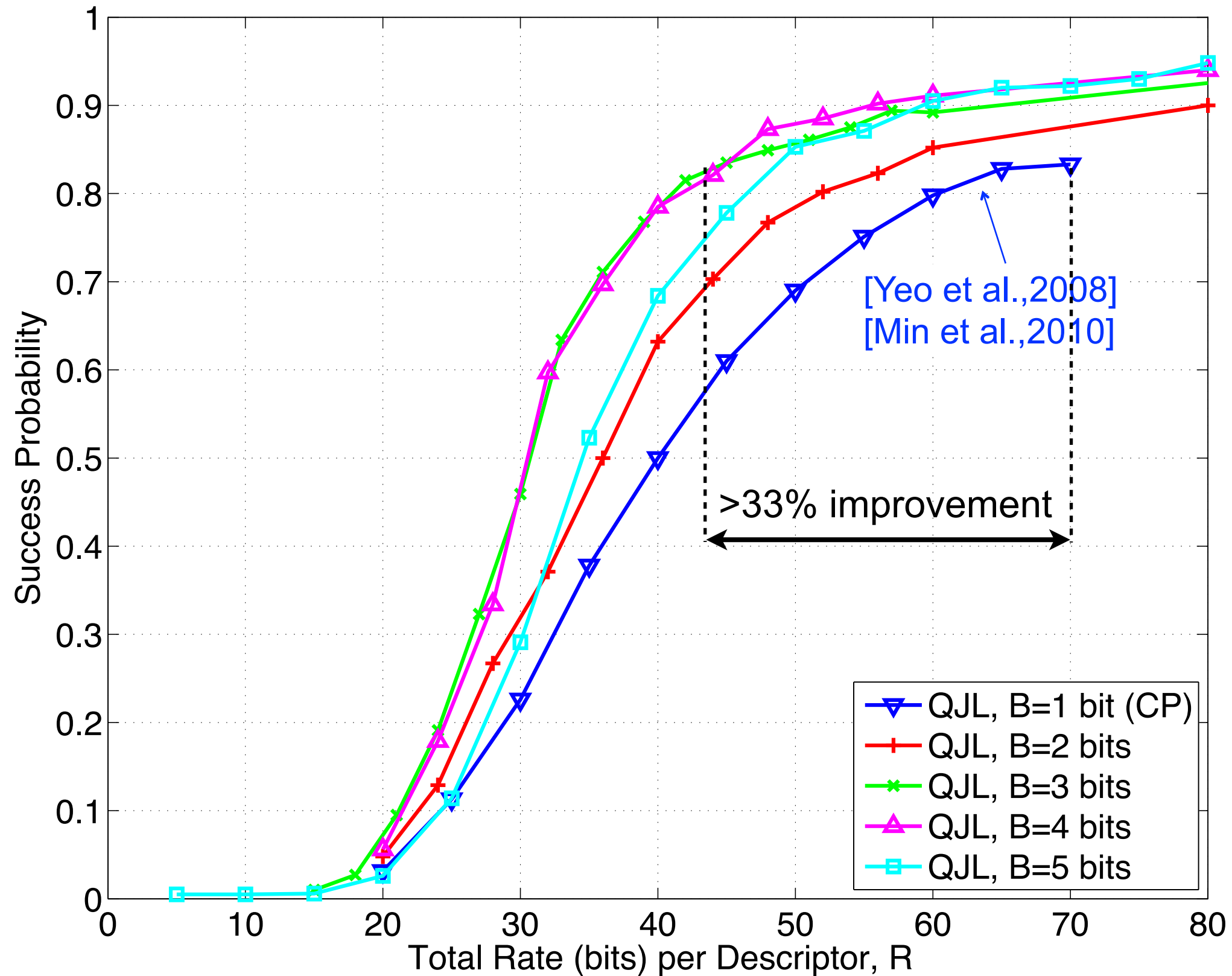


# Success Probability



- Yeo C., Ahammad P., and Ramchandran K., “Coding of image feature descriptors for distributed rate-efficient visual correspondences,” *International Journal of Computer Vision*, vol. 94, pp. 267–281, 2011, 10.1007/s11263-011-0427-1.
- Min K., Yang L., Wright J., Wu L., Hua X.-S., and Ma Y., “Compact projection: Simple and efficient near neighbor search with practical memory requirements,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2010.

# Success Probability with Fixed Rate

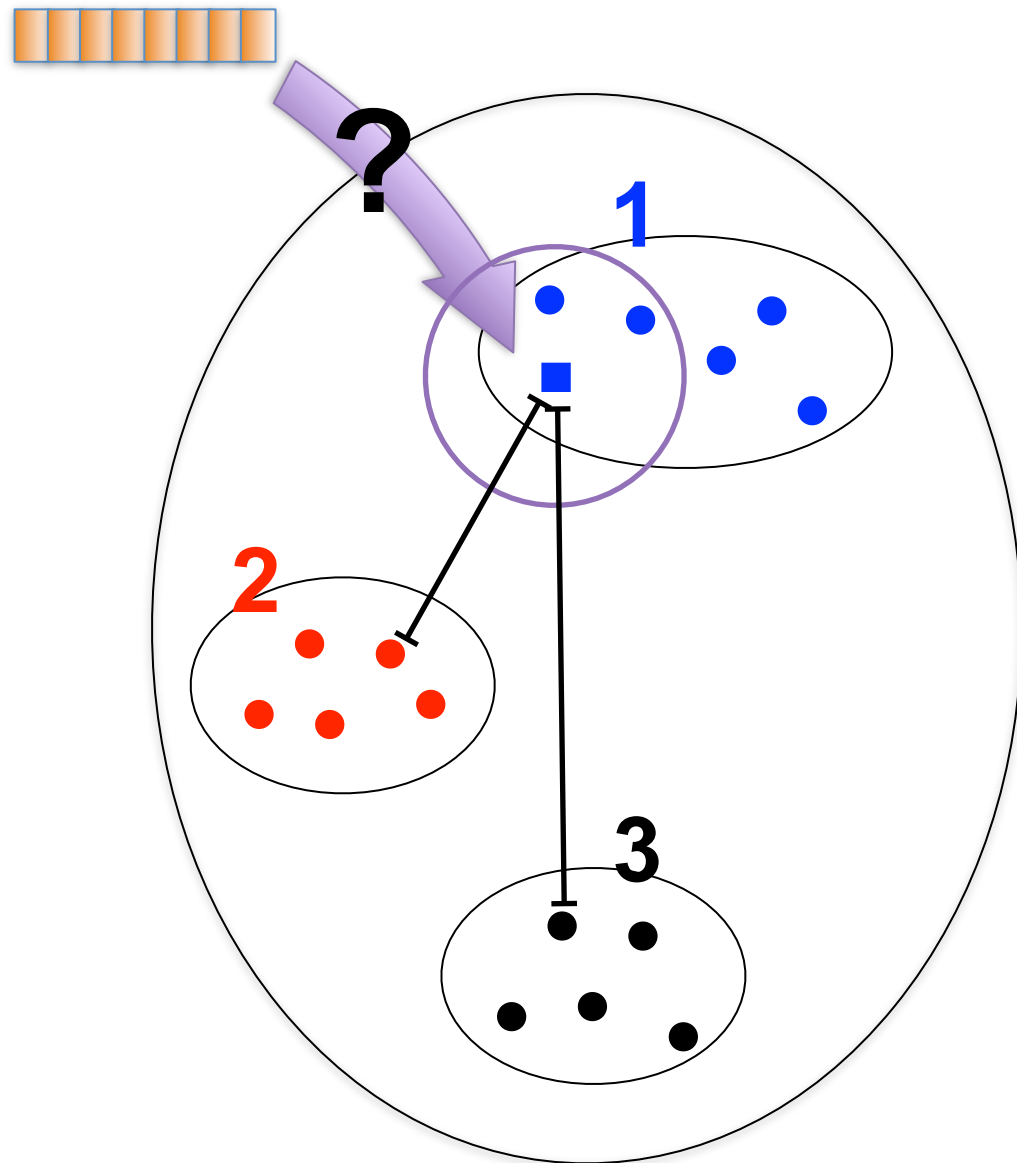


# Information Scalability

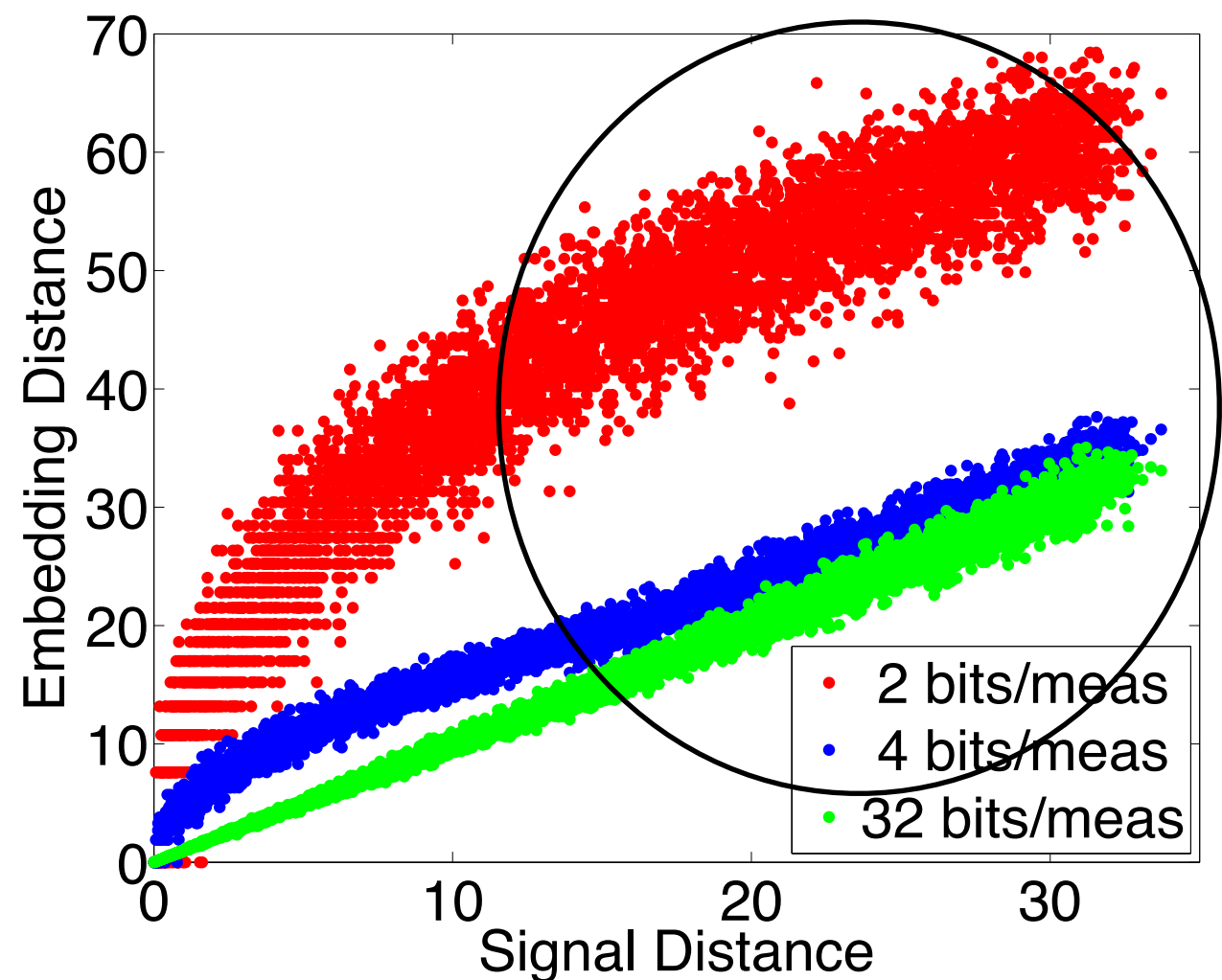
Inference relies on clusters of signals

Large distances not necessary to determine clusters and nearest neighbors

Should not spend bits encoding large distances!



**But how?**



# **UNIVERSAL QUANTIZED EMBEDDINGS**

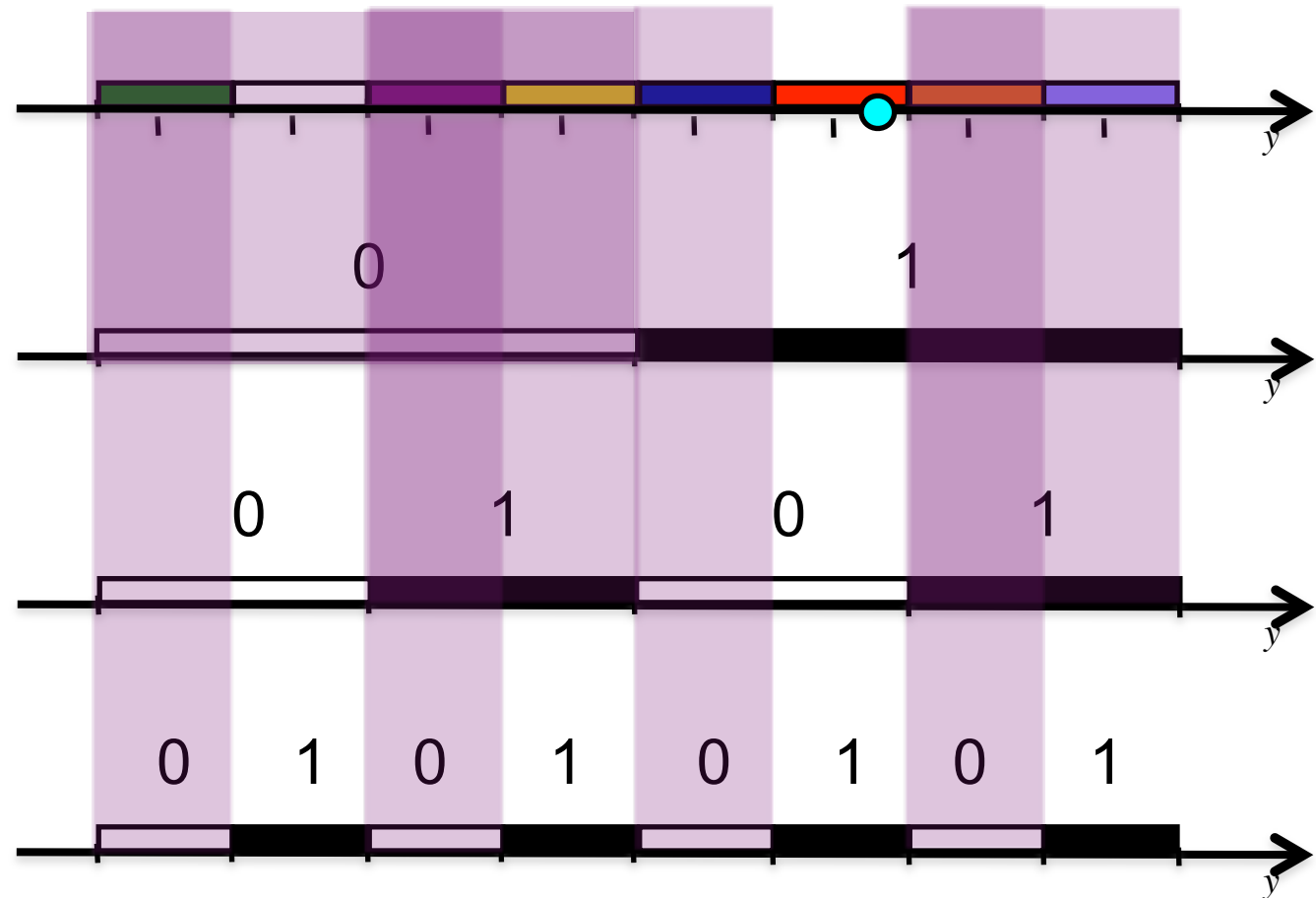
# What can a bit tell us?

3 bit quantization intervals

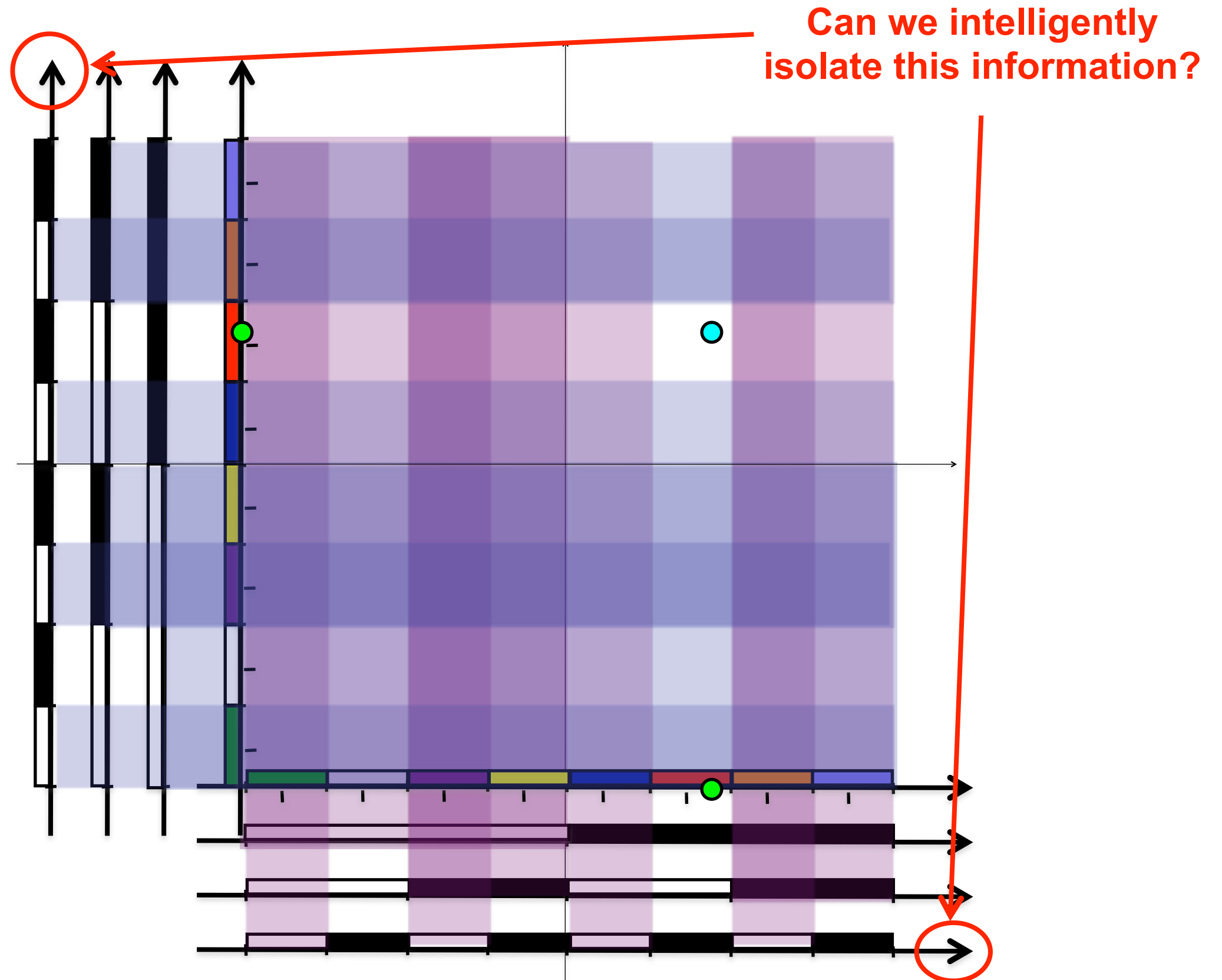
1<sup>st</sup> bit (MSB)

2<sup>nd</sup> bit

3<sup>rd</sup> bit (LSB)



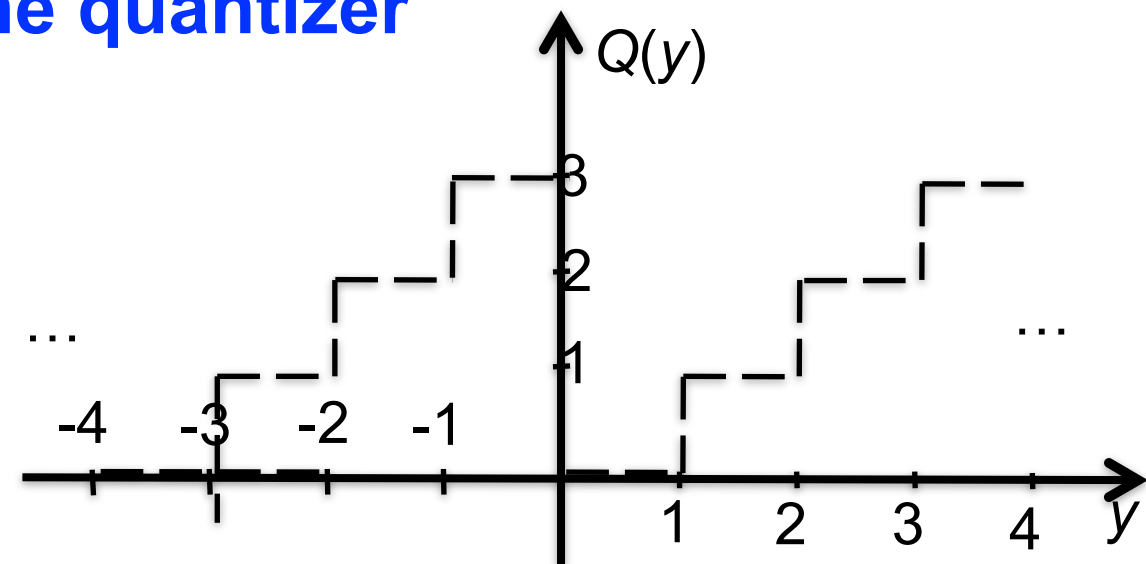
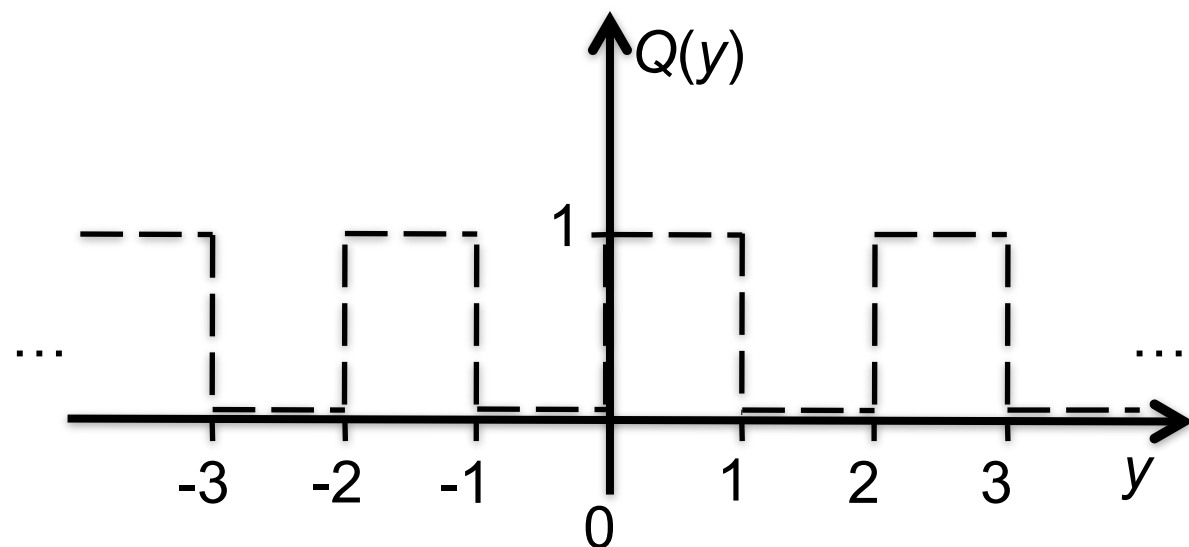
# What can a bit tell us?





# Rate-Efficient Scalar Quantization

Solution: **Modify the quantizer**



Non-monotonic quantizer: Multiple intervals quantize to same value  
(Focus on 1-bit quantizer today)

**measurements**  
(w/ i.i.d. gaussian matrix)

**dither**  
(i.i.d. uniform)

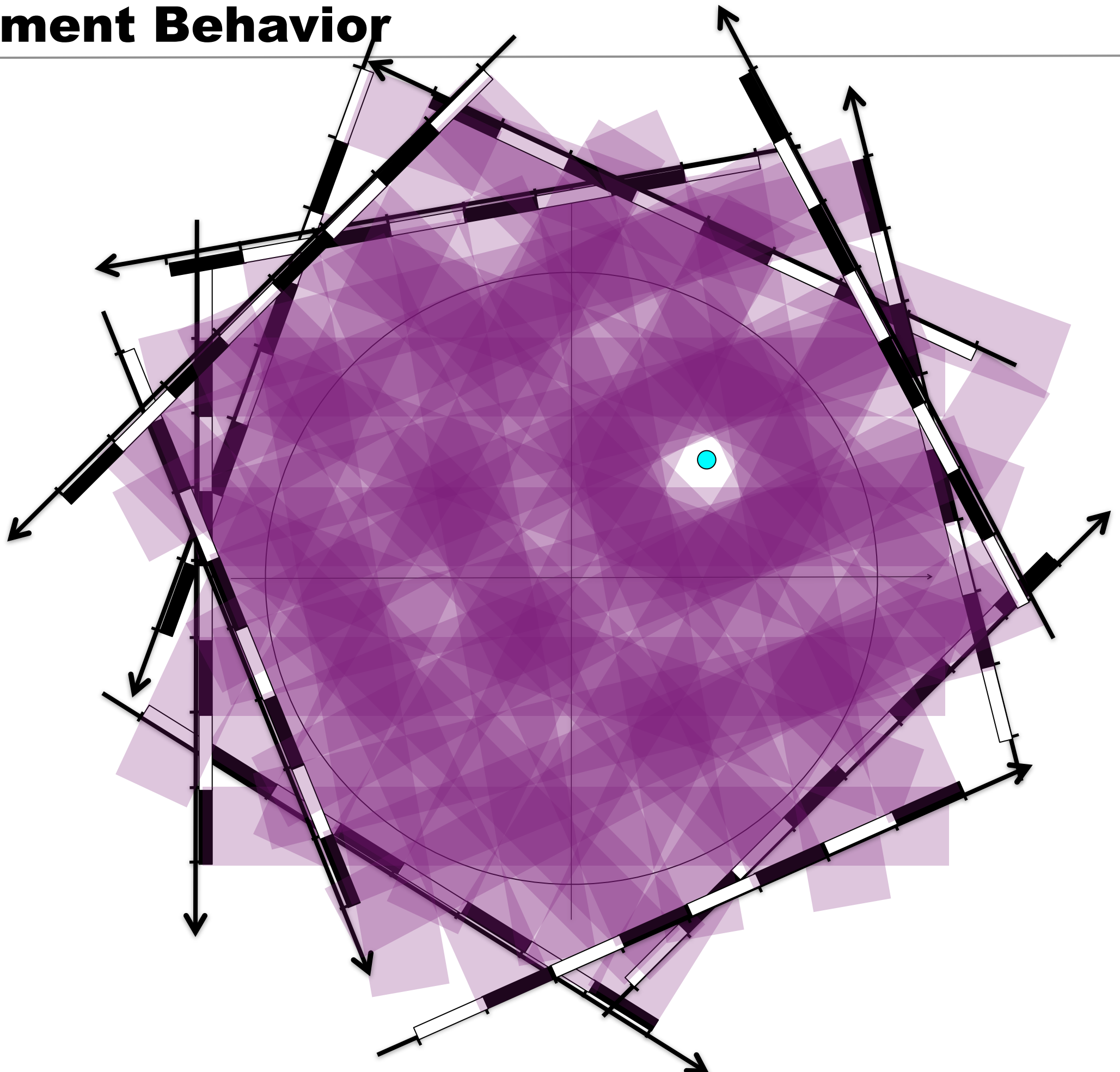
$$q_m = Q \left( \frac{\langle \mathbf{x}, \mathbf{a}_m \rangle + w_m}{\Delta_m} \right), \quad \mathbf{q} = Q(\Delta^{-1}(\mathbf{A}\mathbf{x} + \mathbf{w}))$$

**scalar quantizer**  
(non-monotonic)

**scaling/precision parameter**  
( $\Delta_m = \Delta$ , same for all measurements)

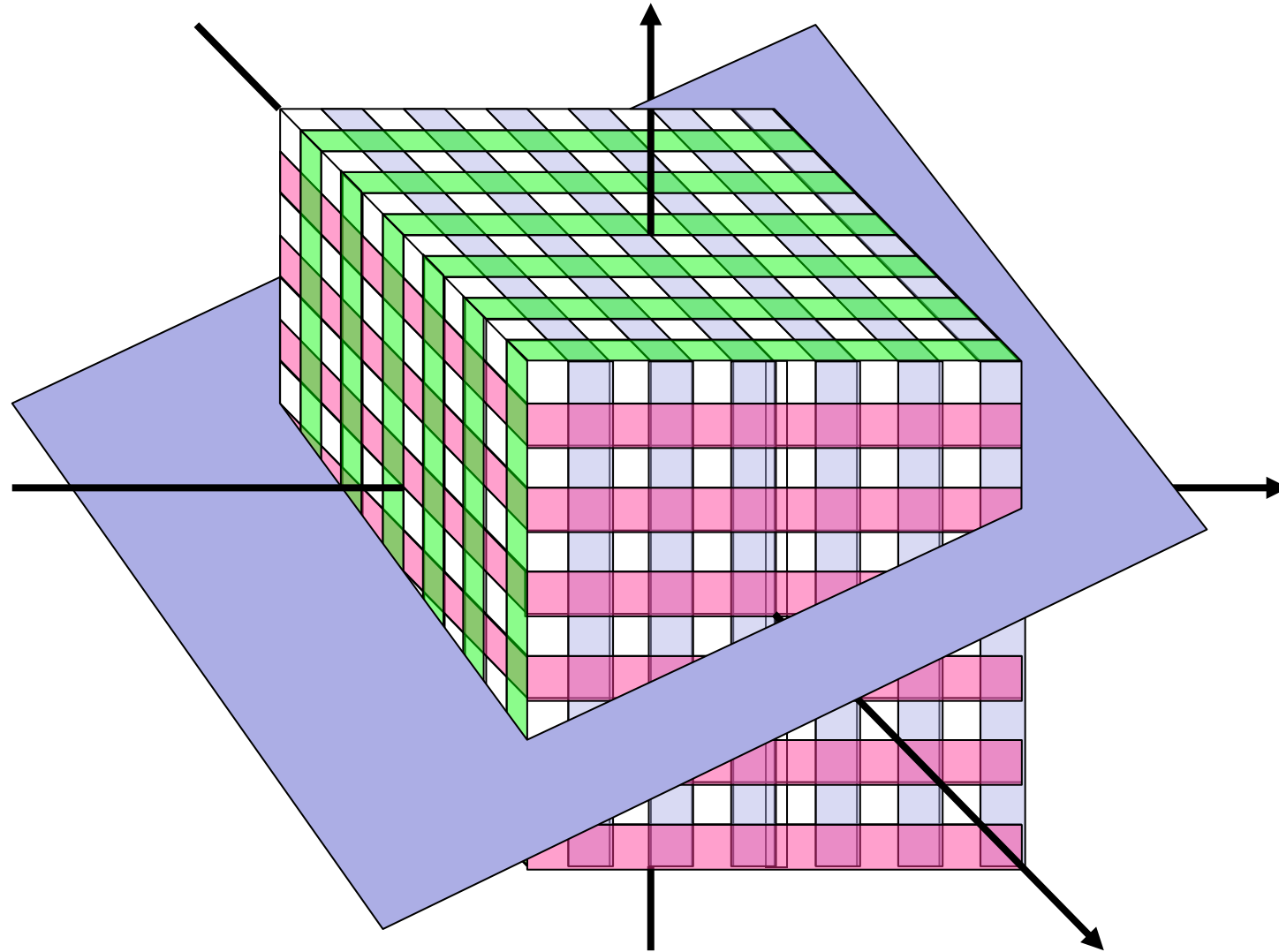
# Measurement Behavior

---



# Quantizer Geometry (1 bit)

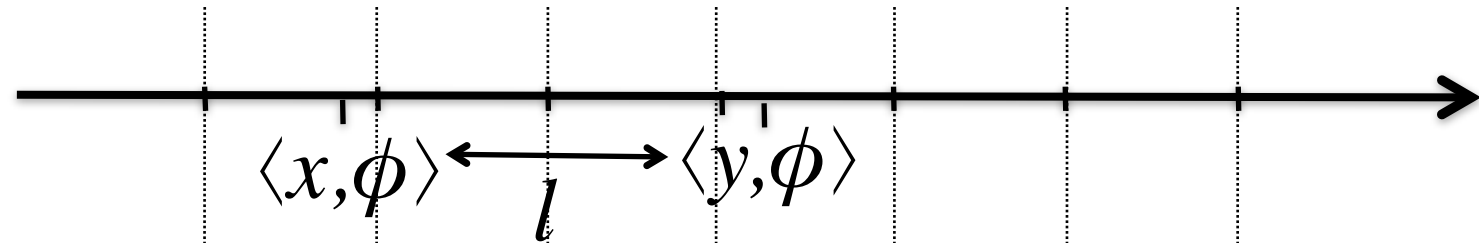
---



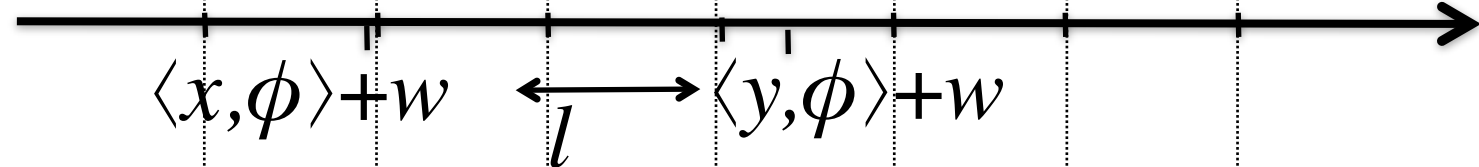
Quantization cells are **not** continuous  
Signal subspace intersects **most** of them

# Pairs of Signals, Single Measurement

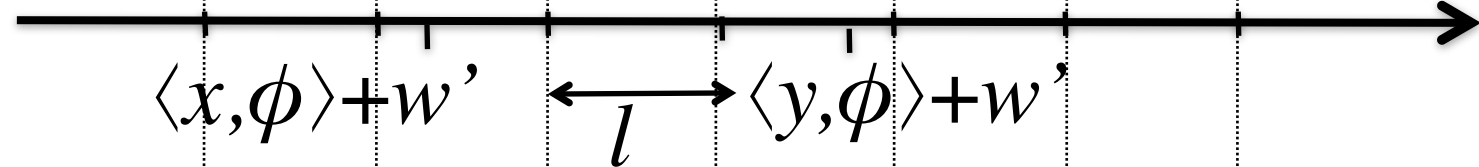
Projection  
(measurement)



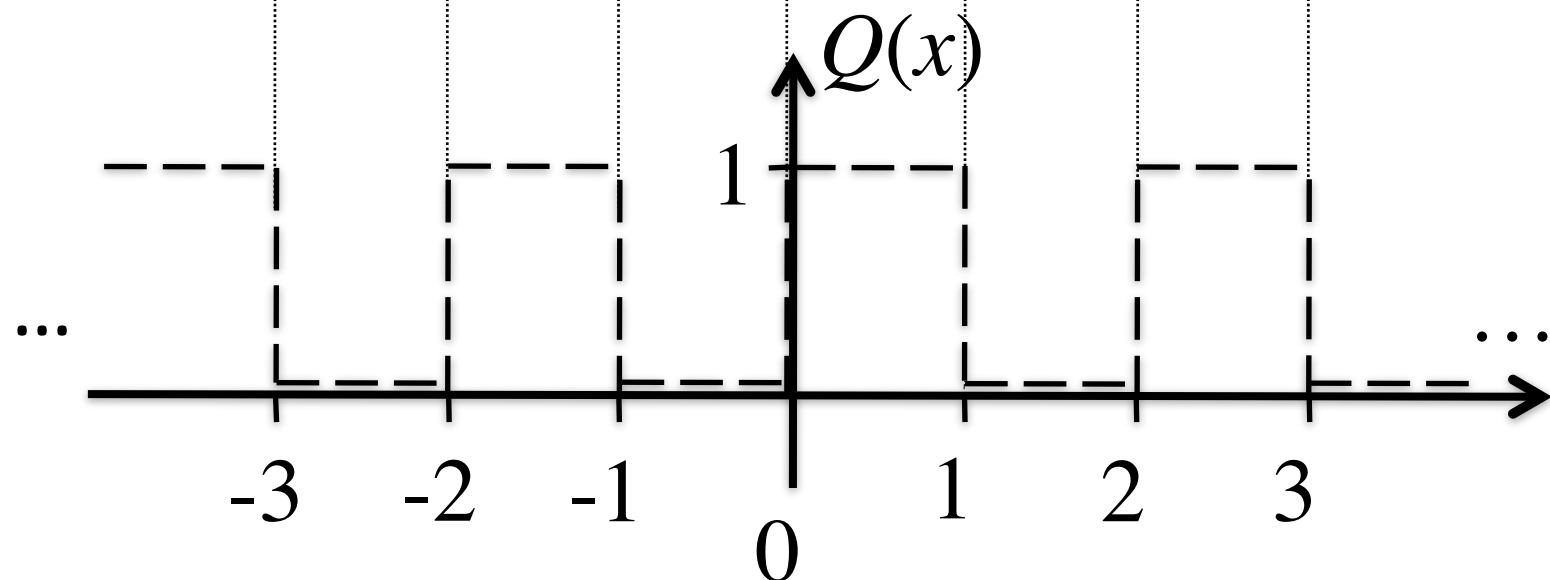
Projection  
with dither



Projection with  
different dither

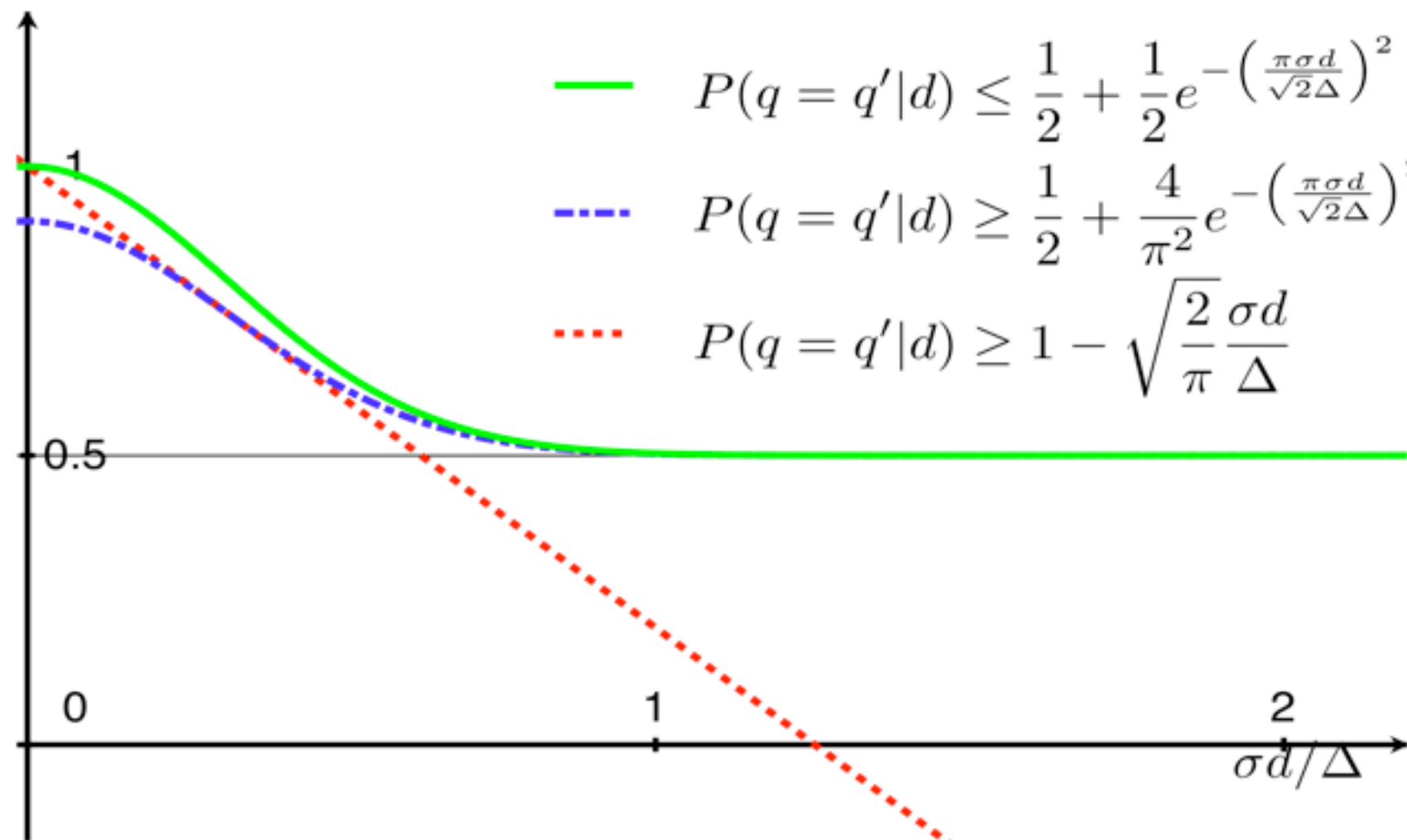


Quantization  
function



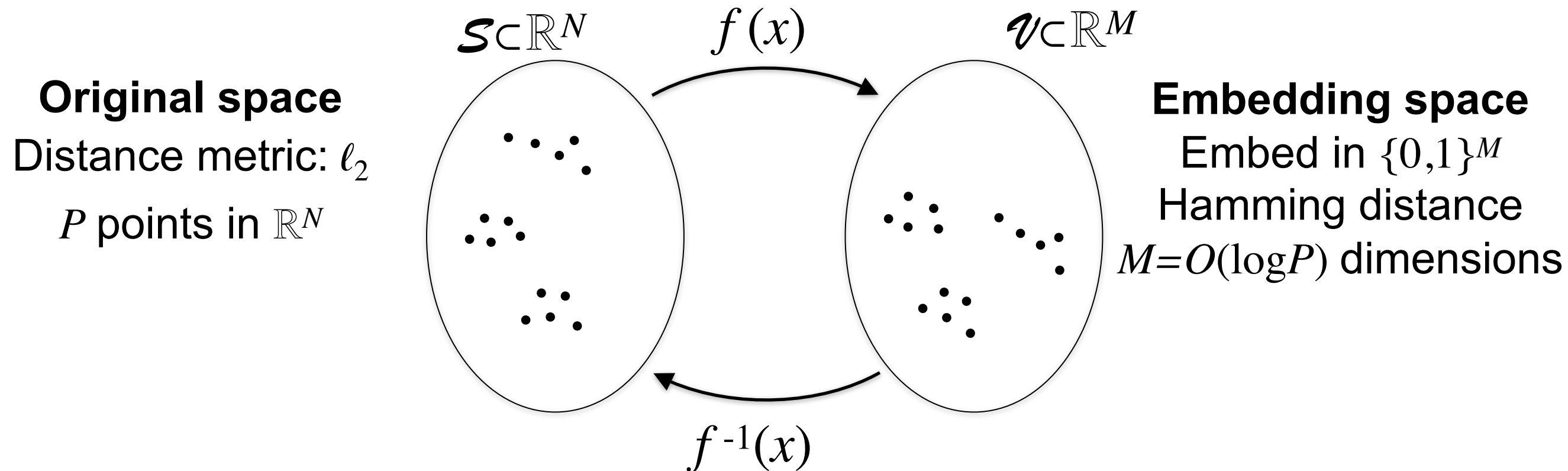
# Pairs of Signals, Single Measurement

$P(q=q')$  : probability that a single measurement is consistent for a pair of signals, given their distance  $d$



**In other words:**  
**Hamming distance** of embedding is  
**proportional to  $\ell_2$  distance**  
**up to a point**

# Embedding Properties

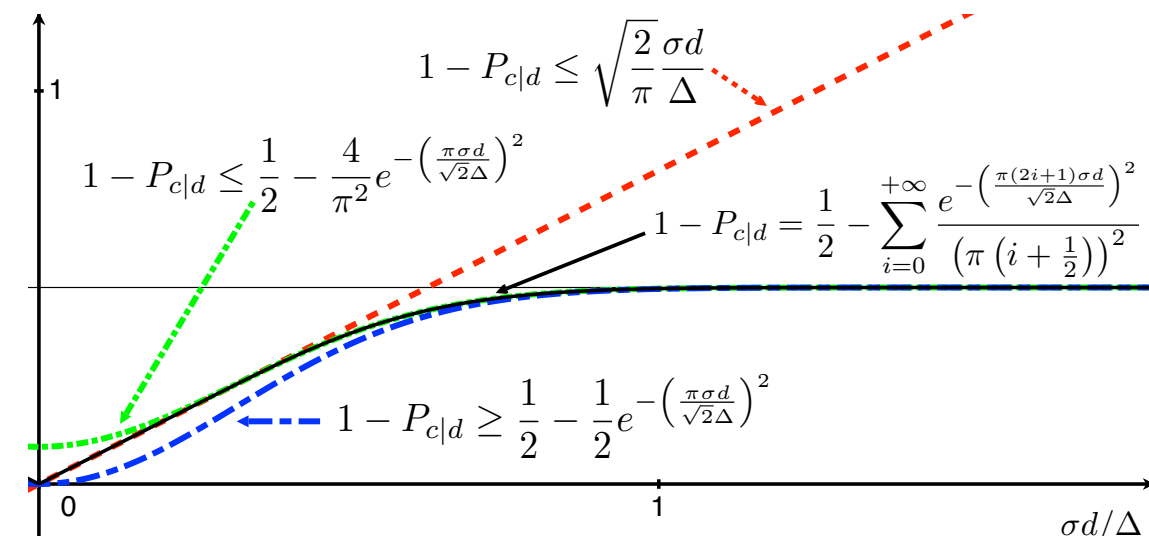


For all  $x, y$  in  $\mathcal{S}$ :

$$g(d) - \epsilon \leq d_H(f(x), f(y)) \leq g(d) + \epsilon$$

$$g(d) = 1 - P_{c|d}$$

as long as  $M = O\left(\frac{1}{\epsilon^2} \log P\right)$

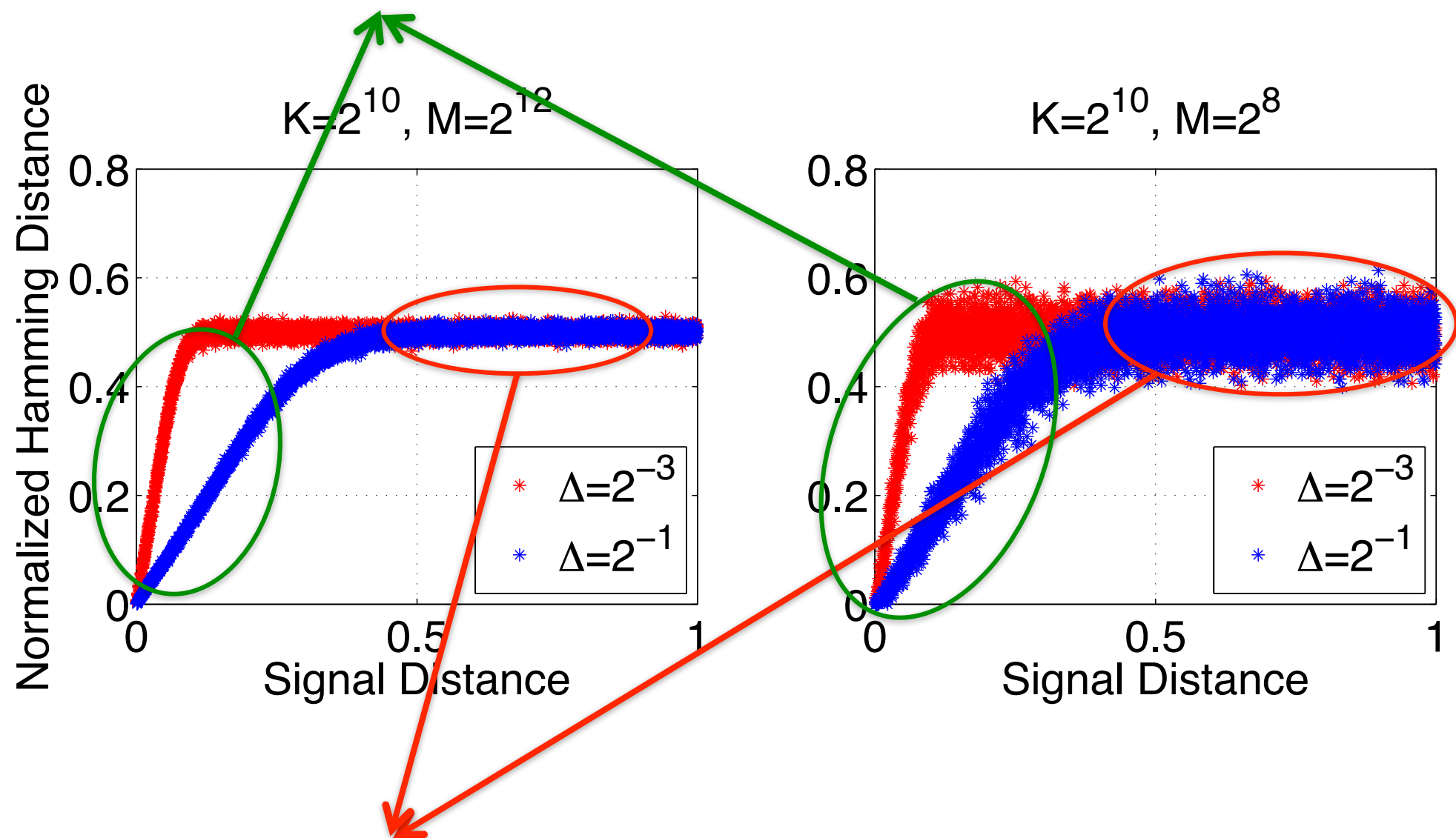


- Boufounos P. T. and Rane S., "Secure Binary Embeddings for Privacy Preserving Nearest Neighbors," *Proc. Workshop on Information Forensics and Security (WIFS)*, Foz do Iguaçu, Brazil, November 29 – December 2, 2011.

# Error Behavior

$$g(d) - \epsilon \leq d_H(f(x), f(y)) \leq g(d) + \epsilon$$

“Linear” region:  $\ell_2 \propto d_H$ , slope controlled by  $\Delta$



“Flat” region: no distance information

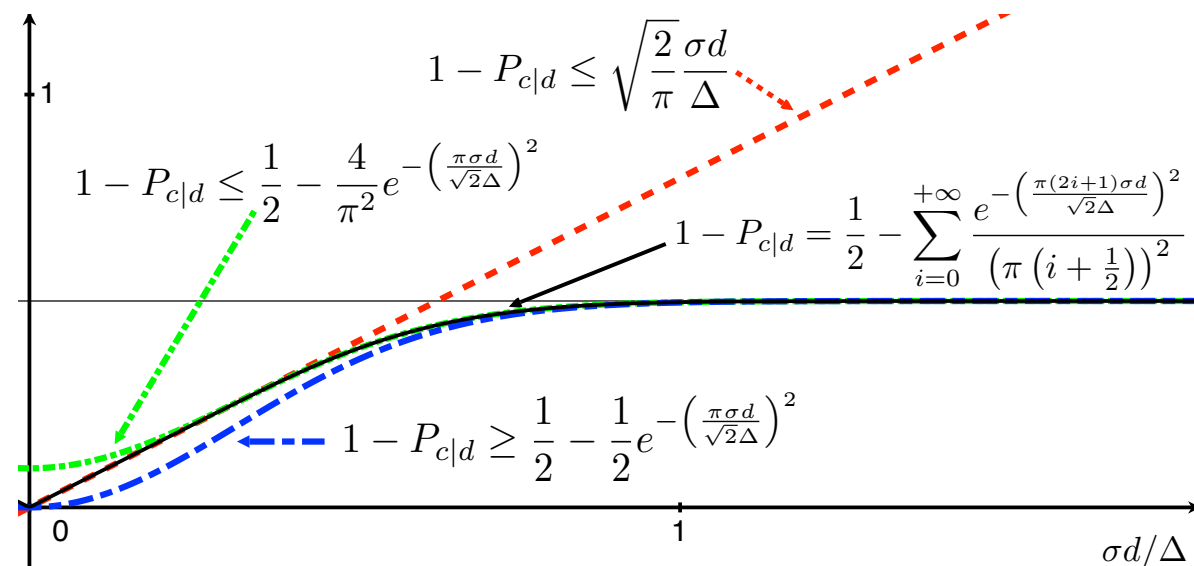


# Error Behavior

$$g(d) - \epsilon \leq d_H(f(x), f(y)) \leq g(d) + \epsilon$$

$$M = O\left(\frac{1}{\epsilon^2} \log P\right)$$

Similar trade-off as J-L but on  $g(d)=1-P_{c|d}$



**Distance estimate:**  $\hat{d} = g^{-1}(d_H(f(x), f(y)))$

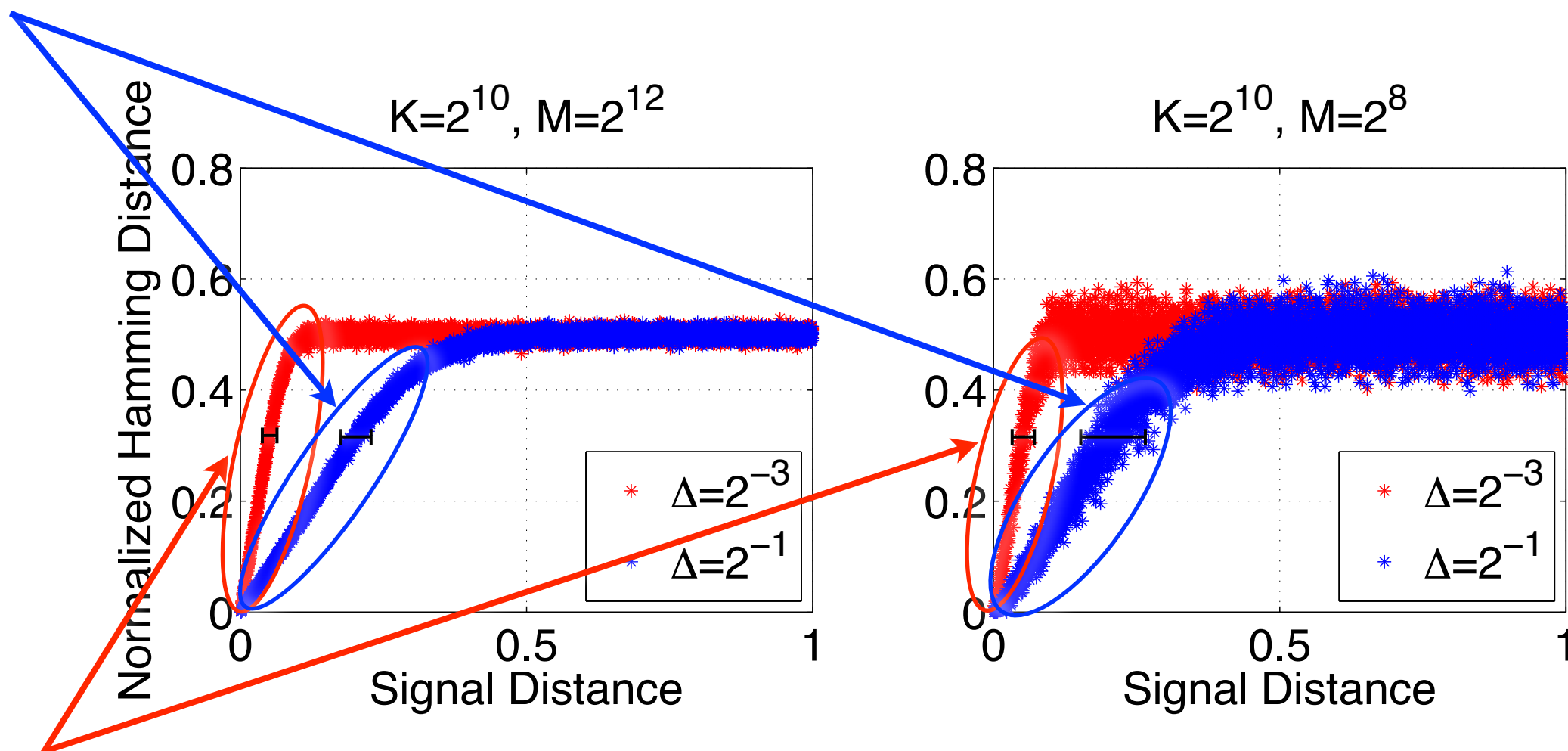
**Estimate ambiguity:**  $\hat{d} - \frac{\epsilon}{g'(\hat{d})} \lesssim d \lesssim \hat{d} + \frac{\epsilon}{g'(\hat{d})}$

**Properties (slope) controlled by choice of  $\Delta$**

# Error Behavior

$$g(d) - \epsilon \leq d_H(f(x), f(y)) \leq g(d) + \epsilon$$

**Large  $\Delta$ :** small slope, more ambiguity, preserves larger distances

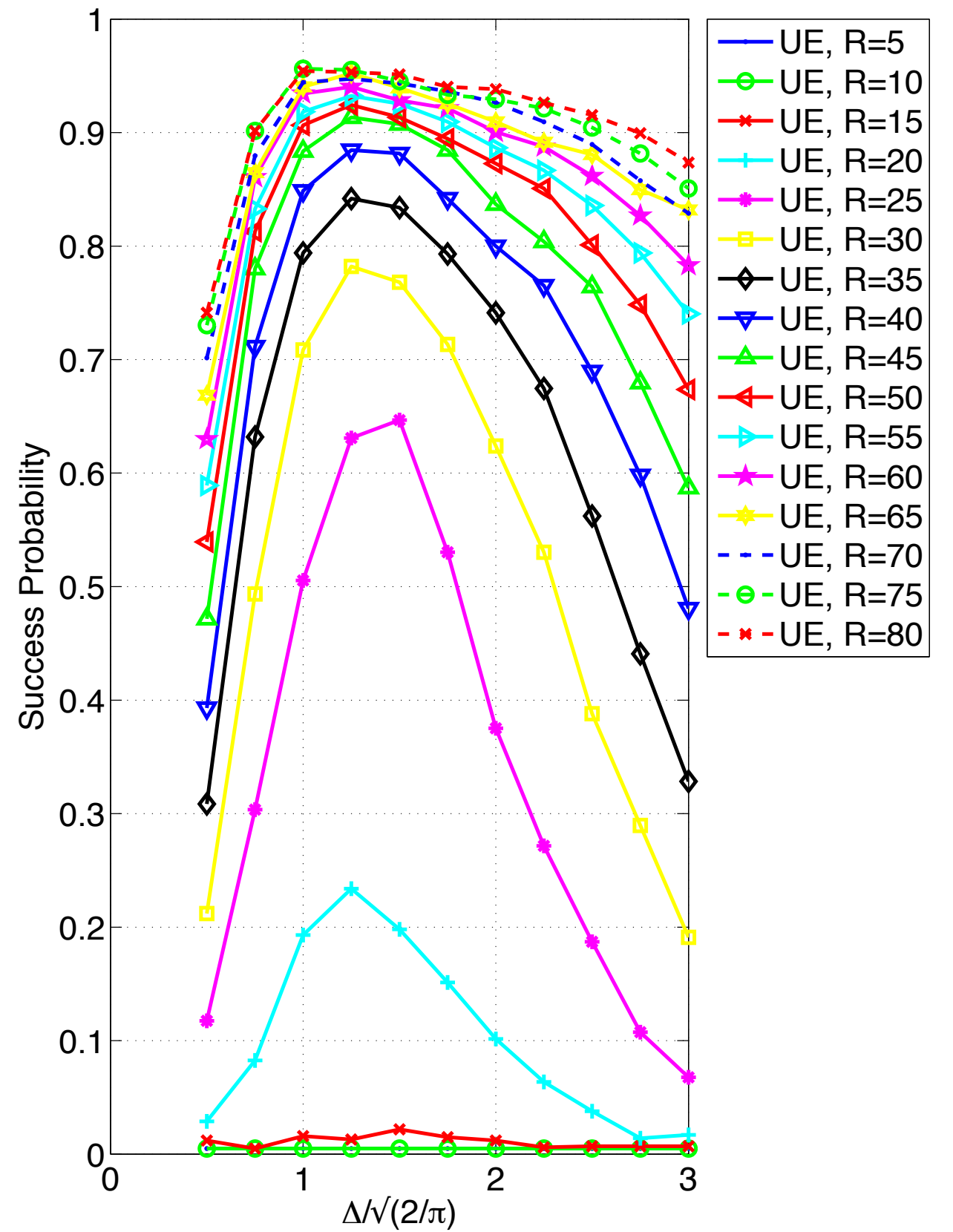
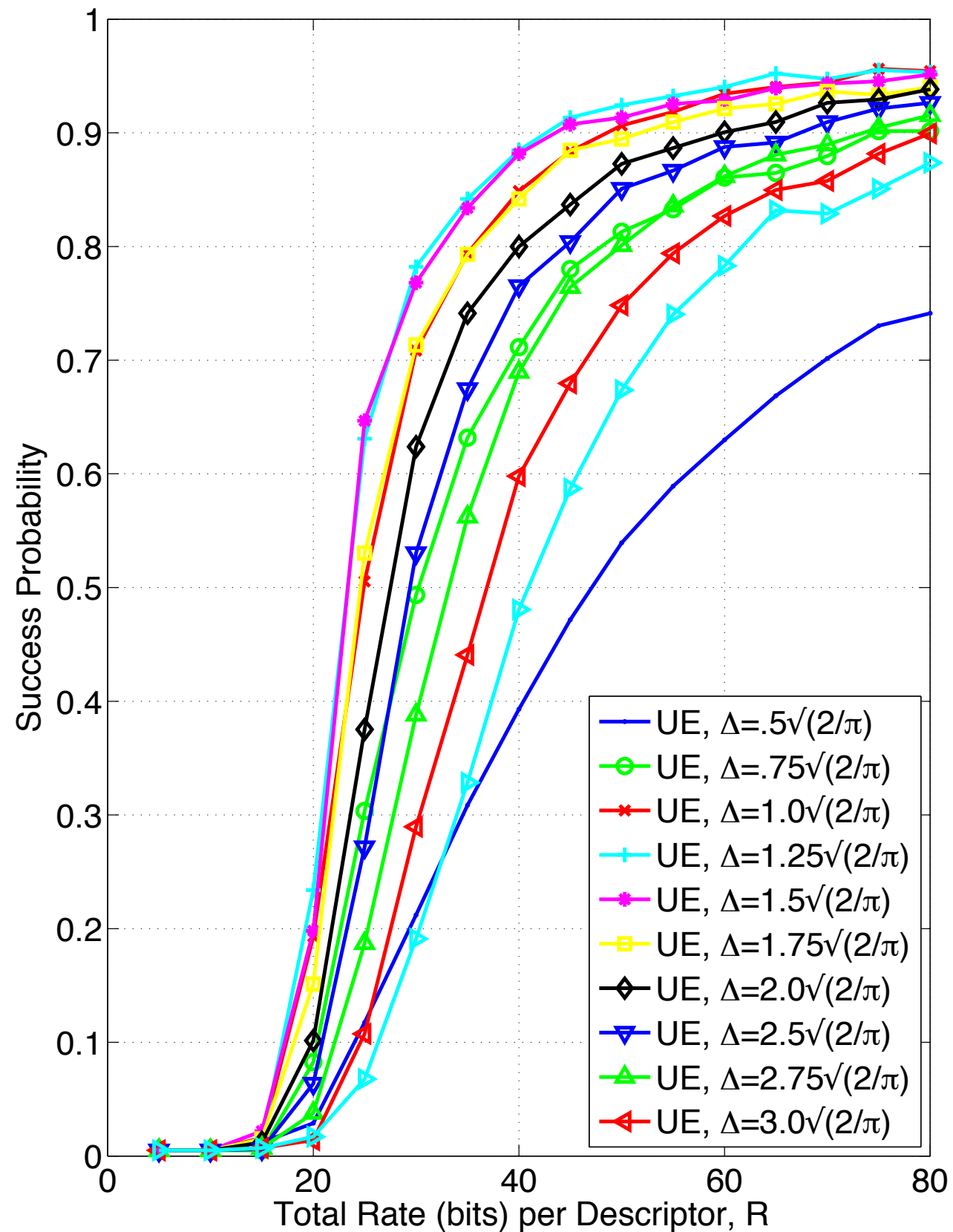


**Small  $\Delta$ :** large slope, less ambiguity, preserves smaller distances

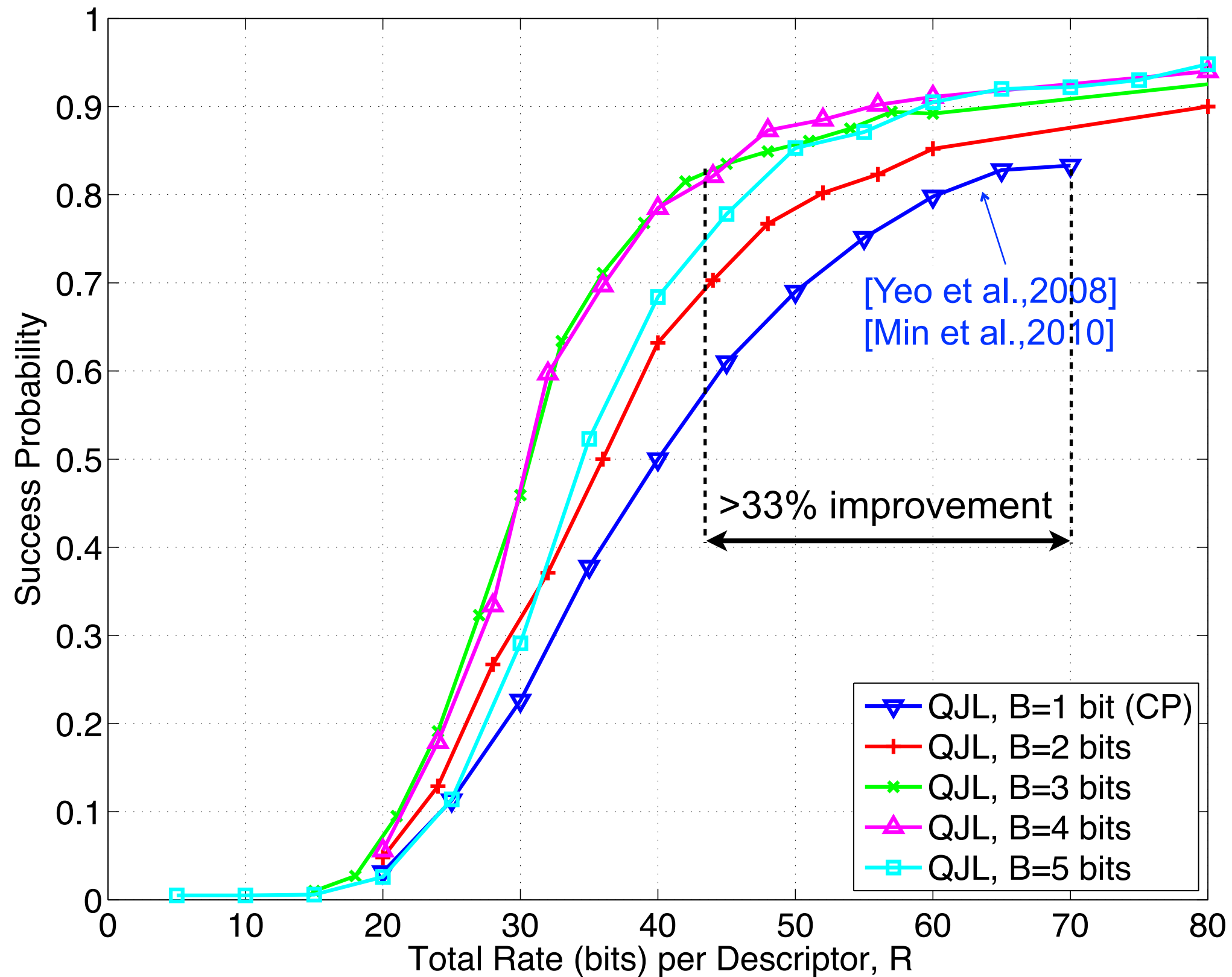
- Boufounos P. T. and Rane S., “Efficient Coding of Signal Distances Using Universal Quantized Embeddings,” *Proc. Data Compression Conference (DCC)*, Snowbird, UT, March 20-22, 2013.

**IN PRACTICE**

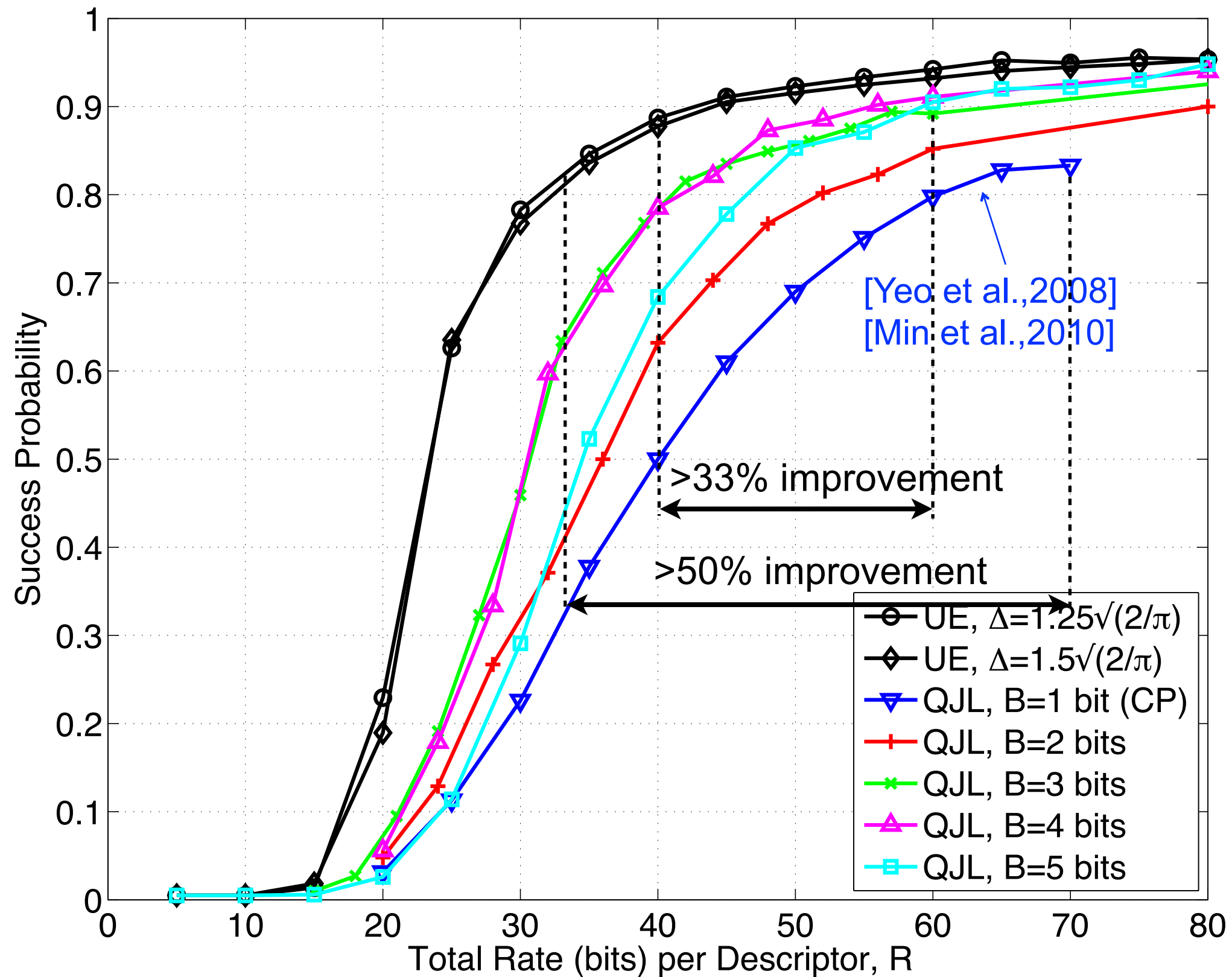
# In practice



# In practice



# In practice



# **BEYOND EMBEDDINGS**



# Reconstruction

---

- **Consistent reconstruction**: find a signal that quantizes to same bits, i.e.,

$$\hat{\mathbf{x}} \text{ s.t. } \mathbf{q} = Q \left( \Delta^{-1} (\Phi \hat{\mathbf{x}} + \mathbf{w}) \right)$$

- Very **good theoretical guarantees**

- **Exponential error decay** with number of bits  $\varepsilon = O(c^{-B})$

- Reconstruction is a **very hard** problem

- Seems to have combinatorial complexity
  - Probably NP

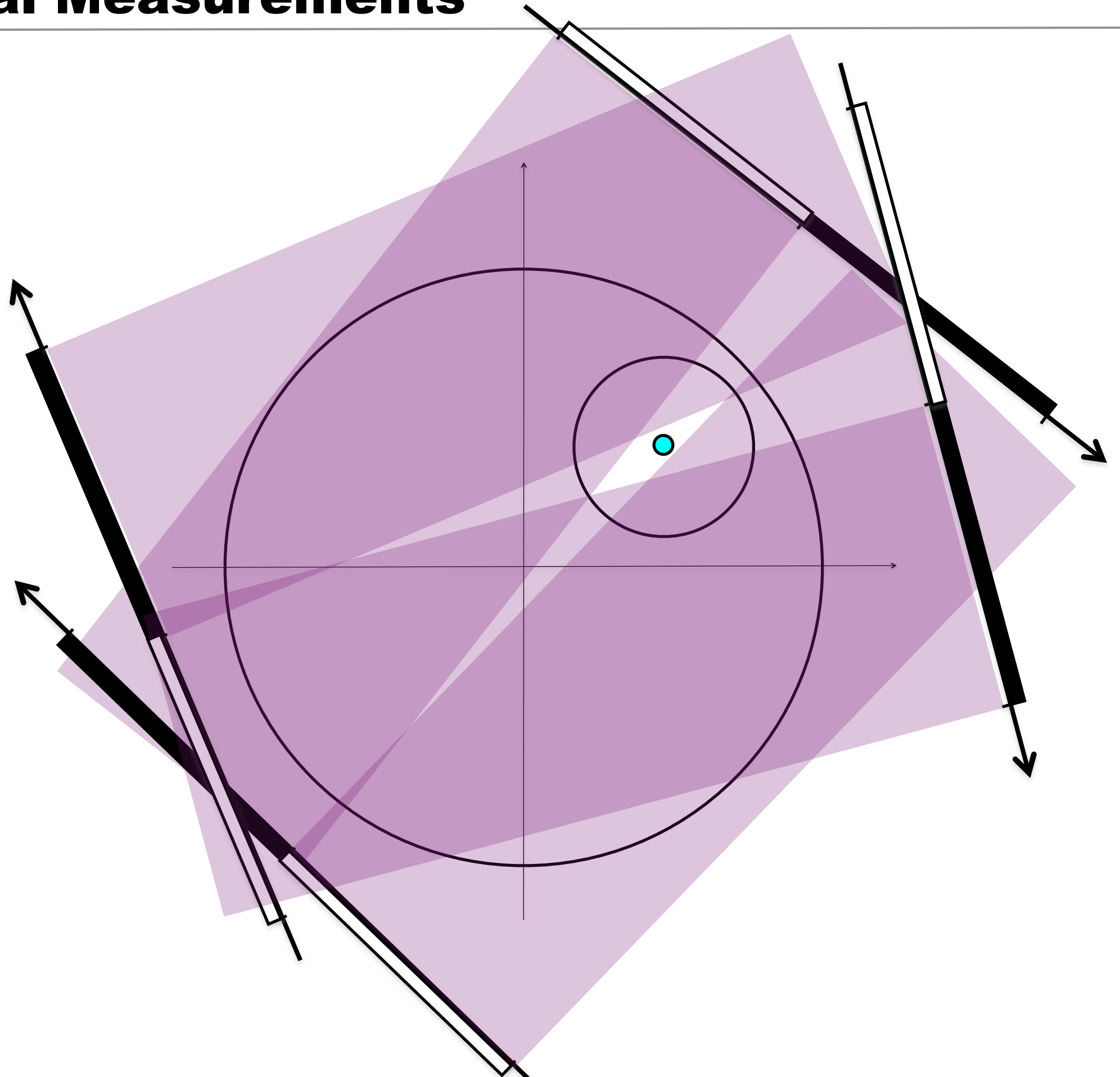
- Need to **enable efficient reconstruction**

- Classical methods **exploit bit hierarchy** to make problem convex
  - Should **maintain theoretical guarantees**

- **Solution**: Construct **bit hierarchy**; **sub-problems** become **convex**

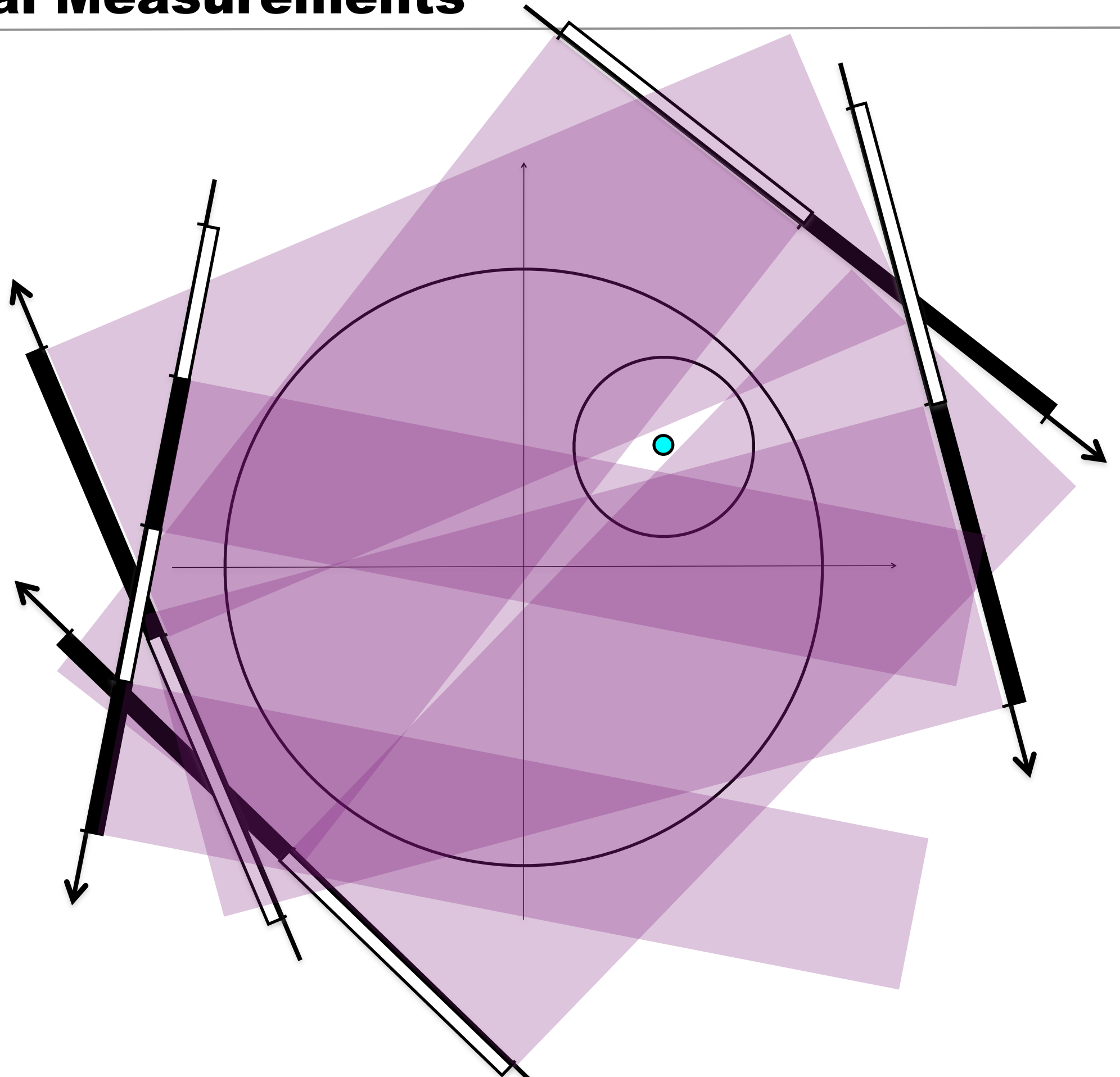
# Hierarchical Measurements

---



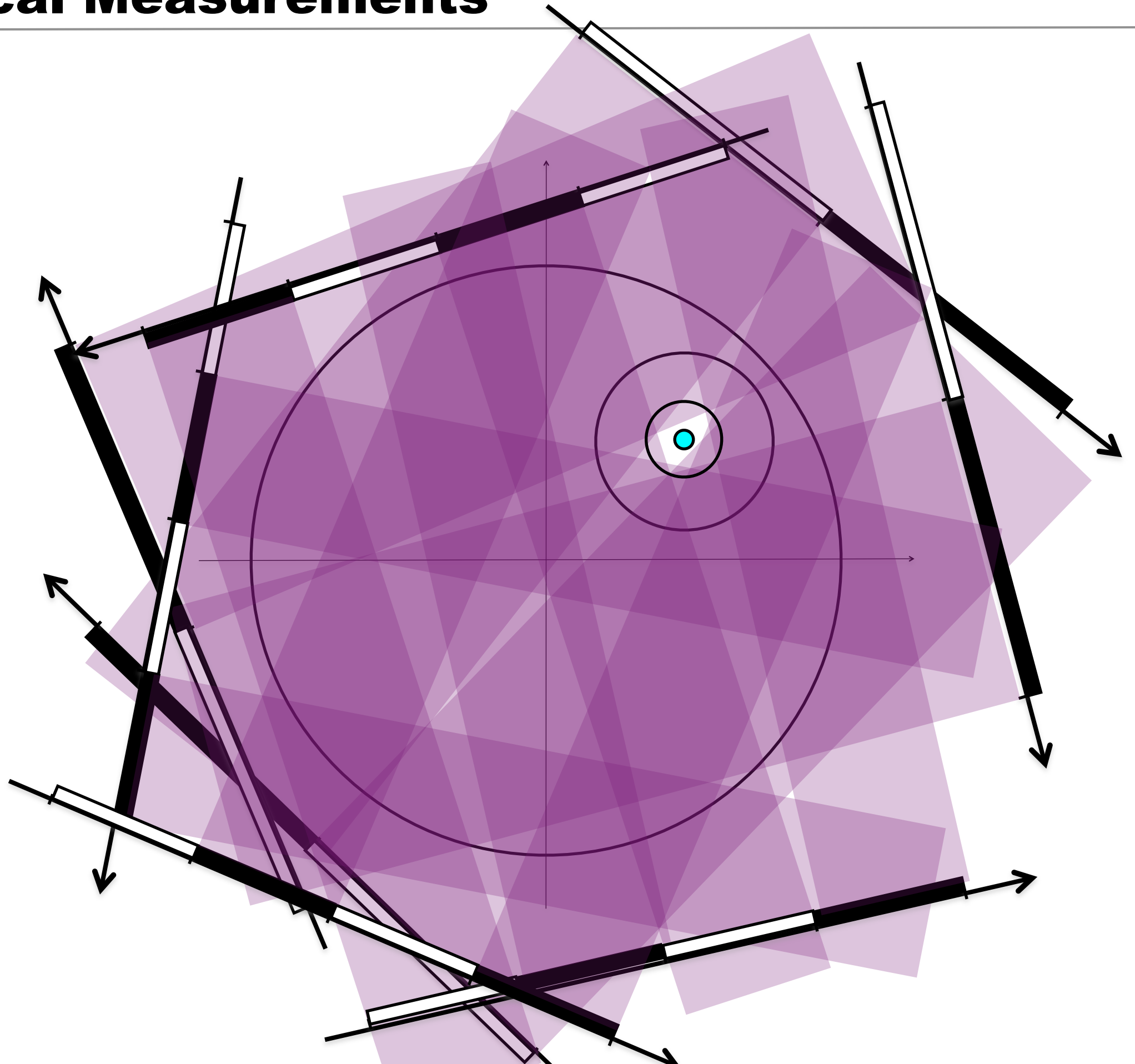
# Hierarchical Measurements

---



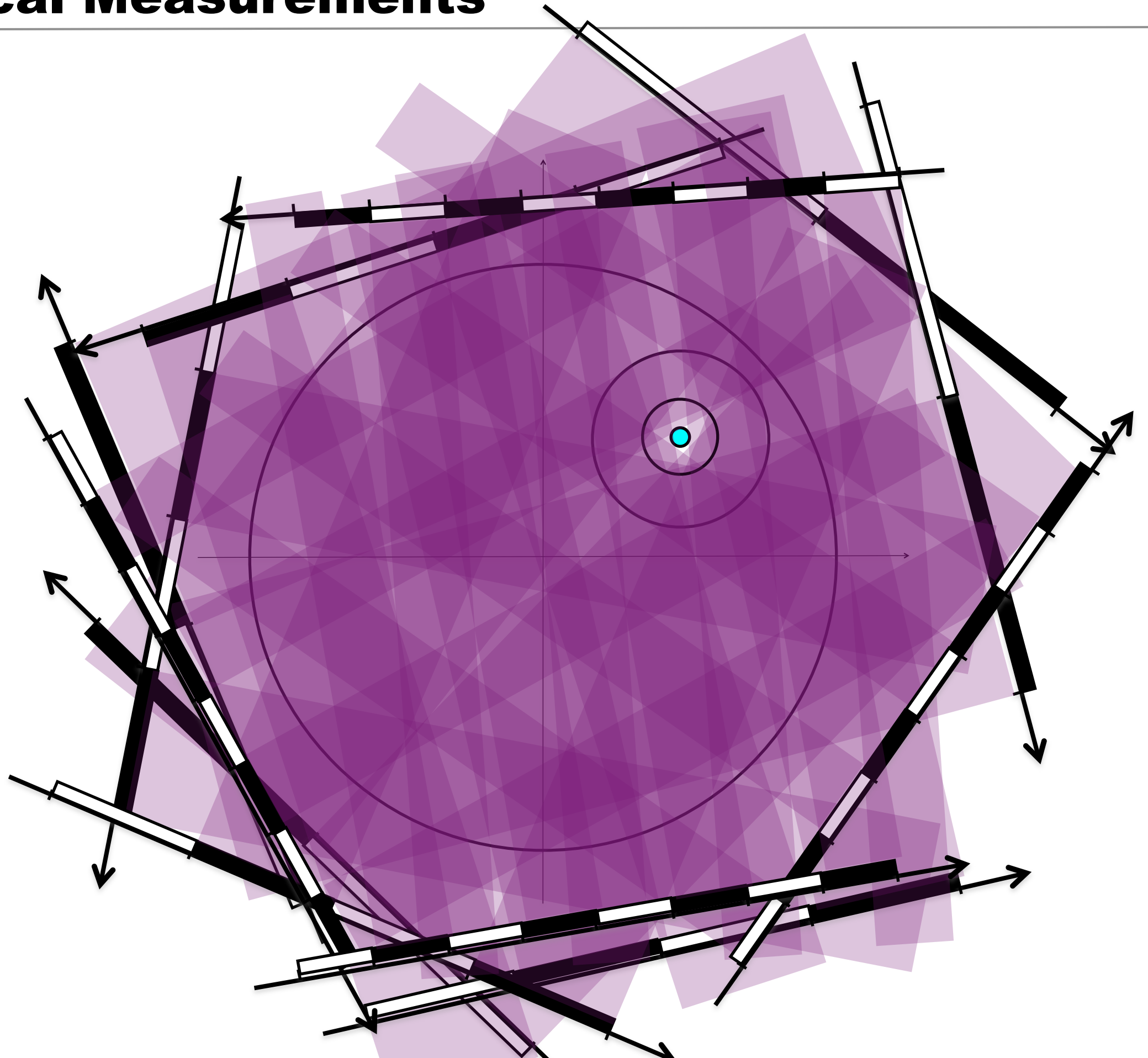
# Hierarchical Measurements

---



# Hierarchical Measurements

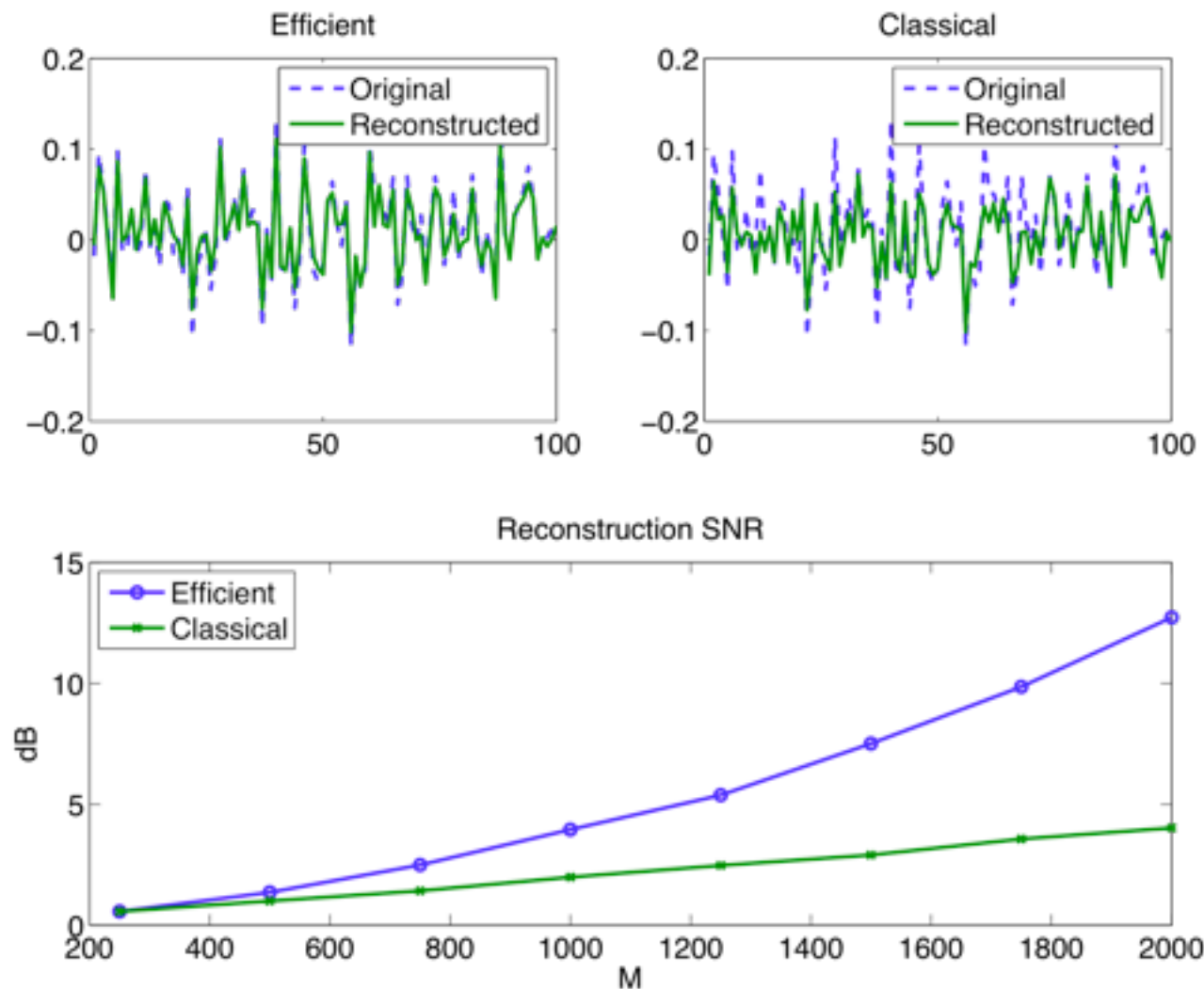
---



# Reconstruction

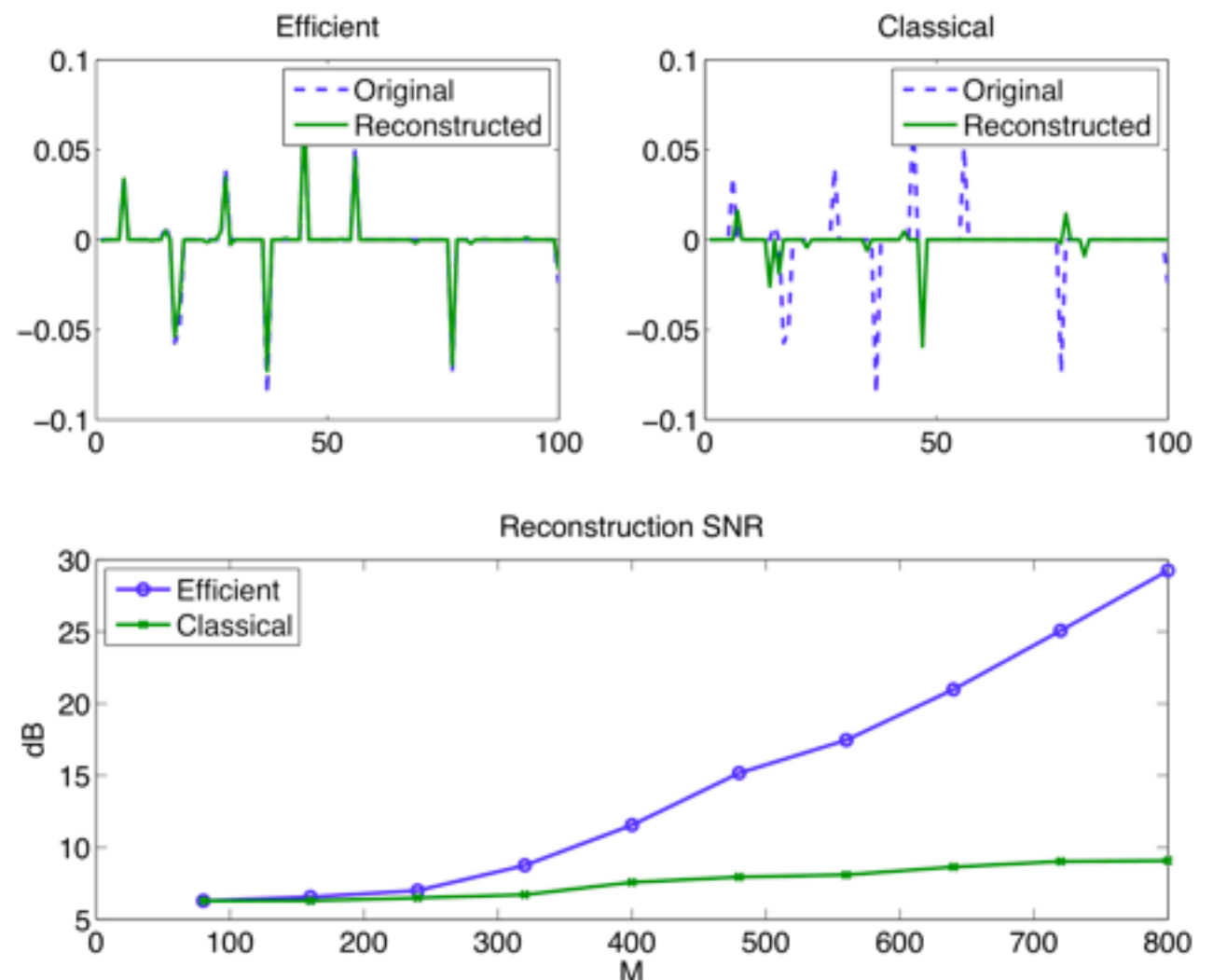
## Sampling

- Given uncertainty pick  $\Delta$  such that reconstruction is convex
- Take enough measurements to scale uncertainty by  $\alpha < 1$
- Scale  $\Delta \leftarrow \alpha\Delta$  for next set of measurements
- Iterate until desired precision



## Reconstruction

- For first set of measurements formulate convex reconstruction
- Solve for consistency
- Use solution to incorporate next set of measurements and determine consistency constraints
- Iterate until all measurement sets are incorporated





# Privacy Preserving Properties

---

Assume we have encoding of two signals  $x, y$ , but not  $\mathbf{A}$  and  $\mathbf{w}$   
What does the encoding reveal about their relationship?

$$I(f(x); f(y) | d) \leq 10M e^{-\left(\frac{\pi \sigma d}{\Delta}\right)^2}$$

**Mutual information decays very fast with  $d$ .**

Information theoretic privacy-preserving guarantee:

**When signals are far apart,  
encoding reveals nothing about their relationship!**

Very useful for security applications  
(e.g., privacy-preserving nearest neighbors,  
secure biometric authentication)

- Boufounos P. T. and Rane S., “Secure Binary Embeddings for Privacy Preserving Nearest Neighbors,” *Proc. Workshop on Information Forensics and Security (WIFS)*, Foz do Iguaçu, Brazil, November 29 – December 2, 2011.

# Further Reading

---

- Johnson W. and Lindenstrauss J., “Extensions of Lipschitz mappings into a Hilbert space,” *Contemporary Mathematics*, vol. 26, pp. 189–206, 1984.
- Andoni, A. and Indyk, P., “Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions,” *Comm. ACM*, vol. 51, no. 1, pp. 117–122, 2008.
- Datar M., Immorlica N., Indyk P., and Mirrokni V., “Locality-Sensitive Hashing Scheme Based on p-Stable Distributions,” *Proc. Symposium on Computational Geometry*, 2004
- Jacques L., Laska J. N., Boufounos P. T., Baraniuk R. J., "Robust 1-Bit Compressive Sensing via Binary Stable Embeddings of Sparse Vectors," *IEEE Trans. Info. Theory*, v. 59, no. 4, April, 2013.
- Plan, Y. and Vershynin, R., “Dimension reduction by random hyperplane tessellations,” preprint, arXiv:1111.4452, 2011.
- Ai, A., Lapanowski, A., Plan, Y., Vershynin, R., “One-bit compressed sensing with non-Gaussian measurements”, *Linear Algebra and Applications*, to appear.
- Boufounos P. T., "Universal Rate-Efficient Scalar Quantization," *IEEE Trans. Info. Theory*, v. 58, no. 3, pp. 1861-1872, March, 2012.
- Boufounos P. T. and Rane S., “Secure Binary Embeddings for Privacy Preserving Nearest Neighbors,” *Proc. Workshop on Information Forensics and Security (WIFS)*, Foz do Iguaçu, Brazil, November 29 – December 2, 2011.
- Li M., Rane S., and Boufounos P. T., “Quantized embeddings of scale-invariant image features for mobile augmented reality,” *IEEE 14th International Workshop on Multimedia Signal Processing (MMSP)*, Banff, Canada, Sept. 17-19, 2012
- Boufounos P. T. and Rane S., “Efficient Coding of Signal Distances Using Universal Quantized Embeddings,” *Proc. Data Compression Conference (DCC)*, Snowbird, UT, March 20-22, 2013.
- Yeo C., Ahammad P., and Ramchandran K., “Coding of image feature descriptors for distributed rate-efficient visual correspondences,” *International Journal of Computer Vision*, vol. 94, pp. 267–281, 2011, 10.1007/s11263-011-0427-1.
- Min K., Yang L., Wright J., Wu L., Hua X.-S., and Ma Y., “Compact projection: Simple and efficient near neighbor search with practical memory requirements,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2010.



# Final Thoughts/Discussion

---

- Quantization very important in signal processing
  - Signal acquisition systems
  - Information embedding/transmission
  - Information hiding, security, privacy
- Not your college-level quantization
  - High-dimensional geometrical problem
  - Additive noise model inadequate
  - Tight bounds with better models
  - Consistency is important
  - Saturation can be useful
  - Non-linear concentration of measure occurs
- Quantization is a very active area of research
  - 1-bit CS/Quantized CS
  - Sigma-Delta for compressive and non-compressive systems
  - Geometry of non-linear inverse problem solving
  - Quantized Embeddings
  - Vector Quantization (a whole other tutorial)

# Still Open and Interesting (a small sampling...)

---

- Oversampling and Quantization
  - Beyond consistency: quantization with additive noise
- Quantized CS
  - Interaction between sparsity/sensing/quantization (signal/measurement model)
  - 1-bit CS algorithmic convergence guarantees (e.g. BIHT, RSS, MSP)
  - Consistent QCS theory for any bitdepth (1-bit to high-res)
  - Optimal quantizer design for non-gaussian measurements
  - Rate-distortion performance: CS for compression
  - Sigma-Delta CS for 1-bit quantization
  - Vector Quantization of CS measurements
- Universal Quantization and Embeddings
  - General reconstruction algorithms: is reconstruction possible?
  - Embedding guarantees for more general embeddings (e.g. multi-bit)
  - Embedding behavior design
  - Tighter connections with LSH
  - Other security/privacy-preserving properties

# Today's Topics

---

1. Modern Scalar Quantization
2. Compressive Sensing Overview
3. Compressive Sensing and Quantization
4. 1-bit Compressive Sensing
5. Locality Sensitive Hashing and Universal Quantization

**For more:**

**Repository:** <http://www.boufounos.com/resources-on-quantization/>

<http://dsp.rice.edu/1bitCS/>

<http://nuit-blanche.blogspot.com>

<http://nuit-blanche.blogspot.com/search/label/1bit>

<http://nuit-blanche.blogspot.com/search/label/QuantCS>

<http://www.boufounos.com/research/quantization/>

**Questions/Comments?**

[petros@boufounos.com](mailto:petros@boufounos.com)

<http://boufounos.com>

[laurent.jacques@uclouvain.be](mailto:laurent.jacques@uclouvain.be)

<http://perso.uclouvain.be/laurent.jacques/>