

Quantized Iterative Hard Thresholding: Bridging 1-bit and High Resolution Quantized Compressed Sensing

Laurent Jacques, Kévin Degraux, Christophe De Vleeschouwer

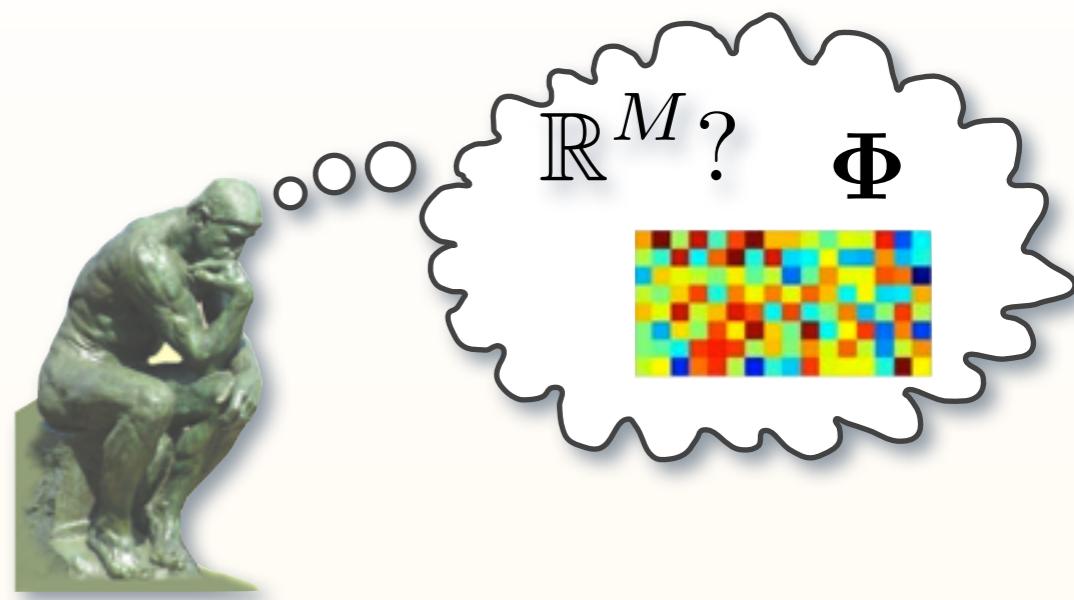


Louvain University (UCL), Louvain-la-Neuve, Belgium

Compressive Sampling and Quantization

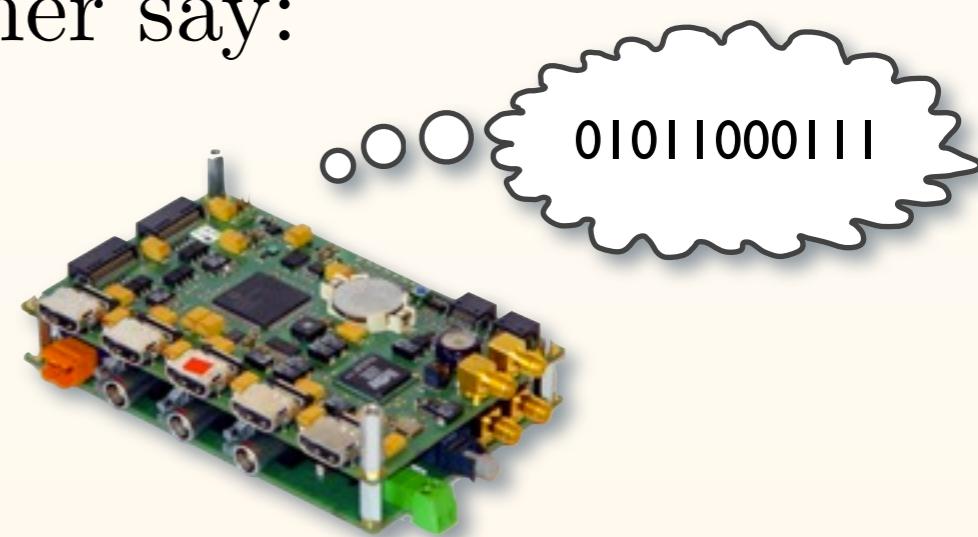
Compressed sensing theory says:

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



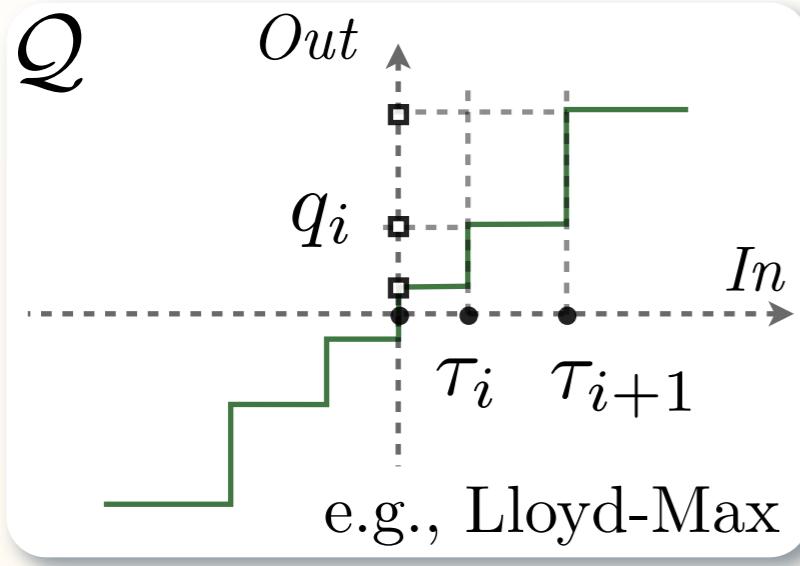
Information theory and sensor designer say:

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



(this talk) Focus on scalar quantization

Turning measurements into bits → scalar quantization



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

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

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

Important conventions:

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

1. First extreme: CS and high-resolution scalar quantization

Former solutions

- ▶ Quantization is like a noise

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

quantization
distortion

Former solutions

- ▶ Quantization is like a noise

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

- ▶ CS is robust, *e.g.*, with ...
 - ▶ Basis Pursuit DeNoise (BPDN):

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

- ▶ Iterative Hard Thresholding (IHT):

$$\hat{x} = \underset{\mathbf{u} \in \mathbb{R}^N}{\operatorname{argmin}} \|\Phi \mathbf{u} - y\|^2 \text{ s.t. } \|\mathbf{u}\|_0 \leq K$$

approximating: $\mathbf{x}^{(n+1)} = \mathcal{H}_K(\mathbf{x}^{(n)} + \mu \Phi^* (\mathbf{y} - \Phi \mathbf{x}^{(n)}))$

Former solutions

- ▶ Quantization is like a noise

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

- ▶ CS is robust (e.g., with BPDN or IHT)

If $\|n\| \leq \epsilon$ and $\frac{1}{\sqrt{M}}\Phi$ is RIP(δ, aK) with $\delta \leq \delta_0$, then

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

for some $C, D > 0$ and $e(K) = \text{deviation to } K\text{-sparsity}$.

$$e(K) = \|x - x_K\|_1/\sqrt{K} \quad \text{or} \quad \|x - x_K\|_2 + \|x - x_K\|_1/\sqrt{K}$$

[Candès, 08] ($a = 2, \delta_0 < \sqrt{2} - 1$)

[Blumensath, Davies, 09] ($a = 3, \delta_0 < 1/32$)

Former solutions

- Quantization is like a noise

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

- CS is robust (e.g., with BPDN or IHT)

If $\|n\| \leq \epsilon$ and $\frac{1}{\sqrt{M}}\Phi$ is RIP(δ, aK) with $\delta \leq \delta_0$, then

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

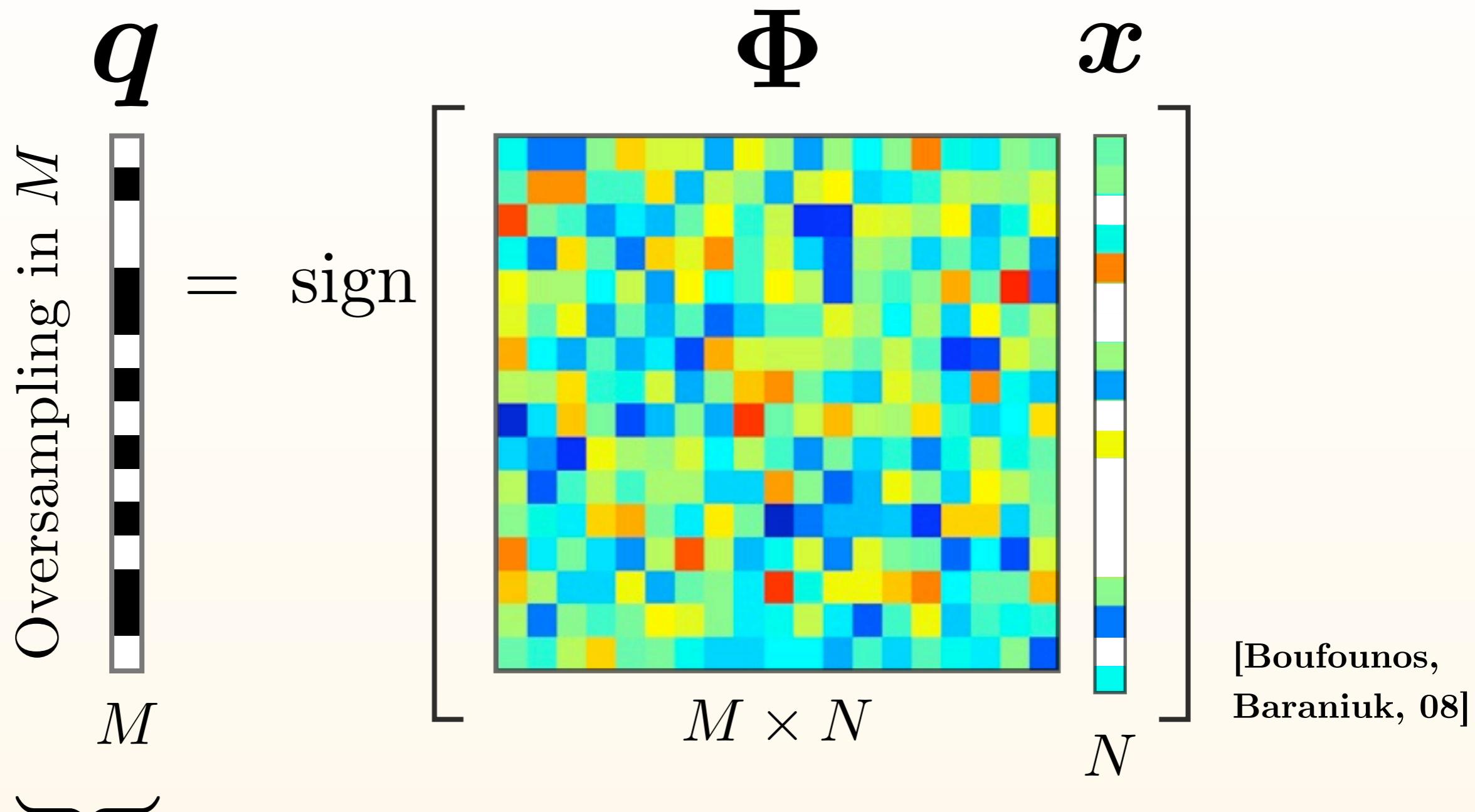
for some $C, D > 0$ and $e(K) = \text{deviation to } K\text{-sparsity}$.

- Question: finding a ϵ for QCS? High-Resolution Bounds!

e.g., $\frac{\epsilon}{\sqrt{M}} = O(\alpha)$ (unif. \mathcal{Q} bin width α) or $O(2^{-B})$ for B -bit \mathcal{Q}

2. Second extreme: CS and 1-bit scalar quantization

1-bit Compressed Sensing

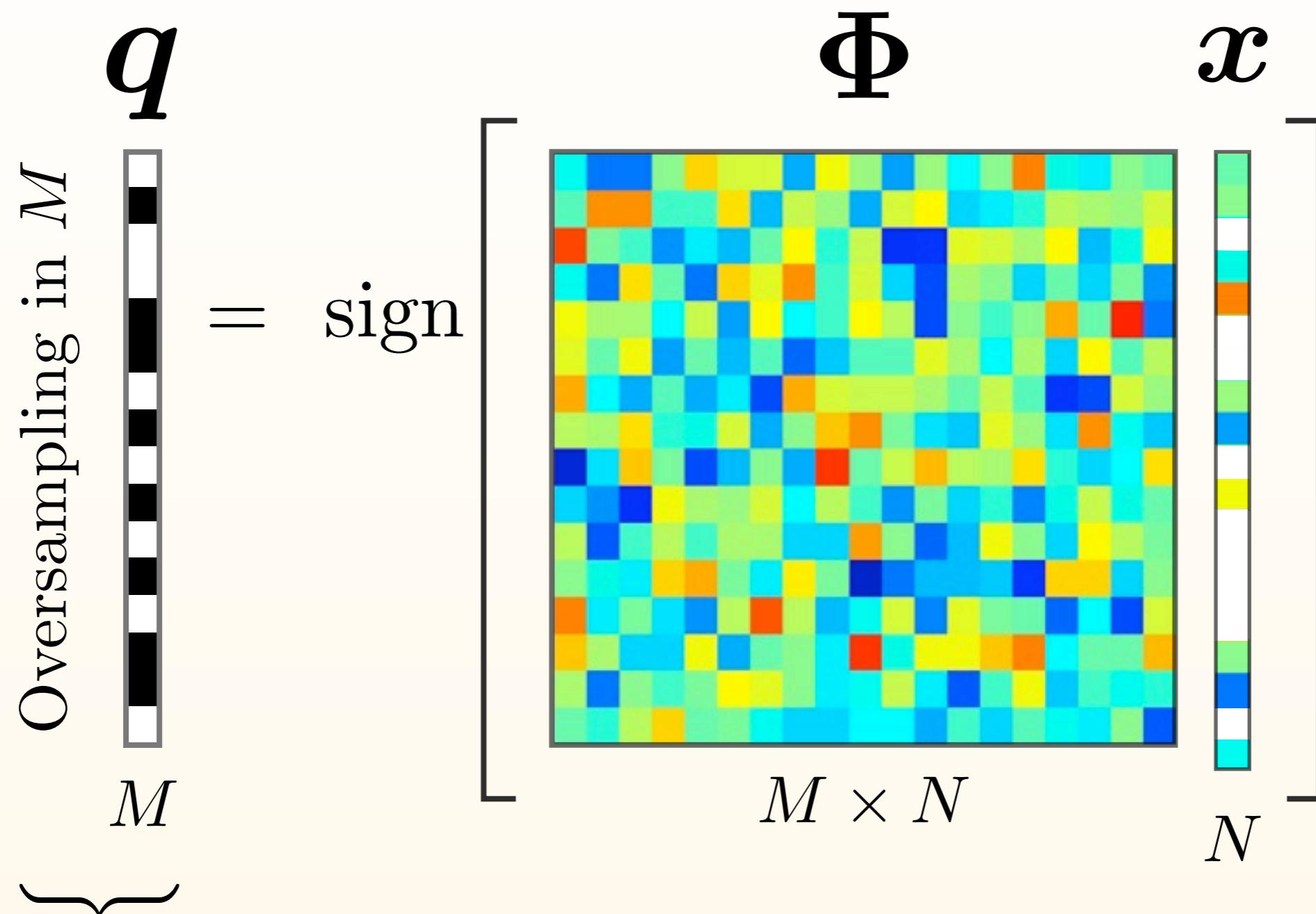


M -bits! But, which information inside q ?

[Boufounos,
Baraniuk, 08]

1-bit Computational Compressed Sensing

bits matter!



M -bits! But, which information inside q ?

1-bit Computational Compressed Sensing

bits matter!

$$q = \text{sign} \begin{bmatrix} \Phi & x \\ M \times N & N \end{bmatrix}$$

Oversampling in M

Warning 1: signal amplitude is lost!

⇒ Amplitude is arbitrarily fixed

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

Warning 2: \exists forbidden sensing!

(e.g. Bernoulli + canonical sparsity)

[Plan, Vershynin, 11]

Why does it work?

\mathbf{x} on S^2

M vectors:

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

iid Gaussian

1-bit Measurements

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

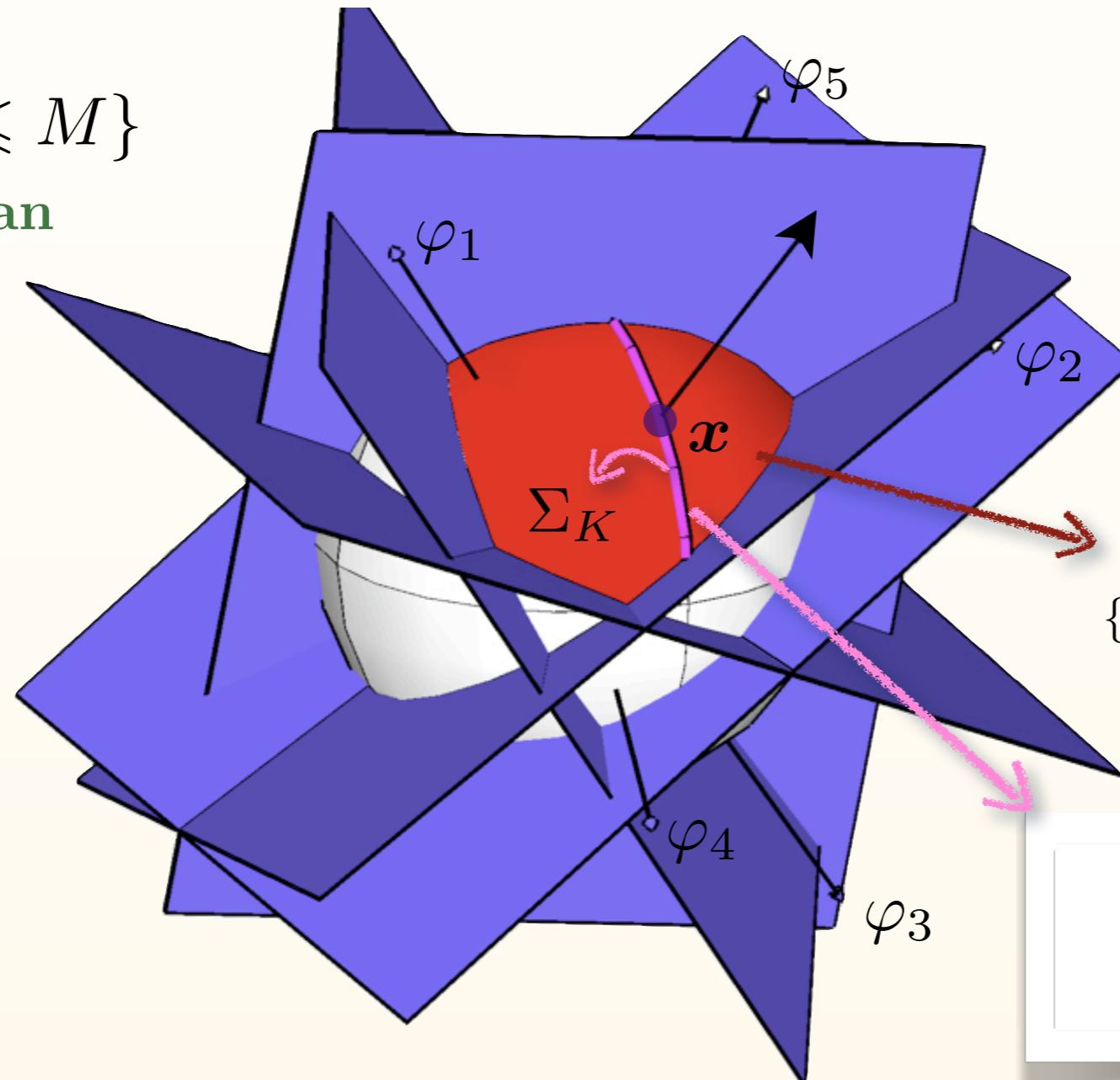
$$\langle \varphi_2, \mathbf{x} \rangle > 0$$

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

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

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

⋮



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

Lower bound on
this width?

1-bit CS bounds

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} ,

$$\begin{array}{l} \text{if only } b \text{ bits are different} \\ \text{between } A(\mathbf{x}) \text{ and } A(\mathbf{s}) \end{array} \Rightarrow \|\mathbf{x} - \mathbf{s}\| \leq \frac{K+b}{K} \epsilon$$

[LJ, Laska, Boufounos, Baraniuk, 13] ($b = 0$)
[LJ, Degraux, De Vleeschouwer, 13] ($b \neq 0$)

If $M = O(\epsilon^{-2}K \log N) \rightarrow$ Binary ϵ Stable Embedding (B ϵ SE):

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

\neq RIP like

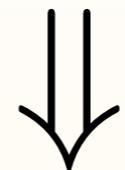
[LJ, Laska, Boufounos, Baraniuk, 13] [Plan, Vershynin, 11, 12]

Lesson: possible recovery should be as consistent as possible.

Easiest 1-bit reconstruction

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

$$\left| \frac{\sqrt{\pi}/2}{M} \langle \text{sign}(\Phi\mathbf{x}), \Phi\mathbf{s} \rangle - \langle \mathbf{x}, \mathbf{s} \rangle \right| \leq \epsilon \quad (\text{w.h.p})$$



[Plan, Vershynin, 12]

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

Compute

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

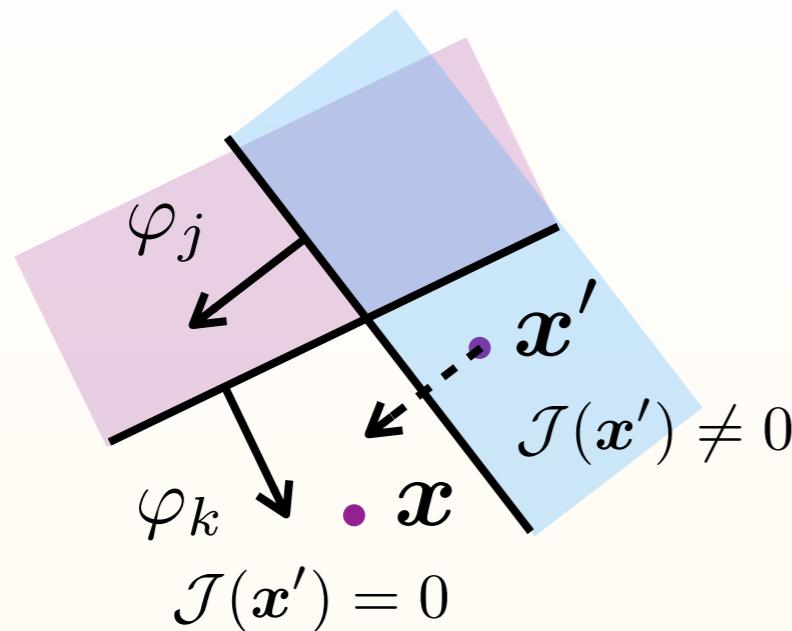
Then, if previous property holds, $\|\mathbf{x} - \hat{\mathbf{x}}\| \leq 2\epsilon$.

[LJ, Degraux, De Vleeschouwer, 13]

$$\propto \arg \max_{\mathbf{u} \in \mathbb{R}^N} \mathbf{q}^T \Phi \mathbf{u} \quad \text{s.t.} \quad \|\mathbf{u}\|_0 \leq K, \quad \|\mathbf{u}\| = 1$$

\equiv “PV-L0 problem”
[Bahmani, Boufounos, Raj, 13]

Binary Iterative Hard Thresholding (BIHT)



Idea: “greedily”

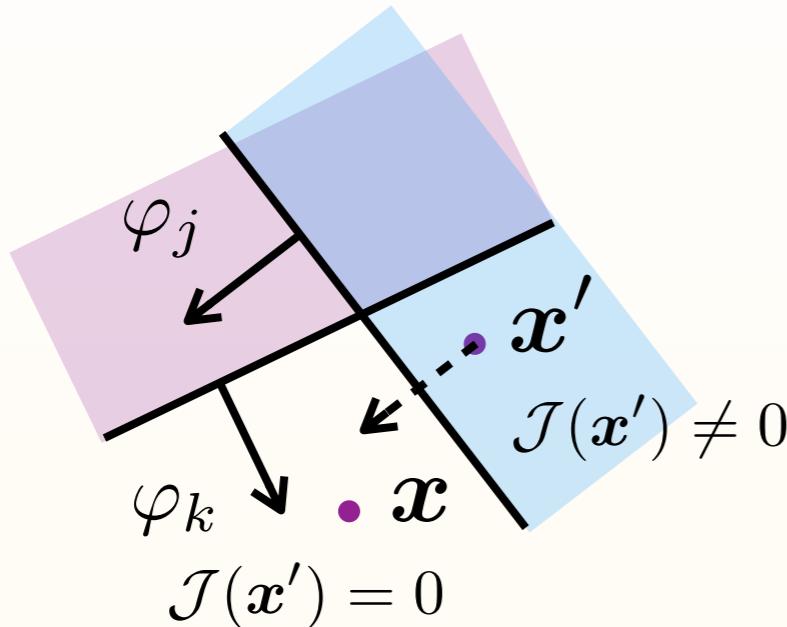
- ▶ forcing consistency
- ▶ forcing sparsity

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

with $(\lambda)_- = (\lambda - |\lambda|)/2$

or
$$\boxed{\mathcal{J}(\mathbf{x}') = \|[\text{diag}(\mathbf{q})(\Phi \mathbf{x}')]_- \|_1}$$

Binary Iterative Hard Thresholding (BIHT)



Idea: “greedily”

- ▶ forcing consistency
- ▶ forcing sparsity

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

with $(\lambda)_- = (\lambda - |\lambda|)/2$

or $\mathcal{J}(\mathbf{x}') = \|[\text{diag}(\mathbf{q})(\Phi \mathbf{x}')]_- \|_1$

Objective: $\text{argmin}_{\mathbf{u}} \mathcal{J}(\mathbf{u})$ s.t. $\|\mathbf{u}\|_0 \leq K$

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



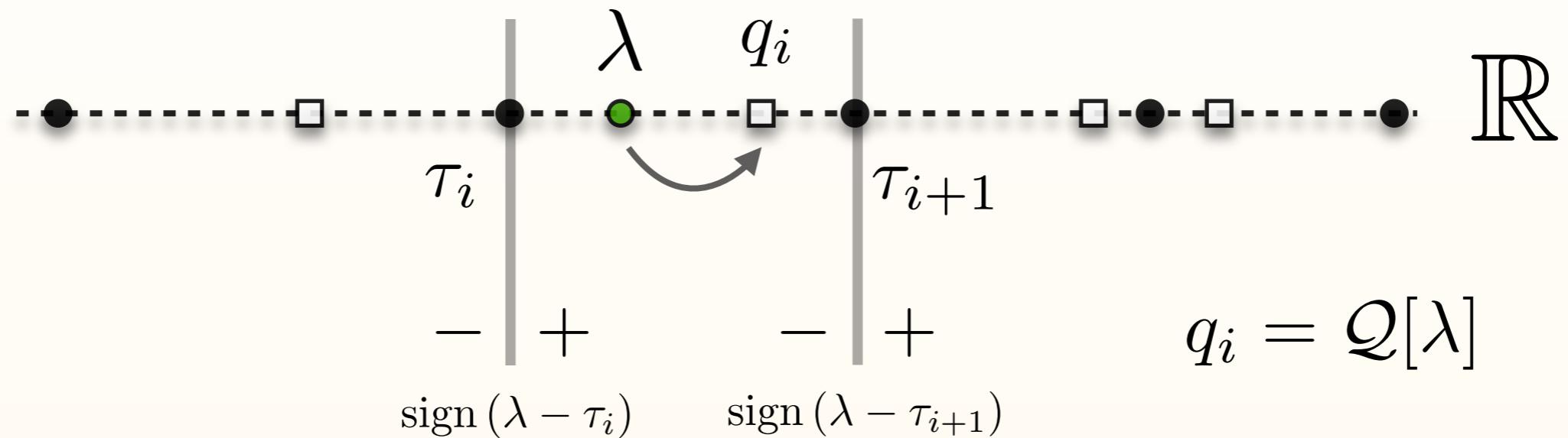
$$\begin{aligned} \mathbf{a}^{l+1} &= \mathbf{x}^l + \frac{\tau}{2} \Phi^T (\mathbf{q} - A(\mathbf{x}^l)), \\ \mathbf{x}^{l+1} &= \mathcal{H}_K(\mathbf{a}^{l+1}), \quad l \leftarrow l + 1 \end{aligned}$$

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

3. Bridging 1-bit & B -bit CS?

Bridging 1-bit & B -bit CS?

- B -bit quantizer defined with thresholds:



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

- 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

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

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

Remark: for symmetric \mathcal{Q} , $\mathcal{Q}(\lambda) = \sum_{j=2}^{2^B} w_j \text{sign}(\lambda - \tau_j)$

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

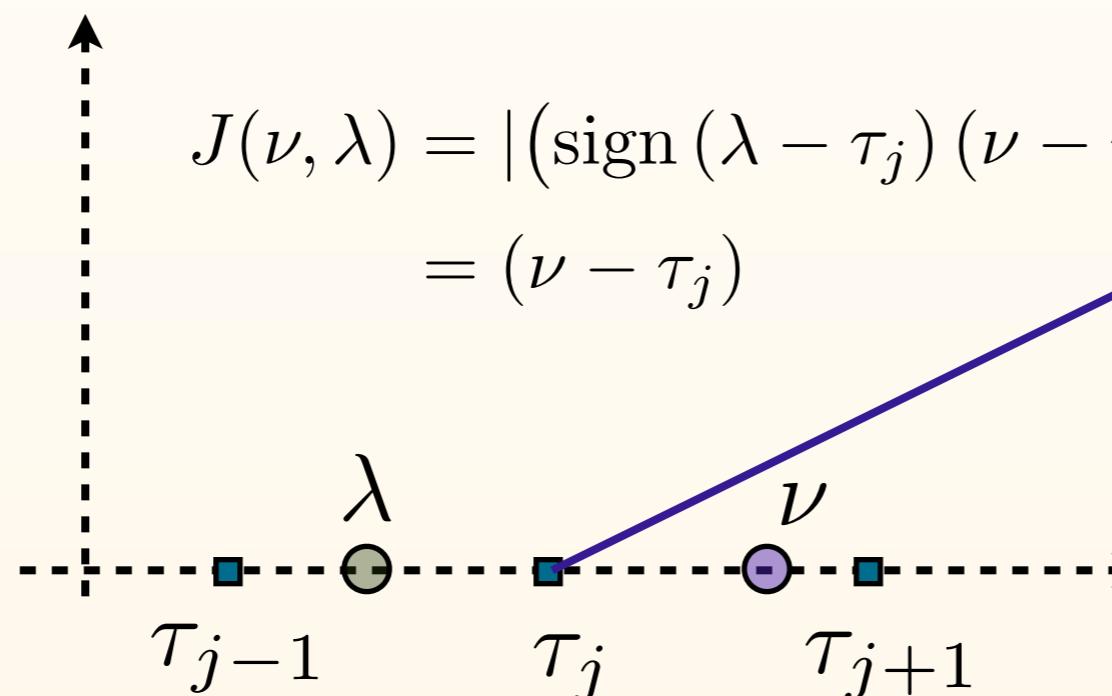
$$J(\nu, \lambda) = \sum_{j=2}^{2^B} w_j \left| (\text{sign}(\lambda - \tau_j)(\nu - \tau_j))_- \right|,$$

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 ℓ_1 -sided norm



(for $w_j = 1$)

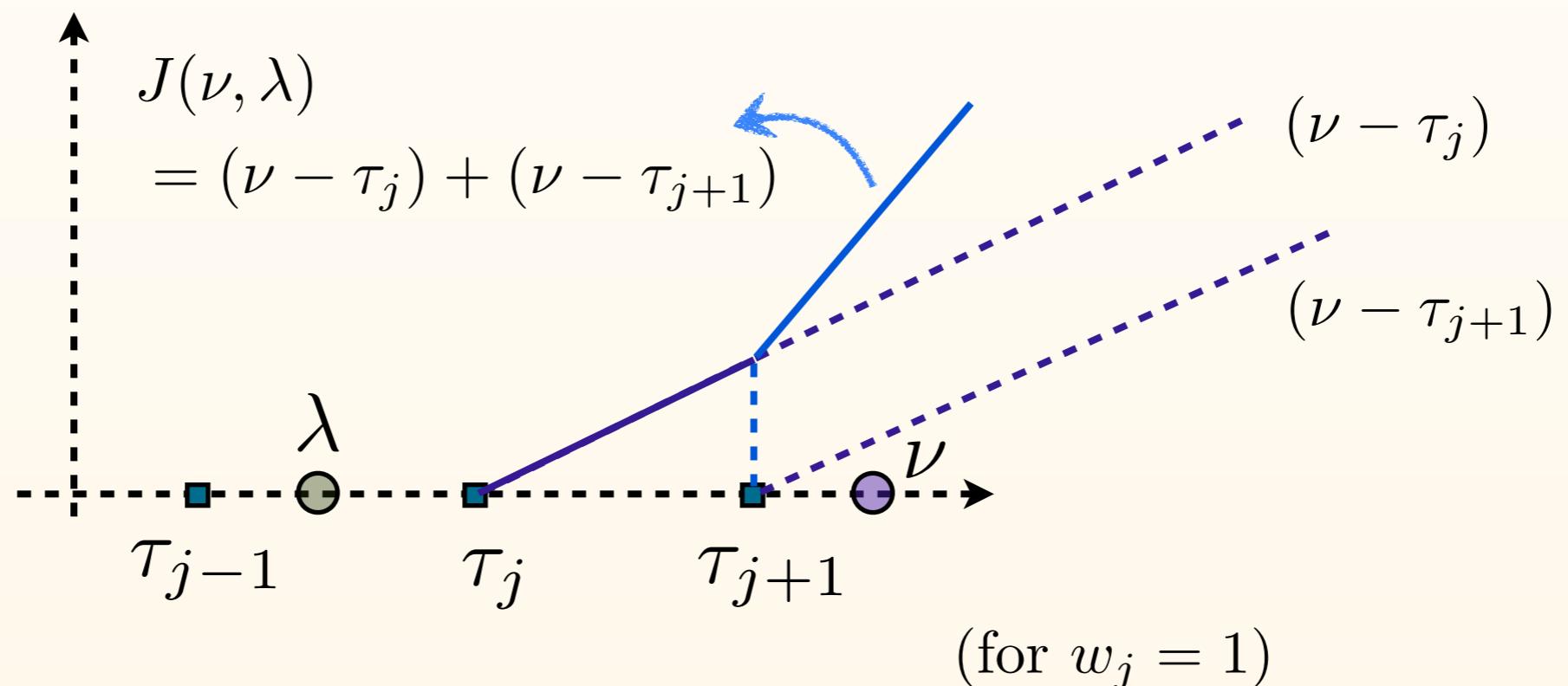
Bridging 1-bit & B -bit CS?

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

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

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

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



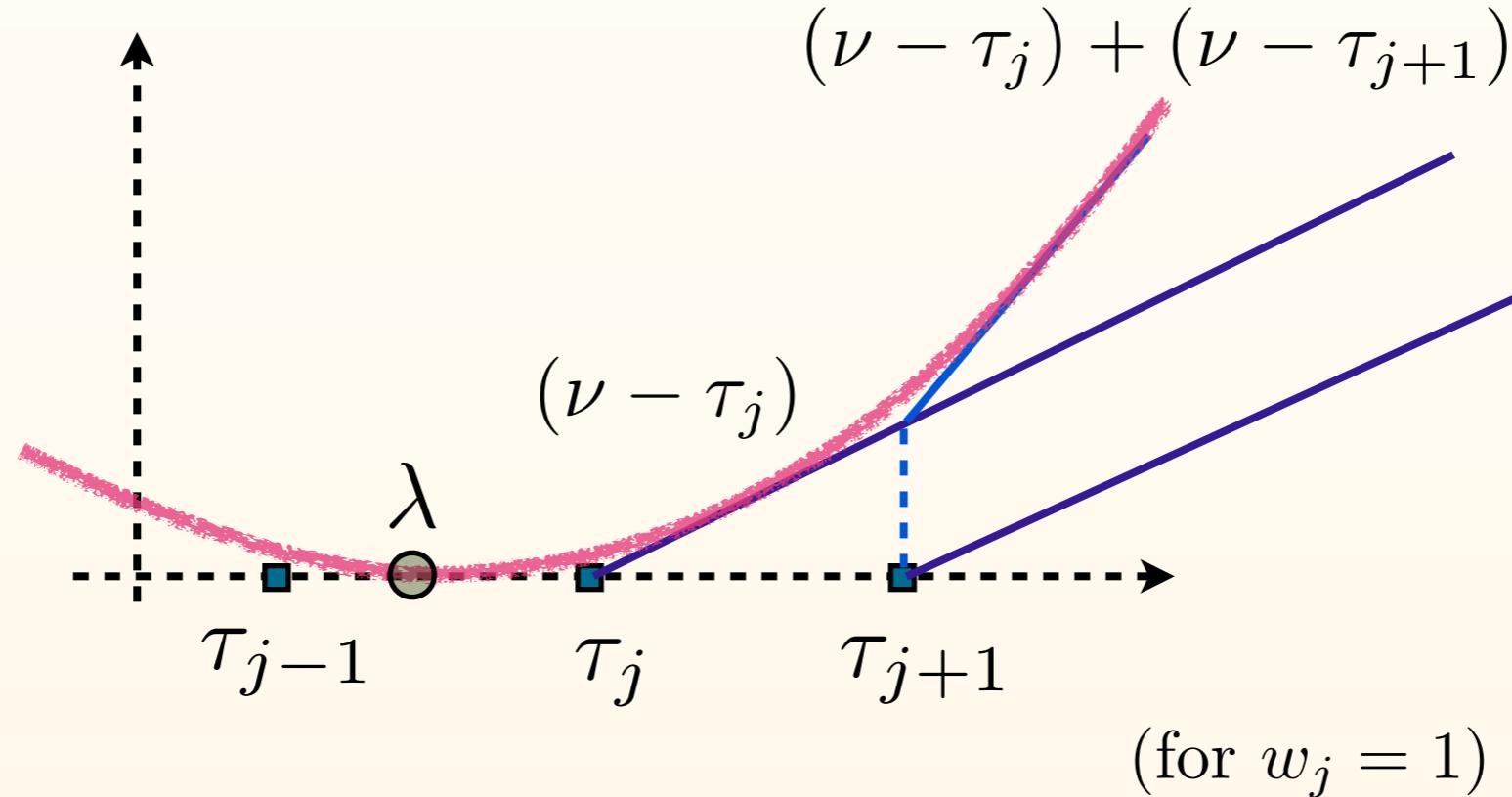
Bridging 1-bit & B -bit CS?

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

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

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

Illustration:



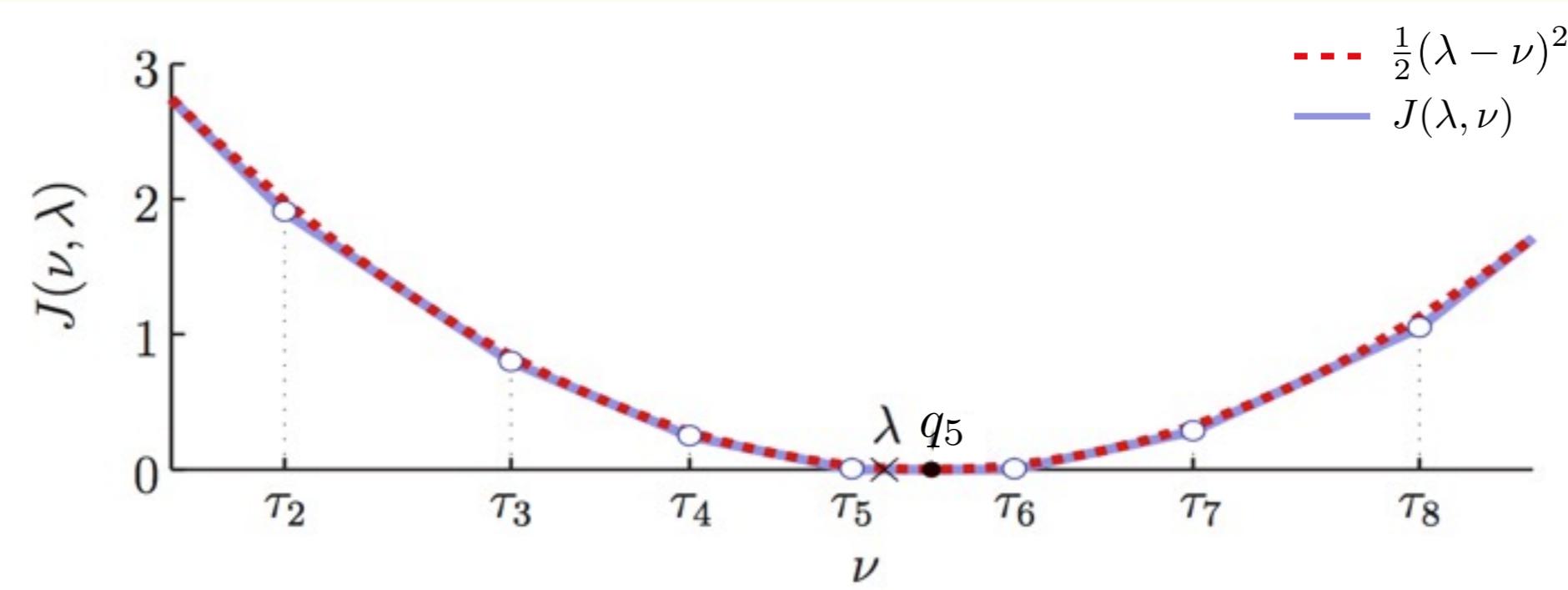
Bridging 1-bit & B -bit CS?

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

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

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

Illustration: more bins, $\lambda \in \mathcal{R}_5$



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

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

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

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

Remarks:

- J is convex in ν
- 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 $B \gg 1$:
 $J(\nu, \lambda) \rightarrow \frac{1}{2}(\nu - \lambda)^2$ and $\mathcal{J}(\mathbf{u}, \mathbf{v}) \rightarrow \frac{1}{2}\|\mathbf{u} - \mathbf{v}\|^2$ (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 Σ_K (as for Iterative Hard Thresholding)

[Blumensath, Davies, 08]

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

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 Σ_K (as for Iterative Hard Thresholding)
[Blumensath, Davies, 08]

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

- combinatorial but greedy solution (as for IHT):

$$\mathbf{x}^{(n+1)} = \mathcal{H}_K[\mathbf{x}^{(n)} - \mu \partial \mathcal{E}_B(\mathbf{x}^{(n)})] \text{ and } \mathbf{x}^{(0)} = 0. \quad (\mu > 0, \text{ see after})$$

$$\Phi^*(\text{sign } (\Phi\mathbf{u}) - \text{sign } (\Phi\mathbf{x})) \xleftarrow[B=1]{} \mathbf{x}^{(n+1)} = \mathcal{H}_K[\mathbf{x}^{(n)} - \mu \partial \mathcal{E}_B(\mathbf{x}^{(n)})] \xrightarrow[B \gg 1]{} \Phi^*(\Phi\mathbf{u} - \mathbf{q})$$

↓

BIHT! Quantized IHT (QIHT) IHT!

“all that.. for this!?” ;-)

Bridging 1-bit & B -bit CS?

- ▶ So, QIHT reads

$$\boldsymbol{x}^{(n+1)} = \mathcal{H}_K \left[\boldsymbol{x}^{(n)} + \mu \boldsymbol{\Phi}^* (\boldsymbol{q} - \mathcal{Q}_B(\boldsymbol{\Phi} \boldsymbol{x}^{(n)})) \right] \text{ and } \boldsymbol{x}^{(0)} = 0.$$

- ▶ Setting μ ?

- ▶ for $B = 1$, μ has no impact

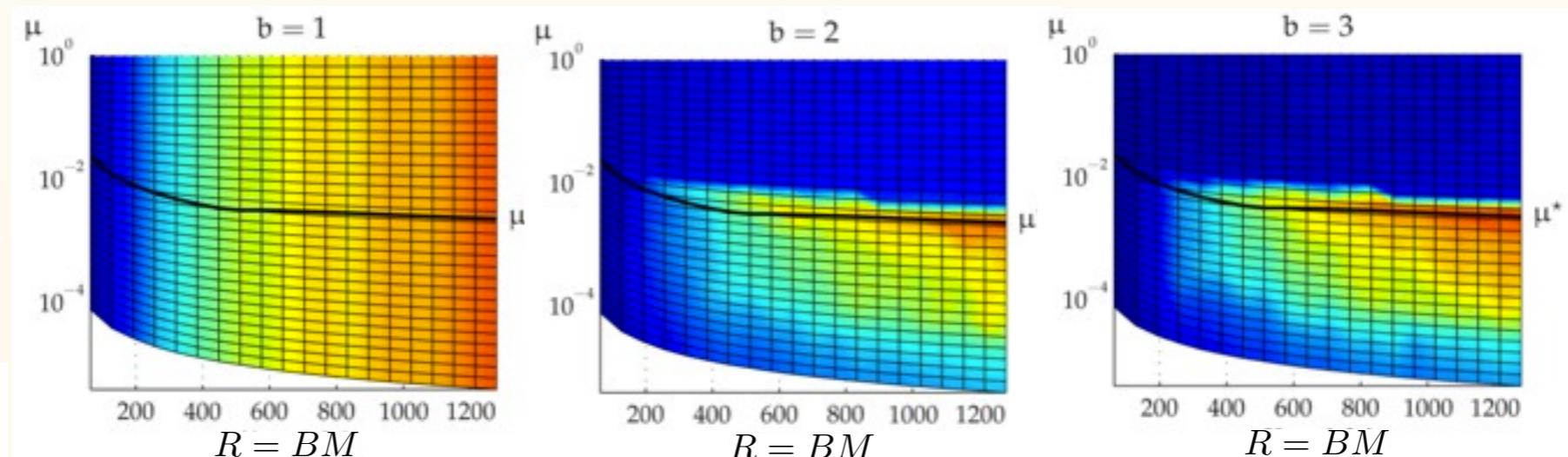
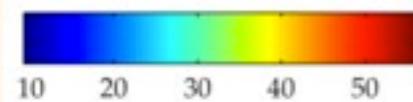
- ▶ for high B , IHT solution: [Blumensath, Davies, 08]

If $\boldsymbol{\Phi}/\sqrt{M}$ is RIP, $\mu < \frac{1}{(1+\delta_{2K})M}$

- ▶ for any B , heuristic: $\mu = \frac{1}{M}(1 - \sqrt{cK/M})$ for some $c > 0$.

Validated by
extensive
simulations
($c = 3$)

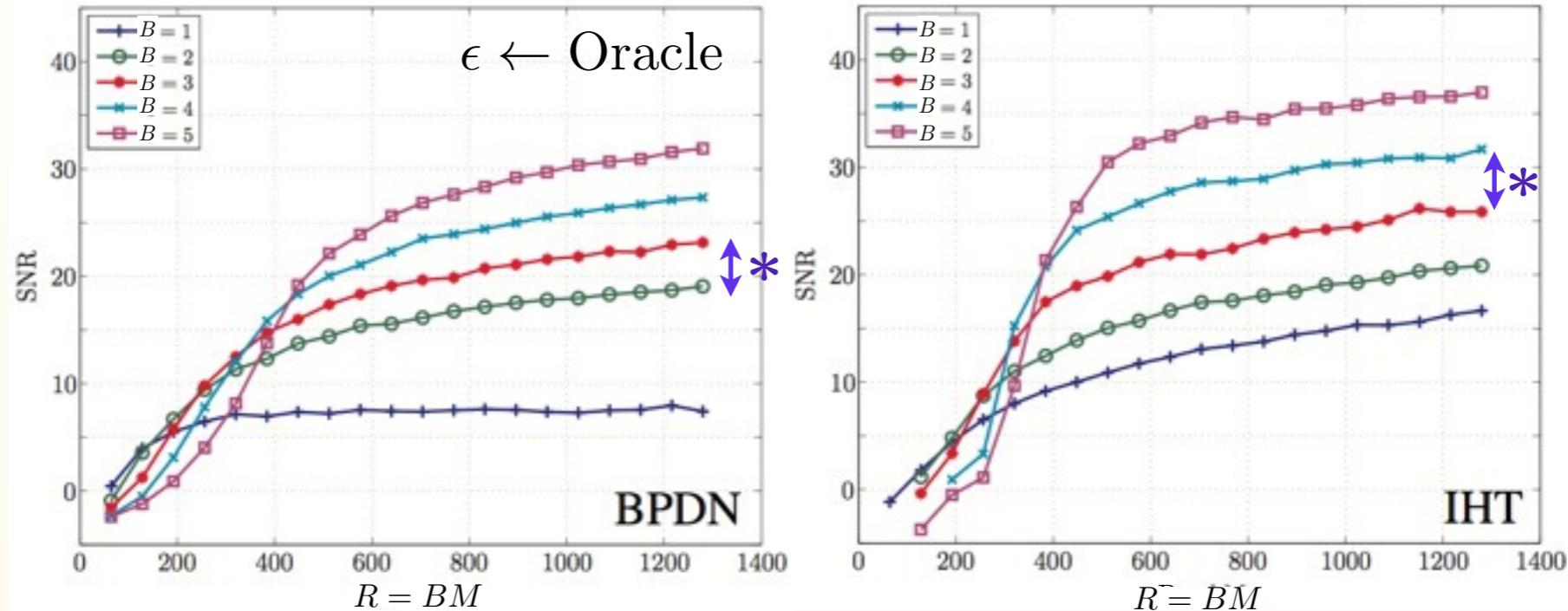
SNR:



QIHT simulations

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

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

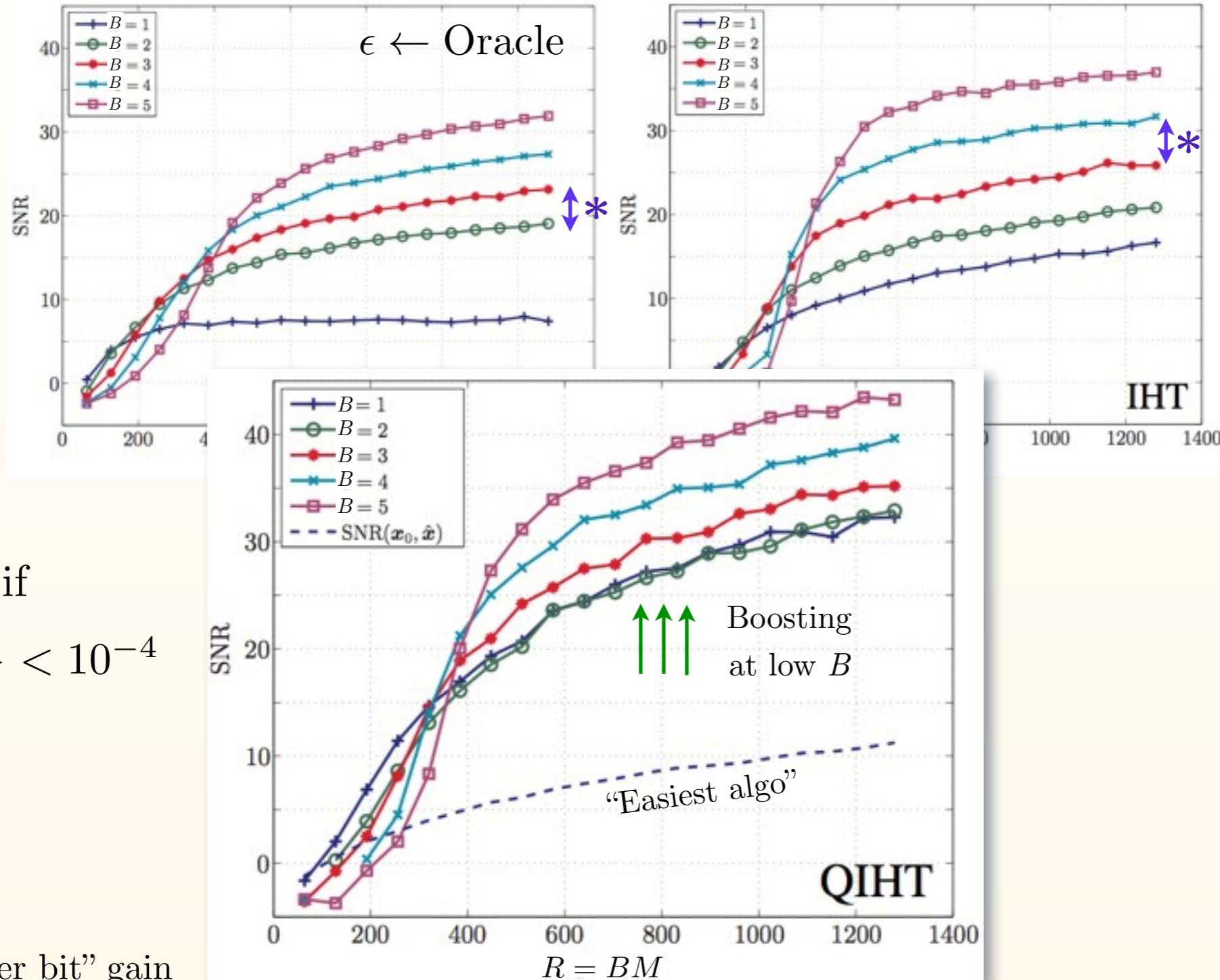


*: almost “6dB per bit” gain

QIHT simulations

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

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

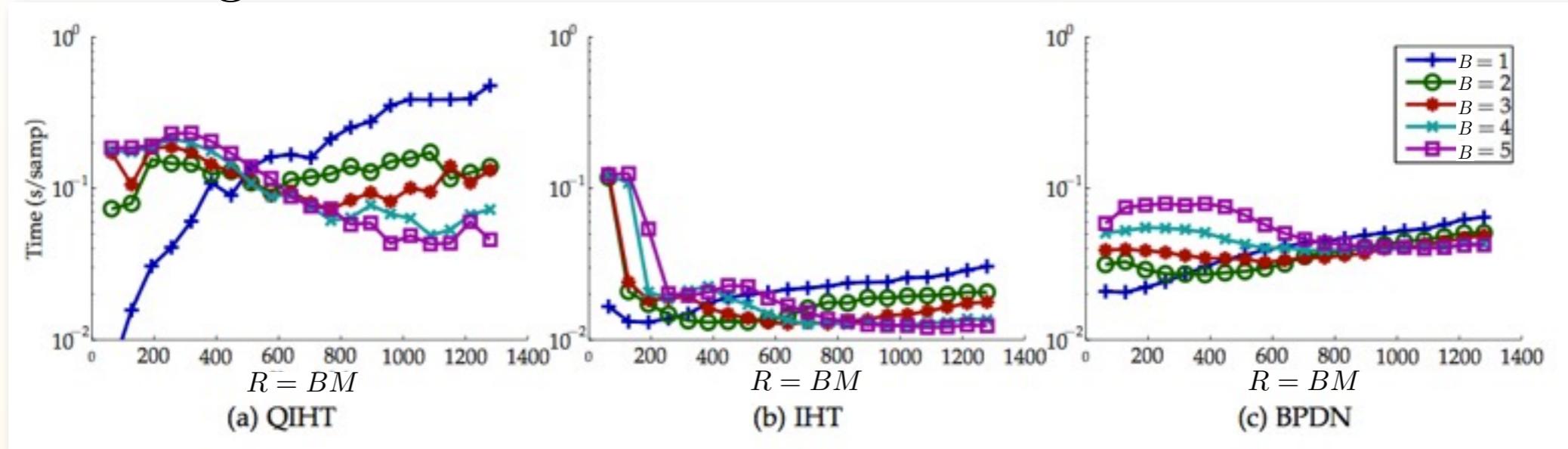


QIHT simulations

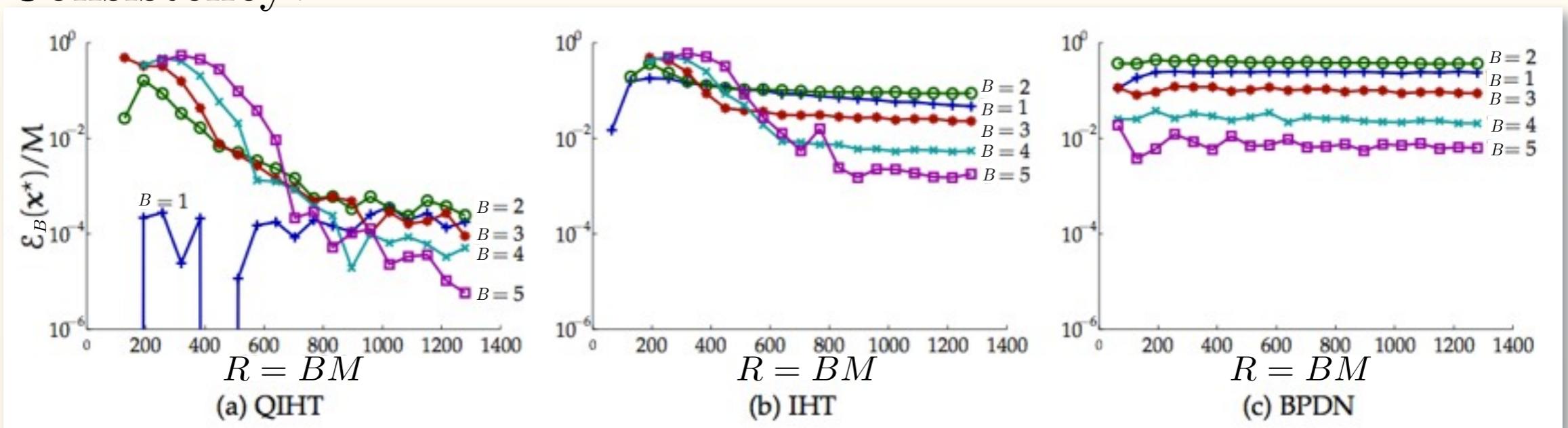
$N = 1024$, $K = 16$, $R = BM \in \{64, 128, \dots, 1280\}$, 100 trials (+ Lloyd-Max Gauss. Q.)

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

Running time:



Consistency?



Conclusions

- IHT framework can integrate quantization
- Lot of theoretical works to do/come...
 - BIHT/QIHT convergence/optimality guarantees?
 - New embedding results?
considering $\mathcal{Q}(\lambda) = \sum_{j=2}^{2^B} w_j \text{sign}(\lambda - \tau_j)$ = sum of 1-bit costs + thresholds?
 - Blend of RIP/BεSE? as in
$$|(1 - \delta)\|\mathbf{x} - \mathbf{x}'\| - \epsilon| \leq \text{dist}(\mathcal{Q}(\Phi\mathbf{x}), \mathcal{Q}(\Phi\mathbf{x}')) \leq (1 + \delta)\|\mathbf{x} - \mathbf{x}'\| + \epsilon ?$$
for sparse \mathbf{x}, \mathbf{x}' and with $\epsilon \rightarrow_{B,M} 0$ and $\delta \rightarrow_M 0$?
- Coming:
Application to B -bit compressive sensors (RMPI)

Further Reading

- (this work) L. Jacques, K. Degraux, C. De Vleeschouwer, “Quantized Iterative Hard Thresholding: Bridging 1-bit and High-Resolution Quantized Compressed Sensing”, SAMPTA 2013, <http://arxiv.org/abs/1305.1786>.
- 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.
- 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.
- S. Bahmani, P.T. Boufounos, B. Raj, “Robust 1-bit Compressive Sensing via Gradient Support Pursuit”, <http://arxiv.org/abs/1304.6626>