

# Sequential Projection Learning for Hashing with Compact Codes

Jun Wang<sup>†</sup>, Sanjiv Kumar<sup>‡</sup>, Shih-Fu Chang<sup>†</sup>

<sup>†</sup> Columbia University, New York, USA

<sup>‡</sup> Google Research, New York, USA

## Background and Related Work

### Nearest neighbor search for large databases

- Exhaustive nearest neighbor search not practical for large scale databases with high dimensional points
- Approximate nearest neighbor (ANN) search achieves sub-linear query time
- Popular ANN search techniques include tree-based approaches and hashing approaches

### Related work

- Different choices of projections:
  - random projections
  - locality sensitive hashing (**LSH**, Indyk et al. 98)
  - Shift Invariant Kernel based Hashing (**SIKH**, Raginsky et al. 09)
  - principal projections
  - spectral hashing (**SH**, Weiss, et al. 08)
- Different forms of hash functions
  - identity function: LSH & Boosted SSC (**BSSC**, Shakhnarovich, 05)
  - sinusoidal function: SH & SIKH
- Other related work
  - Restricted boltzmann machines (**RBM**s, Hinton et al. 06)
  - Jian et al. 08 (metric learning)
  - Kulis et al. 09 & Mu et al. 10 (kernelized)

### Main issues

- Existing hashing methods mostly rely on random or principal projections, which are not very compact or accurate
- Simple metrics are usually not enough to express semantic similarity

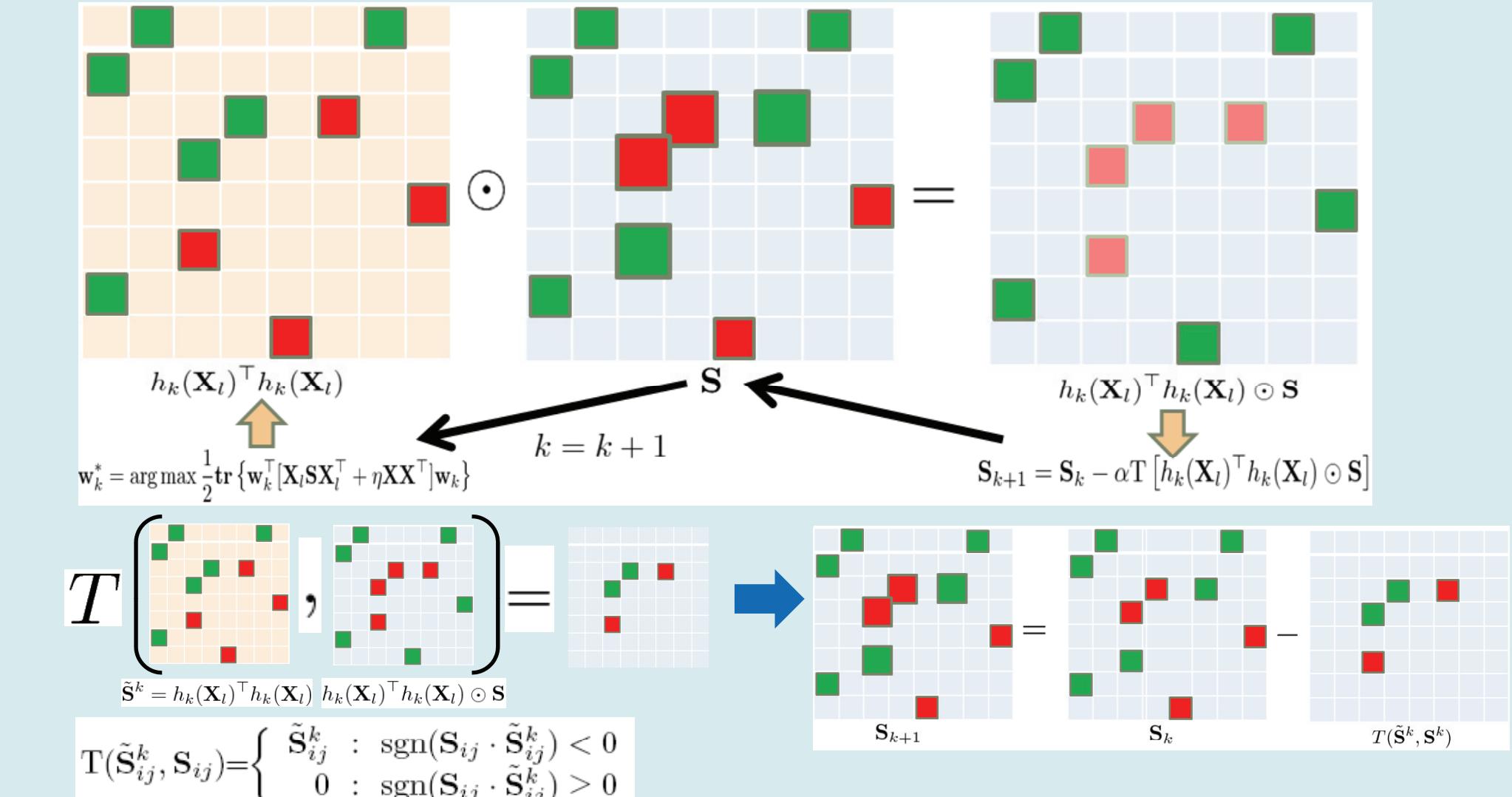
## Semi-Supervised Sequential Projection Learning (S3PLH) - Solution

### Final objective function after relaxation

$$J(\mathbf{W}) = \frac{1}{2} \text{tr} \{ \mathbf{W}^\top \mathbf{X}_l \mathbf{S} \mathbf{X}_l^\top \mathbf{W} \} + \frac{\eta}{2} \text{tr} \{ \mathbf{W}^\top \mathbf{X} \mathbf{X}^\top \mathbf{W} \}$$

$\mathbf{M} = \mathbf{X}_l \mathbf{S} \mathbf{X}_l^\top + \eta \mathbf{X} \mathbf{X}^\top$  “adjusted” covariance matrix  
The orthogonal hash codes can be obtained by doing eigen-decomposition on  $\mathbf{M}$  - PCAH

### Sequential solution via updating label matrix



## Unsupervised Sequential Projection Learning (USPLH) - Extension

- Generate pseudo pair-wise label data from existing hash bits

$$h_k(\mathbf{x}) = -1 \quad h_k(\mathbf{x}) = 1$$

$\mathcal{M} = \{(x_i, x_j)\} : h(x_i) \cdot h(x_j) = -1, |\mathbf{w}^\top (\mathbf{x}_i - \mathbf{x}_j)| \leq \epsilon$   
 $\mathcal{C} = \{(x_i, x_j)\} : h(x_i) \cdot h(x_j) = 1, |\mathbf{w}^\top (\mathbf{x}_i - \mathbf{x}_j)| \geq \zeta$

- Use pseudo label matrix to learn projections

$$\mathbf{S}_{MC}^k : \begin{cases} (x_i, x_i) \in \mathcal{M} \Rightarrow \mathbf{S}_{MC}^k(i, i) = 1 \\ (x_i, x_i) \in \mathcal{C} \Rightarrow \mathbf{S}_{MC}^k(i, i) = -1 \end{cases}$$

Pseudo label matrix from kth hash bit

$$\mathbf{M}_k = \sum_{i=0}^{k-1} \lambda^{k-i} \mathbf{X}_{MC}^i \mathbf{S}_{MC}^i \mathbf{X}_{MC}^i {}^\top + \eta \mathbf{X} \mathbf{X}^\top$$

“adjusted” covariance matrix with pseudo labels

## Semi-Supervised Hashing - Formulation

### Setting and notations

$$h_k(\mathbf{x}) = \text{sgn}(\mathbf{w}_k^\top \mathbf{x} + b_k), b_k = \text{mean}(\mathbf{w}_k^\top \mathbf{x})$$

In the setting of SSH, one is given the data  $\mathbf{X}$ , and a subset  $\mathbf{X}_l \in \mathbb{R}^{D \times l}$  with pair-wise labels:

$$\begin{aligned} (\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{M} &\text{ neighbor pair} \\ (\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{C} &\text{ nonneighbor-pair} \end{aligned}$$

Define the pairwise label matrix  $\mathbf{S} \in \mathbb{B}^{l \times l}$  as:

$$S_{ij} = \begin{cases} 1 & : (\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{M} \\ -1 & : (\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{C} \\ 0 & : \text{otherwise.} \end{cases}$$

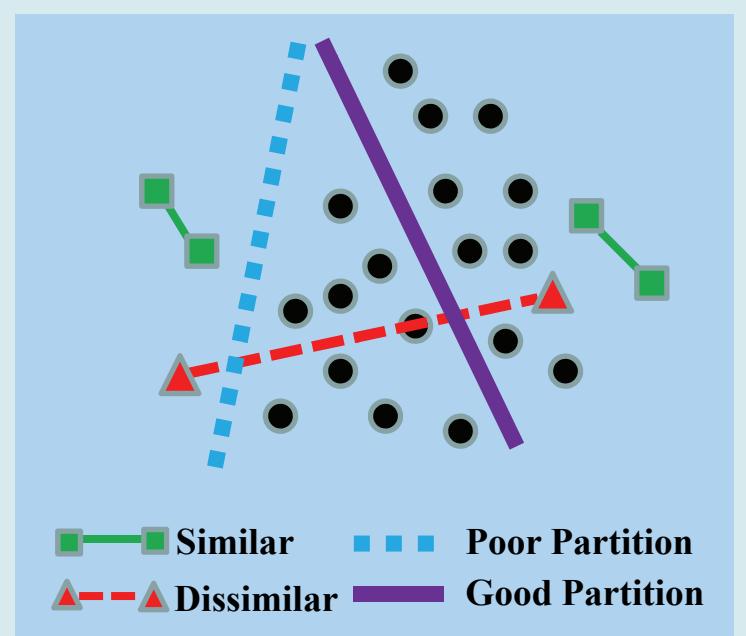


Illustration of semi-supervised hashing to achieve partitioning with maximum empirical fitness and entropy.

### Empirical Fitness

$$\begin{aligned} h_k(\mathbf{X}_l)^\top h_k(\mathbf{X}_l) &= \mathbf{S} = h_k(\mathbf{X}_l)^\top h_k(\mathbf{X}_l) \odot \mathbf{S} \\ J(\mathbf{H}) &= \sum_k \left\{ \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{M}} h_k(\mathbf{x}_i) h_k(\mathbf{x}_j) - \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{C}} h_k(\mathbf{x}_i) h_k(\mathbf{x}_j) \right\} \end{aligned}$$

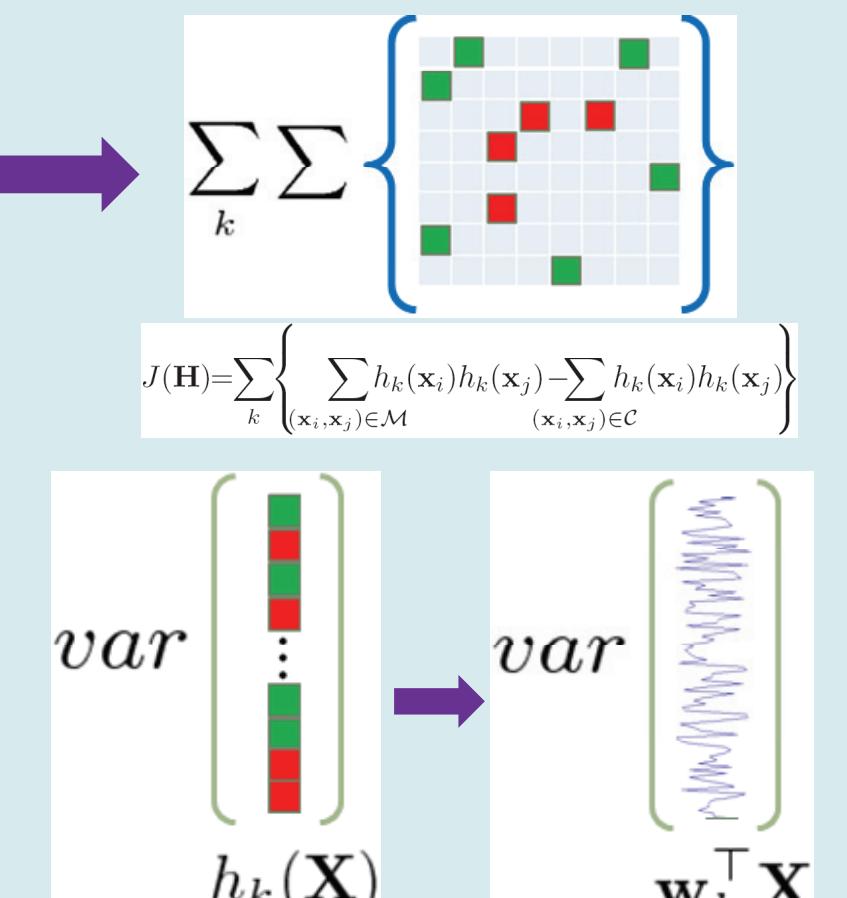
### Variance as Regularizer

$$p[h_k(x) = 1] = p[h_k(x) = -1] \iff \max \text{var}[h_k(x)]$$

**Information theoretic term:** maximum variance partition equals to balanced partition

$$\max \text{var}[h_k(\mathbf{x})] \geq \frac{1}{\beta} \text{var}[\mathbf{w}_k^\top \mathbf{x}]$$

Lower bound of maximum bit variance



## Experiment and Results

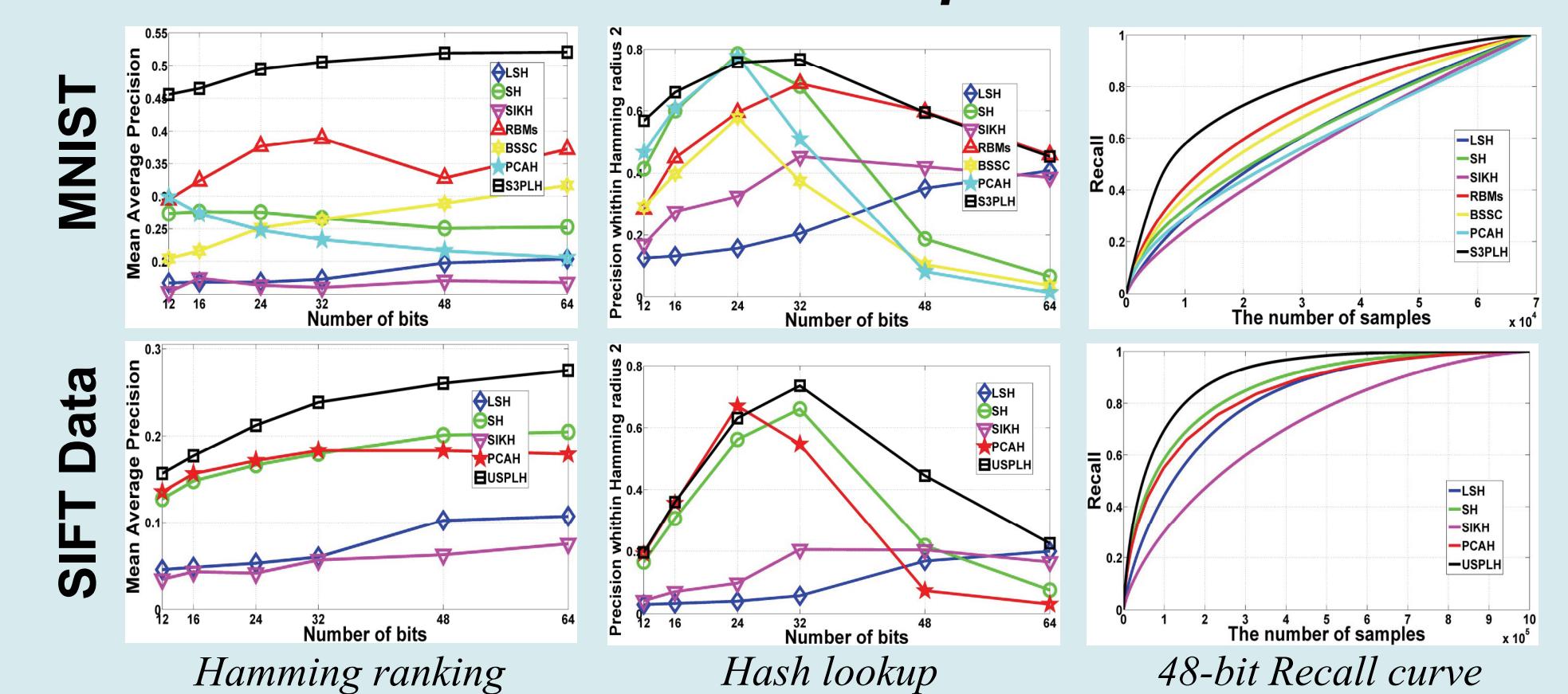
### Datasets

- MNIST Digits: 70K, 1K training labels, 1K query tests (semi-supervised)
- SIFT Feature: 1 Million, 2K pseudo labels, 10K query tests (unsupervised)

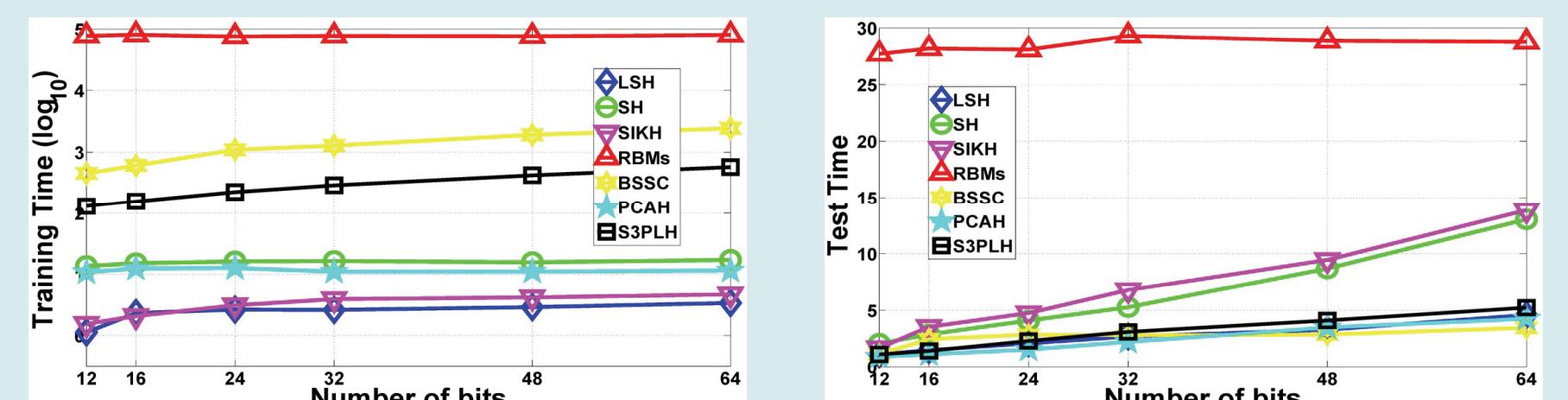
### Experimental protocol

- Hash lookup (precision within hamming radius 2)  
Empty return treated as failed query with zero precision
- Hamming ranking (Mean Average Precision)
- Compared methods: LSH, SH, SIKH, RBMs, BSSC, PCAH

### Performance evaluation and comparison



### Computational evaluation



### Conclusion:

- A semi-supervised paradigm for hashing learning (empirical fitness with information theoretic regularization)
- Sequential learning idea for error correction
- Extension of unsupervised case
- Easy implementation and highly scalable

### Future work:

- Theoretical analysis of performance guarantee
- Weighted hamming embedding