

Approximate Query Processing

1

- Space Partitioning-based
 - ▣ Tree
 - ▣ Encoding
 - ▣ Locality Sensitive Hashing
- Graph-based Methods

Notes:

- Recent works mainly in the Database area
- Prefer ease of exposition over rigor
- Categorization is not fixed/unique

Locality Sensitive Hashing (LSH)

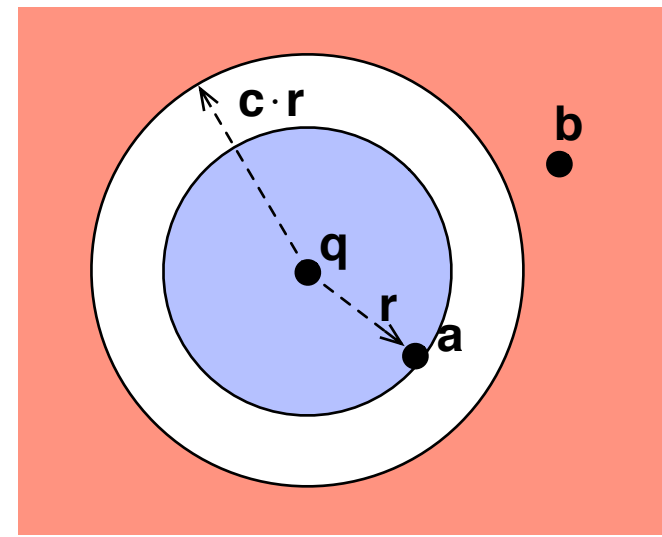
2

□ From the perspective of collision probability

□ (Ordinary) hash function h :

■ $\Pr[\underline{h(x) = h(y)}] = \varepsilon, \text{ if } x \neq y$

c.f., Cryptographic hash functions
 $\Pr [h(x) = h(y)] = 2^{-m}, \text{ if}$
 $\text{Hamming}(x, y) = 1$



Locality Sensitive Hashing (LSH)

3

□ From the perspective of collision probability

▣ (Ordinary) hash function h :

■ $\Pr[\underline{h(x) = h(y)}] = \varepsilon, \text{ if } x \neq y$

▣ LSH

■ $\Pr[\underline{h(x) = h(y)}]$ increases with **locality**

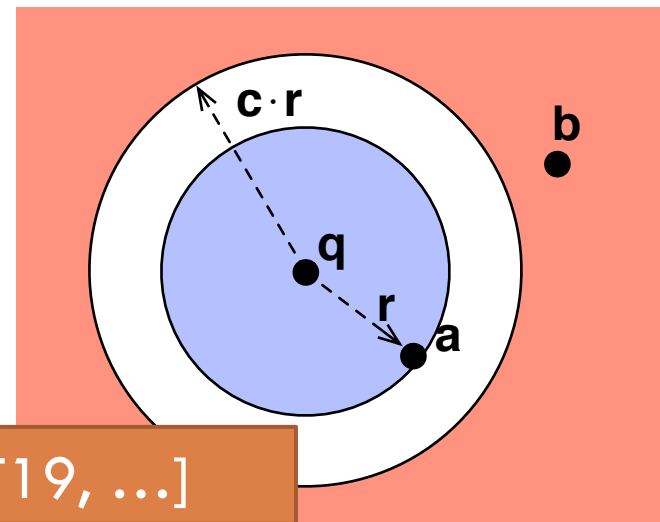
■ Randomness comes from r.v. $h \in H$

c.f., Cryptographic hash functions
 $\Pr [h(x) = h(y)] = 2^{-m}$, if
 $\text{Hamming}(x, y) = 1$

(r_1, r_2, p_1, p_2) -sensitive [IM98]

- $\Pr[h(x) = h(y)] \geq p_1$, if $\text{dist}(x, y) \leq r_1$
- $\Pr[h(x) = h(y)] \leq p_2$, if $\text{dist}(x, y) \geq r_2$

$\Pr[h(x) = h(y)] = \text{sim}(x, y)$ [C02] **too narrow**



Can be generalized [SWQZ+14, ACPS18, CKPT19, ...]

Locality Sensitive Hashing (LSH)

4

□ Equality search

▣ Index: store o into bucket $h(o)$

▣ Query:

■ retrieve every o_i in the bucket $h(q)$

■ verify if $o_i = q$

$$\Pr[h(q) = h(o)] = 1/B$$

Locality Sensitive Hashing (LSH)

5

□ Equality search

▣ Index: store o into bucket $h(o)$

$$\Pr[h(q) = h(o)] = 1/B$$

▣ Query:

▣ **retrieve** every o_i in the bucket $h(q)$

▣ **verify** if $o_i = q$

□ LSH

c.f., [PIM12] for the rigorous QP procedure

▣ $\forall h \in \text{LSH-family}, \Pr[Q(h(q)) = Q(h(o))] = f(\text{Dist}(q, o))$

▣ $Q(\cdot)$: quantization (**not essential**)

▣ “Near-by” points have more chance of colliding with q than “far-away” points

▣ Similar index & query procedures, with a **weak** probabilistic guarantee

➔ Repeat to boost the guarantee

LSH Families

6

- Many are known
 - ▣ L_p ($0 < p \leq 2$): use p -stable distribution to generate the projection vector
 - For L_2 , just use random Gaussian vector
 - Other families exists, e.g., sparse random projection
 - ▣ Angular distance (arccos): SimHash
 - ▣ Jaccard: minhash (based on random permutation)
 - ▣ Hamming:
 - random projection
 - covering LSH

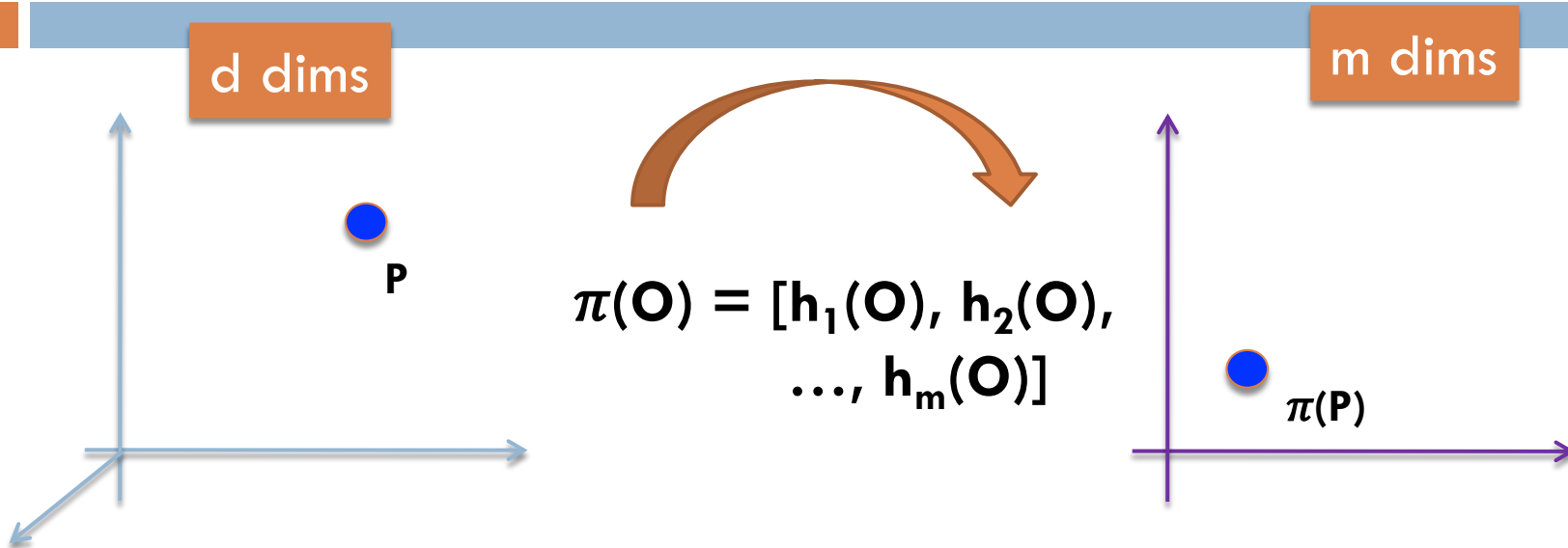
Comments

7

- New queries can be reduced to known LSH cases
 - ▣ Maximum inner product search (MIPS)
 - ▣ Set containment
 - ▣ Group aggregated query
- Related to various distortion-bounded embedding
 - ▣ Edit distance: CGK-embedding to Hamming with $O(K)$ distortion

Probabilistic Mapping

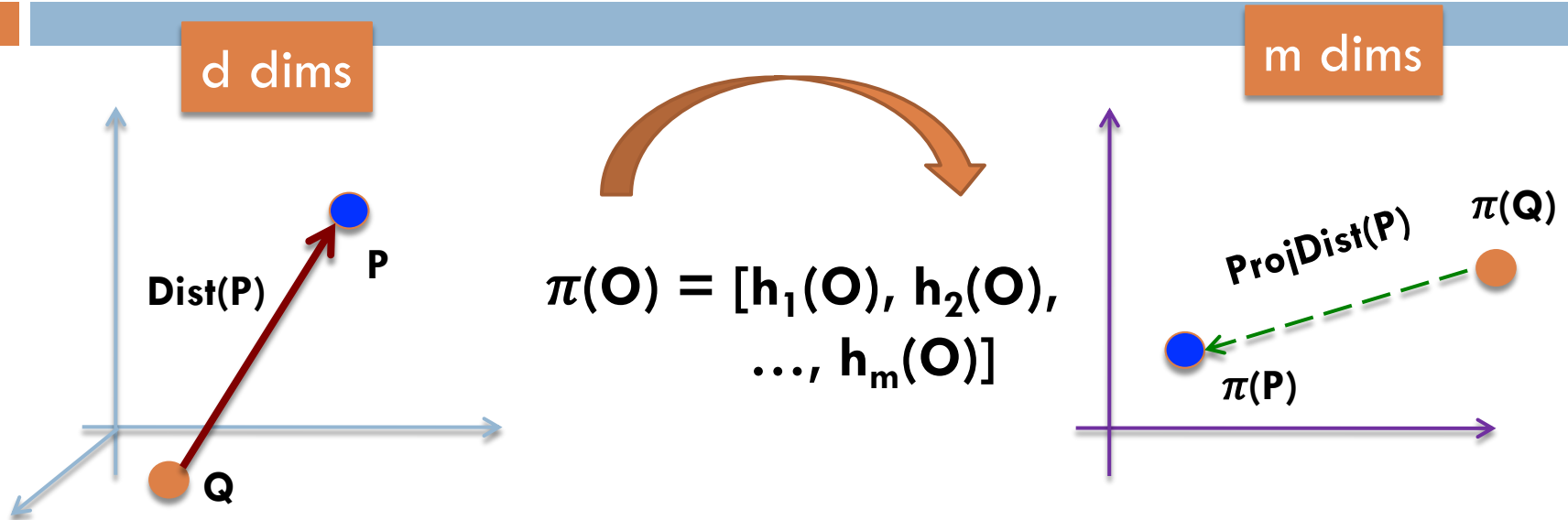
8



- Probabilistic, linear mapping from the **original space** to the **projected space**

Probabilistic Mapping

9

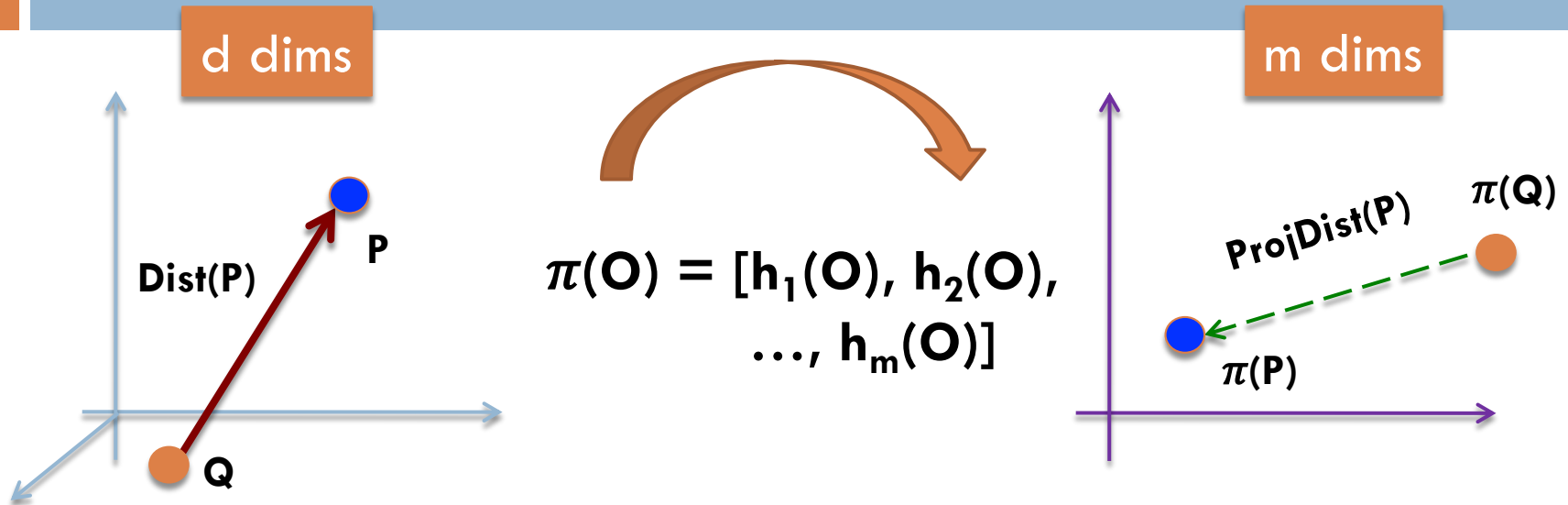


- Probabilistic, linear mapping from the **original space** to the **projected space**
- What about the **distances** (wrt Q or $\pi(Q)$) in these two spaces?

Probabilistic Distance Tracking

Property of the Mapping

10

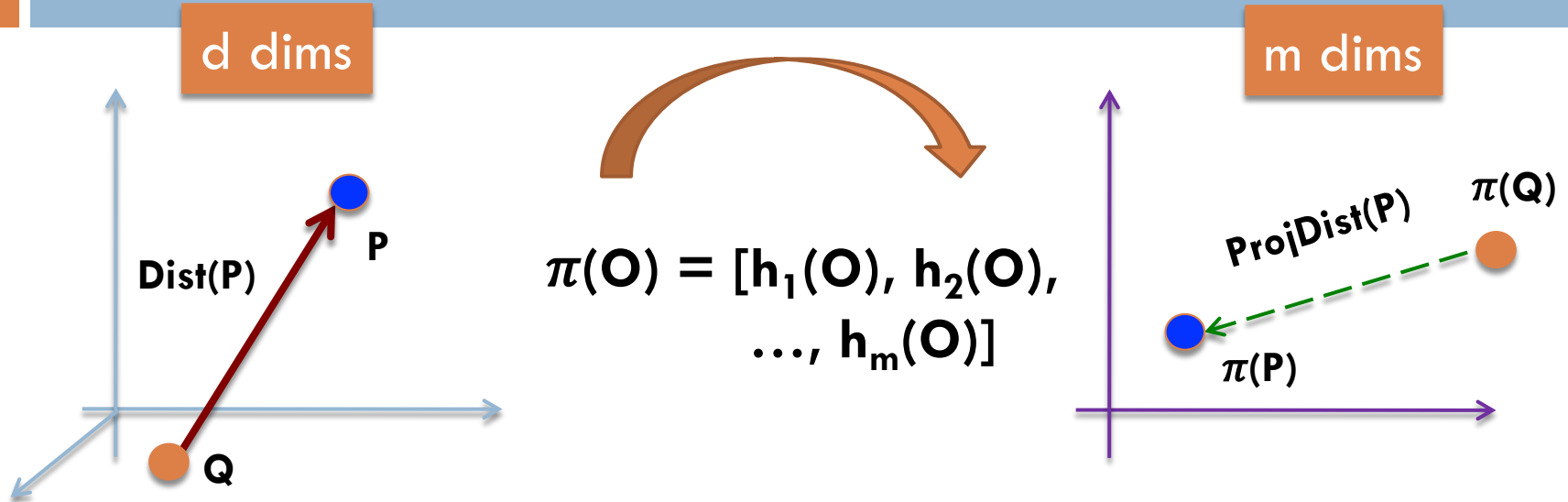


- $\text{ProjDist}(P)^2 \sim \text{Dist}(P)^2 * \chi_m^2$ [SWQZ+14]
 - ▣ $\text{ProjDist}(P)^2$ can be computed (incrementally) from $h_i(P)$ and $h_i(Q)$ due to the linearity of the hash function
 - ▣ Can be generalized to other p-stable LSH functions

Probabilistic Distance Tracking

Property of the Mapping

11



LSH provides a probabilistic distance-preserving mapping between the two spaces

Johnson & Lindenstrauss Lemma [JL84] only works for L2 and induces a method that requires more space than LSH [AIR18]

Roadmap

12

- Roadmap
 - ▣ Practical LSH methods (i.e., linear index complexity)
 - ▣ Data-dependent LSH methods

Practical Variations of LSH

13

- Easy to relax the LSH method in practice at the cost of no worst-case guarantees
 - **E2LSH**: use fewer number of random projections
 - **Multiprobe LSH** (entropyLSH and other variants): space-time tradeoff
 - **LSH in practice**: use empirically tuned parameters (k, l)
 - **SRS**: a space-saving LSH index with in-memory or disk implementations
 - **QALSH**: use multiple B-trees for the index and use collision counting strategy based on LSH
 - **HD-index**: space filling curves as pseudo-LSH functions
 - **SK-LSH**: Replace LSB-tree/forest by a dimension-wise linear mapping
 - ...

Data-sensitive Hashing

14

- LSH is data-insensitive
 - ▣ Indexing hyper-parameters determined by the shape of the data only
 - ▣ Indexing parameters are randomly generated
- Efforts to make data-sensitive, LSH-like methods
 - ▣ [AR15]
 - Aim: break the lower bounds of ρ
 - ▣ DSH
 - ▣ OPFA / NeOPFA
- Learning-to-hash methods
 - ▣ NSH [PCM15]
 - ▣ [LYZX+18] and many in the ML/CV communities


c.f., [AIR18]

c.f., <https://learning2hash.github.io> and <https://cs.nju.edu.cn/lwj/L2H.html>

DSH [GJLO14]

15

- Learn a family of (hash) functions, H , that preserves k NN of queries

1. Training data:  $\mathbf{W}_{ij} = \begin{cases} 1 & , \text{ if } o_j \in kNN(q_i) \\ -1 & , \text{ if } o_j \notin kNN(q_i) \wedge o_j \text{ is sampled} \\ 0 & \text{ otherwise.} \end{cases}$

- sampled queries

- their k -NN objects (+ve)

- samples non- $c*k$ -NN objects (-ve)

2. Function family:

- Thresholded linear functions $h(\mathbf{x}; \mathbf{a}) = \text{sgn}(\mathbf{a}^\top \mathbf{x})$

3. Learn one hash function

$$\arg \min_h \sum_i \sum_j \ell(q_i, o_j) \mathbf{W}_{ij} \quad , \text{ where } \ell(q, o) = (h(q) - h(o))^2$$

DSH [GJLO14]

16

- Learn a family of (hash) functions, H , that preserves k NN of queries
- 4. Learn multiple hash functions
 - Multiplicative updates on W_{ij}
 - Increase W_{ij} if incorrectly classified
 - Decrease W_{ij} if correctly classified
 - (Under some assumptions) obtain H that satisfies the (k, ck, p_1, p_2) -sensitive property for the training data
 - if $o \in \text{NN}(q, k)$, then collision probability from $H \geq p_1$
 - if $o \notin \text{NN}(q, ck)$, then collision probability from $H \leq p_2$

Key difference with AdaBoost: High recall and low precision

Learned ANN Index [LZSW+20]

17

- Focus on external I/O

- Use B-trees and maximize the use of sequential I/Os

- Scheme:

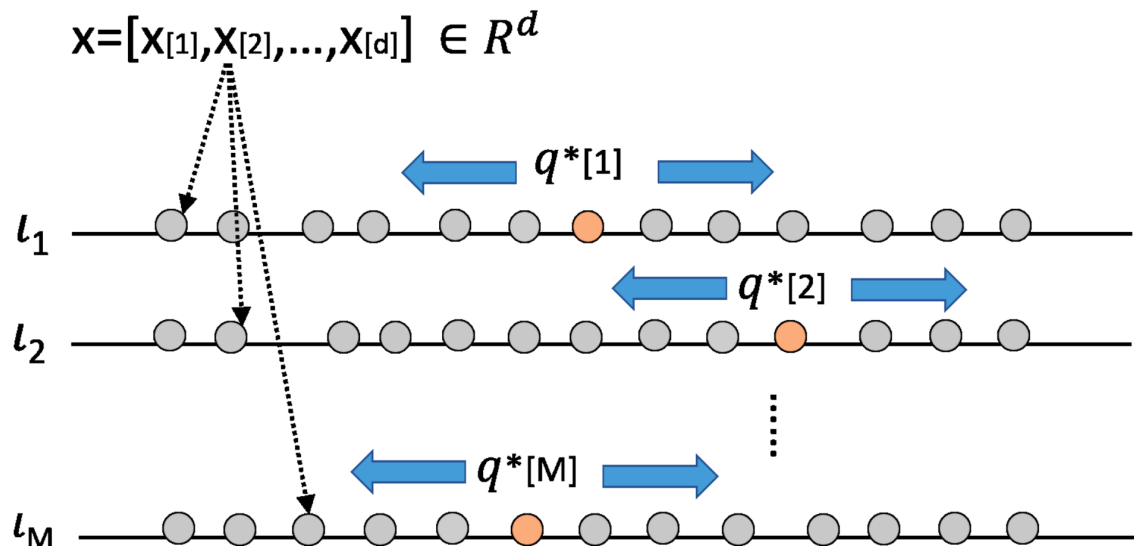
- $H: \mathbb{R}^d \rightarrow \mathbb{R}^M$
 - Index each dimension of $H(X)$ in a clustered B-tree

- Query processing

- Collect candidates on each of the M projected dimensions
 - When T candidates are seen on all M lists, rerank them and return top- k

$\approx \text{MedRank}$
(minfreq = 1)
[FKS03]

must c.f., [AFKPS08]



Function family

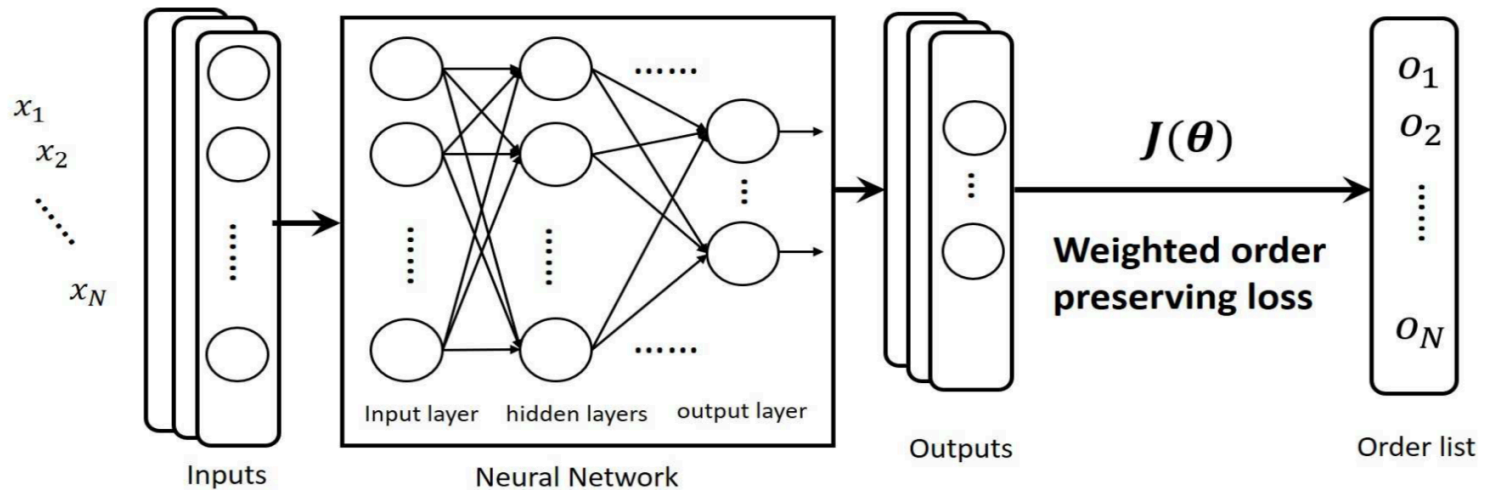
18

□ Consider

▣ linear functions

$$\blacksquare H(x)[m] = \mathbf{w}_m^T \mathbf{x}$$

▣ non-linear functions



W

How to learn the parameters?

19

Consider the linear functions:
 $H(x)[m] = \mathbf{w}_m^\top \mathbf{x}$

- Goal:
 - ▣ Encourage **segment**-order preserving mappings

Part of
the Loss
function

$$J^*(\mathbf{w}_m) = \sum_{i=1}^L \sum_{\tilde{x} \in l_i^o} \mathbf{1}_{r(\tilde{x}) \in [t \cdot (i-1), t \cdot i]}$$

Continuous
Relaxation

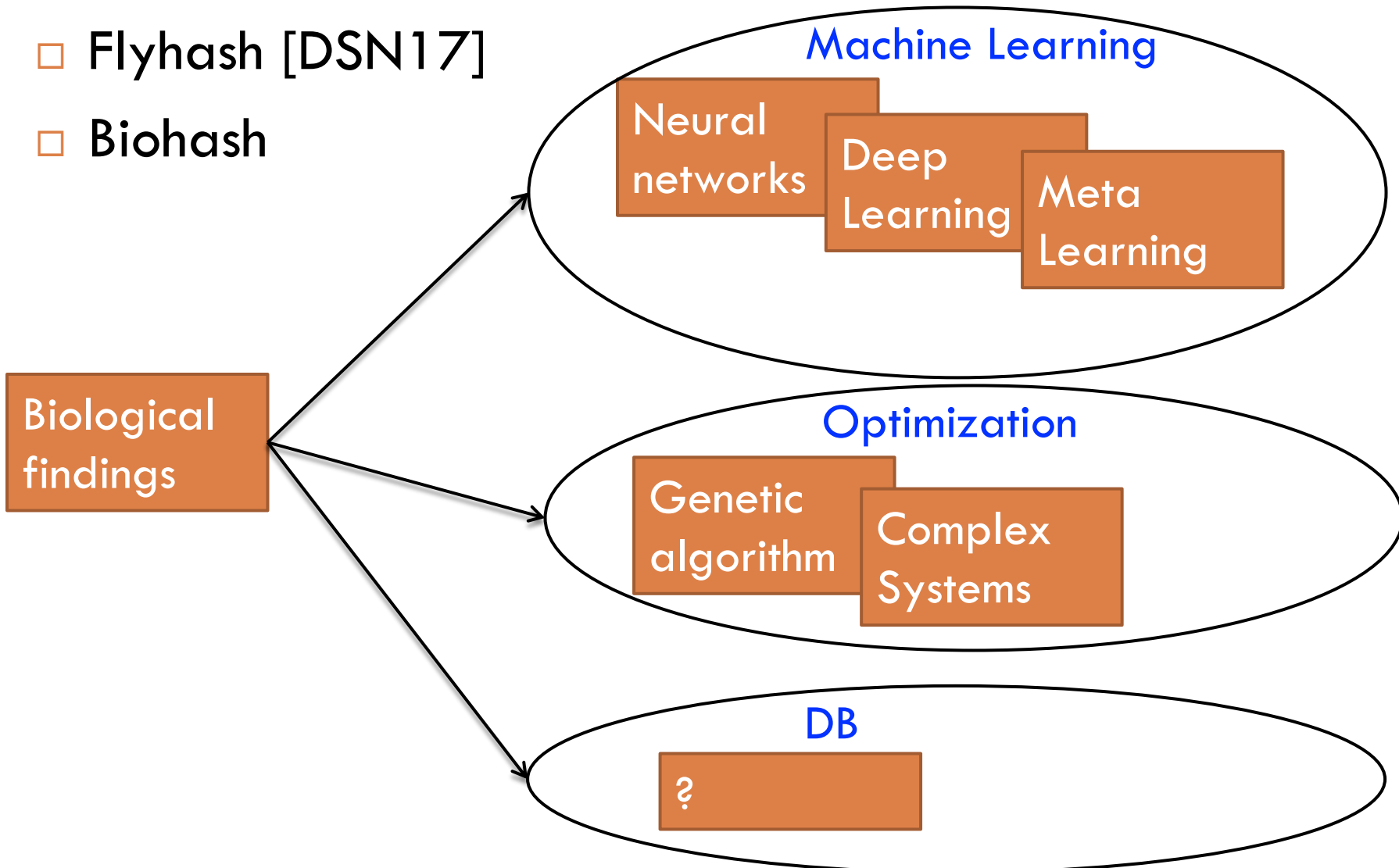
mapped \mathbf{x} in the i -th segment

\mathbf{x} in the i -th segment

Biology Inspired Hashing

20

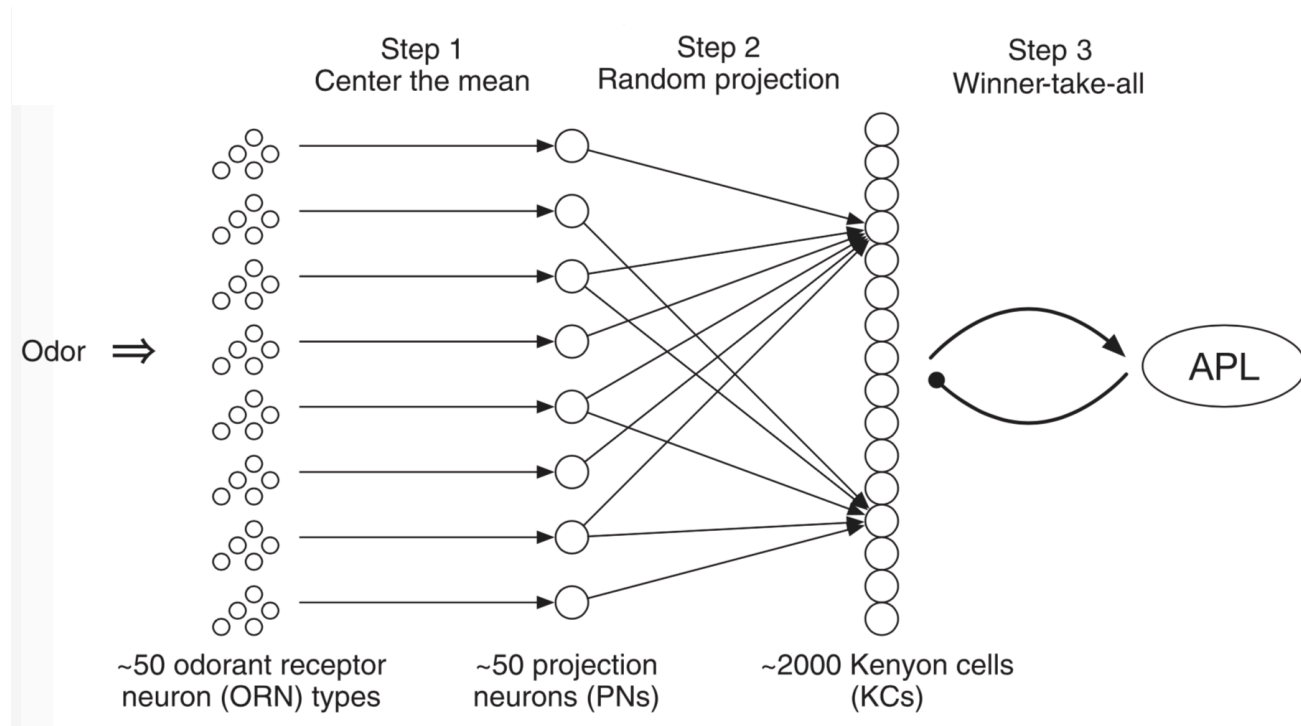
- Flyhash [DSN17]
- Biohash



FlyHash

21

- The fly olfactory circuit generates a **low-overlapping, sparse** neuron activation pattern when an odor is presented



FlyHash

22

□ Difference with LSH

- \mathbf{W} is a sparse binary random matrix

- Dimensionality expansion !!

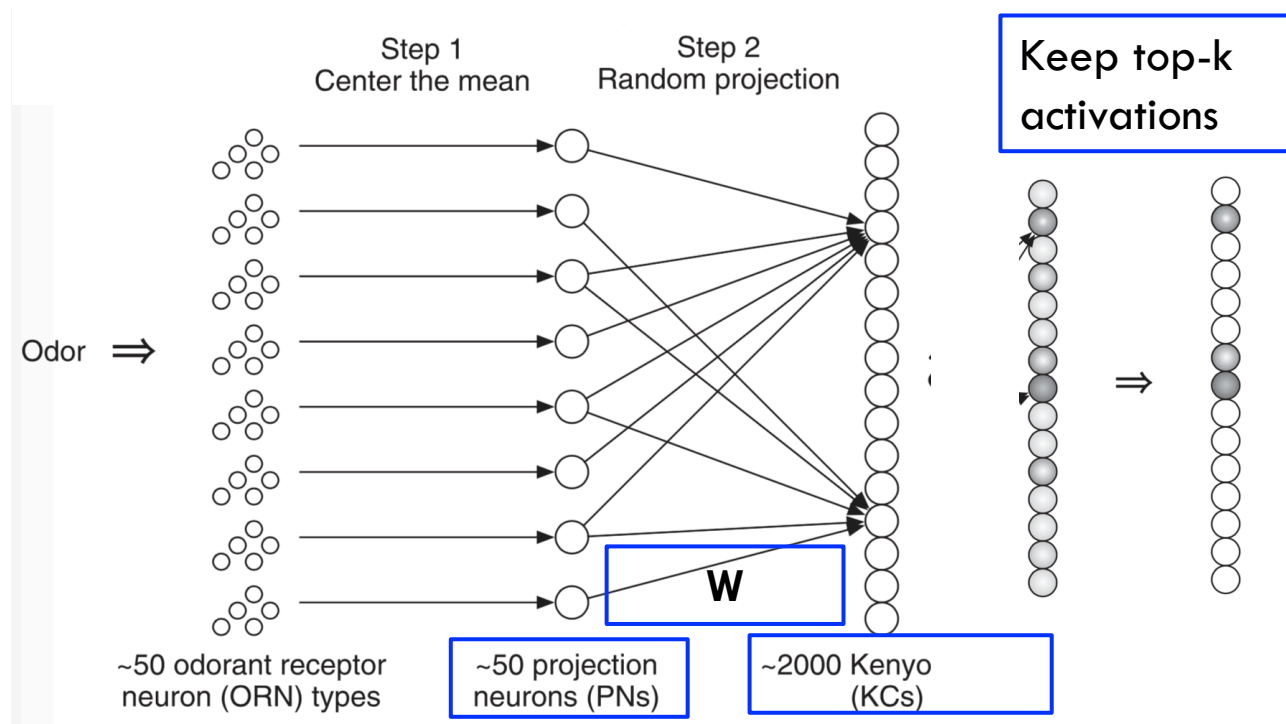
- Sparsification

- L2 distance approximately preserved in expectation

$$\mathbf{z} = \sigma(\mathbf{W}\mathbf{x})$$



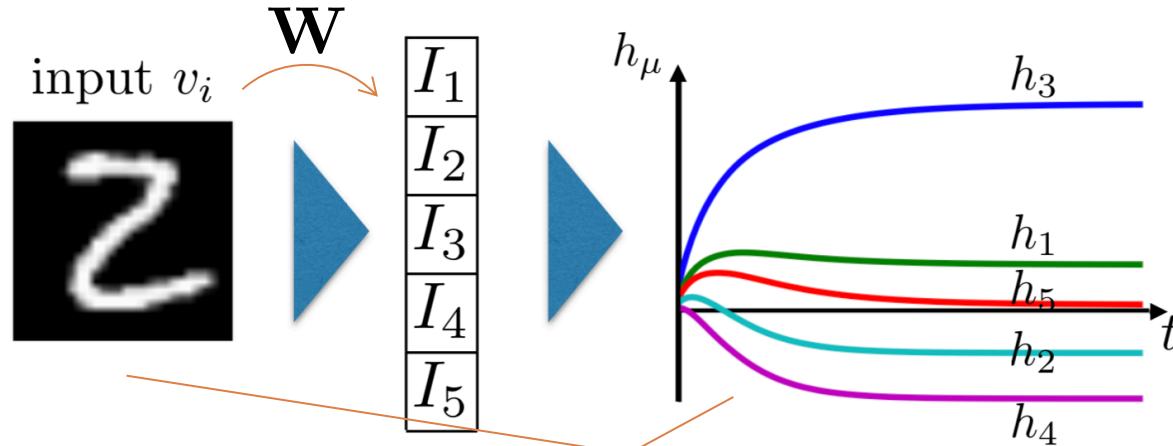
Enable



- Unsupervised learning inspired by biological synaptic plasticity rules
- Overview
 - ▣ (Given W) Stabilizing the hidden **competing** neurons
 - ▣ Learning the projection matrix W

Learning h

24



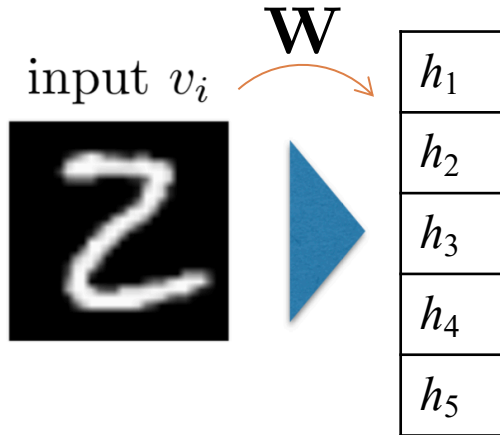
$$\tau \frac{d}{dt} (h_\mu) = I_\mu - \gamma \sum_{\nu \neq \mu} \text{RELU}(h_\nu) - h_\mu, \quad \text{where} \quad I_\mu = \langle \mathbf{W}_\mu, \mathbf{v} \rangle$$

Competing for activation

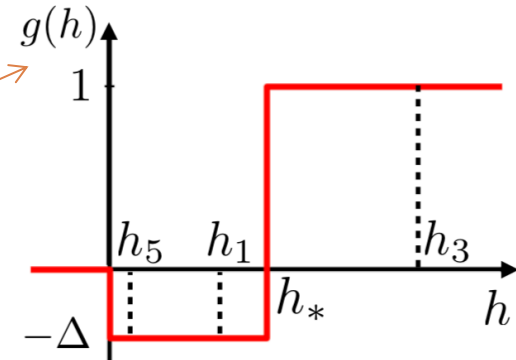
- Fixing the \mathbf{W} , the dynamical equation will converge to a stable hidden vector \mathbf{h}

Learning W

25



h_* is the activation threshold



$$\tau_L \frac{d}{dt}(\mathbf{W}_i) = g(Q) (\mathbf{v}_i - \langle \mathbf{W}, \mathbf{v} \rangle \mathbf{W}_i), \quad \text{where } Q = \frac{\langle \mathbf{W}, \mathbf{v} \rangle}{\langle \mathbf{W}, \mathbf{W} \rangle^{\frac{1}{2}}}$$

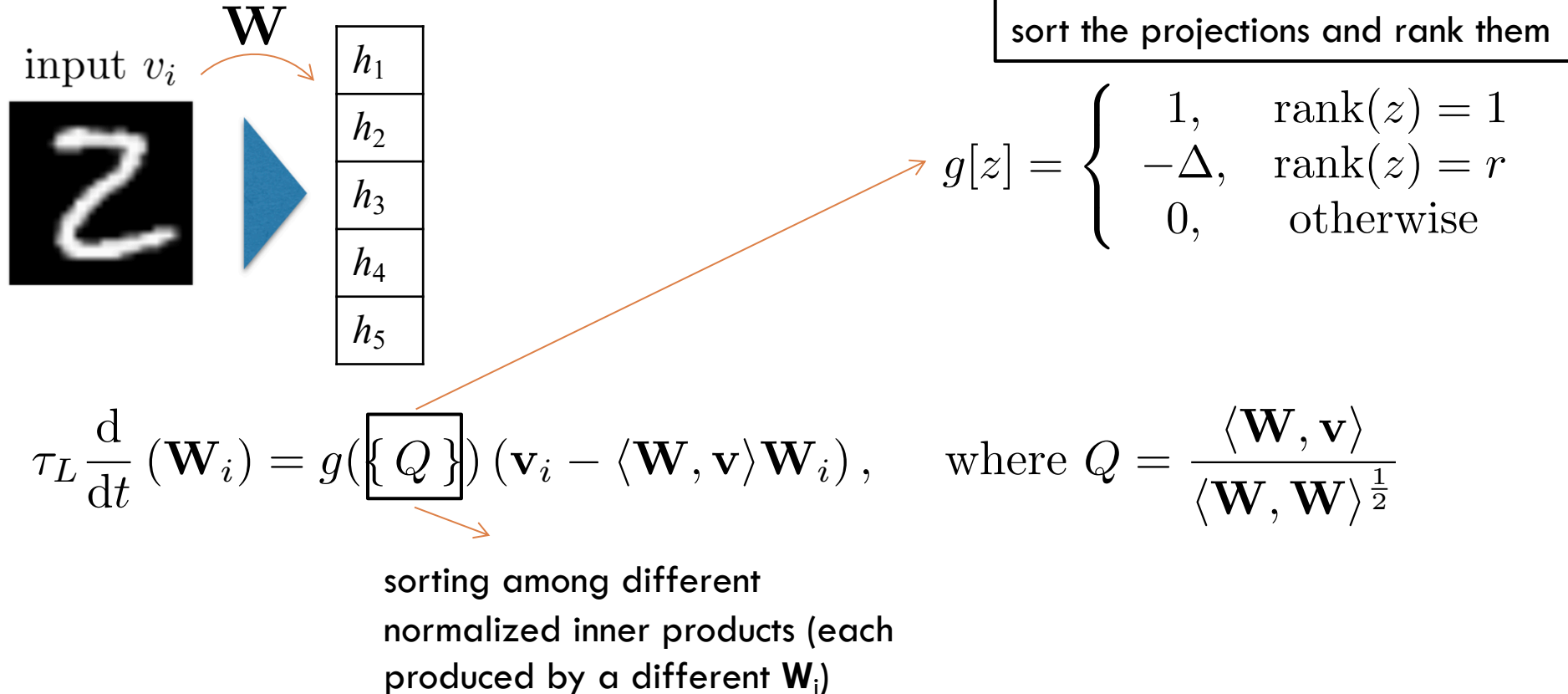
$p = 2$, for any one hidden neuron
 \mathbf{W} is its corresponding weight vector
 $R = 1$

Force $|\mathbf{W}|_p$ to converge to R^p

- Fixing the h , the dynamical equation will converge to a final weight matrix W

BioHash – Learning \mathbf{W}

26



- The rest is the same as FlyHash (i.e., WTA sparsification)