

# **BLOOM FILTERS - RECAP**

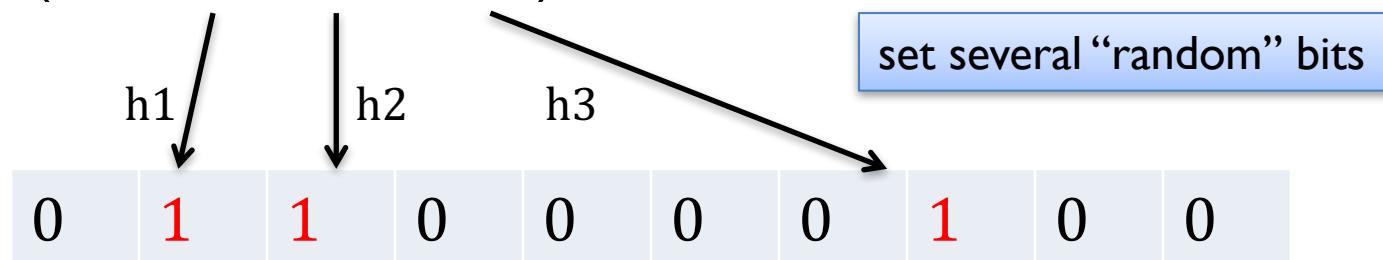
# Bloom filters

- Interface to a Bloom filter
  - `BloomFilter(int maxSize, double p);`
  - `void bf.add(String s); // insert s`
  - `bool bd.contains(String s);`
    - // If s was added return true;
    - // else with probability at least  $1-p$  return false;
    - // else with probability at most  $p$  return true;
  - I.e., a noisy “set” where you can test membership (and that’s it)

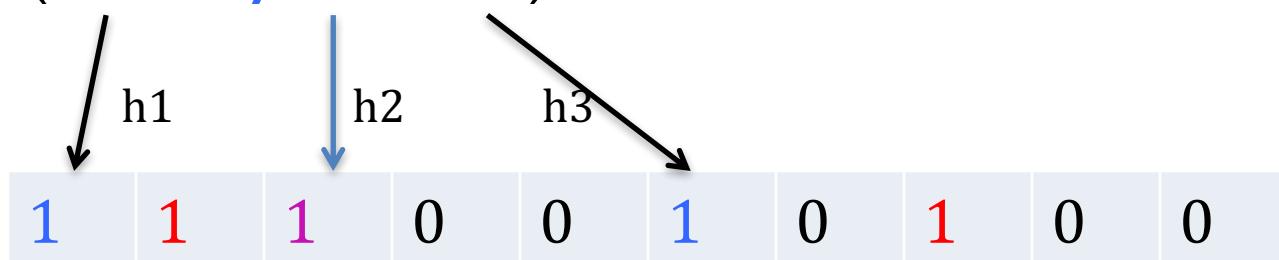
# Bloom filters



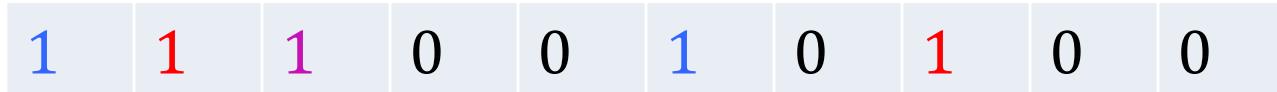
`bf.add("fred flintstone"):`



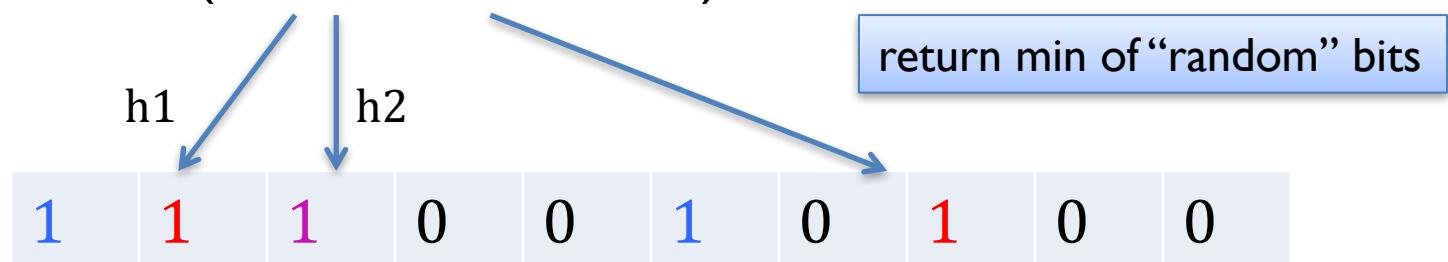
`bf.add("barney rubble"):`



# Bloom filters

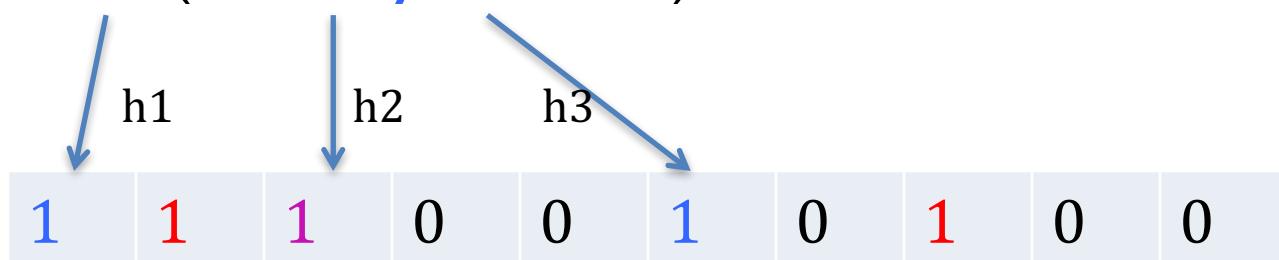


bf.contains (“fred flintstone”):

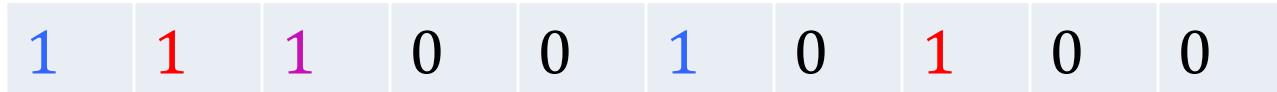


return min of “random” bits

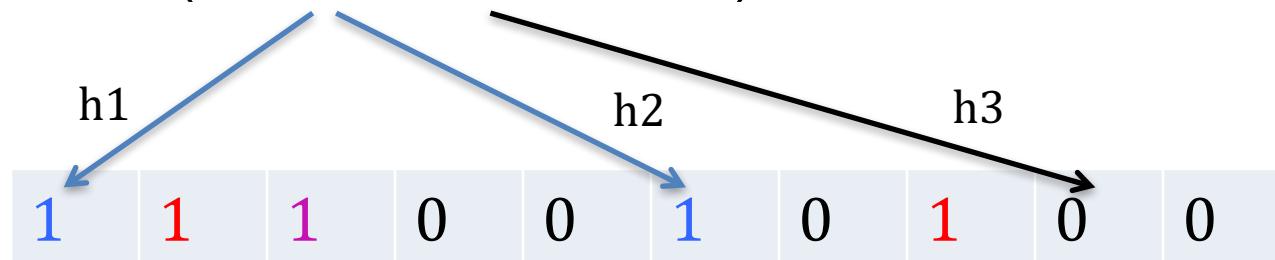
bf.contains(“barney rubble”):



# Bloom filters

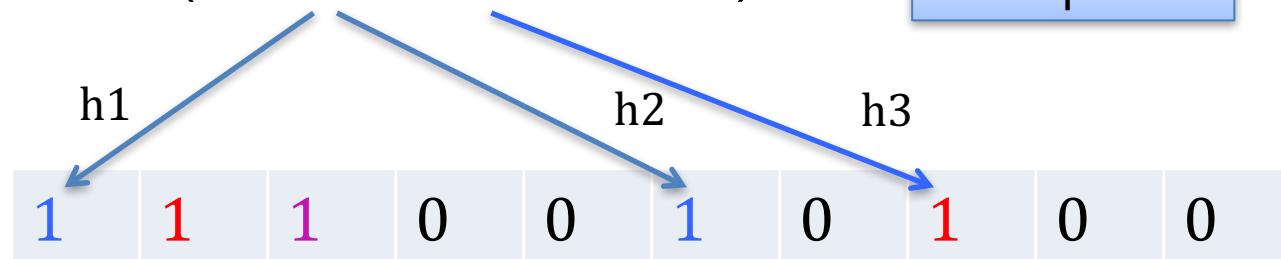


`bf.contains("wilma flintstone"):`



`bf.contains("wilma flintstone"):`

a false positive



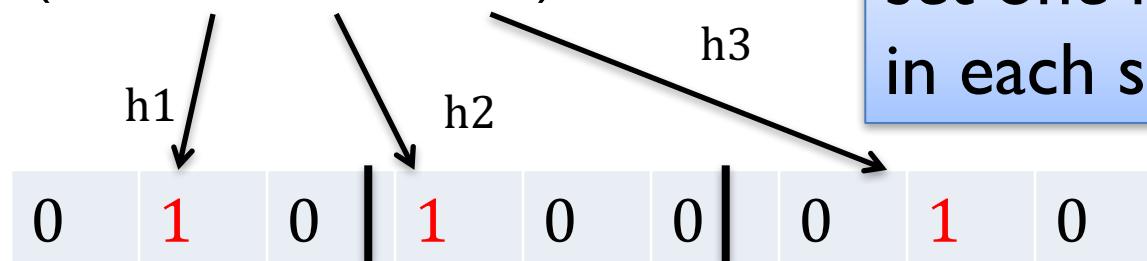
# **BLOOM FILTERS VS COUNT-MIN SKETCHES**

# Bloom filters – a variant

split the bit vector into  $k$  ranges, one for each hash function

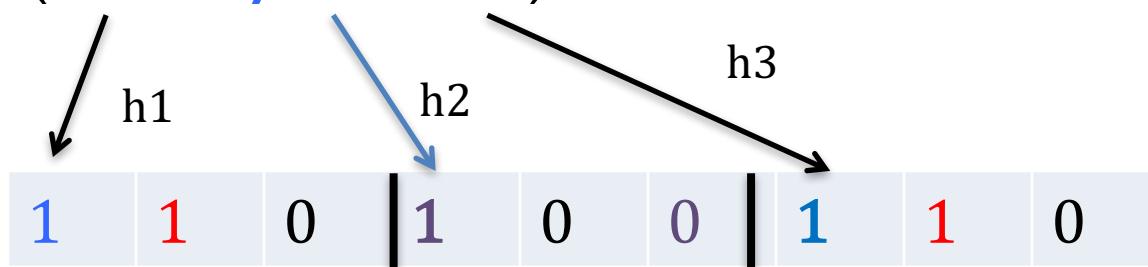


bf.add("fred flintstone"):



set one random bit  
in each subrange

bf.add("barney rubble"):

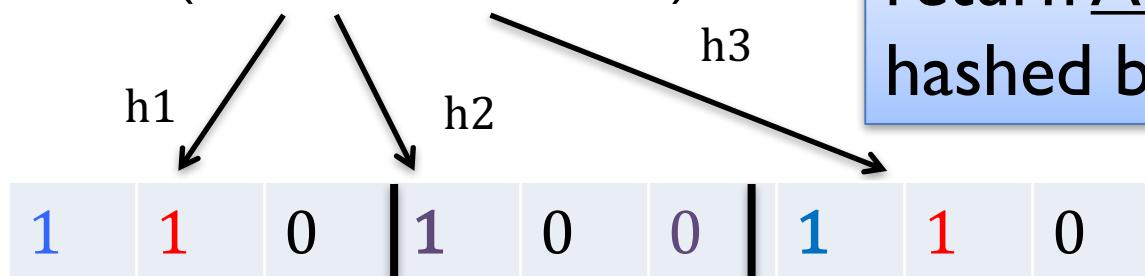


# Bloom filters – a variant

split the bit vector into  $k$  ranges, one for each hash function

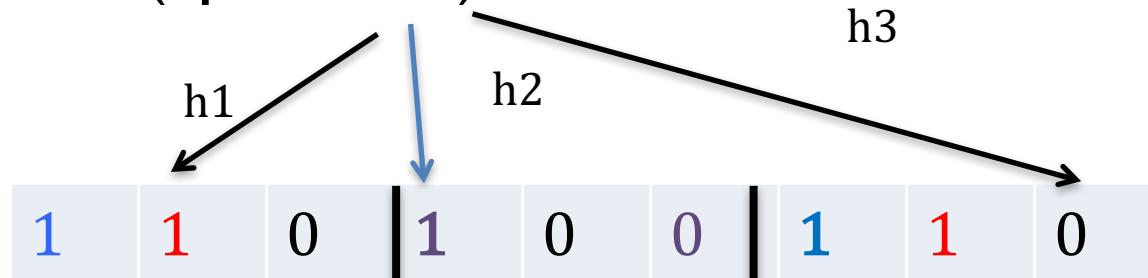


bf.contains("fred flintstore"):



return AND of all  
hashed bits

bf.contains("pebbles"):



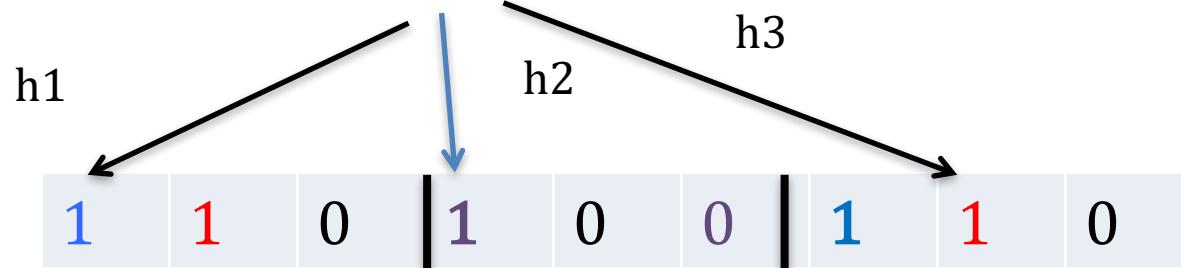
# Bloom filters – a variant

split the bit vector into  $k$  ranges, one for each hash function



bf.contains("pebbles"):

a false positive!

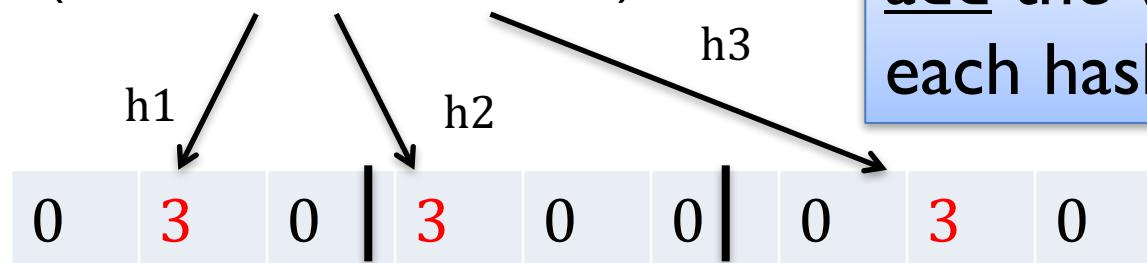


# Count-min sketches

split a real vector into  $k$  ranges, one for each hash function

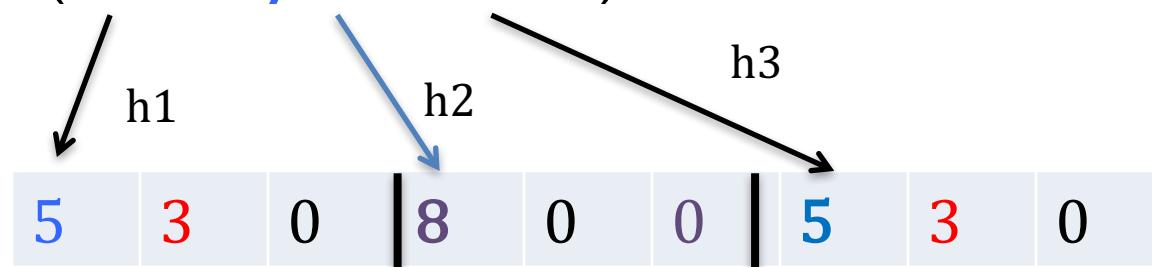


cm.inc("fred flintstone", 3):



add the value to  
each hash location

cm.inc("barney rubble", 5):

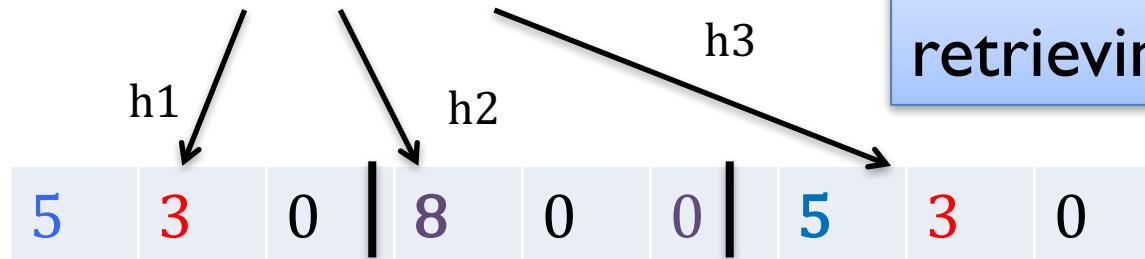


# Count-min sketches

split a real vector into  $k$  ranges, one for each hash function

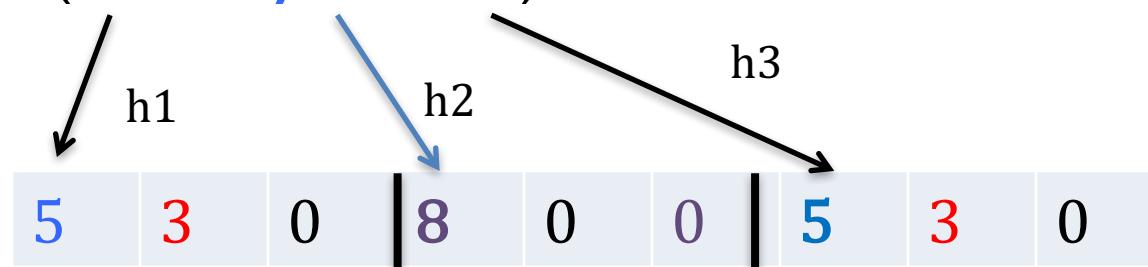


cm.get("fred flintstone"): 3



take min when retrieving a value

cm.get("barney rubble"): 5

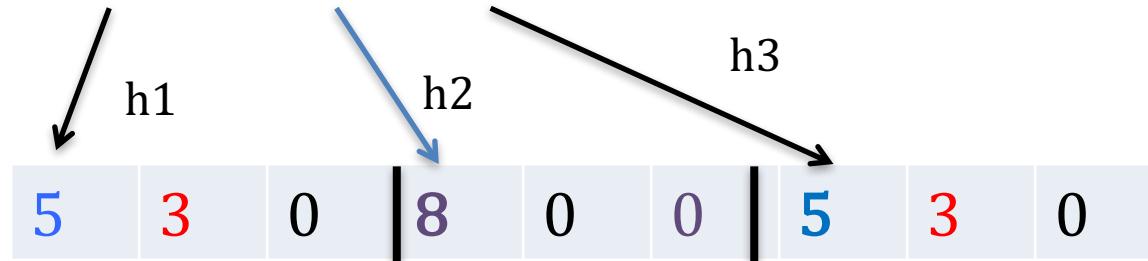


# Count-min sketches

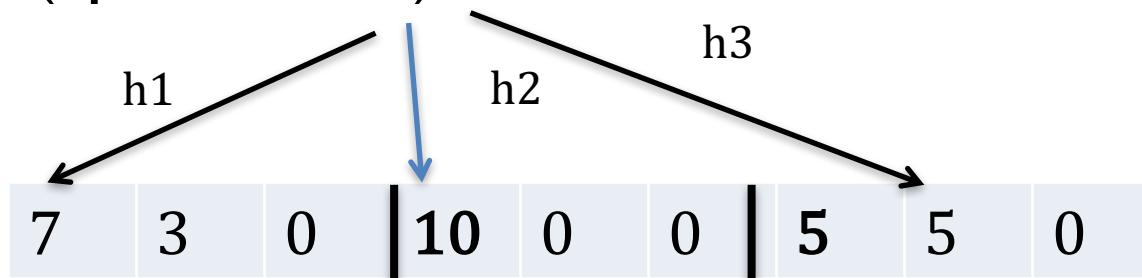
split a real vector into  $k$  ranges, one for each hash function



`cm.get("barney rubble"):` 5



`cm.add("pebbles", 2):`



# Count-min sketches

Equivalently, use a matrix, and each hash leads to a different row

cm.inc("fred flintstone", 3):

0	3	0
3	0	0
0	3	0

cm.inc("barney rubble",5):

5	3	0
8	0	0
5	3	0

# **LOCALITY SENSITIVE HASHING (LSH)**

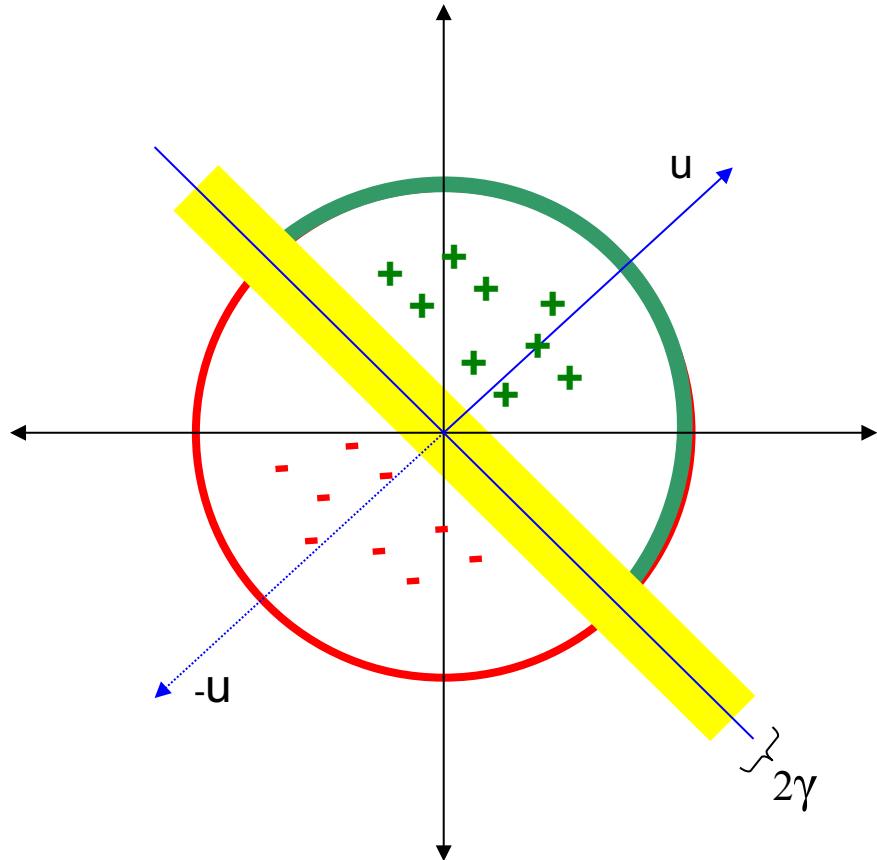
# LSH: key ideas

- Goal:
  - map feature vector  $\mathbf{x}$  to bit vector  $\mathbf{bx}$
  - ensure that  $\mathbf{bx}$  preserves “similarity”

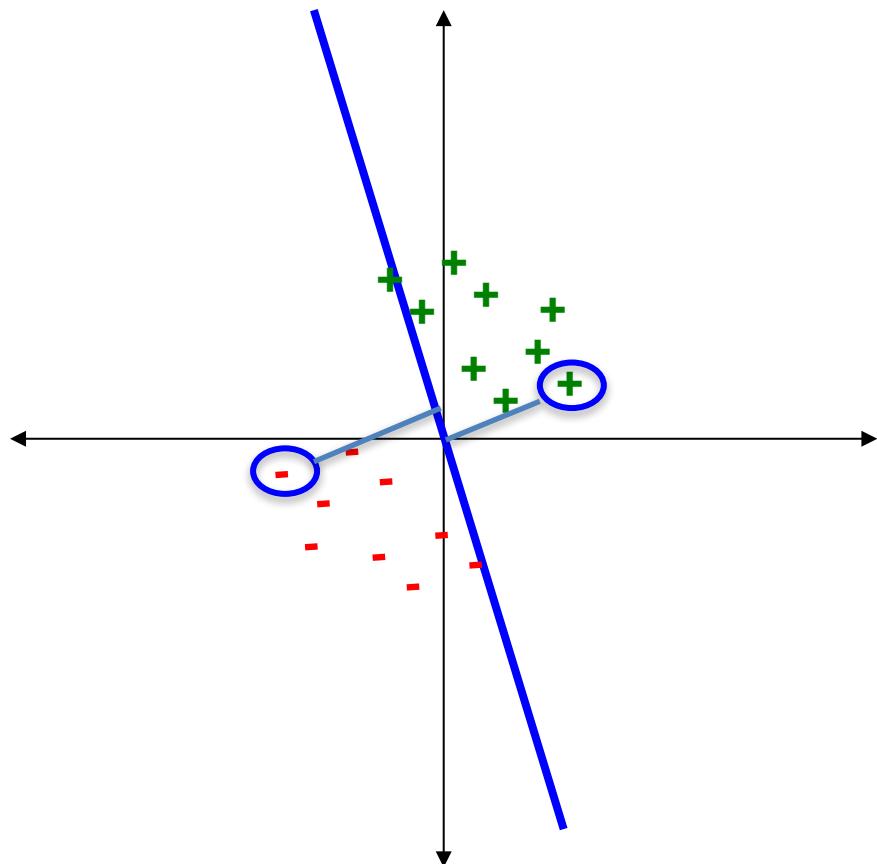
# Random Projections



# Random projections

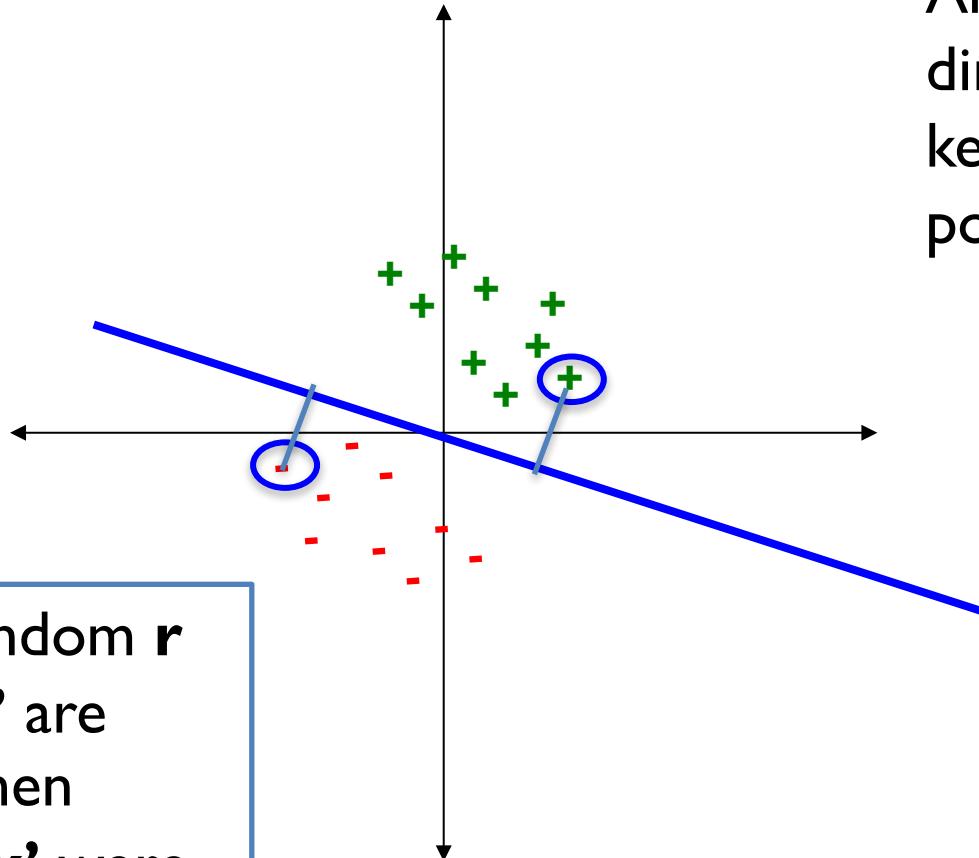


# Random projections



To make those points “close” we need to project to a direction orthogonal to the line between them

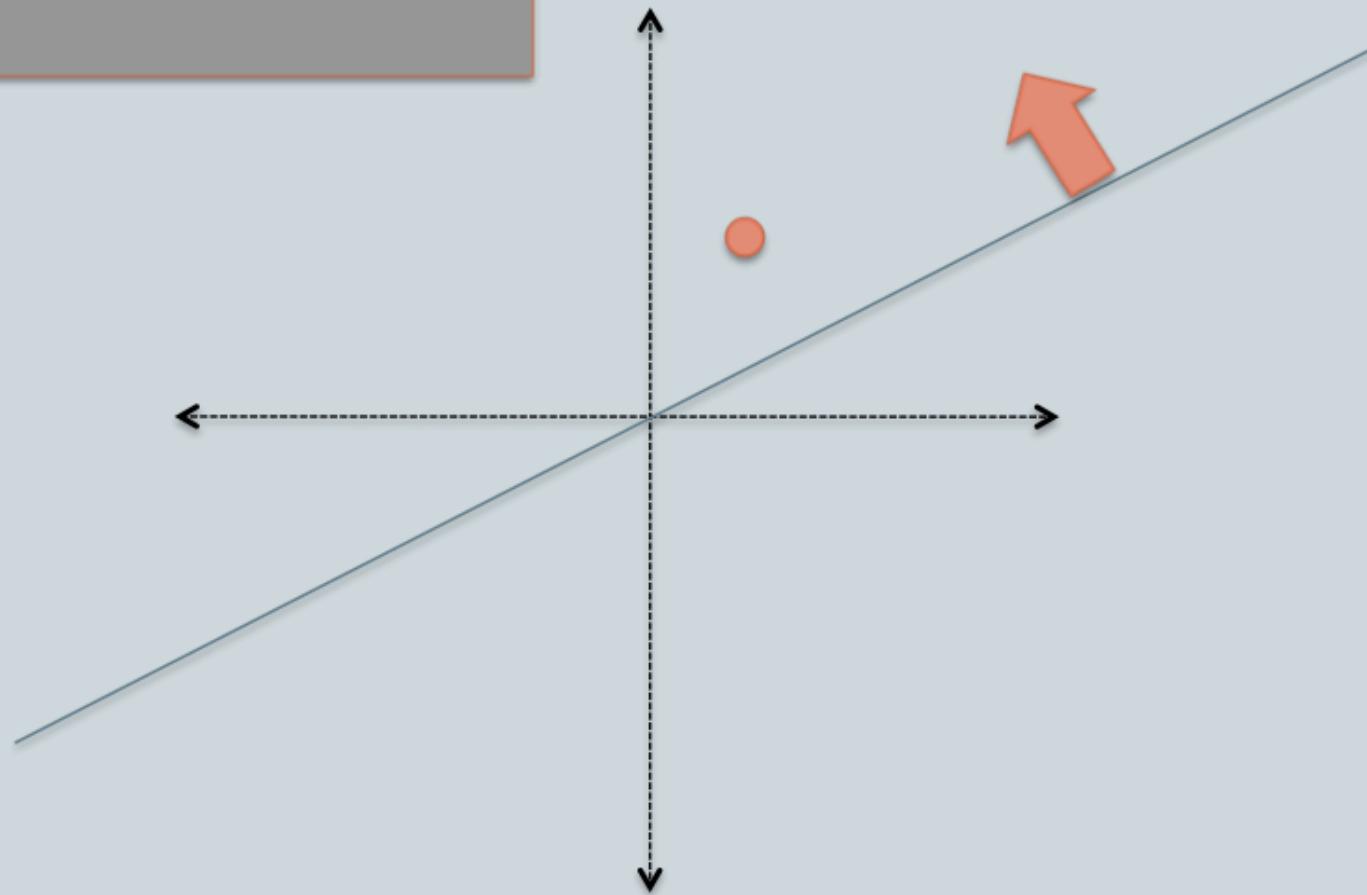
# Random projections



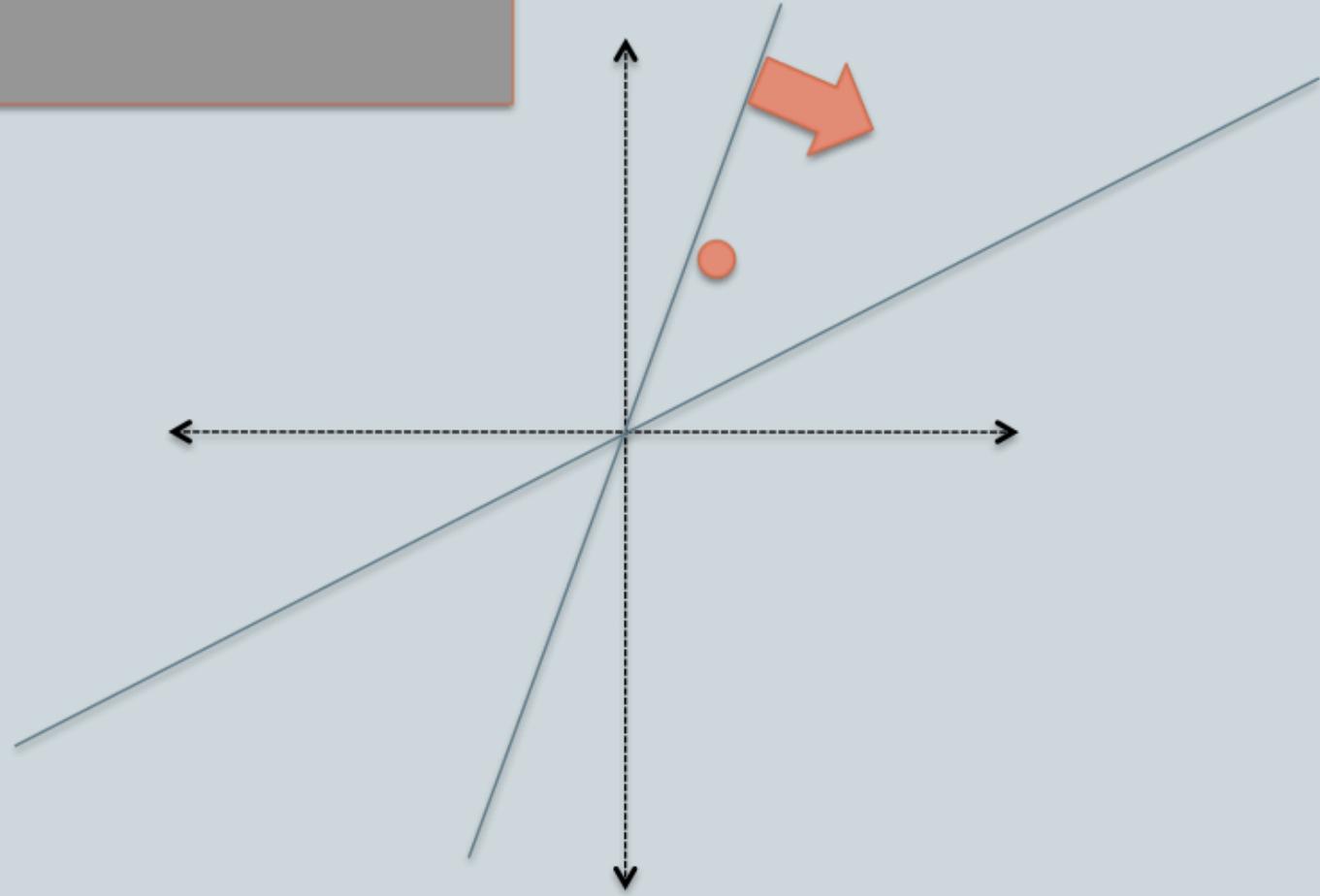
So if I pick a random  $r$  and  $r \cdot x$  and  $r \cdot x'$  are closer than  $\gamma$  then *probably*  $x$  and  $x'$  were close to start with.

# LSH: key ideas

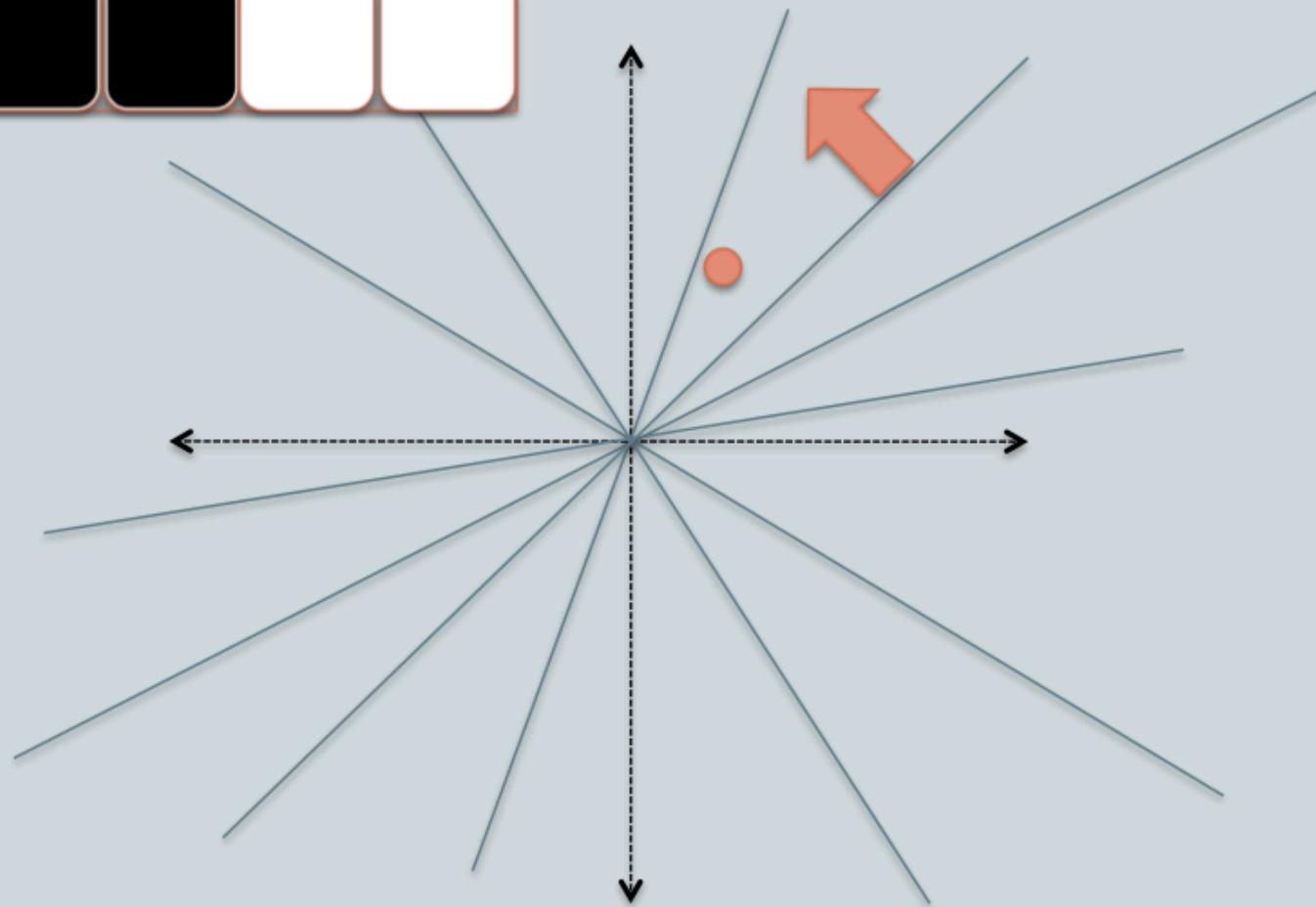
- Goal:
  - map feature vector  $\mathbf{x}$  to bit vector  $\mathbf{bx}$
  - ensure that  $\mathbf{bx}$  preserves “similarity”
- Basic idea: use *random projections* of  $\mathbf{x}$ 
  - Repeat many times:
    - Pick a random hyperplane  $\mathbf{r}$  by picking random weights for each feature (say from a Gaussian)
    - Compute the inner product of  $\mathbf{r}$  with  $\mathbf{x}$
    - Record if  $\mathbf{x}$  is “close to”  $\mathbf{r}$  ( $\mathbf{r} \cdot \mathbf{x} >= 0$ )
      - the next bit in  $\mathbf{bx}$
    - Theory says that if  $\mathbf{x}'$  and  $\mathbf{x}$  have small cosine distance then  $\mathbf{bx}$  and  $\mathbf{bx}'$  will have small Hamming distance



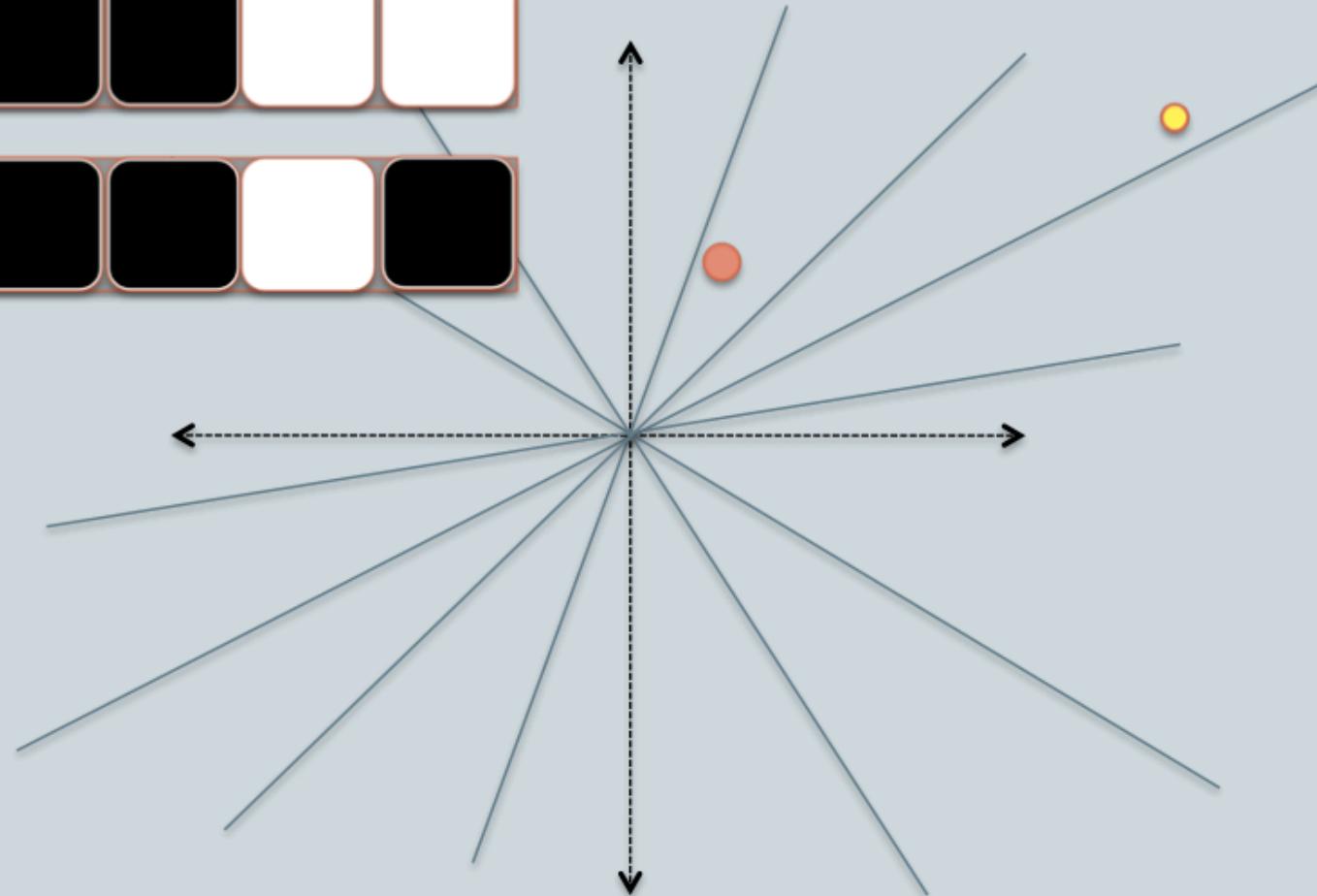
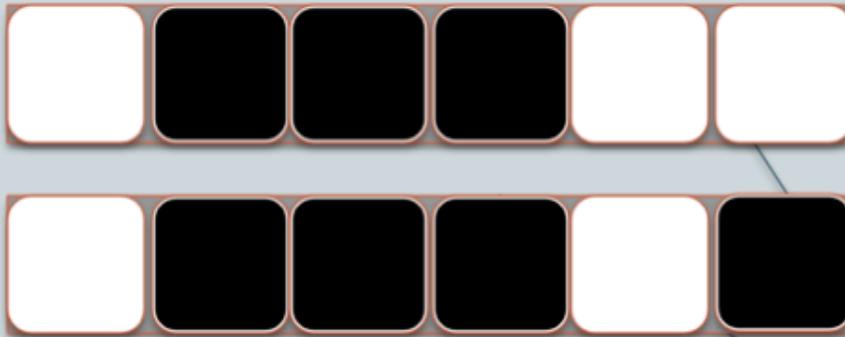
[Slides: Ben van Durme]



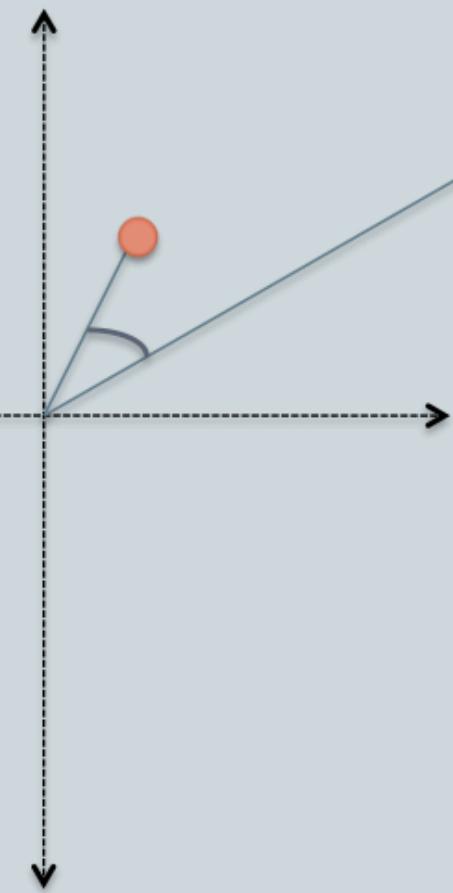
[Slides: Ben van Durme]



[Slides: Ben van Durme]



[Slides: Ben van Durme]



Hamming Distance :=  $h = 1$

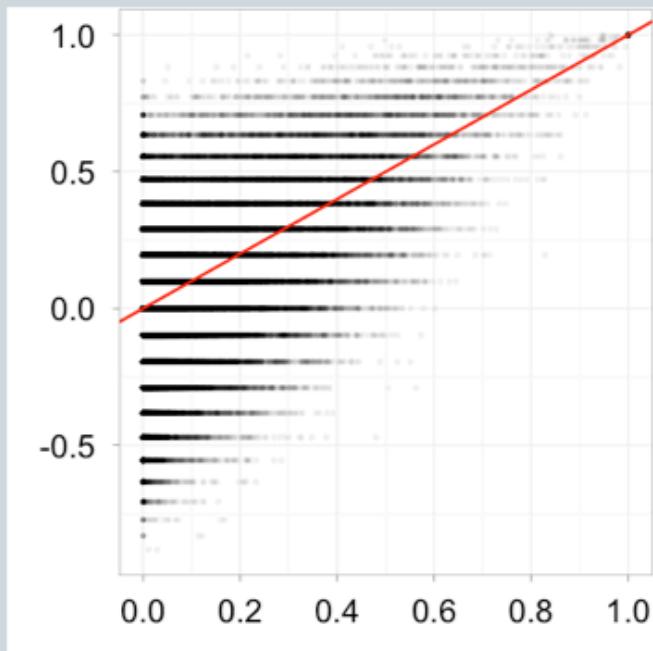
Signature Length :=  $b = 6$

$$\cos(\theta) \approx \cos\left(\frac{h}{b}\pi\right)$$

$$= \cos\left(\frac{1}{6}\pi\right)$$

32 bit signatures

Approximate Cosine

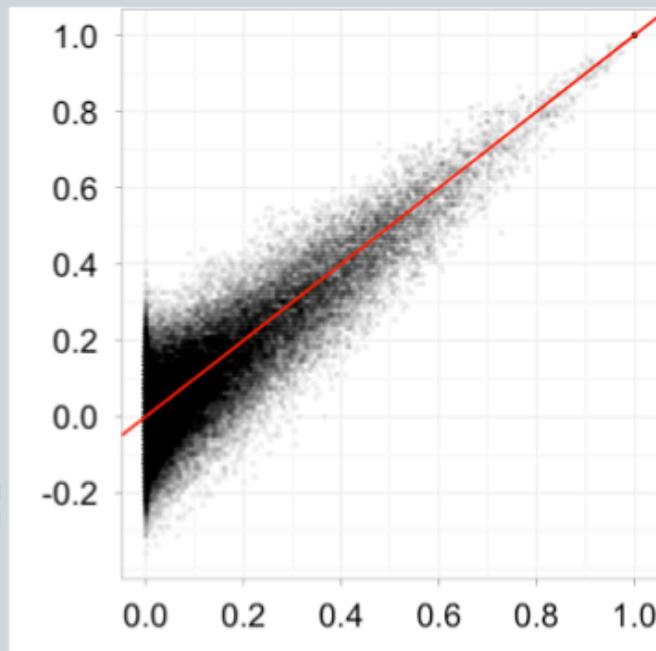


True Cosine

Cheap

256 bit signatures

Approximate Cosine



True Cosine

Accurate

[Slides: Ben van Durme]

# LSH applications

- Compact storage of data
  - and we can still compute similarities
- LSH also gives very fast approximations:
  - approx nearest neighbor method
    - just look at other items with  $\mathbf{bx}' = \mathbf{bx}$
    - also very fast nearest-neighbor methods for Hamming distance
  - very fast clustering
    - cluster = all things with same  $\mathbf{bx}$  vector

# **Online LSH and Pooling**

# Online Generation of Locality Sensitive Hash Signatures

Benjamin Van Durme and Ashwin Lall



human language technology  
center of excellence

JOHNS HOPKINS  
UNIVERSITY

DENISON  
UNIVERSITY

# LSH algorithm

- Naïve algorithm:
  - Initialization:
    - For  $i=1$  to outputBits:
      - For each feature  $f$ :
        - » Draw  $r(f,i) \sim \text{Normal}(0,1)$
    - Given an instance  $\mathbf{x}$ 
      - For  $i=1$  to outputBits:  
 $\text{LSH}[i] =$   
 $\text{sum}(\mathbf{x}[f]*r[i,f] \text{ for } f \text{ with non-zero weight in } \mathbf{x}) > 0 ?$   
 $1 : 0$
    - Return the bit-vector LSH

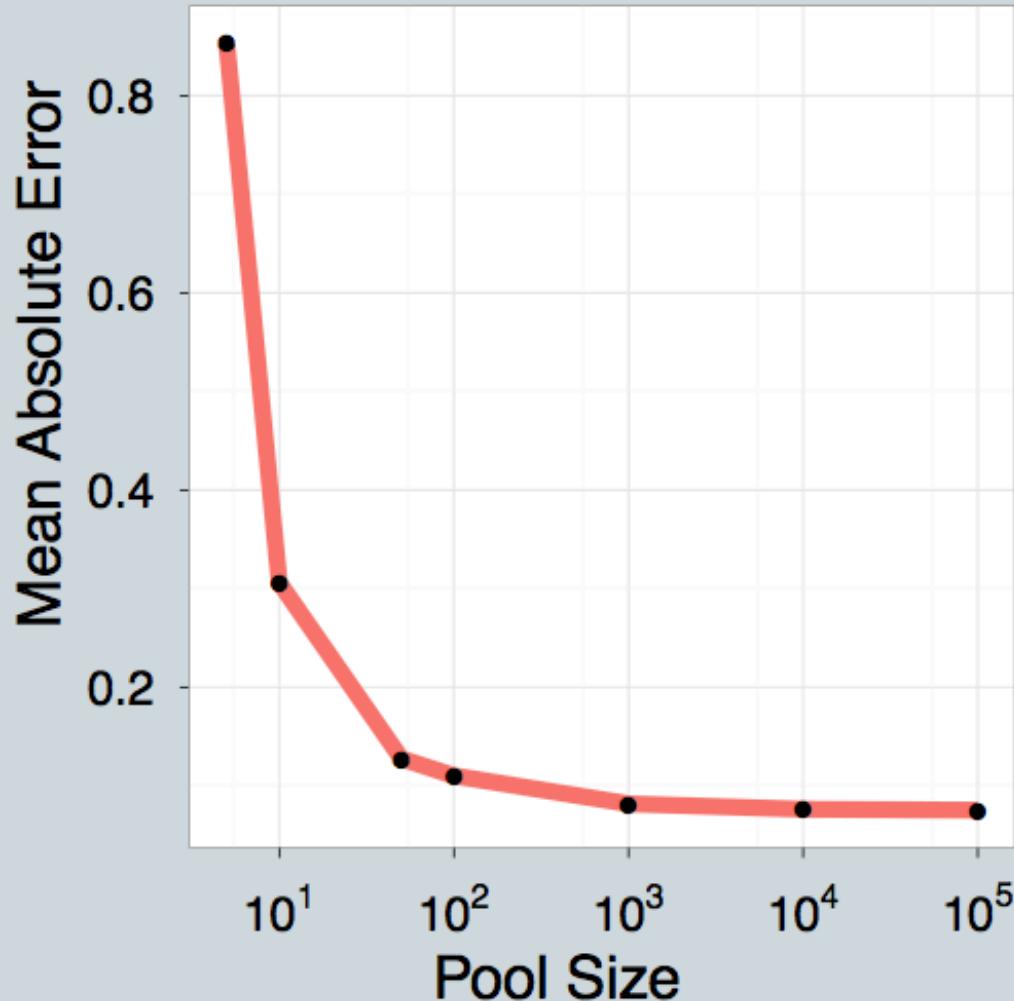
# LSH algorithm

- But: storing the  $k$  classifiers is expensive in high dimensions
  - For each of 256 bits, a dense vector of weights for every feature in the vocabulary
- Storing seeds and random number generators:
  - Possible but somewhat fragile

# LSH: “pooling” (van Durme)

- Better algorithm:
  - Initialization:
    - Create a pool:
      - Pick a random seed  $s$
      - For  $i=1$  to  $\text{poolSize}$ :
        - » Draw  $\text{pool}[i] \sim \text{Normal}(0,1)$
    - For  $i=1$  to  $\text{outputBits}$ :
      - Devise a random hash function  $\text{hash}(i,f)$ :
        - » E.g.:  $\text{hash}(i,f) = \text{hashcode}(f) \text{ XOR } \text{randomBitString}[i]$
  - Given an instance  $\mathbf{x}$ 
    - For  $i=1$  to  $\text{outputBits}$ :
$$\text{LSH}[i] = \sum(\mathbf{x}[f] * \text{pool}[\text{hash}(i,f) \% \text{poolSize}] \text{ for } f \in \mathbf{x}) > 0 ? 1 : 0$$
    - Return the bit-vector LSH

# The Pooling Trick



# LSH: key ideas: pooling

- Advantages:
  - with pooling, this is a compact re-encoding of the data
    - you don't need to store the  $r$ 's, just the pool

# Locality Sensitive Hashing (LSH) in an On-line Setting

# LSH: key ideas: online computation

- Common task: **distributional clustering**
  - for a word  $w$ ,  $v(w)$  is sparse vector of words that co-occur with  $w$
  - cluster the  $v(w)$ 's

...guards at Pentonville prison in North London discovered that an escape attempt...

An American Werewolf in London is to be remade by the son of the original director...

...UK pop up shop on Monmouth Street in London today and on Friday the brand...

$v(\text{London})$ : Pentonville, prison, in, North, .... and, on Friday

# LSH: key ideas: online computation

- Common task: distributional clustering
  - for a word  $w$ ,  $v(w)$  is sparse vector of words that co-occur with  $w$
  - cluster the  $v(w)$ 's

London is similar to:

**Milan**.<sup>.97</sup>, **Madrid**.<sup>.96</sup>, **Stockholm**.<sup>.96</sup>, **Manila**.<sup>.95</sup>, **Moscow**.<sup>.95</sup>

in is similar to:

**during**.<sup>.99</sup>, **on**.<sup>.98</sup>, **beneath**.<sup>.98</sup>, **from**.<sup>.98</sup>, **onto**.<sup>.97</sup>

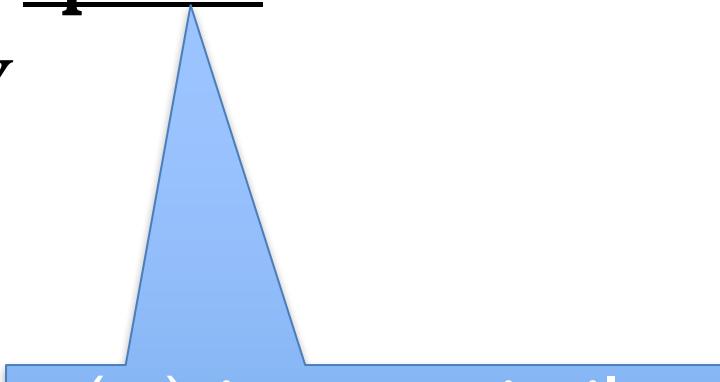
sold is similar to:

**deployed**.<sup>.84</sup>, **presented**.<sup>.83</sup>, **sacrificed**.<sup>.82</sup>, **held**.<sup>.82</sup>, **installed**.<sup>.82</sup>

# LSH: key ideas: online computation

- Common task: distributional clustering
  - for a word  $w$ ,  $v(w)$  is sparse vector of words that co-occur with  $w$
  - cluster the  $v(w)$ 's

Levy, Omer, Yoav Goldberg, and Ido Dagan. "Improving distributional similarity with lessons learned from word embeddings." *Transactions of the Association for Computational Linguistics* 3 (2015): 211-225.



$v(w)$  is very similar to a word embedding (eg, from word2vec or GloVe)

$$\vec{v} \in \mathbb{R}^d$$

$v$  is context vector;  $d$  is vocab size

$$\vec{r}_i \sim N(0, 1)^d \quad r_i \text{ is } i\text{-th random projection}$$

$$h_i(\vec{v}) = \begin{cases} 1 & \text{if } \vec{v} \cdot \vec{r}_i \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

...guards at Pentonville prison in North London ...

An American Werewolf in London is to be remade

...UK pop up shop on Monmouth Street in London

$v(\text{London})$ : Pentonville, prison, in, North, .... and

$h_i(v)$   $i$ -th bit of LSH encoding

if  $\vec{v} = \sum_j \vec{v}_j$  because context vector is sum of mention contexts

then  $\vec{v} \cdot \vec{r}_i = \sum_j \vec{v}_j \cdot \vec{r}_i$

Break into local products

these come one by one as we stream thru the corpus

online

$$h_{it}(\vec{v}) = \begin{cases} 1 & \text{if } \sum_j^t \vec{v}_j \cdot \vec{r}_i \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

---

## Algorithm 1 STREAMING LSH ALGORITHM

---

### Parameters:

$m$  : size of pool

$d$  : number of bits (size of resultant signature)

$s$  : a random seed

$h_1, \dots, h_d$  : hash functions mapping  $\langle s, f_i \rangle$  to  $\{0, \dots, m-1\}$

### INITIALIZATION:

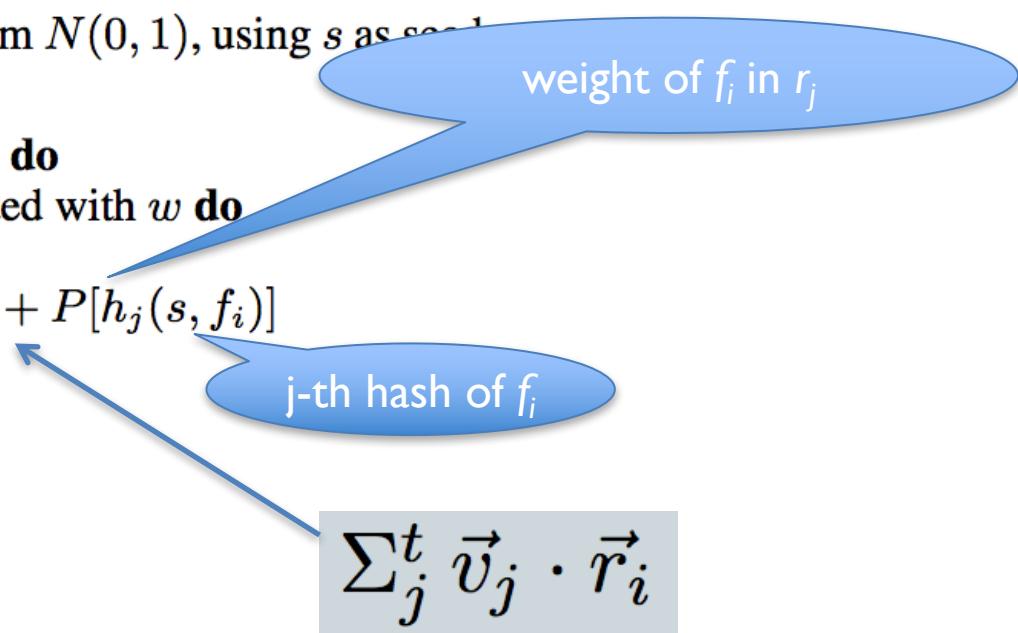
- 1: Initialize floating point array  $P[0, \dots, m-1]$
- 2: Initialize  $H$ , a hashtable mapping words to floating point arrays of size  $d$
- 3: **for**  $i := 0 \dots m-1$  **do**
- 4:      $P[i] :=$  random sample from  $N(0, 1)$ , using  $s$  as seed

### ONLINE:

- 1: **for** each word  $w$  in the stream **do**
- 2:     **for** each feature  $f_i$  associated with  $w$  **do**
- 3:         **for**  $j := 1 \dots d$  **do**
- 4:              $H[w][j] := H[w][j] + P[h_j(s, f_i)]$

### SIGNATURE COMPUTATION:

- 1: **for** each  $w \in H$  **do**
- 2:     **for**  $i := 1 \dots d$  **do**
- 3:         **if**  $H[w][i] > 0$  **then**
- 4:              $S[w][i] := 1$
- 5:         **else**
- 6:              $S[w][i] := 0$



weight of  $f_i$  in  $r_j$

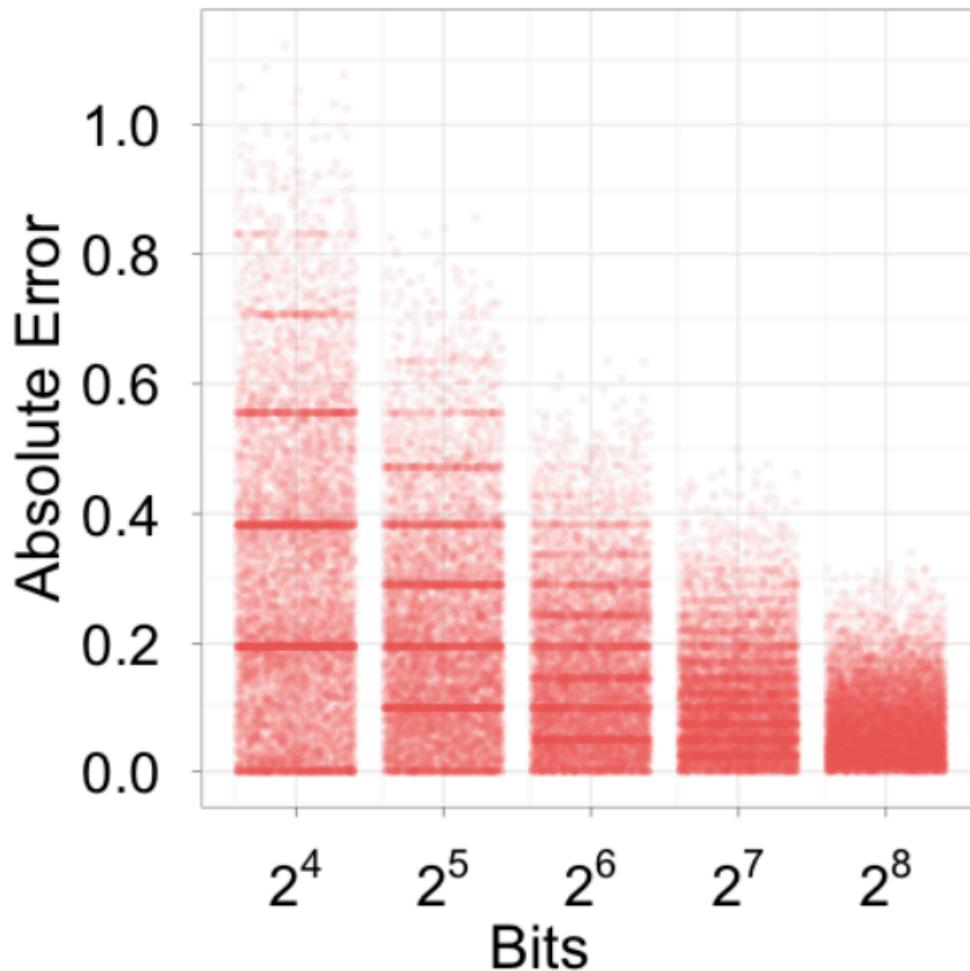
j-th hash of  $f_i$

$$\sum_j^t \vec{v}_j \cdot \vec{r}_i$$

# Experiment

- Corpus: 700M+ tokens, 1.1M distinct bigrams
- For each, build a feature vector of words that co-occur near it, using on-line LSH
- Check results with 50,000 actual vectors

# Experiment



**London**

**Milan<sub>.97</sub>, Madrid<sub>.96</sub>, Stockholm<sub>.96</sub>, Manila<sub>.95</sub>, Moscow<sub>.95</sub>**

ASHER<sub>0</sub>, Champaign<sub>0</sub>, MANS<sub>0</sub>, NOBLE<sub>0</sub>, come<sub>0</sub>

Prague<sub>1</sub>, Vienna<sub>1</sub>, suburban<sub>1</sub>, synchronism<sub>1</sub>, Copenhagen<sub>2</sub>

Frankfurt<sub>4</sub>, Prague<sub>4</sub>, Taszar<sub>5</sub>, Brussels<sub>6</sub>, Copenhagen<sub>6</sub>

Prague<sub>12</sub>, Stockholm<sub>12</sub>, Frankfurt<sub>14</sub>, Madrid<sub>14</sub>, Manila<sub>14</sub>

Stockholm<sub>20</sub>, Milan<sub>22</sub>, Madrid<sub>24</sub>, Taipei<sub>24</sub>, Frankfurt<sub>25</sub>

Closest based on true cosine

**London**

Milan<sub>.97</sub>, Madrid<sub>.96</sub>, Stockholm<sub>.96</sub>, Manila<sub>.95</sub>, Moscow<sub>.95</sub>  
ASHER<sub>0</sub>, Champaign<sub>0</sub>, MANS<sub>0</sub>, NOBLE<sub>0</sub>, come<sub>0</sub>  
Prague<sub>1</sub>, Vienna<sub>1</sub>, suburban<sub>1</sub>, synchronism<sub>1</sub>, Copenhagen<sub>2</sub>

**London**

Milan<sub>.97</sub>, Madrid<sub>.96</sub>, Stockholm<sub>.96</sub>, Manila<sub>.95</sub>, Moscow<sub>.95</sub>  
ASHER<sub>0</sub>, Champaign<sub>0</sub>, MANS<sub>0</sub>, NOBLE<sub>0</sub>, come<sub>0</sub>  
Prague<sub>1</sub>, Vienna<sub>1</sub>, suburban<sub>1</sub>, synchronism<sub>1</sub>, Copenhagen<sub>2</sub>  
Frankfurt<sub>4</sub>, Prague<sub>4</sub>, Taszar<sub>5</sub>, Brussels<sub>6</sub>, Copenhagen<sub>6</sub>  
Prague<sub>12</sub>, Stockholm<sub>12</sub>, Frankfurt<sub>14</sub>, Madrid<sub>14</sub>, Manila<sub>14</sub>  
Stockholm<sub>20</sub>, Milan<sub>22</sub>, Madrid<sub>24</sub>, Taipei<sub>24</sub>, Frankfurt<sub>25</sub>

Closest based on 32 bit sig.'s

**Cheap**

# Points to review

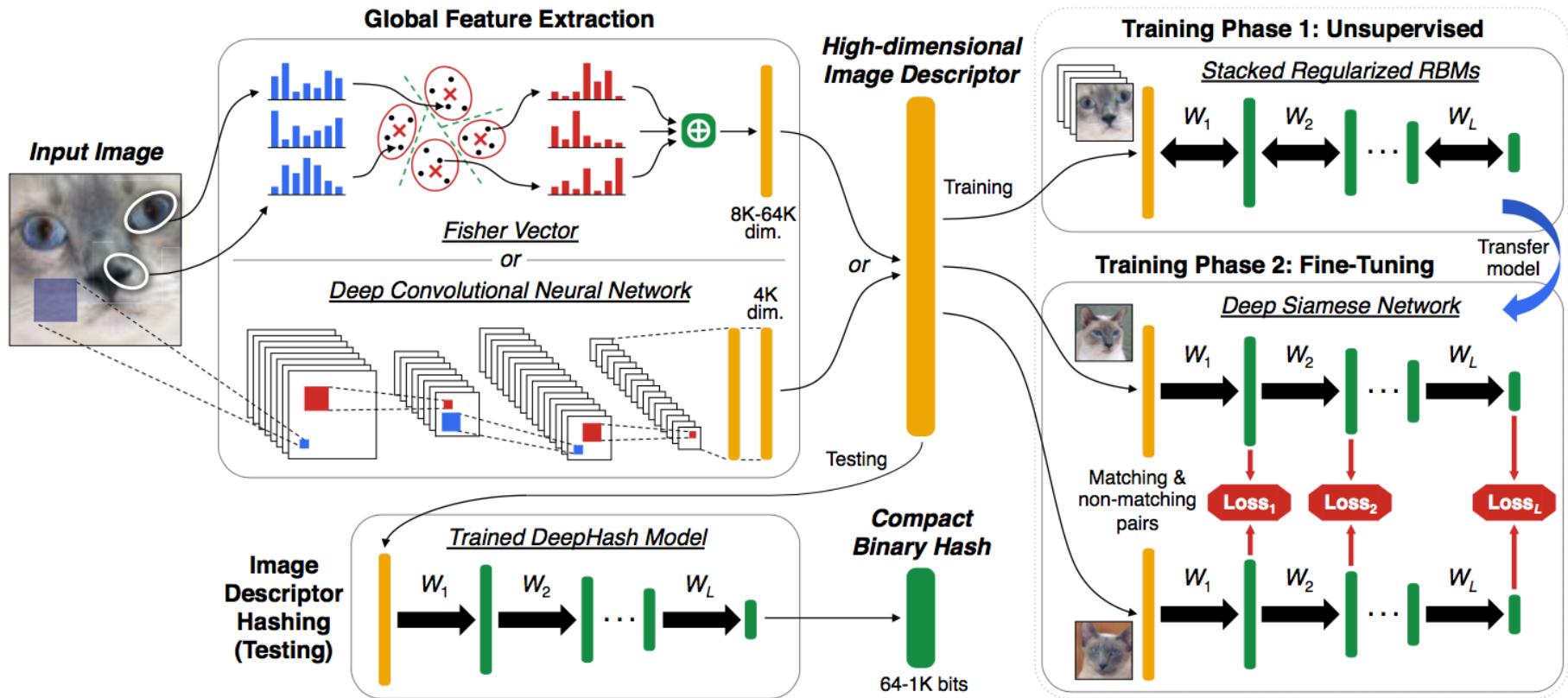
- APIs for:
  - Bloom filters, CM sketch, LSH
- Key applications of:
  - Very compact noisy sets
  - Efficient counters accurate for *large* counts
  - Fast approximate cosine distance
- Key ideas:
  - Uses of hashing that allow collisions
  - Random projection
  - Multiple hashes to control  $\text{Pr}(\text{collision})$
  - Pooling to compress a lot of random draws

# A DEEP-LEARNING VARIANT OF LSH

# DeepHash: Getting Regularization, Depth and Fine-Tuning Right

Jie Lin<sup>\*,1,3</sup>, Olivier Morère<sup>\*,1,2,3</sup>, Vijay Chandrasekhar<sup>1,3</sup>, Antoine Veillard<sup>2,3</sup>, Hanlin Goh<sup>1,3</sup>  
I2R<sup>1</sup>, UPMC<sup>2</sup>, IPAL<sup>3</sup>

ICMR 2017



# DeepHash

Image

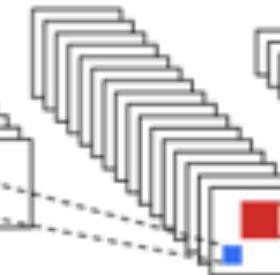
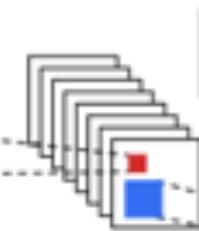
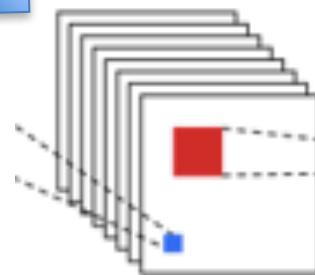


Compact Bit Vector  
64-1000 bits long

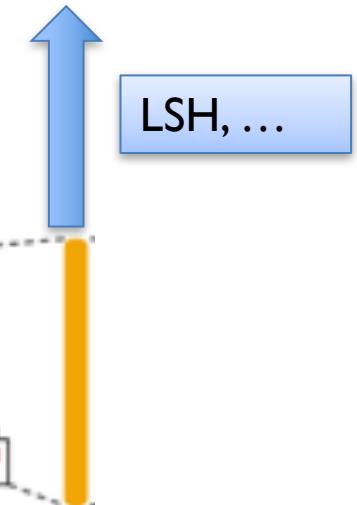


# DeepHash

Image



Compact Bit Vector  
64-1000 bits long



LSH, ...

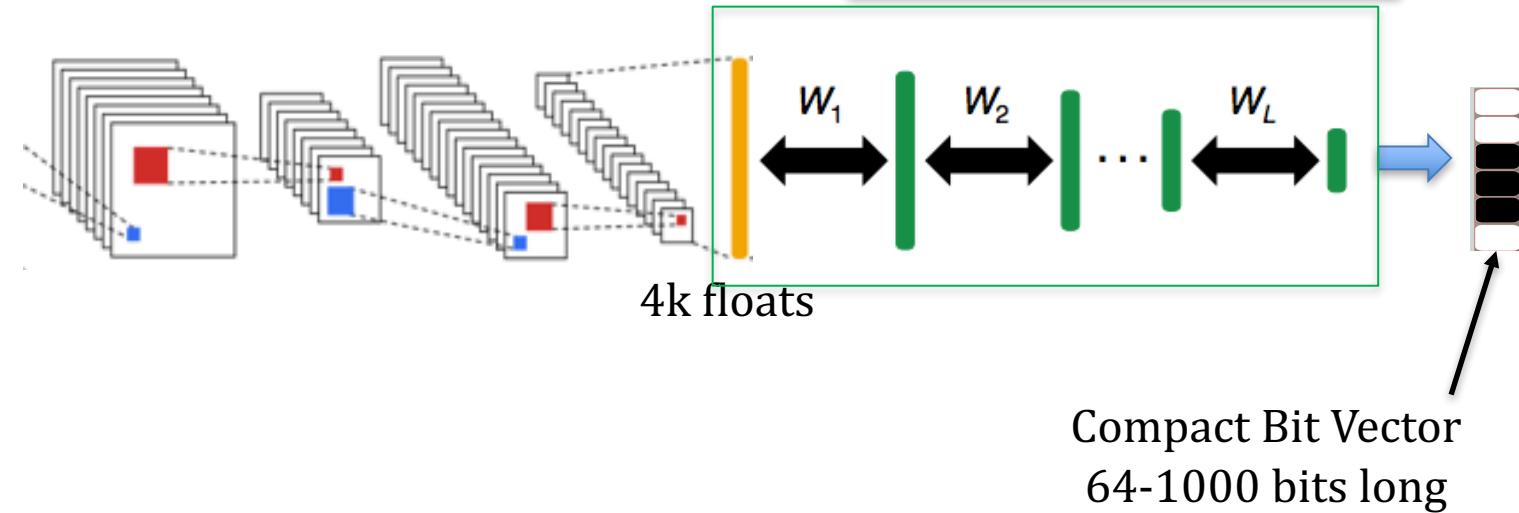
4k floats

# DeepHash

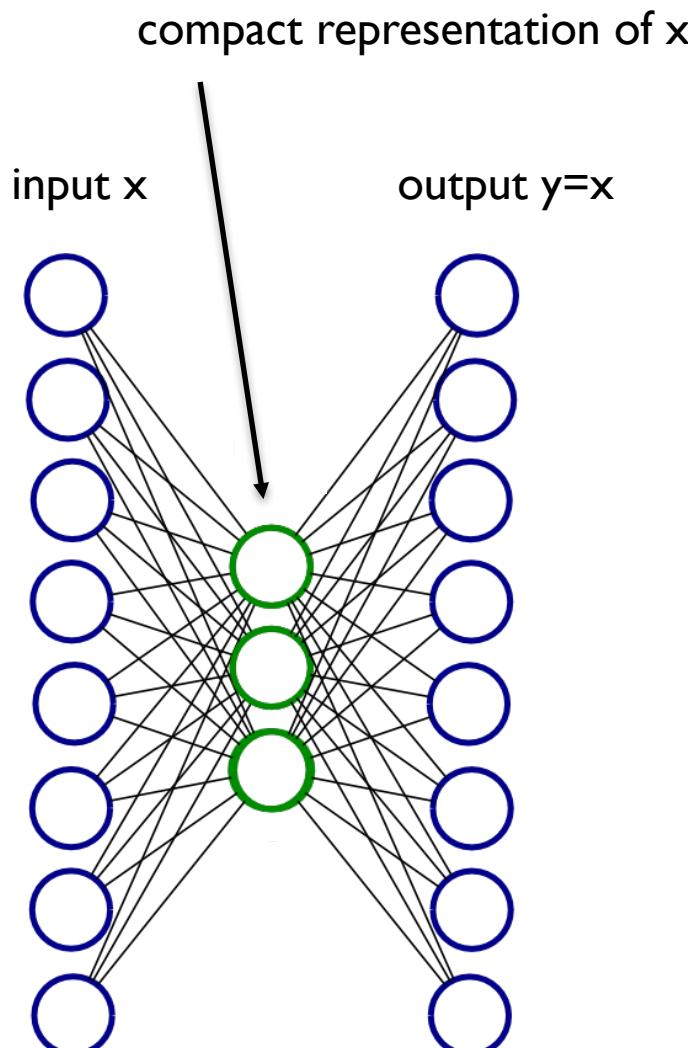
Image



Deep Restricted  
Boltzmann Machine

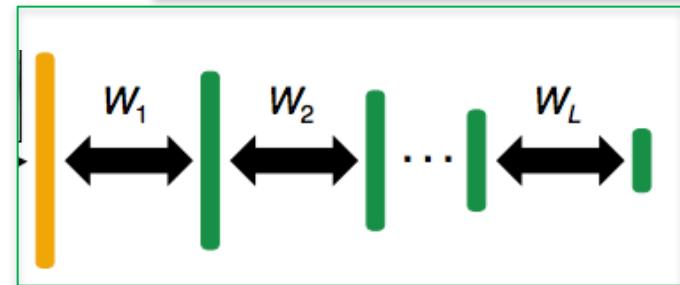


# DeepHash



An autoencoder

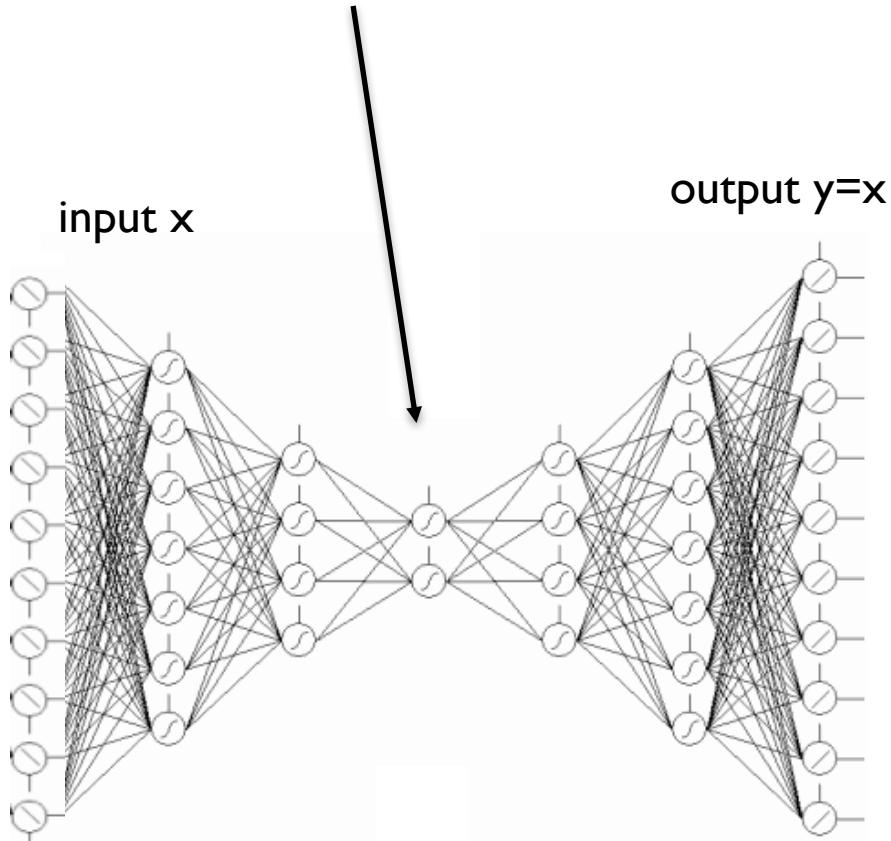
Deep Restricted Boltzmann Machine



Restricted Boltzmann Machine  
is closely related to  
an autoencoder

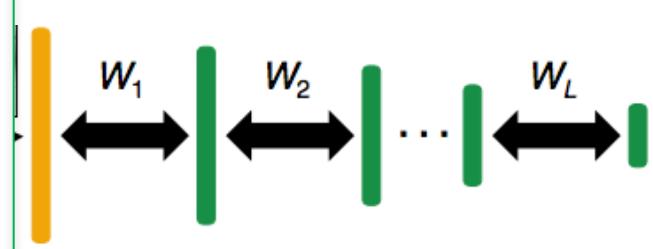
# DeepHash

compact representation of  $x$



A deeper autoencoder

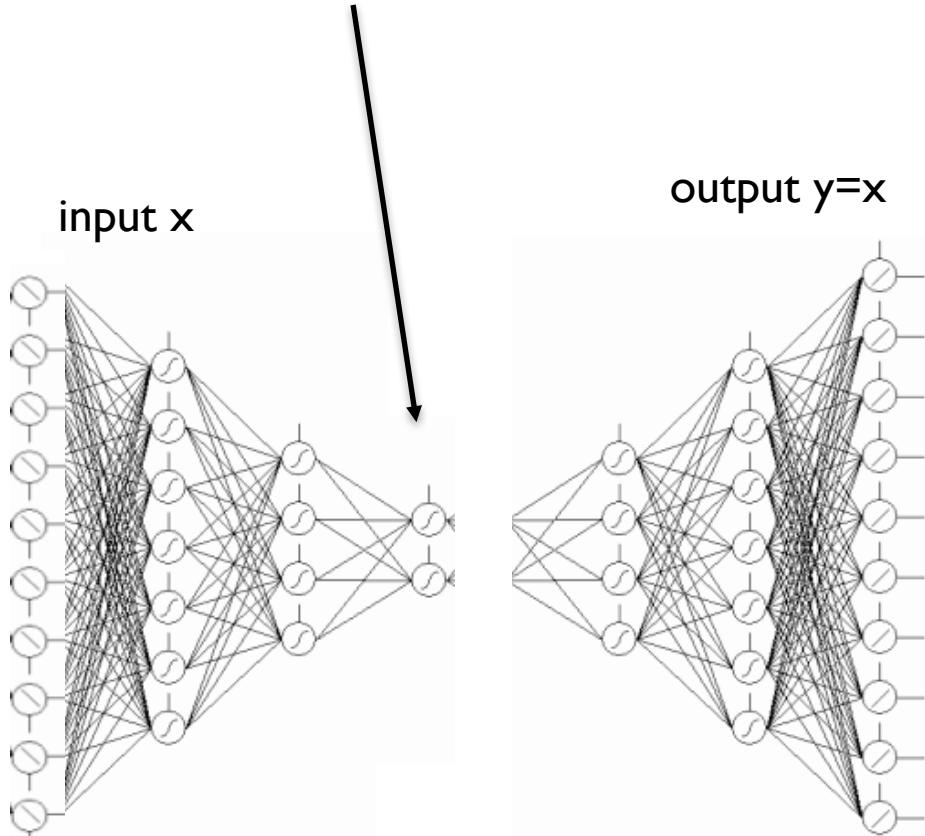
Deep Restricted Boltzmann Machine



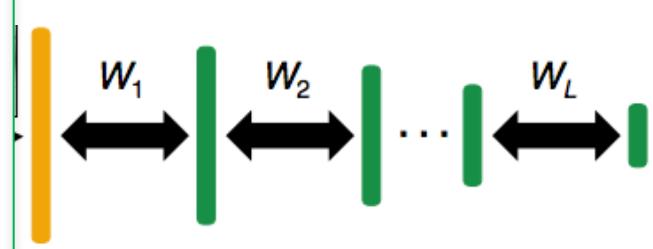
Deep Restricted Boltzmann Machine is closely related to a deep autoencoder

# DeepHash

compact representation of  $x$



Deep Restricted  
Boltzmann Machine



Deep Restricted  
Boltzmann Machine is

But the RBM is symmetric:  
weights  $W$  from layer  $j$  to  
 $j+1$  are the transpose of  
weights from  $j+1$  back to  $j$

RBM is also stochastic: compute  $\text{Pr}(\text{hidden}|\text{visible})$ , sample  
from that distribution, compute  $\text{Pr}(\text{visible}|\text{hidden})$ , sample, ...  
53

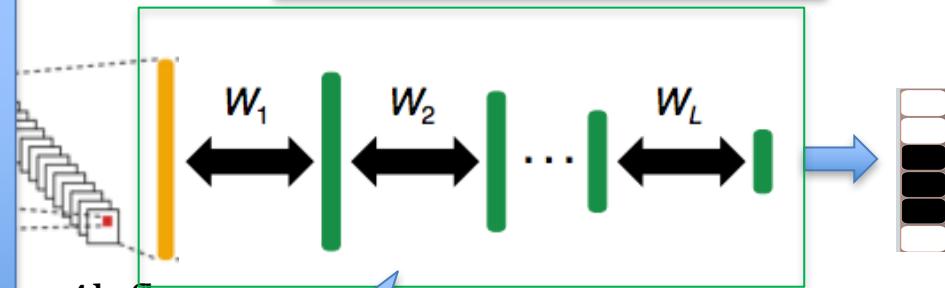
# DeepHash

Image

Another trick: regularize so that the representations are dense (about 50-50 “on” and “off” for an image ) and each bit has 50-50 chance of being “on”

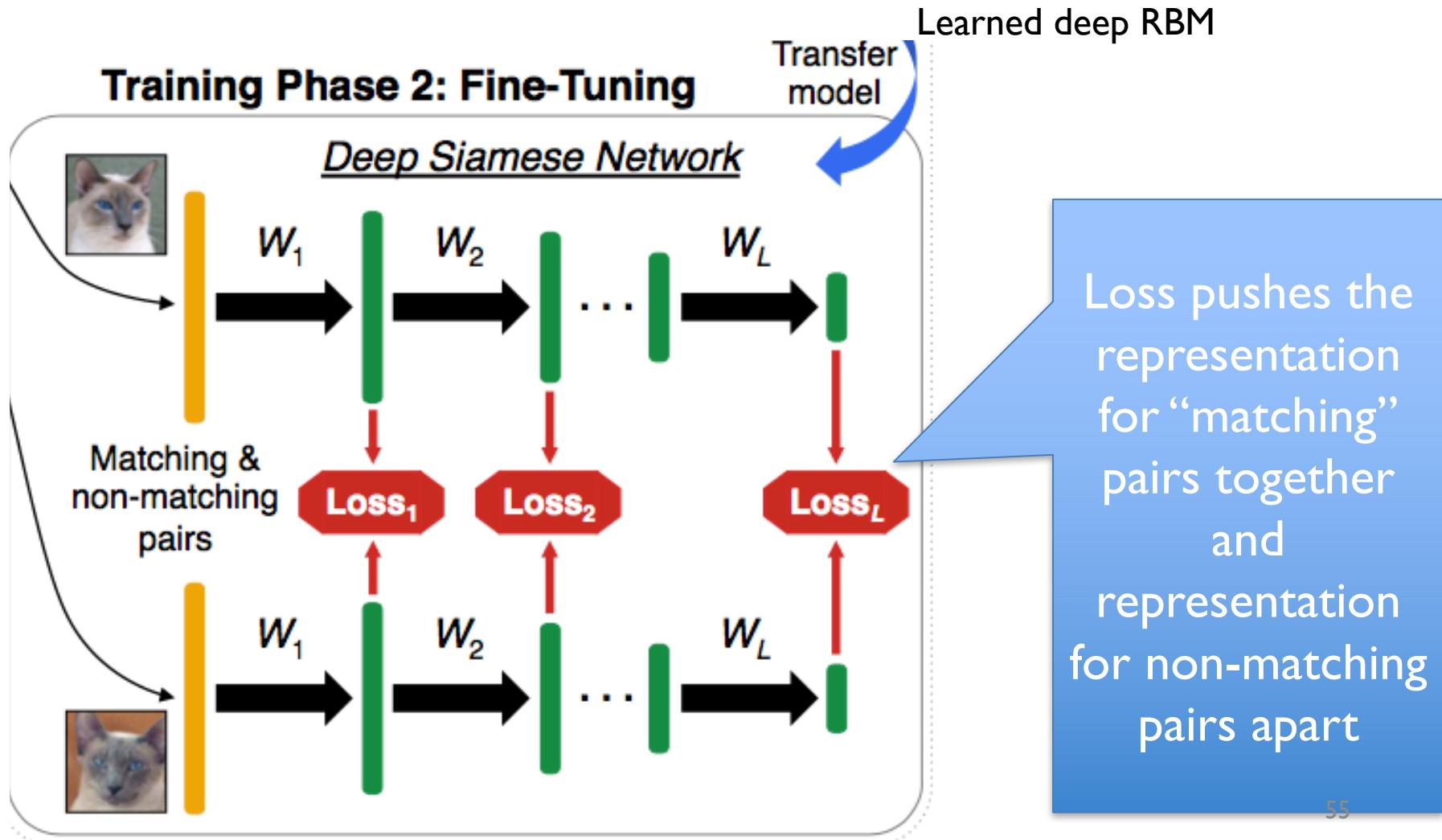
And then more training....

Deep Restricted Boltzmann Machine

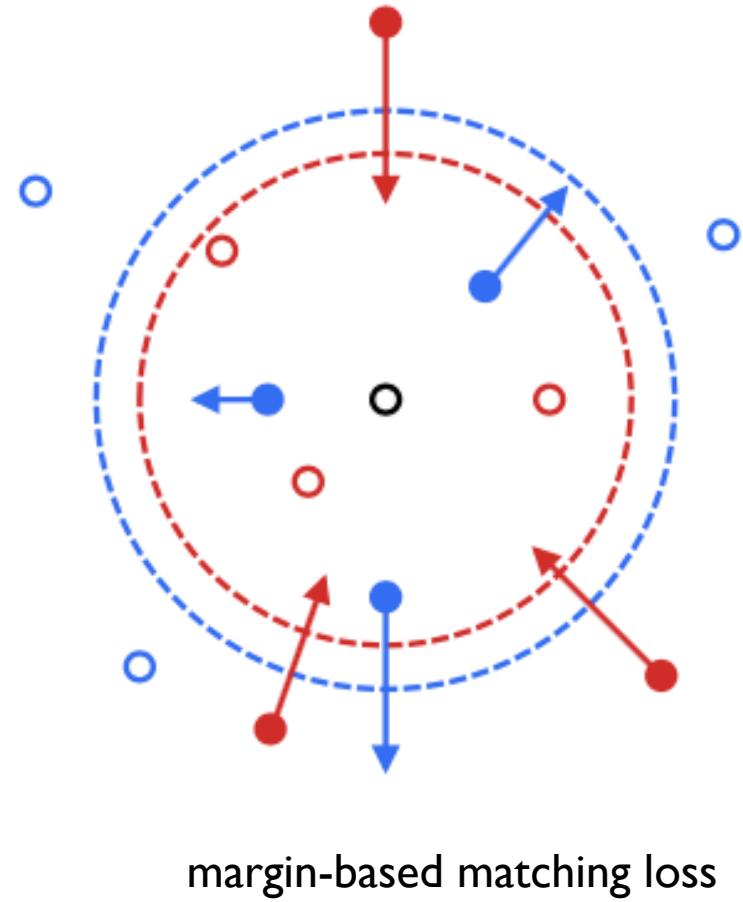
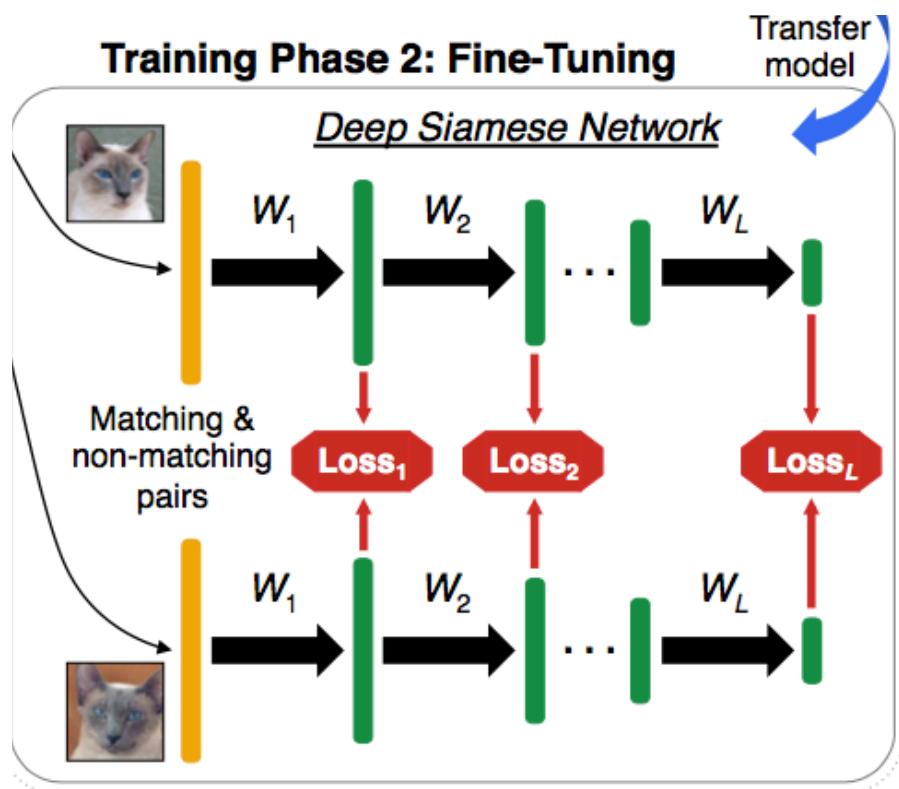


This model is trained to compress image features, then reconstruct the image features from the reconstructions.

# Training on matching vs non-matching pairs



# Training on matching vs non-matching pairs



# Training on matching vs non-matching pairs

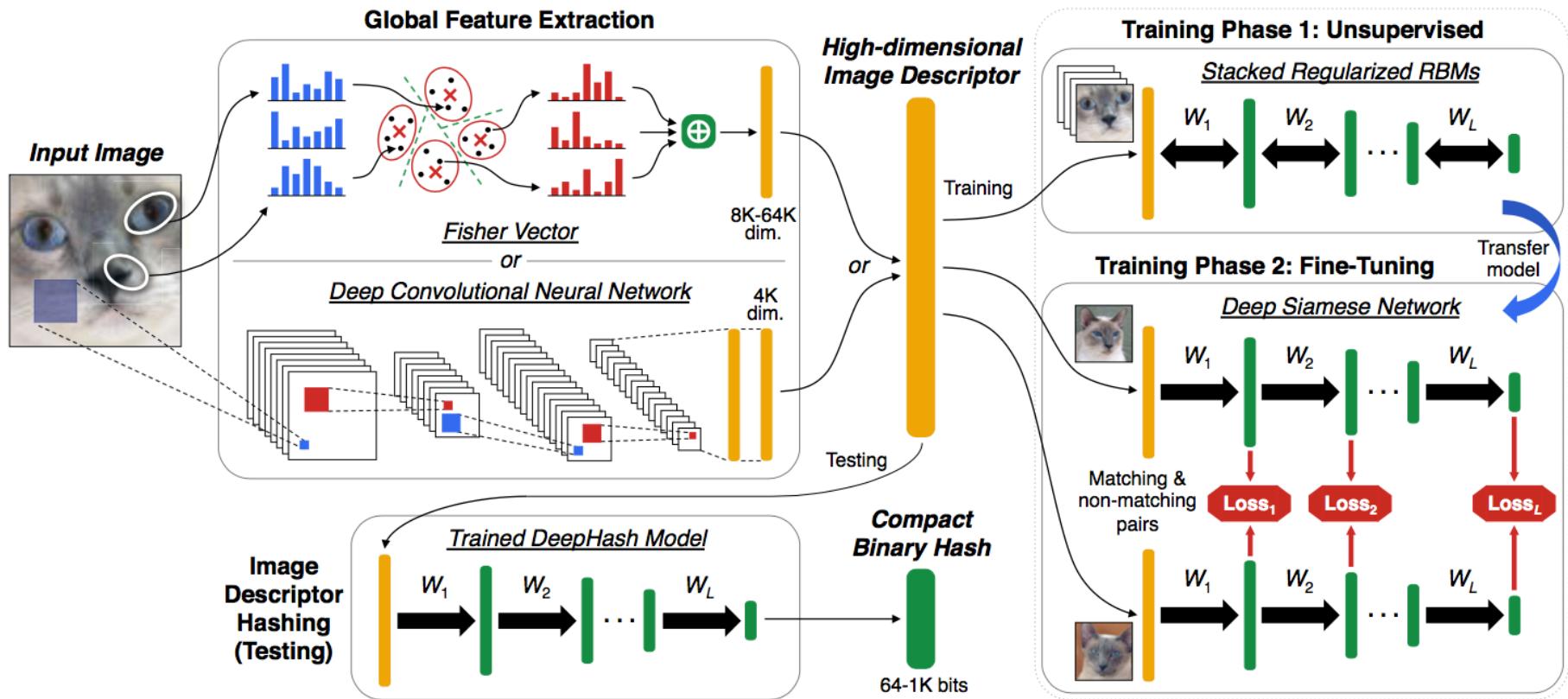
matching pairs were photos of the same “landmark”



# DeepHash: Getting Regularization, Depth and Fine-Tuning Right

Jie Lin<sup>\*,1,3</sup>, Olivier Morère<sup>\*,1,2,3</sup>, Vijay Chandrasekhar<sup>1,3</sup>, Antoine Veillard<sup>2,3</sup>, Hanlin Goh<sup>1,3</sup>  
I2R<sup>1</sup>, UPMC<sup>2</sup>, IPAL<sup>3</sup>

ICML 2017



# DeepHash: Getting Regularization, Depth and Fine-Tuning Right

Jie Lin<sup>\*,1,3</sup>, Olivier Morère<sup>\*,1,2,3</sup>, Vijay Chandrasekhar<sup>1,3</sup>, Antoine Veillard<sup>2,3</sup>, Hanlin Goh<sup>1,3</sup>  
I2R<sup>1</sup>, UPMC<sup>2</sup>, IPAL<sup>3</sup>

ICML 2017

