

N L P R

多媒体内容索引和排序

程 健 研究员

jcheng@nlpr.ia.ac.cn

中国科学院自动化研究所

2022. 11. 7





大纲

- 背景介绍
- 多媒体内容索引
- 多媒体内容排序
- 近似近邻搜索
- 总结与展望

多媒体内容索引与排序



浩如烟海
无从下手

唯有索引
方得真经





多媒体内容索引与排序



中国科学院统一自动化系统

登录 检索首页 我的图书馆 参数设置 结果列表 最近检索 读

所有字段 三国演义

检索 二次检索

排序: 年(降序)/著者 格式: 封面视图 选中记录 整合集合 重新查询 分类浏览 标签浏览

记录 1 - 10 of 186 (最大显示记录 20000 条)

1 2 3 4 5 6 7 8 9 10 下一页>

1 ☐



《三国演义》试论 SFX

作者: 董每戡
出版社: 北京出版社
格式: 图书
评级: ☆☆☆

索书号: I207.413/31
年份: 2020
链接:

馆藏复本: 2, 已出借复本: 0

2 ☐



三国演义 SFX

作者: 罗贯中
出版社: 人民文学出版社
格式: 图书
评级: ☆☆☆

索书号: I242.43/94/1
年份: 2020
链接:

馆藏复本: 6, 已出借复本: 0

3 ☐



《三国志演义》史话 SFX

作者: 陈翔华
出版社: 国家图书馆出版社
格式: 图书
评级: ☆☆☆

索书号: I207.413/30
年份: 2019
链接: 在线查询
url

馆藏复本: 0, 已出借复本: 0

4 ☐



凯叔三国演义. 三分天下 SFX

作者: 凯叔
出版社: 湖南少年儿童出版社
格式: 图书
评级: ☆☆☆

索书号: 22.224/K365/42001CB1409468
年份: 2019
链接:

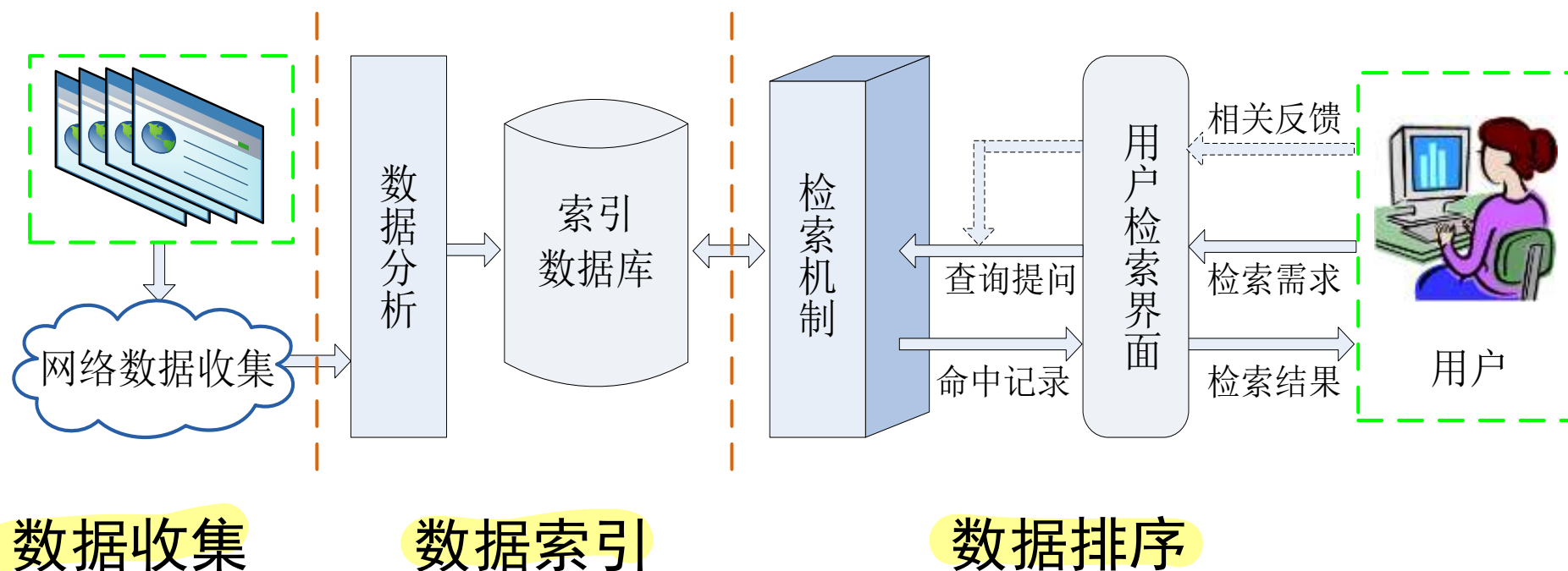
馆藏复本: 4, 已出借复本: 0

5 ☐



凯叔三国演义. 三分天下 SFX

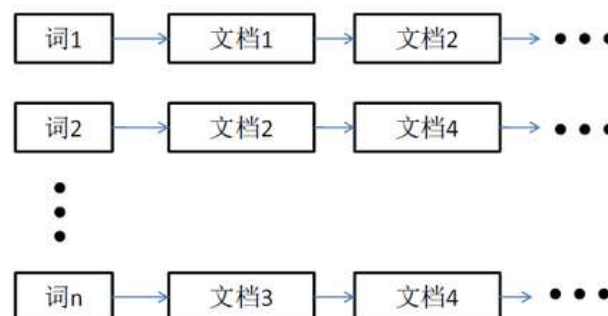
网络多媒体检索的基本流程



数据索引

- 索引信息获取
 - 网页内容解析
 - 基于学习算法的自动标注

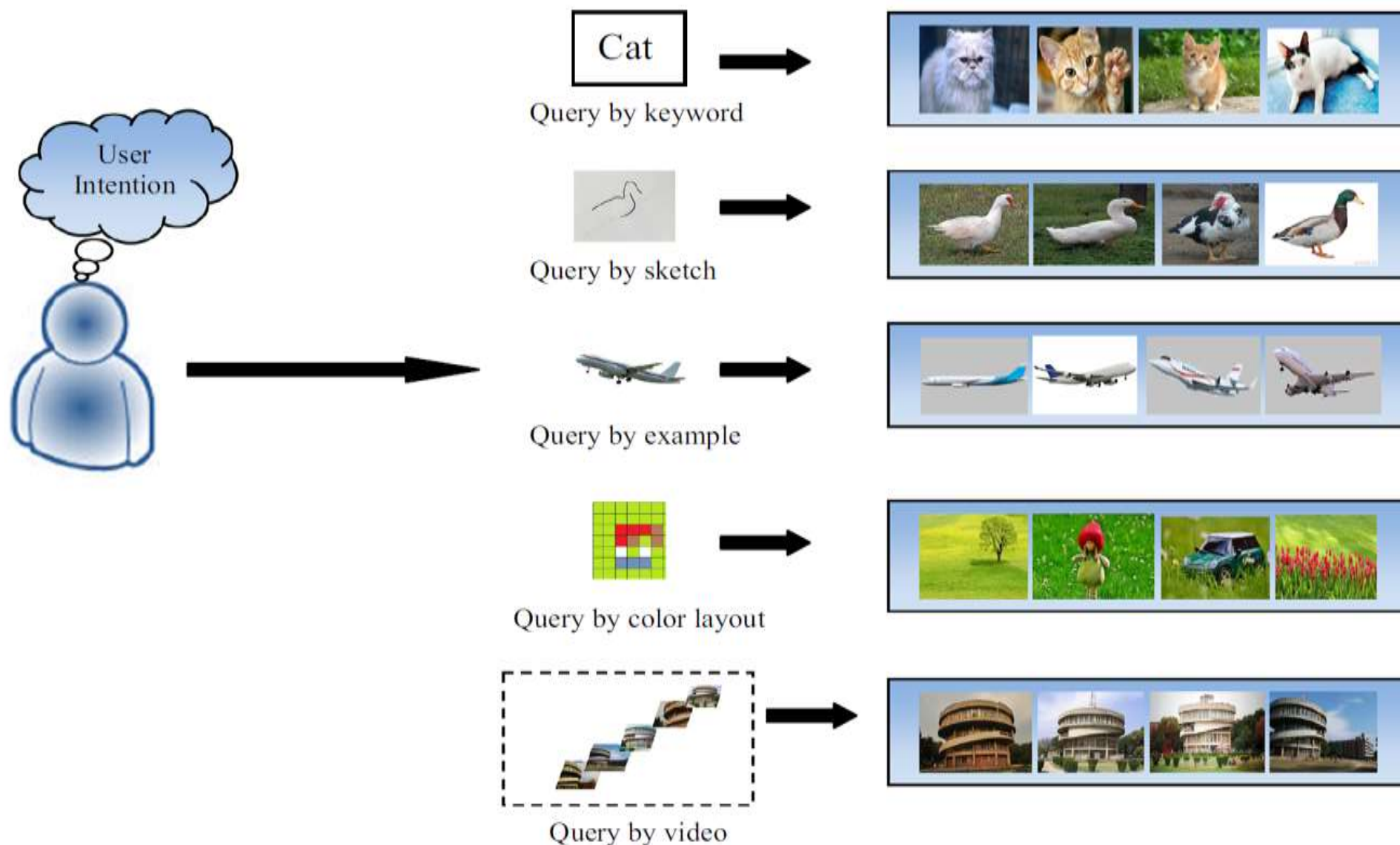
- 建立索引文件
 - 索引倒排表



```
T0 = "it is what it is"  
T1 = "what is it"  
T2 = "it is a banana"  
  
"a": {2}  
"banana": {2}  
"is": {0, 1, 2}  
"it": {0, 1, 2}  
"what": {0, 1}
```

- 关键问题
 - 索引信息的准确性
 - 索引结构的高效性 vs. 有限的内存而巨大的数据量

网络多媒体检索的基本流程



评测指标

- 准确率 (Precision)
- 召回率 (Recall)
- F1 Score
- 平均准确率均值 (Mean Average Precision, MAP)
- NDCG (Normalized discounted cumulative gain)



评测指标

- True Positive : 做出Positive的判定, 而且判定是正确的
- False Positive : 做出Positive的判定, 而且判定是错误的
- True Negative : 正确的Negative判定
- False Negative: 错误的Negative判定

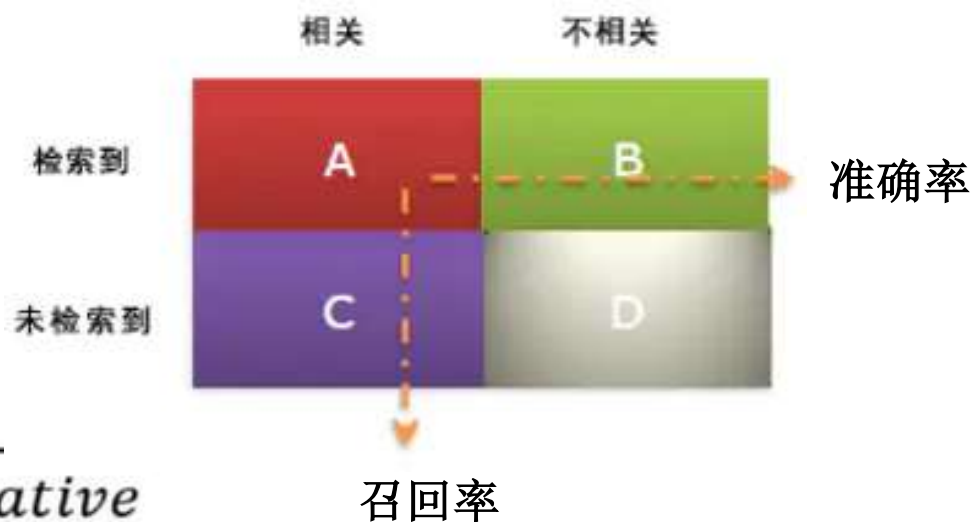
评测指标

➤ 准确率 (Precision)

$$\text{precision} = \frac{\text{true positive}}{\text{true positive} + \text{false positive}}$$

➤ 召回率 (Recall)

$$\text{recall} = \frac{\text{true positive}}{\text{true positive} + \text{false negative}}$$



评测指标

➤ F1 Score

$$F1Score = 2 \times \frac{precision \times recall}{precision + recall}$$

➤ 平均准确率均值 (Mean Average Precision, MAP)

$$AP = \frac{1}{N} \times \sum_{i=1}^N \frac{i}{position(i)}$$

其中，N 表示相关文档总数，position(i) 表示第 i 个相关文档在检索结果列表中的位置。

MAP (Mean Average Precision) 即多个查询的平均正确率 (AP) 的均值，从整体上反映模型的检索性能。

评测指标

- 归一化折损累计增益 (NDCG, Normalized discounted cumulative gain)

$$DCG_p = \sum_{i=1}^p \frac{rel_i}{\log_2(i+1)} = rel_1 + \sum_{i=2}^p \frac{rel_i}{\log_2(i+1)}$$

或: $DCG_p = \sum_{i=1}^p \frac{2^{rel_i} - 1}{\log_2(i+1)}$ 增加相关度比重

$$NDCG_p = \frac{DCG_p}{IDCG_p}$$

← Ideal DCG: 理想中最大的DCG值

大纲

- 背景介绍
- **多媒体内容索引**
- 多媒体内容排序
- 近似近邻搜索
- 总结与展望

索引方法分类

- 传统方法 (NoAI-based)

Tree-based, Hash, Inverted Index...

- 智能方法 (AI-based)

Soft Computing, Machine Learning, KRR...

- 协同智能方法 (Collaborative AI-based)

Collaborative Filtering...



大纲

- 背景介绍
- 多媒体内容索引
- **多媒体内容排序**
- 近似近邻搜索
- 总结与展望

搜索排序

Google



PageRank

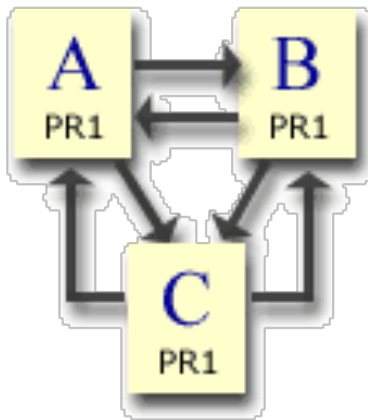
- 如果存在一个网页，它被许多其它链接链接到，则说明这个网页比较重要，则此网页的PageRank值比较高。
- 如果存在一个网页，它本身的PageRank值比较高，且此网页又链接了一个网页，则这个被链接的网页比较重要，其PageRank值较高。

• PageRank的基本公式

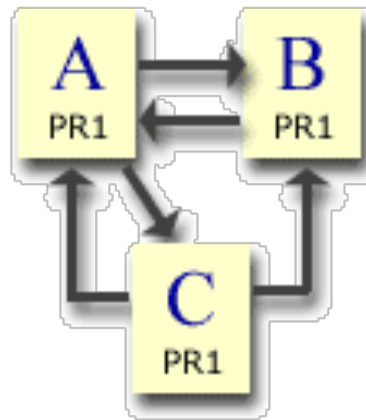
$$r(p) = \alpha \times \sum_{q:(q,p) \in \varepsilon} \frac{r(q)}{\omega(q)} + (1 - \alpha) \times \frac{1}{N}$$

- q – 网页 p 的后向链接, i.e., q 指向 p
- $\omega(q)$ – 网页 q 的前向链接数目.
- $r(q)$ – 网页 q 的PageRank.
- N – 整个网络中网页的总数

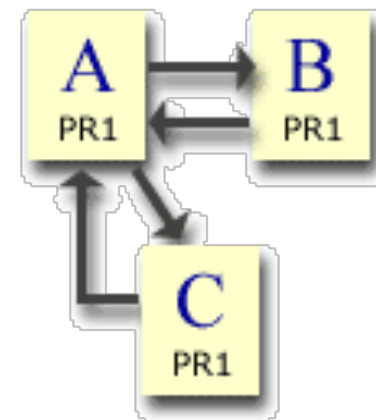
PageRank



Page A = 1
Page B = 1
Page C = 1



Page A = 1.425
Page B = 1
Page C = 0.575



Page A = 1.4594
Page B = 0.7702
Page C = 0.7702

PageRank



优点：与查询无关的静态算法；

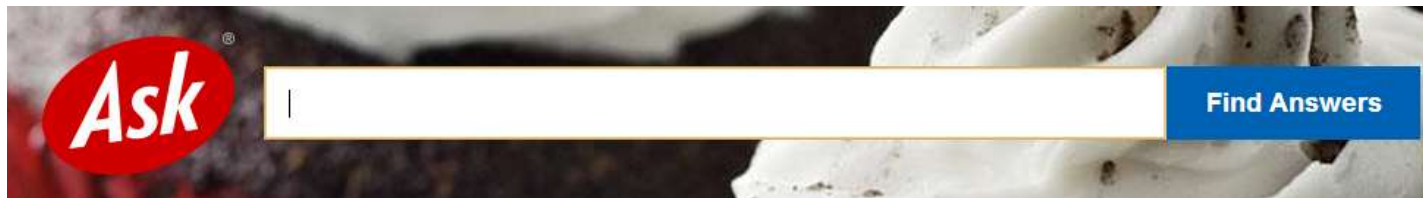
缺点：与主题无关，旧网页比新网页排名高

HITS: Hyperlink Induced Topic Search



基本假设:

- ✓ 一个好的“Authority”页面会被很多好的“Hub”页面指向;
- ✓ 一个好的“Hub”页面会指向很多好的“Authority”页面;



优点: 自然语言、社交网络取得很好效果;

缺点: 计算复杂、主题漂移、易作弊



PageRank vs. HITS

➤ PageRank:

Query-independent, 静态/离线, 全局, 不容易作弊

➤ HITS:

Query-dependent, 在线, 局部（只与Query相关）, 容易作弊

排序学习 (Learning to Rank)

- 单文档方法 (PointWise Rank)
- 文档对方法 (PairWise Rank)
- 文档列表方法 (ListWise Rank)

PointWise Rank

➤ **Classification:**

Discriminative model for IR (SIGIR 2004),

McRank (NIPS 2007)

➤ **Regression:**

Least Square Retrieval Function (TOIS 1989),

Subset Ranking using Regression (COLT 2006)

➤ **Ordinal Classification:**

Pranking (NIPS 2002)

Constraint Ordinal Regression (ICML 2005)

优点：速度快，复杂度低；

缺点：没有考虑文档间的关系，效果一般

PairWise Rank

➤ PairWise Rank:

Ranking SVM (ICANN 1999),

RankBoost (JMLR 2003),

RankNet (ICML 2005),

LambdaRank (NIPS 2006)

.....

优点: 将排序问题转为了文档对顺序的判断;

缺点: 没有考虑文档出现在搜索列表中的位置; 不同的查询, 其相关文档数量差异很大。

ListWise Rank

➤ **Measure-specific:**

AdaRank (SIGIR 2007)

SoftRank (LR4IR 2007)

LambdaMART (inf.retr 2010)

➤ **Non-Measure specific:**

ListNet (ICML 2007),

ListMLE (ICML 2008),

BoltzRank (ICML 2009)

优点：速度快，复杂度低；

?

缺点：没有考虑文档间的关系，效果一般

排序学习 (Learning to Rank)

	pointwise	pairwise	listwise
输入信息的完整度	不完全	部分完全	完全
输入	(x, y)	(x_1, x_2, y)	$(x_1, x_2 \dots x_n, \pi)$
输出	$f(x)$	$f(x)$	$f(x)$
样本复杂度	$O(n)$	$O(n^2)$	$O(n!)$
表现	差	中	好

https://blog.csdn.net/anshuai_3w/1

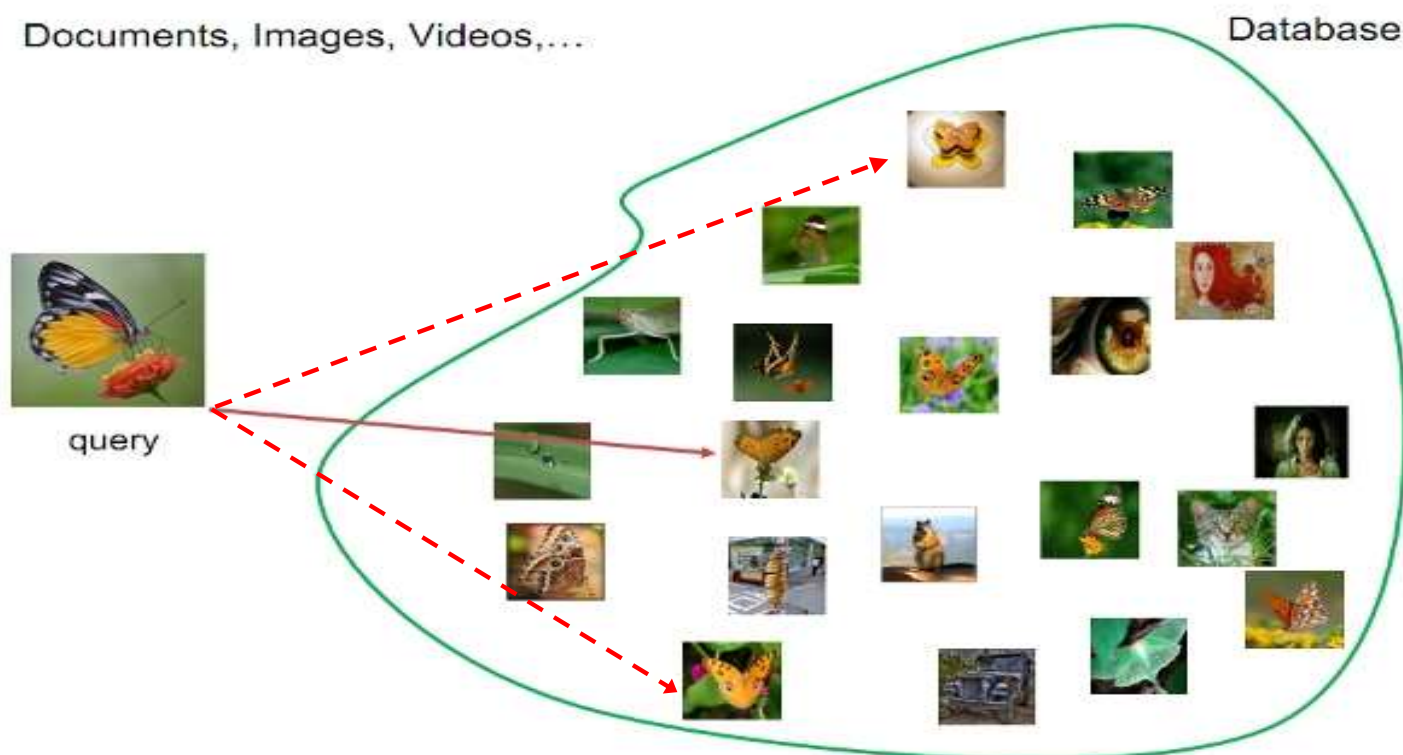


大纲

- 背景介绍
- 多媒体内容索引
- 多媒体内容排序
- **近似近邻搜索**
- 总结与展望

近似近邻搜索

- **近似近邻：检索精度和检索时间的平衡**
- 对大多数应用，近似近邻足够满足需求



基于哈希的表示

How to index data?

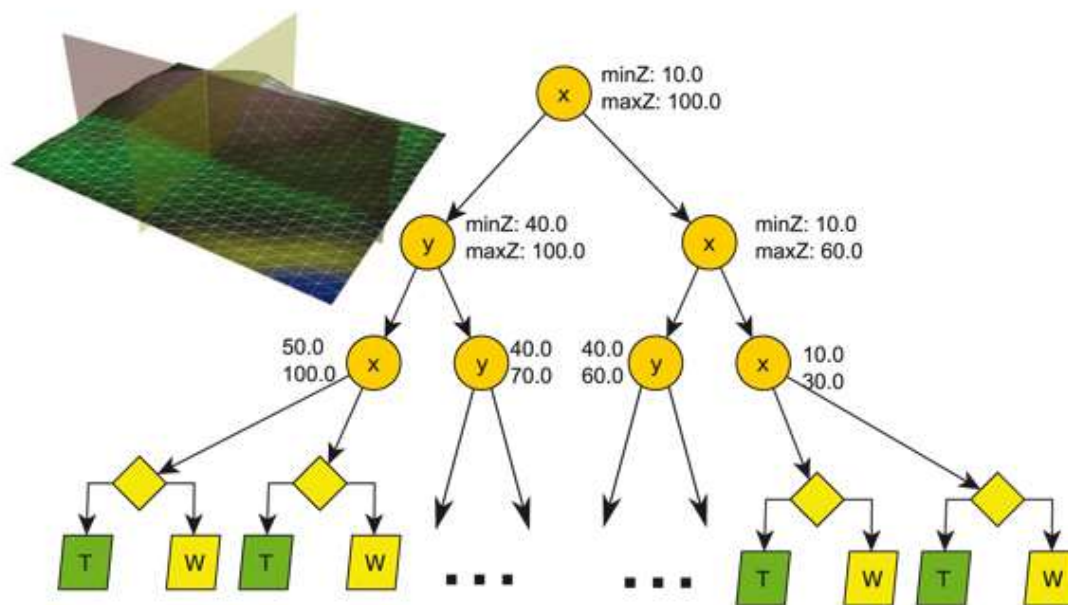
TF-IDF

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

$tf_{i,j}$ = number of occurrences of i in j
 df_i = number of documents containing i
 N = total number of documents

Document analysis

Tree-based

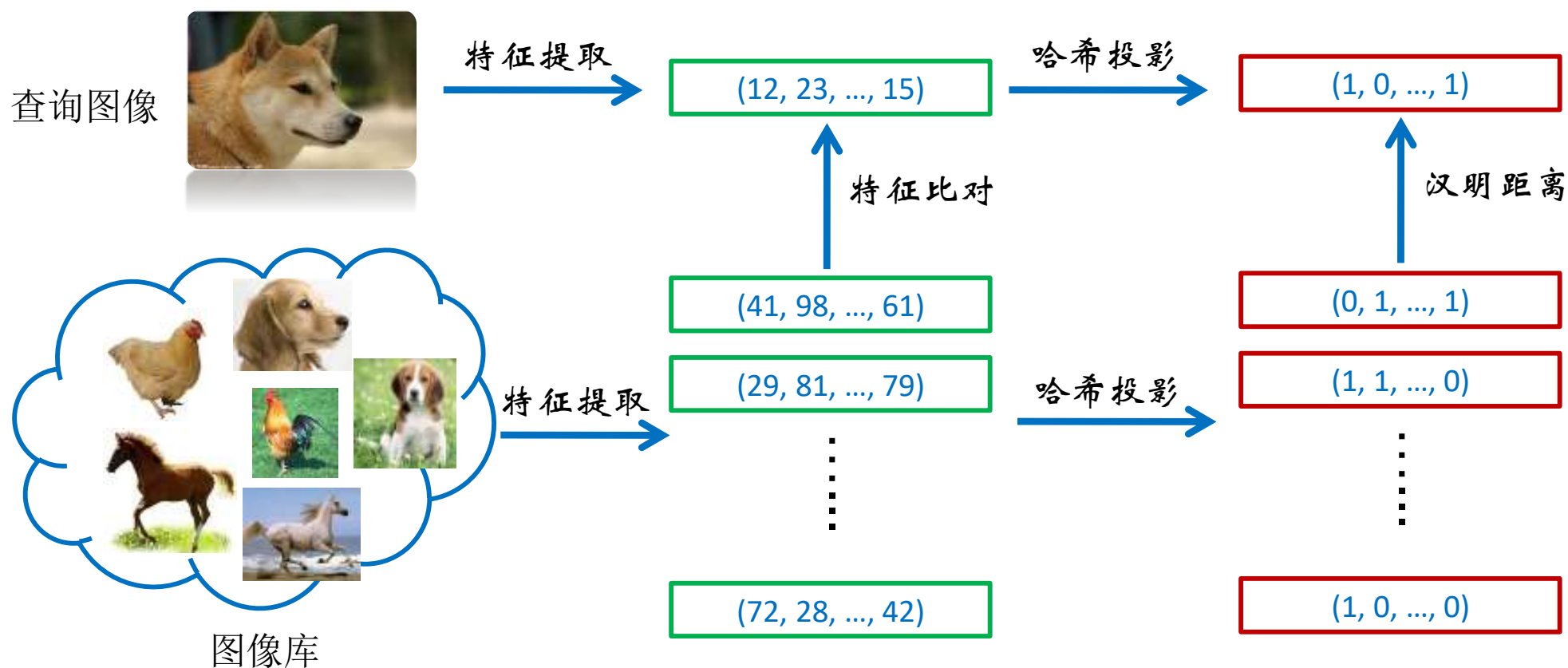


Curse of dimensionality

Pictures from web

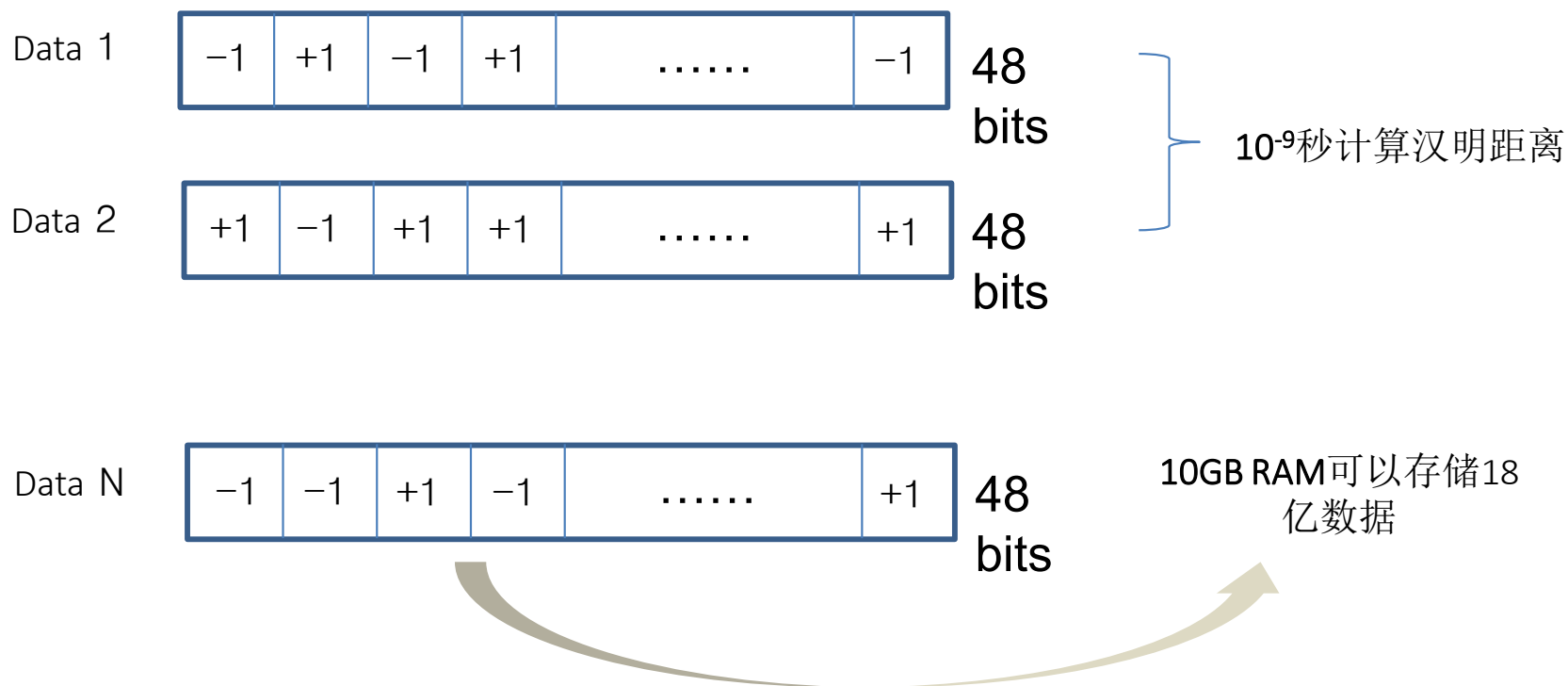
基于哈希的表示

• What is hashing?



基于哈希的表示

- 时间上高效：基于XOR操作的快速计算
- 存储上高效：基于位存储的紧致表达

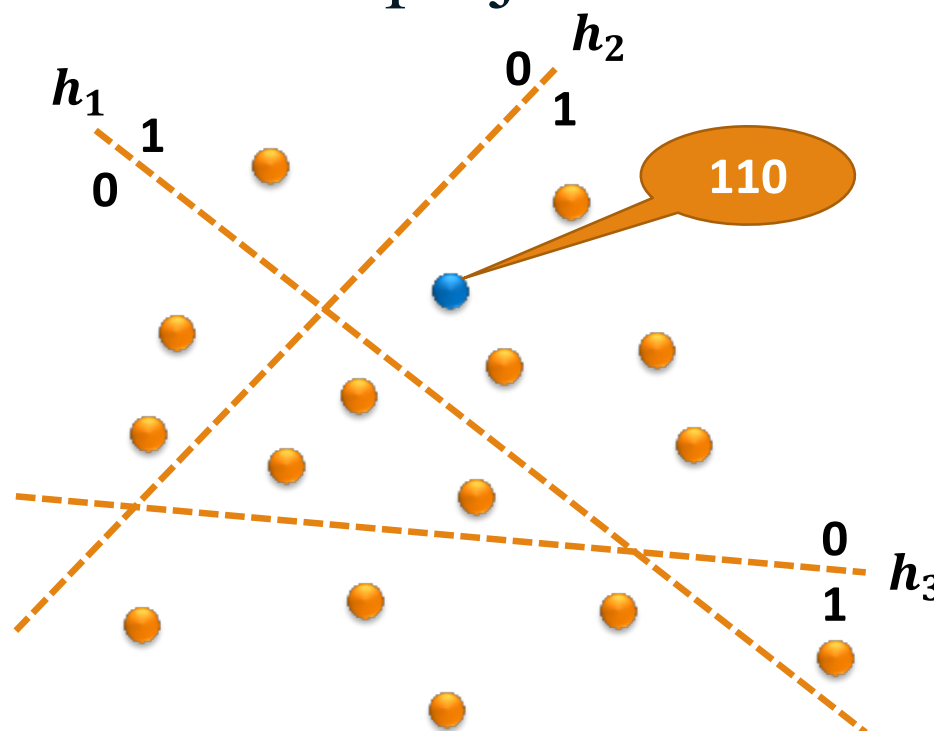


基于哈希的表示

□ Locality Sensitive Hashing (LSH)

- ✓ Date-independent, unsupervised
- ✓ Map the data points with random projections

Hash Function: $B = \text{sgn}(XW)$



LSH与神经科学



REPORT

A neural algorithm for a fundamental computing problem

Sanjoy Dasgupta¹, Charles F. Stevens^{2,3}, Saket Navlakha^{4,*}

+ See all authors and affiliations

Science 10 Nov 2017:
Vol. 358, Issue 6364, pp. 793-796
DOI: 10.1126/science.aam9868

研究发现果蝇的嗅觉环路对相近的嗅觉产生相近的神经活动模式，可将一种味觉习得的行为应用于接触相似味觉时。研究者将其中的三种全新计算策略应用于计算领域，提出一种**局部敏感哈希算法**，该方法可有效改善近似检索的计算表现

LSH与神经科学

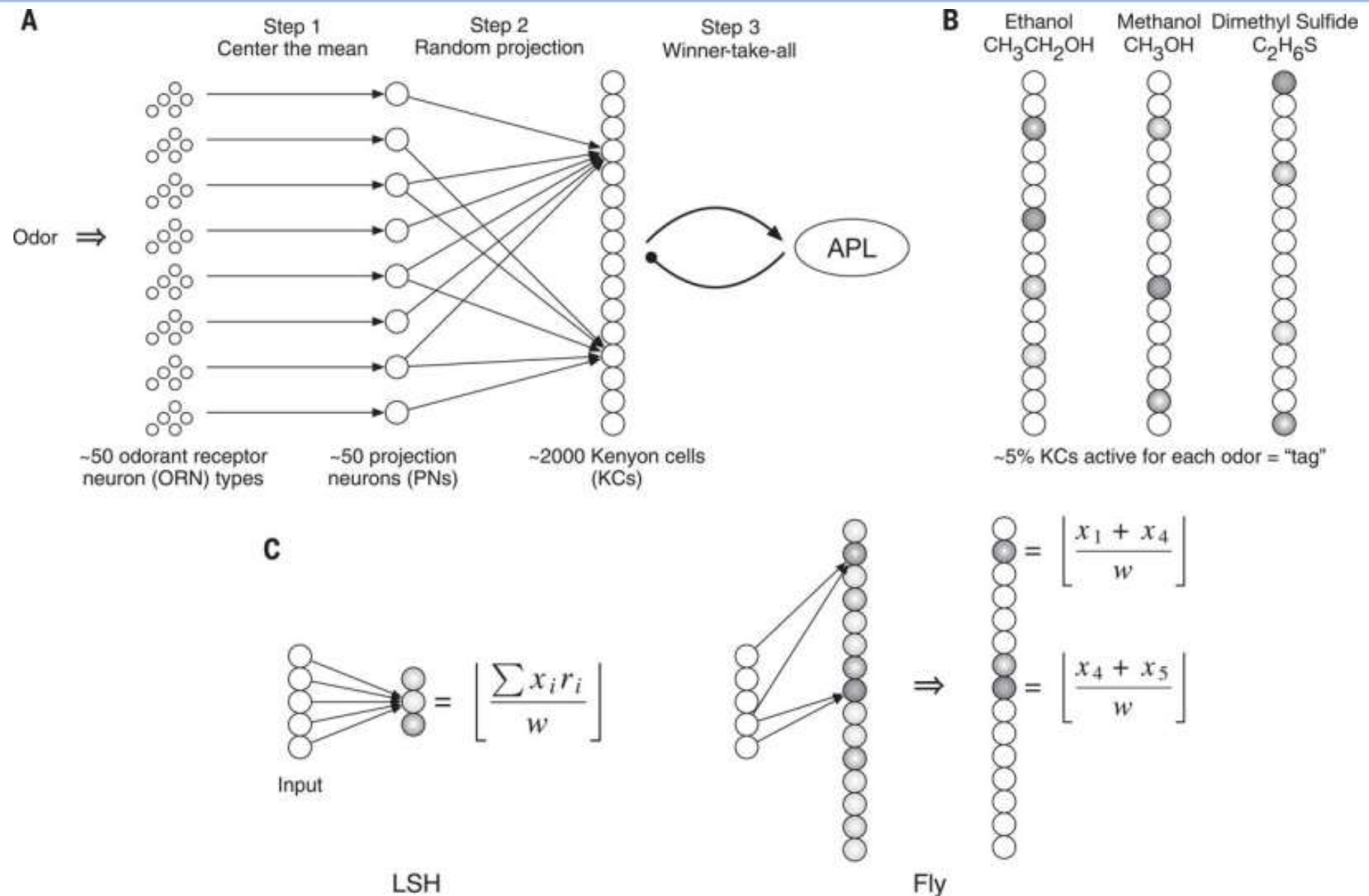


Fig. 1 Mapping between the fly olfactory circuit and locality-sensitive hashing (LSH).

LSH与神经科学

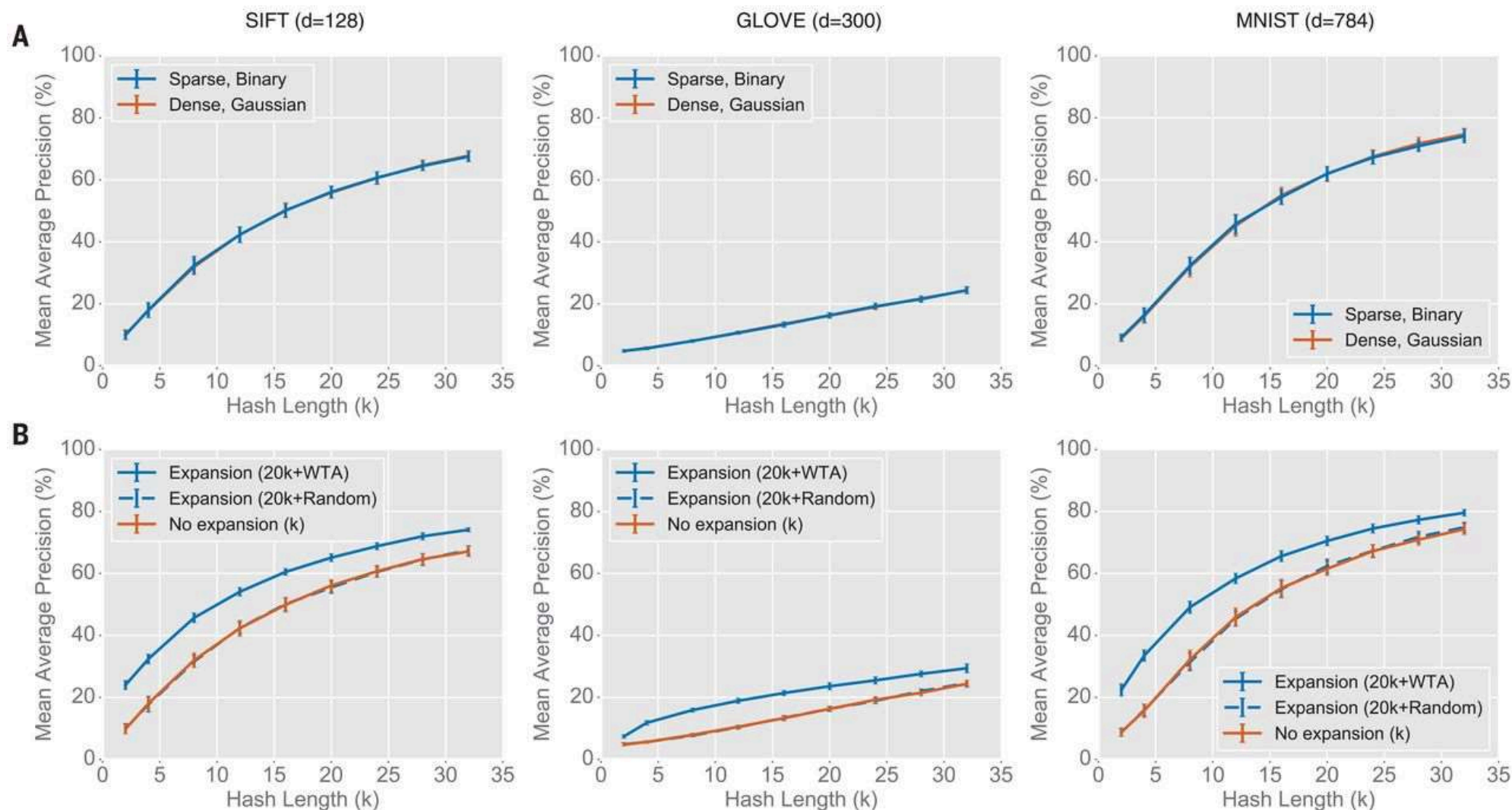


Fig. 2 Empirical comparison of different random projection types and tag-selection methods.

LSH与神经科学

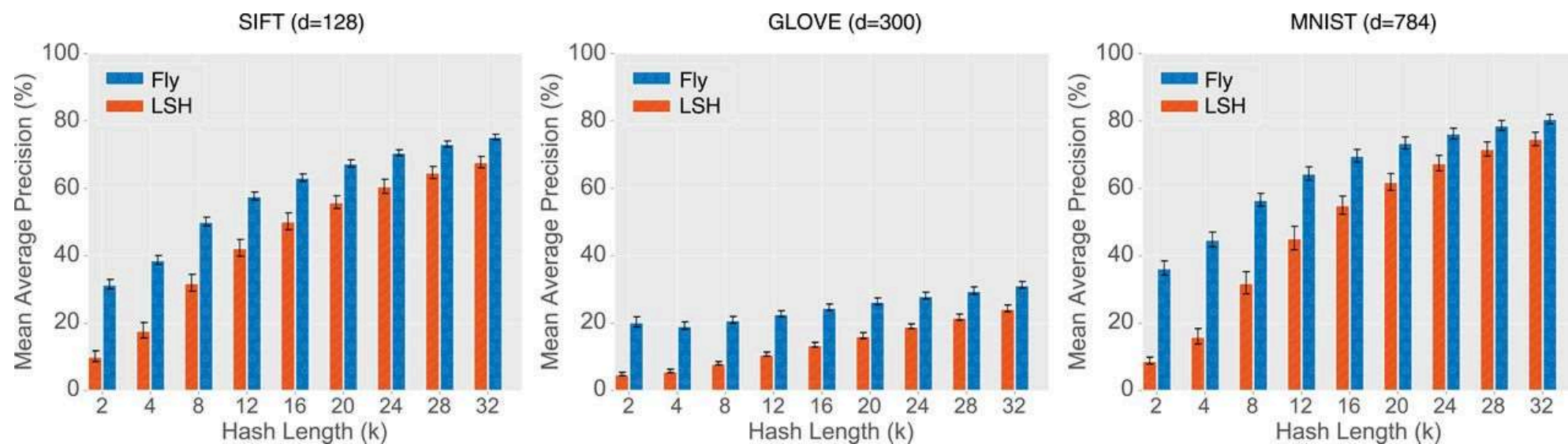
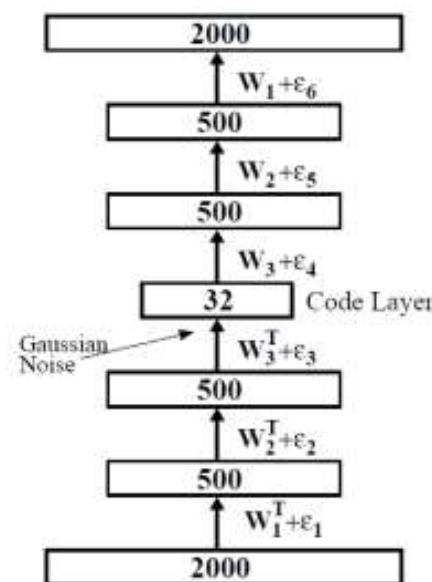
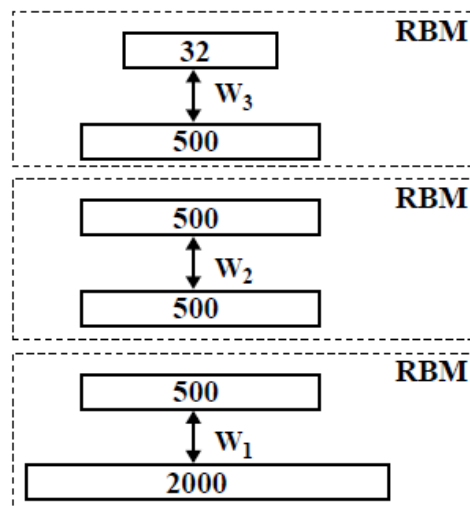


Fig. 3 Overall comparison between the fly algorithm and LSH.

基于哈希的表示

- Restricted Boltzmann Machines (RBMs)
 - ✓ Date-dependent, unsupervised/supervised
 - ✓ Learn the binary codes layer by layer with deep networks

Output layer: final binary codes



基于哈希的表示

□ Spectral Hashing (SH)

- ✓ Date-dependent, unsupervised
- ✓ Design compact binary codes based on spectral graph partitioning

$$\text{minimize : } \sum_{ij} W_{ij} \|y_i - y_j\|^2$$

$$\text{subject to : } y_i \in \{-1, 1\}^k$$

$$\sum_i y_i = 0$$

$$\frac{1}{n} \sum_i y_i y_i^T = I$$



$$\text{minimize : } \text{trace}(Y^T (D - W) Y)$$

$$\text{subject to : } Y(i, j) \in \{-1, 1\}$$

$$Y^T \mathbf{1} = 0$$

$$Y^T Y = I$$

基于哈希的表示

□ Semi-Supervised Hashing (SSH)

- ✓ Date-dependent, semi-supervised
- ✓ Combines the empirical loss over the labeled data with other desirable constraints over both labeled and unlabeled data.

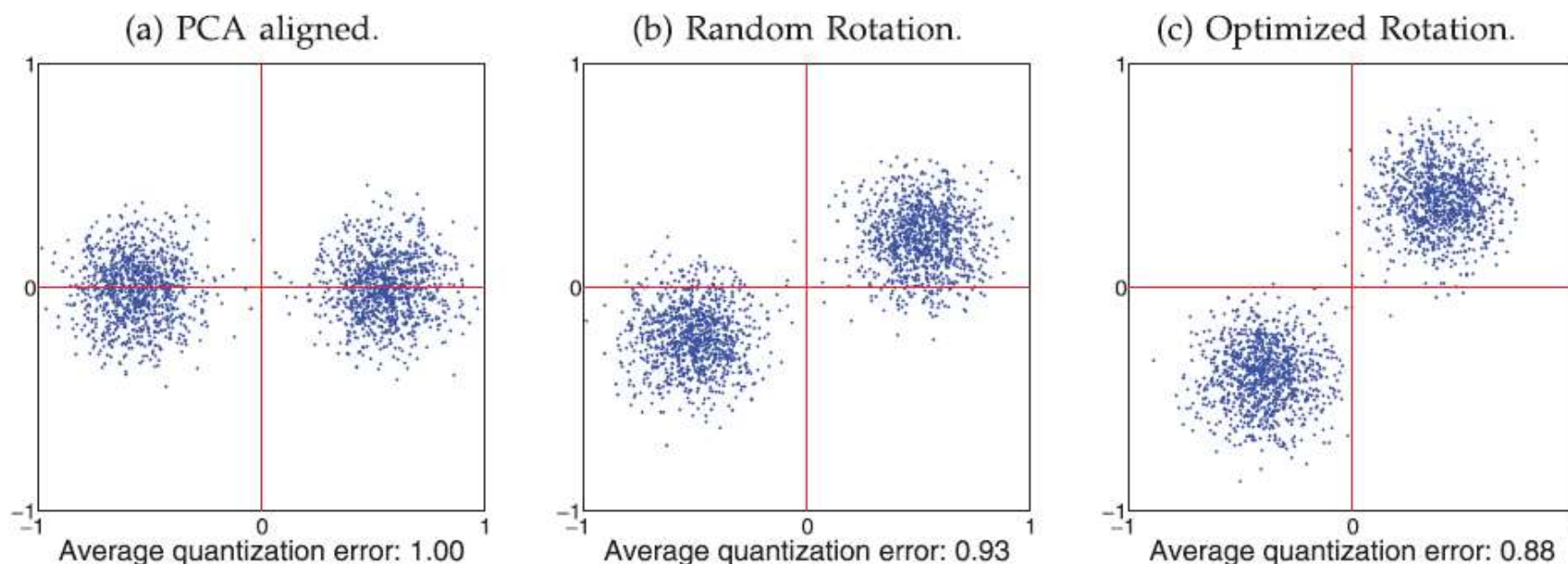
$$\max J(W) = \text{tr} (W^T X_l S X_l^T W) + \eta * \text{tr} (W^T X X^T W)$$

- ✓ A sequential learning scheme (SPLH) is also developed for hash function learning.

基于哈希的表示

- 迭代量化 (Iterative Quantization, ITQ)

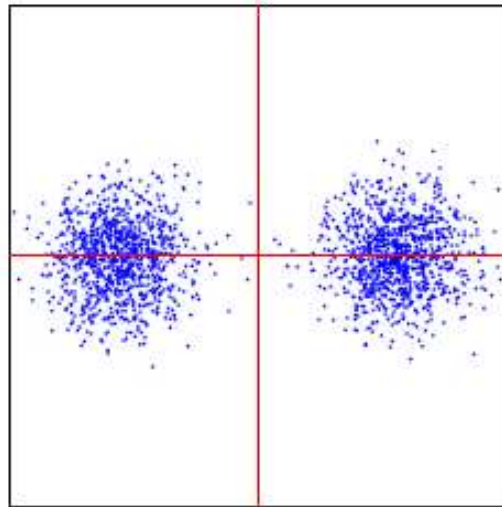
$$\min \|B - XWR\|_F^2, \quad s.t. R^T R = I$$



Yunchao Gong, et.al. Iterative Quantization: A Procrustean Approach to Learning Binary Codes for Large-scale Large Retrieval. CVPR 2011, IEEE TPAMI 2013.

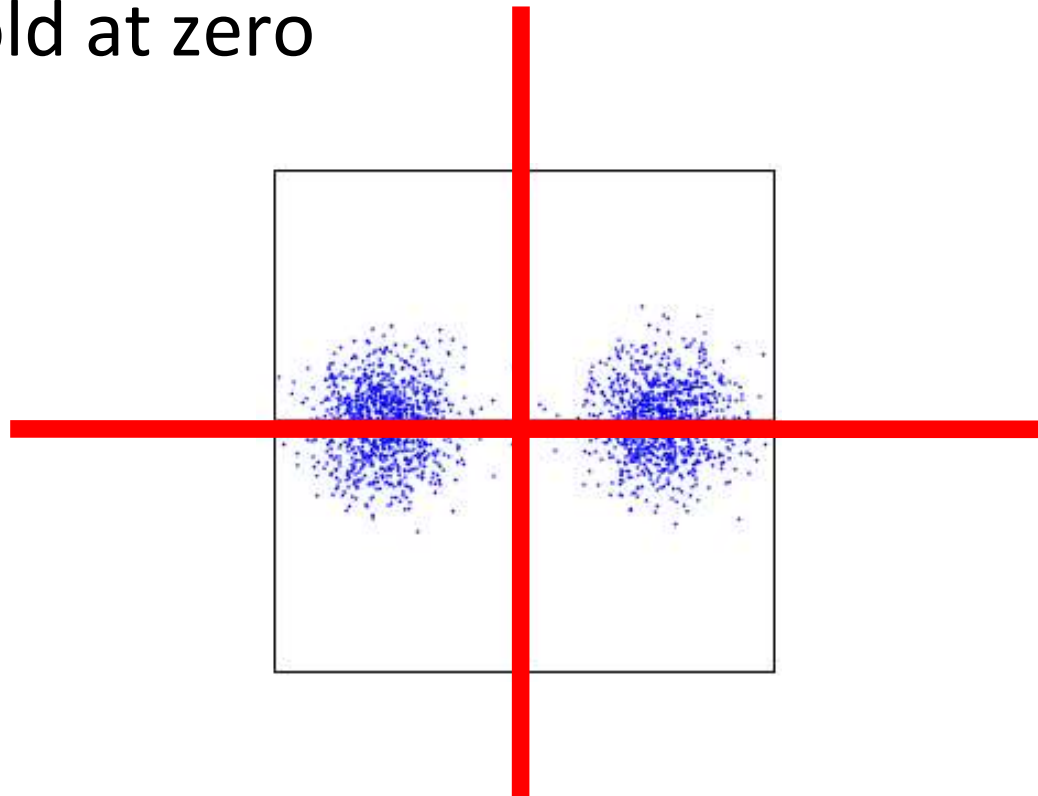
Iterative Quantization(ITQ)

- Baseline scheme:
 - Find PCA embedding of the data
 - For a c -bit code, take top c PCA directions and threshold at zero



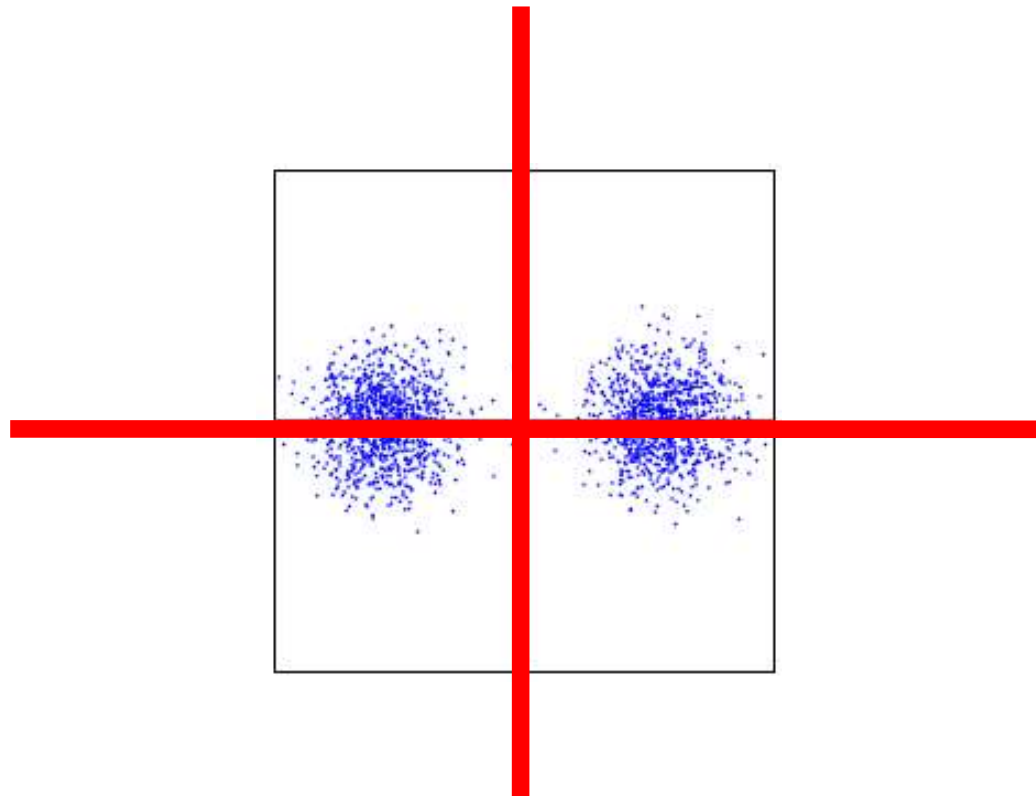
Iterative Quantization(ITQ)

- Baseline scheme:
 - Find PCA embedding of the data
 - For a c -bit code, take top c PCA directions and threshold at zero



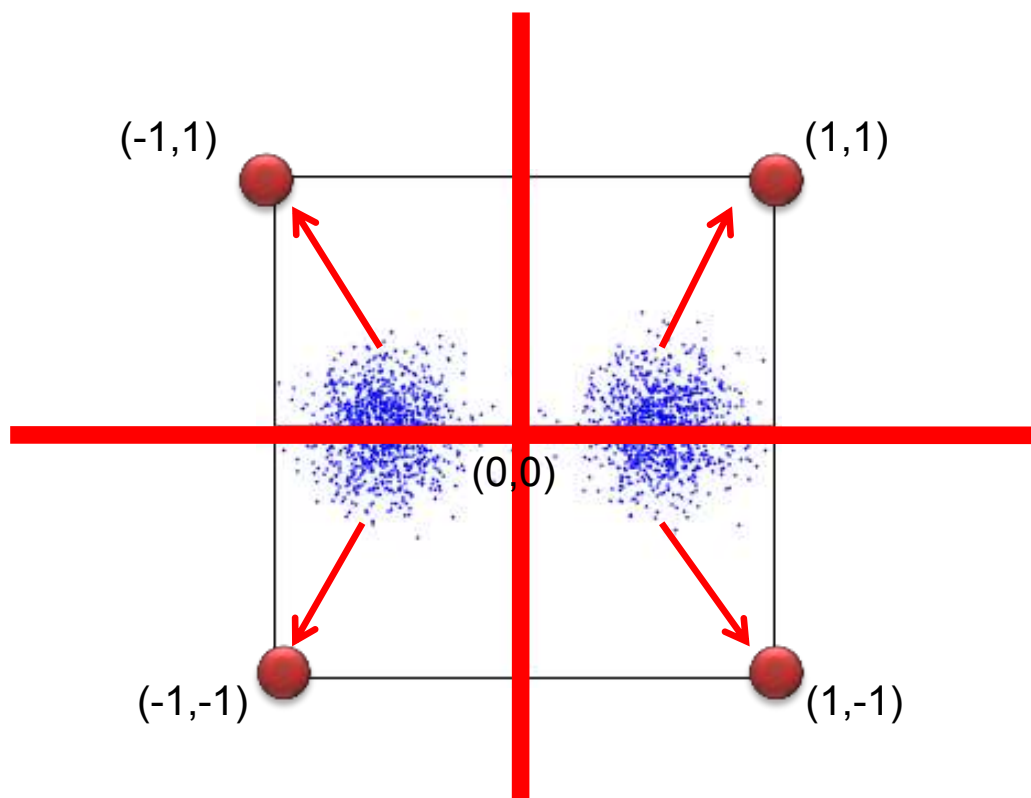
Problem with PCA

- Variance of different dimensions is not balanced



Binary coding as quantization

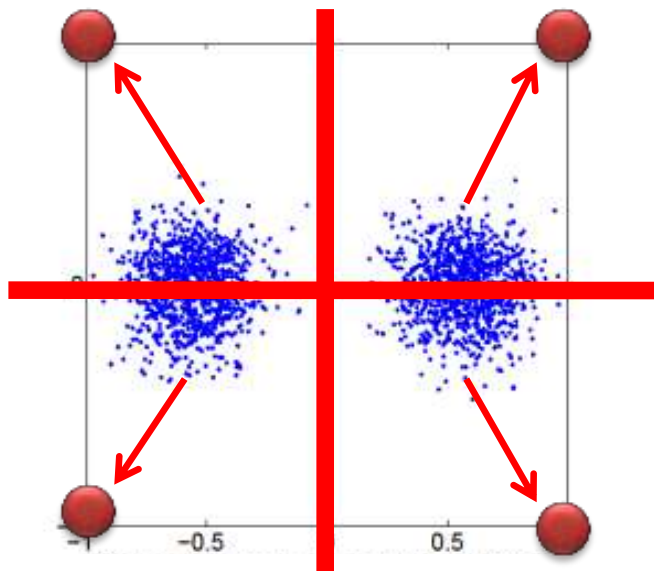
- **Key idea:** similarity-preserving codes should have low quantization error



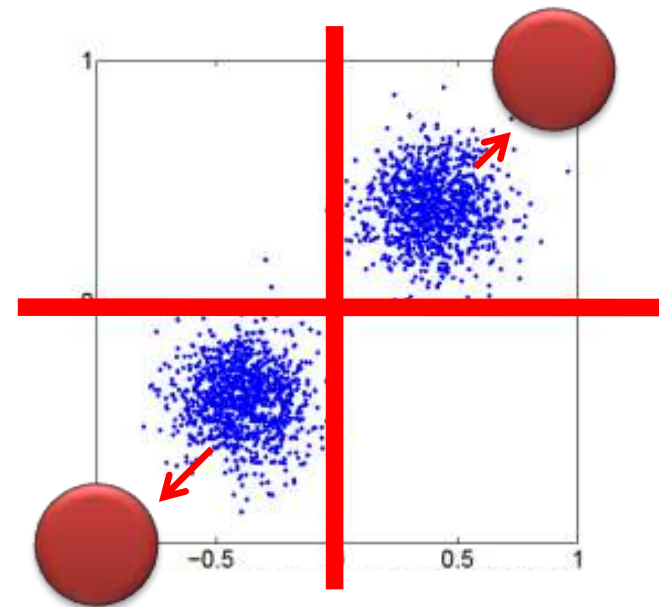
Iterative Quantization(ITQ)

- Rotate the PCA-projected data to minimize quantization error

PCA-projected data



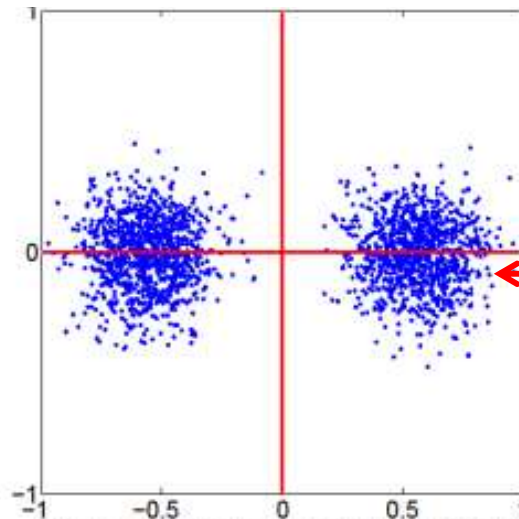
Optimized rotation



Iterative Quantization(ITQ)

- Minimize quantization error of PCA-projected data with respect to the Frobenius norm:

$$\min \mathcal{J}(B, R) = \|B - VR\|_F^2$$

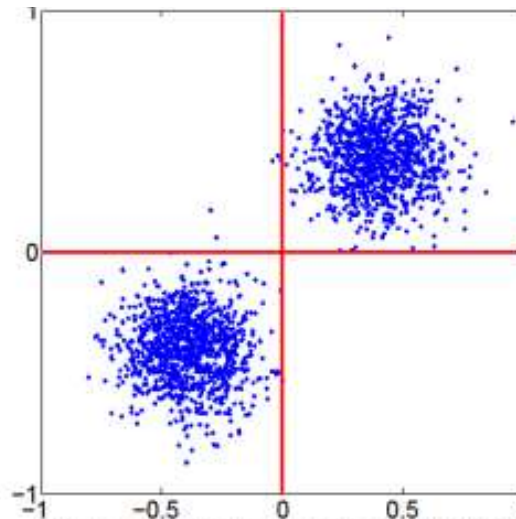


PCA-projected data

Iterative Quantization(ITQ)

- Minimize quantization error of PCA-projected data with respect to the Frobenius norm:

$$\min \mathcal{J}(B, R) = \| \underbrace{B}_{\substack{\text{binary} \\ \text{encoding} \\ \text{matrix}}} - V \underbrace{R}_{\text{rotation}} \|_F^2$$



Iterative Quantization(ITQ)

- Minimize quantization error of PCA-projected data with respect to the Frobenius norm:

$$\min \mathcal{J}(B, R) = \|B - VR\|_F^2$$

subject to $B \in \{-1, 1\}^{n \times c}, \quad R^T R = I$

- **B** is an $n \times c$ matrix where each row is the binary string encoding a data point
- **V** is a matrix of PCA-projected data
- **R** is a $c \times c$ rotation matrix

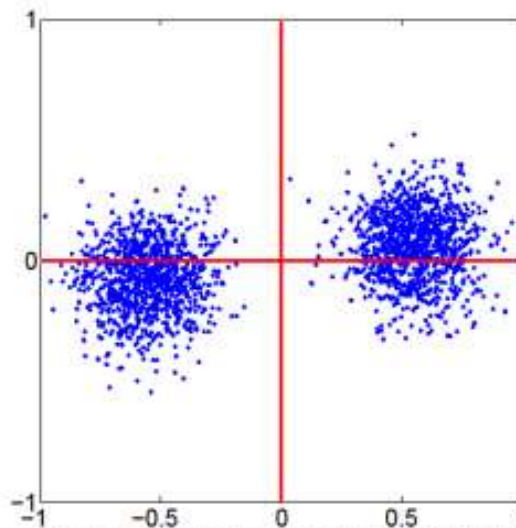
Iterative Quantization(ITQ)

- Minimize quantization error of PCA-projected data with respect to the Frobenius norm:

$$\begin{aligned} \min \mathcal{J}(B, R) &= \|B - VR\|_F^2 \\ \text{subject to } B &\in \{-1, 1\}^{n \times c}, \quad R^T R = I \end{aligned}$$

- Alternating minimization:
 - Initialize ***R*** to a random rotation
 - Fix ***R***, solve for ***B***
 - Fix ***B***, solve for ***R***
 - Iterate until convergence

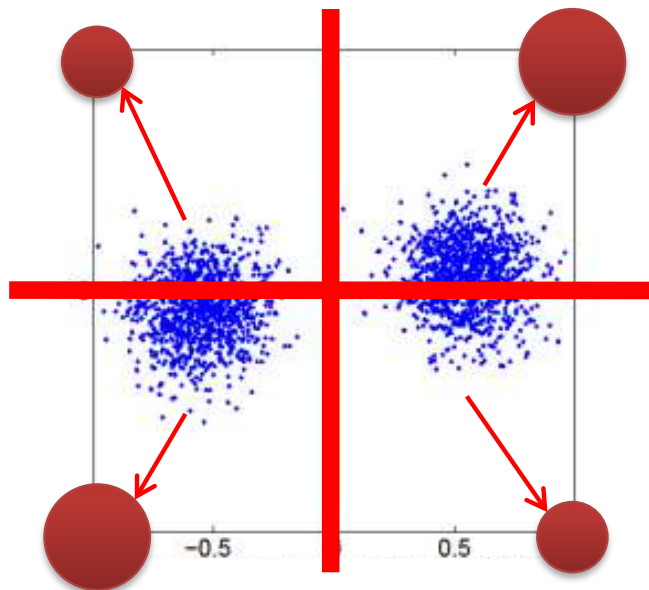
Initialization



randomly rotated data

$$\min \mathcal{J}(R) = \|B - VR\|_F^2,$$

Fix R , find B



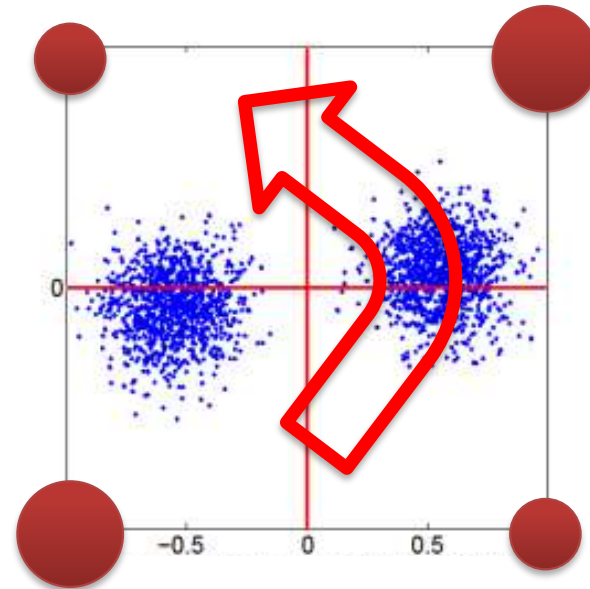
known

$$\min \mathcal{J}(R) = \|B - VR\|_F^2,$$

- Optimal B is given by thresholding rotated coordinates:

$$T = VR \quad \mathcal{J}(B) = \sum_{i=1}^n \sum_{j=1}^c B_{i,j} T_{i,j}. \quad B_{ij} = \begin{cases} 1, & \text{if } T_{i,j} \geq 0; \\ -1, & \text{otherwise.} \end{cases}$$

Fix B , find R



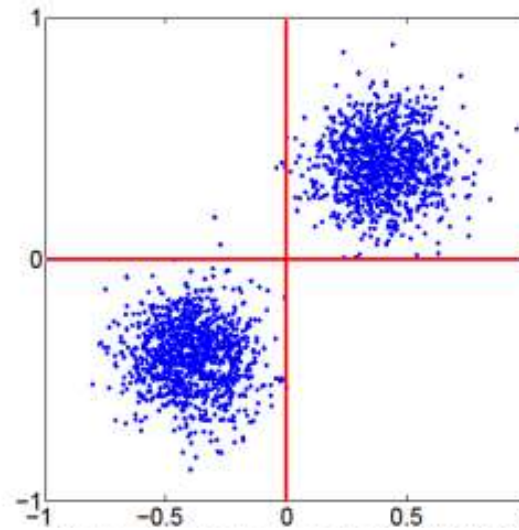
known

$$\min \mathcal{J}(R) = \|B - VR\|_F^2,$$

- Optimal R is found by solving the orthogonal Procrustes problem:

$$B^T V = S \Omega \hat{S}^T \quad R = \hat{S} S^T.$$

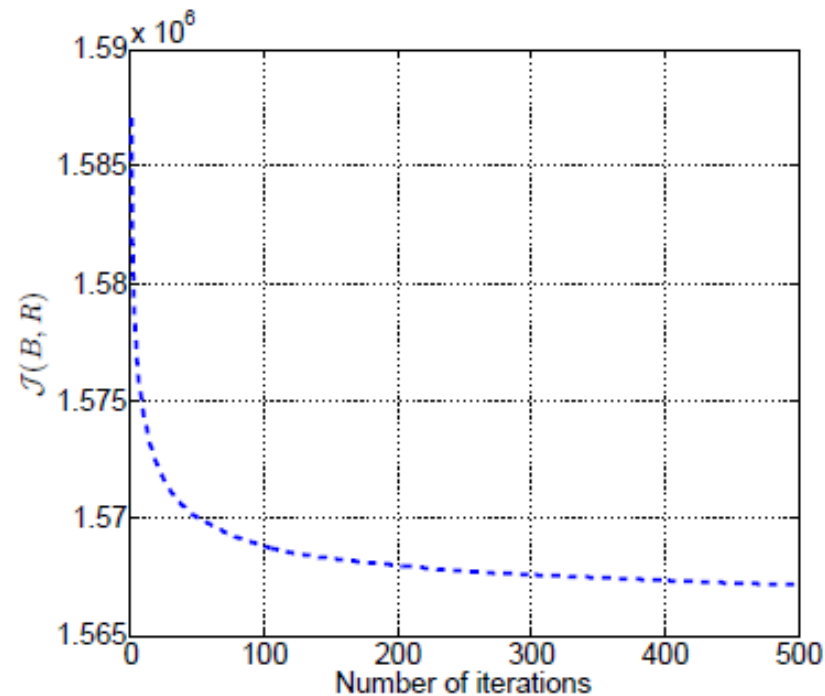
Iterate until convergence



Locally optimal rotation

Iterate until convergence

- Behavior of the objective function:



Evaluation

- CIFAR dataset: ~60,000 images, 11 categories
 - “Tiny images” converted to 320-dimensional gist feature vectors

airplane



automobile



bird



cat



deer



dog



frog



horse



ship

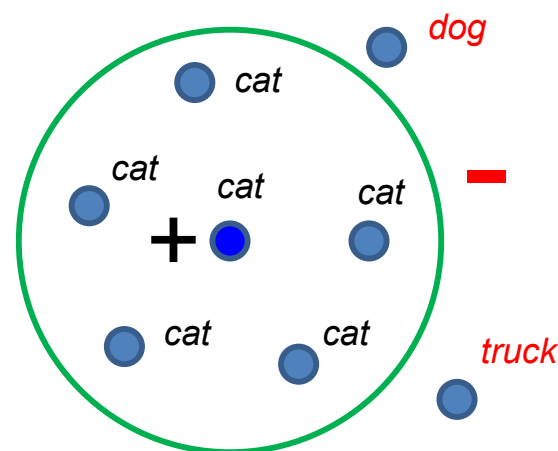
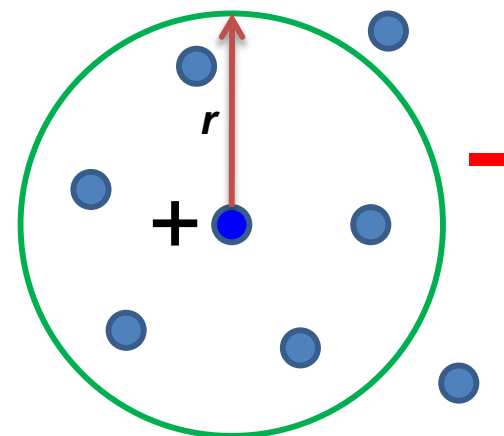


truck



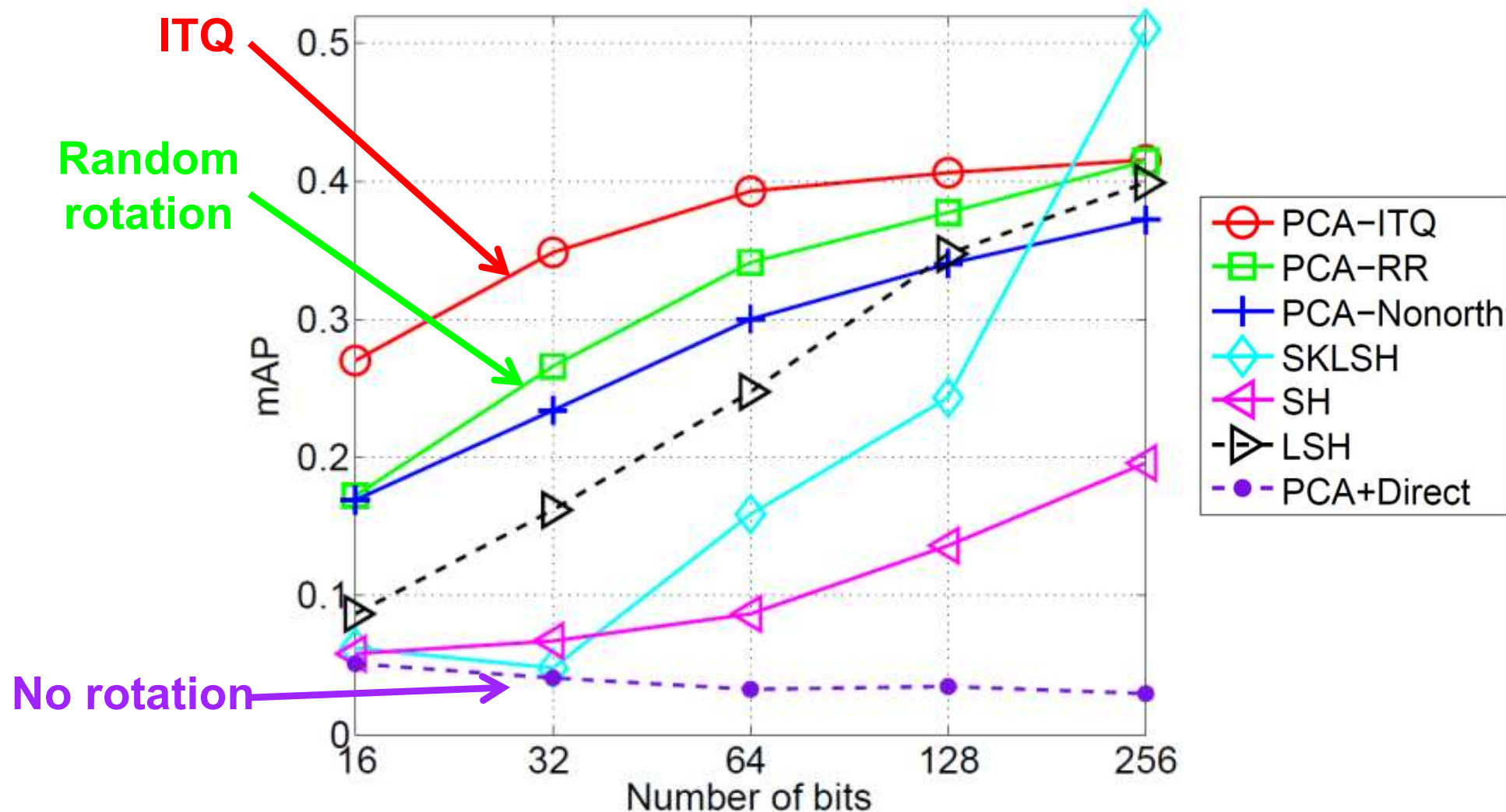
Evaluation

- Euclidean neighbor retrieval
 - Ground truth neighborhood radius defined by average distance to 50th nearest neighbor
 - Performance measured by area under the recall-precision curve (mAP)
- Semantic neighbor retrieval
 - Ground truth defined by class label
 - Performance measured by average precision of top 500 retrieved matches

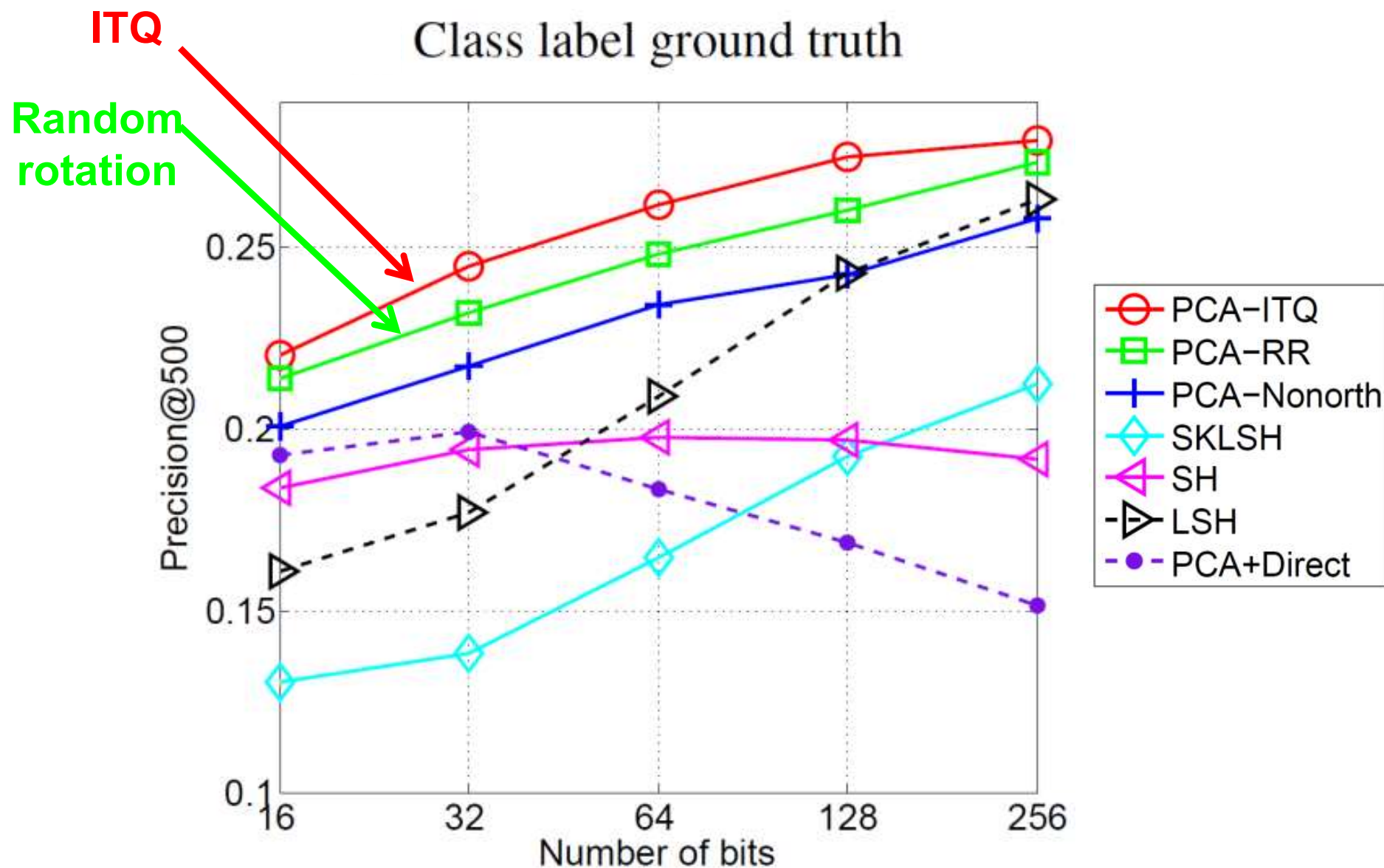


Evaluation

Euclidean ground truth



Evaluation



在线学习哈希方法

在实际检索系统中
数据是流式的、不
断更新的

在线更新

大规模

对于超大规模数据，
很难将其载入内存
进行哈希函数学习

在线学习哈希方法

全部数据 = 之前数据 + 新数据

- 综合所有新老数据，重新学习
- 超高的时间和计算代价

在线学习: Passive-Aggressive (PA) 算法 (JMLR'06)

$$W^t = \arg \min_W \frac{1}{2} \|W - W^{t-1}\|_F^2 + C\xi$$
$$\text{s.t. } L(W; \langle (x_i^t, x_j^t), s_{ij} \rangle) \leq \xi \text{ and } \xi \geq 0$$

Online Hashing
(IJCAI'13)

在线学习哈希方法

■ 矩阵素描(sketching or coreseset)

- ✓ 一个比原始矩阵小的多的矩阵
- ✓ 保留原始矩阵的大多数特性

$$\forall x, \|x\| = 1 \quad |||P^T x|^2 - |Q^T x|^2| \leq \varepsilon \|P\|_F^2$$

在线学习哈希方法

Frequent Directions (FD)

Algorithm 1 Frequent Directions (Liberty [21])

Input: Data matrix $P \in \mathbb{R}^{d \times n}$, sketch matrix $Q \in \mathbb{R}^{d \times l}$.

Output: Sketch matrix Q .

if Q not exists **then**

$Q \leftarrow$ all zeros $d \times l$ matrix

end if

for each column P_i in P , **do**

Insert P_i into a zero valued column of Q

if Q has no zero valued columns **then**

$[U, S, V] = \text{SVD}(Q)$

$C = US$ [just for notation]

Set $\delta = s_{l/2}^2$ [the squared $(l/2)^{th}$ entry of S]

Set $\hat{S} = \sqrt{\max(S^2 - I_l \delta, 0)}$

$Q = U\hat{S}$

end if

end for

Lemma 1. (Liberty [21]) Apply Algorithm 1 to matrix P to obtain a sketch Q with prescribed l , then

$$\forall x, \|x\| = 1 \quad 0 \leq \|P^T x\|^2 - \|Q^T x\|^2 \leq \frac{2}{l} \|P\|_F^2$$

or

$$0 \leq \|PP^T - QQ^T\|_2 \leq \frac{2}{l} \|P\|_F^2$$

在线学习哈希方法

广泛意义上的PCA Hashing:

$$\begin{aligned} \max_{W \in \mathbb{R}^{d \times r}} \quad & tr(W^T (X - \mu)(X - \mu)^T W) \\ \text{s.t.} \quad & W^T W = I_r \end{aligned}$$

其中 μ 是全体数据的均值向量.

能否为矩阵 $X - \mu$ 在线维护一个素描矩阵 Y , 以至于

$$YY^T \approx (X - \mu)(X - \mu)^T$$

在线学习哈希方法

均值漂移问题:

- 在流式问题中, 数据在不断更新, 因此数据的均值 μ 也会不断变化.

给每个数据块加一个虚拟的样本:

- 对于流式数据 $X_t = [D_1, D_2, \dots, D_t]$, 重新设计数据矩阵 E_t :

$$E_t = [D_1 - \overline{D}_1, \quad D_2 - \overline{D}_2, \sqrt{\frac{n_1 m_2}{n_1 + m_2}} (\overline{D}_2 - \mu_1), \dots, \\ D_i - \overline{D}_i, \sqrt{\frac{n_{i-1} m_i}{n_{i-1} + m_i}} (\overline{D}_i - \mu_{i-1}), \dots, \\ D_t - \overline{D}_t, \sqrt{\frac{n_{t-1} m_t}{n_{t-1} + m_t}} (\overline{D}_t - \mu_{t-1})]$$

在线学习哈希方法

基于以上的设计，在任意适合 t ，可以证明：

$$E_t E_t^T = \text{cov}(X_t)$$

通过为新设计的矩阵 E_t 在线维护素描矩阵 Y ，我们可以基于这个新的小矩阵 Y 学习哈希函数。

Algorithm 2 Zero Mean Sketching

Input: Streaming data chunk $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$,
All zeros matrix Y of size $d \times l$.

- 1: Sketch $\mathcal{D}_1 - \overline{\mathcal{D}_1}$ into Y with Algorithm 1
 - 2: $n \leftarrow m_1$ and $\mu \leftarrow \overline{\mathcal{D}_1}$
 - 3: **for** $i = 2 : k$, **do**
 - 4: Sketch $[\mathcal{D}_i - \overline{\mathcal{D}_i}, \sqrt{\frac{nm_i}{n+m_i}}(\overline{\mathcal{D}_i} - \mu)]$ into Y
 - 5: $\mu \leftarrow \frac{n\mu}{n+m_i} + \frac{m_i \overline{\mathcal{D}_i}}{n+m_i}$ [update the data mean]
 - 6: $n \leftarrow n + m_i$ [update the data size]
 - 7: **end for**
-

对比实验

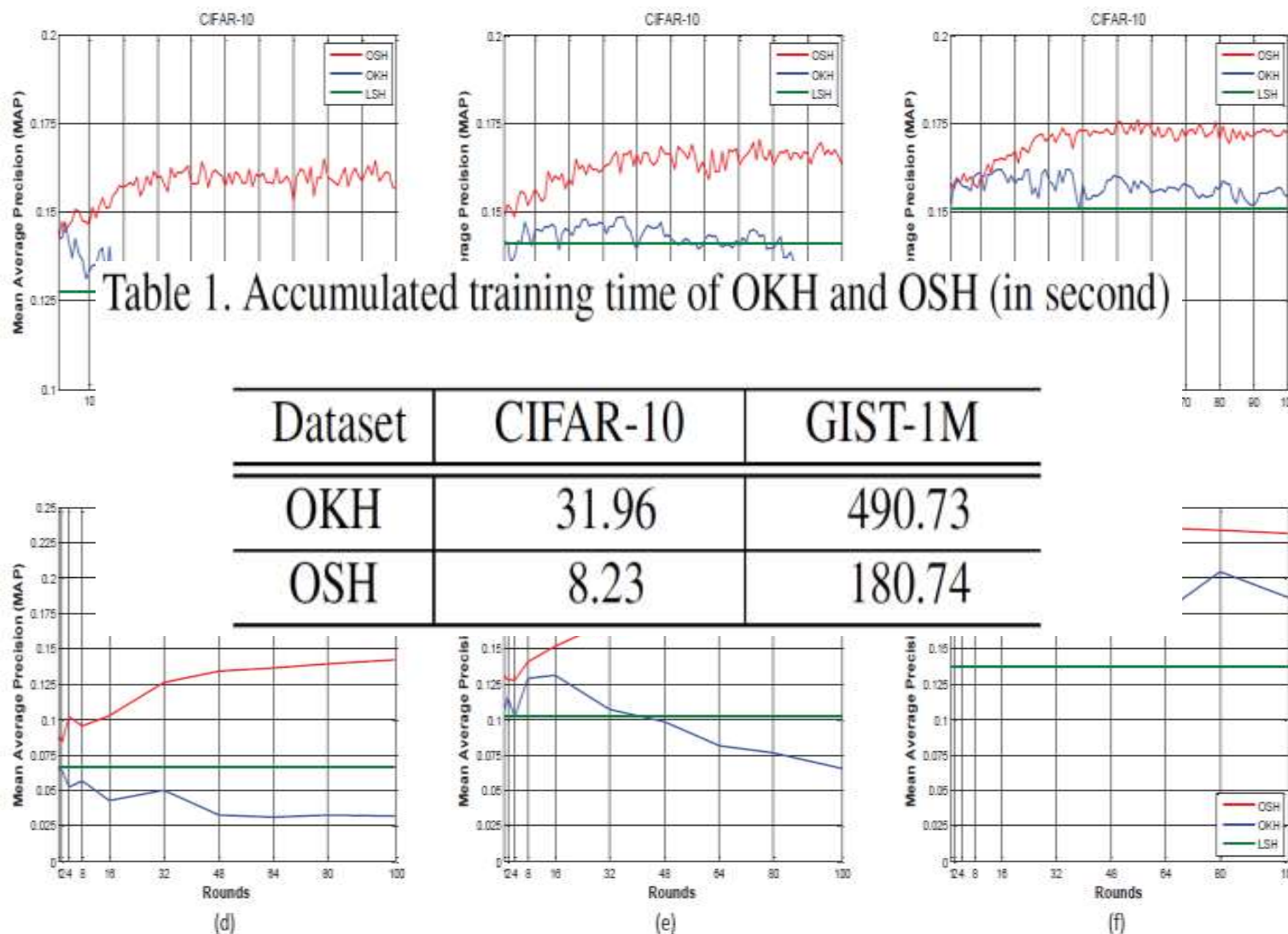
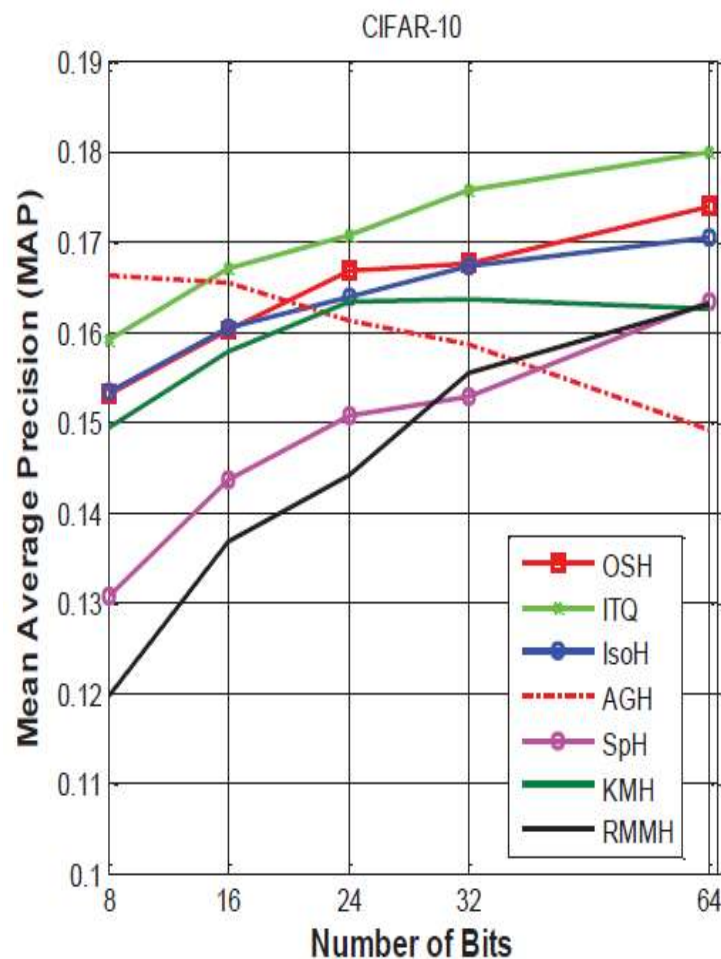
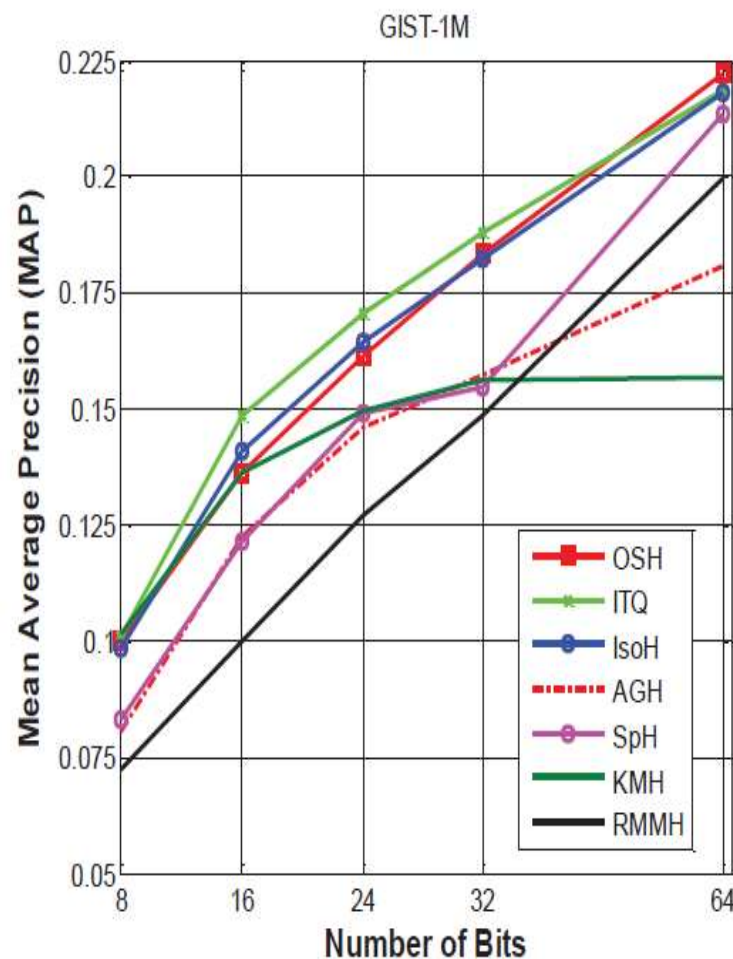


Figure 3. (a)(b)(c) Mean average precision (MAP) on CIFAR-10 dataset at each round with 16, 32, 64 bits. (d)(e)(f) MAP on GIST-1M dataset at each round with 16, 32, 64 bits. (Best viewed in color)

对比实验



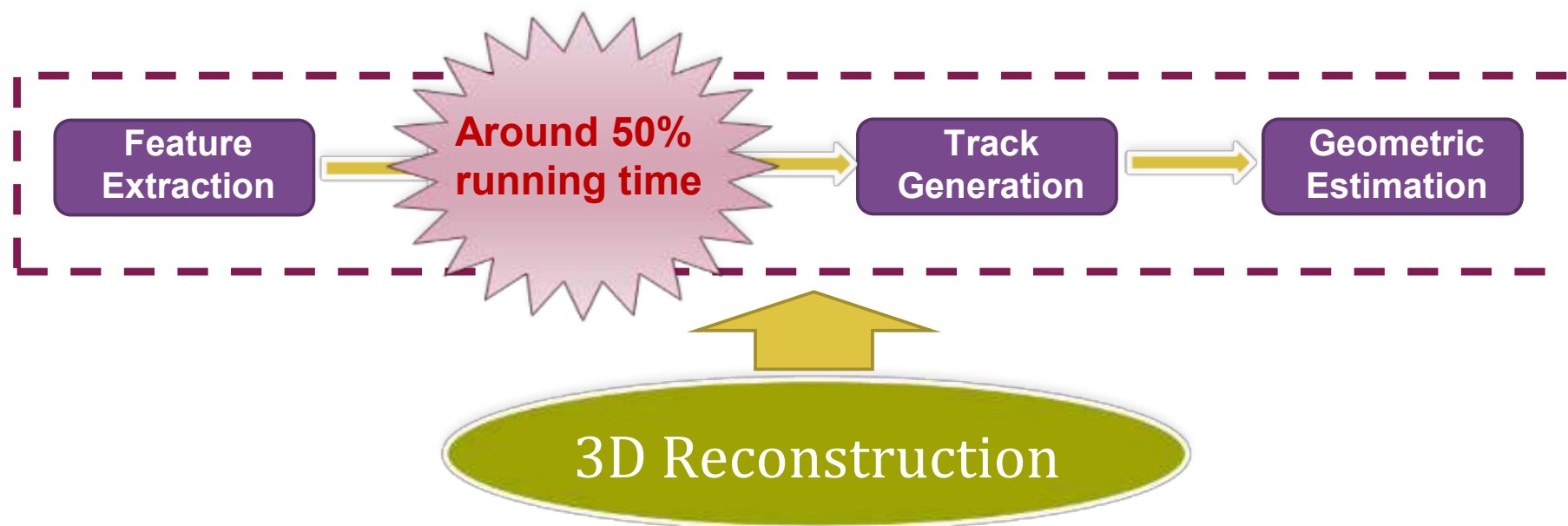
(a)



(b)

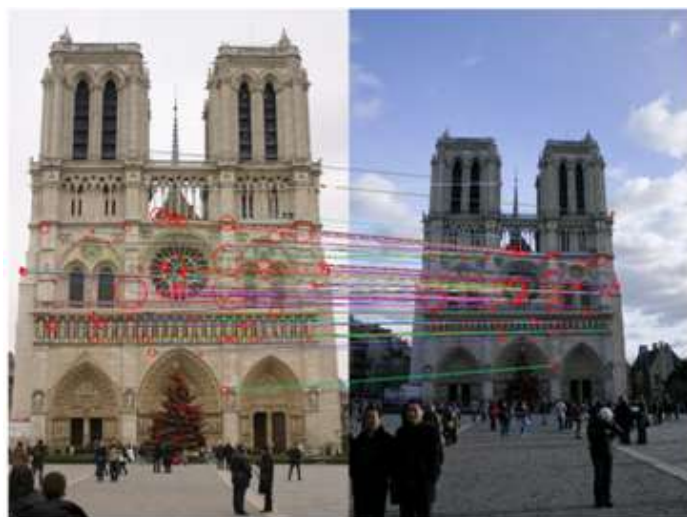
基于哈希的特征匹配

- 三维重建主要包括以下四步:



# Images	# Cores	Match Time	Reconstruction Time	Largest Component
150,000	496	13 Hours	8 Hours	2,106

基于哈希的特征匹配



特征匹配



- 图像搜索
- 三维重建
- 图像分类
- 目标识别

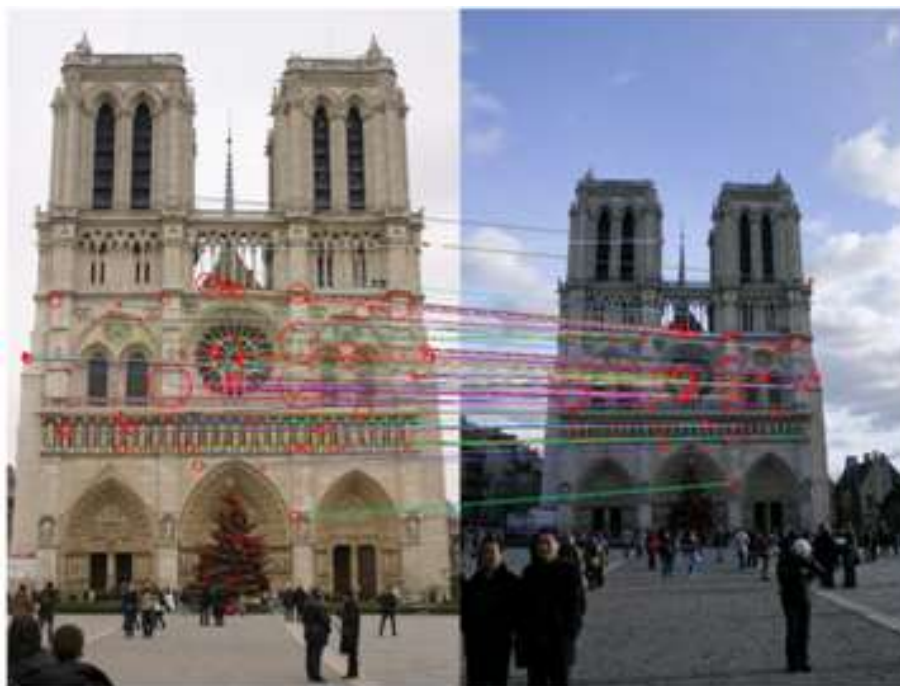
典型应用

- **匹配复杂度太高**: $O(N*(N-1) * M^2)$
- **硬件加速成本高**: GPU、并行

基于哈希的特征匹配

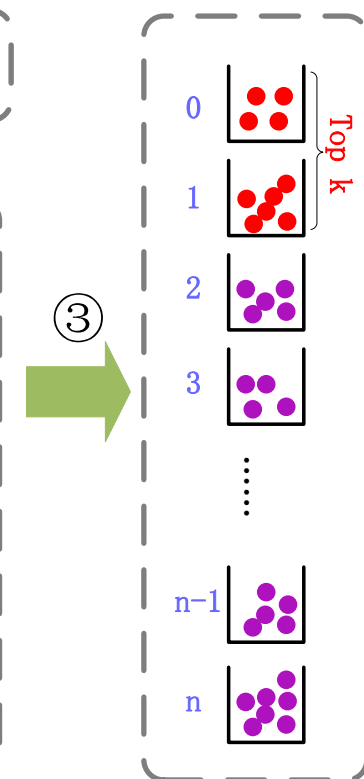
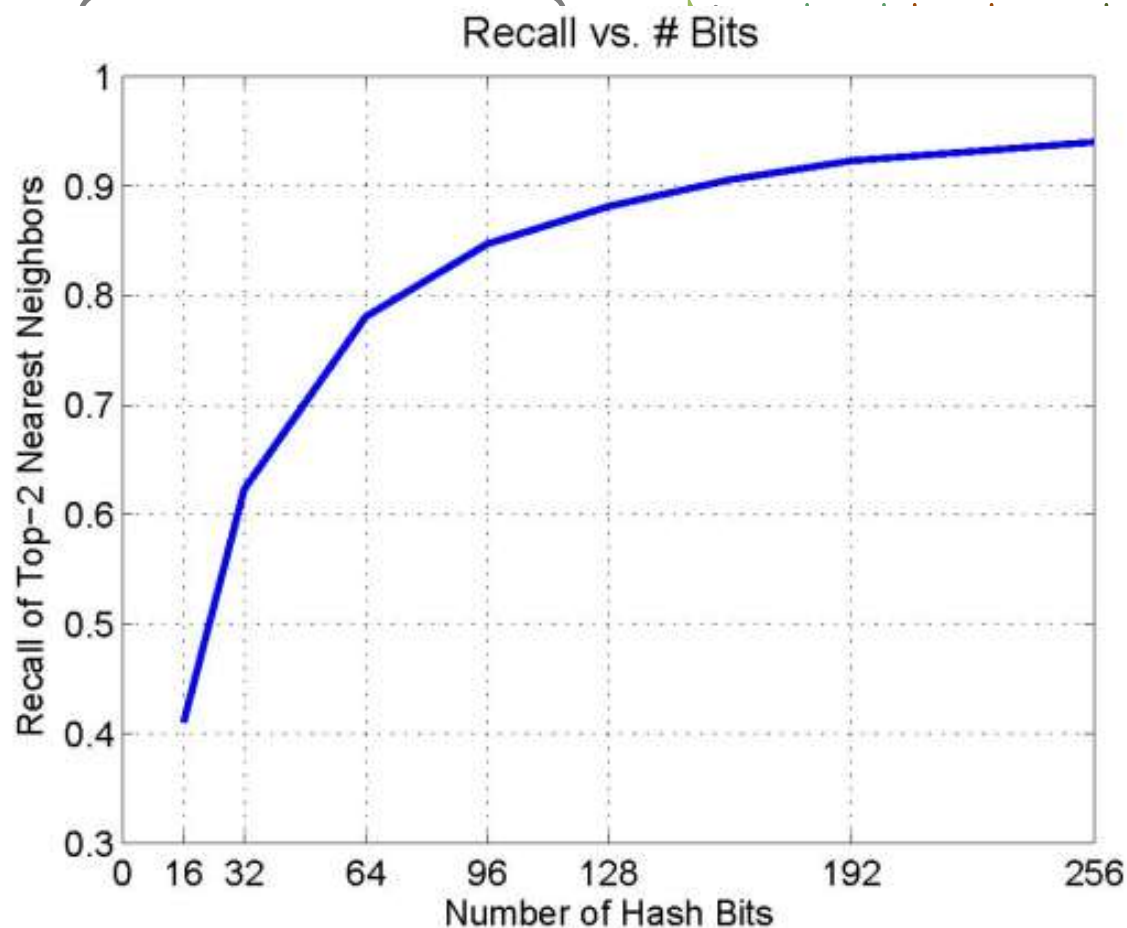
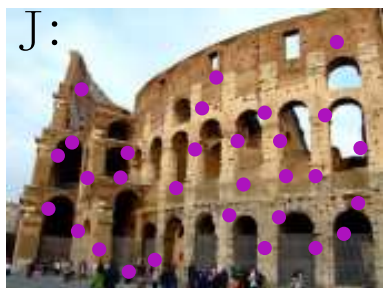
■ 现有特征匹配方法可以分为三大类:

- Point matching
- Line matching
- Region matching



Point matching is searching in essence!

基于哈希的特征匹配



基于哈希的特征匹配



(a) Taj.Mahal: 189 images.



(b) Statue.of.Liberty: 674 images.



(c) Notre.Dame: 715 images.



(d) Colosseum: 1357 images.

基于哈希的特征匹配

比brute force方法快上百倍

比通用的k-d tree方法快10-40倍

Method	Taj.Mahal			Statue.of.Liberty		
	T_{Match}	Speed-Up	Points	T_{Match}	Speed-Up	Points
Brute	52575s	1.00×	124038	30608s	1.00×	88001
KD						3674
LDA						3274
CasH						3587
CasHa						6206
Me						
	T_{Match}	Speed-Up	Points	T_{Match}	Speed-Up	Points
Brute	396729s	1.00×	358121	12307s	1.00×	540308
KDTree	60663s	6.54×	347056	2430s	5.06×	445774
LDAHash	13136s	30.20×	413348	851s	14.46×	492040
CasHash-8Bit	2266s	175.08×	484960	222s	55.44×	393408
CasHash-10Bit	1354s	293.01×	400673	196s	62.79×	512508

基于哈希的特征匹配

OpenMVG (open Multiple View Geometry)



通过测试确信它在不同数据集上都非常有效

I'm actually testing it to **ensure it works well on different datasets** and that I have a code that works as good as your original version.

<https://github.com/openMVG/openMVG/issues/194>



Theia Vision Library

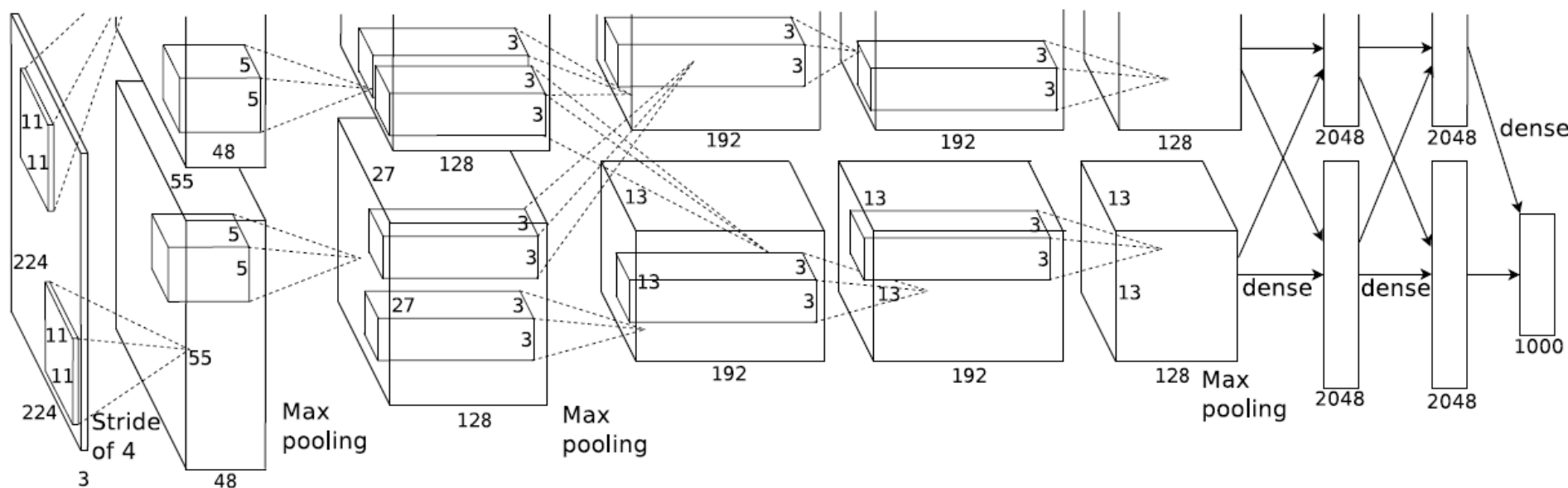
这是非常有用的见解和足够简单的想法

it's **a really useful insight** and is **a simple enough idea**.

<https://github.com/kip622/Theia>

基于量化的CNN加速与压缩

- 卷积神经网络（convolutional neural network）
 - 特征表征能力更强



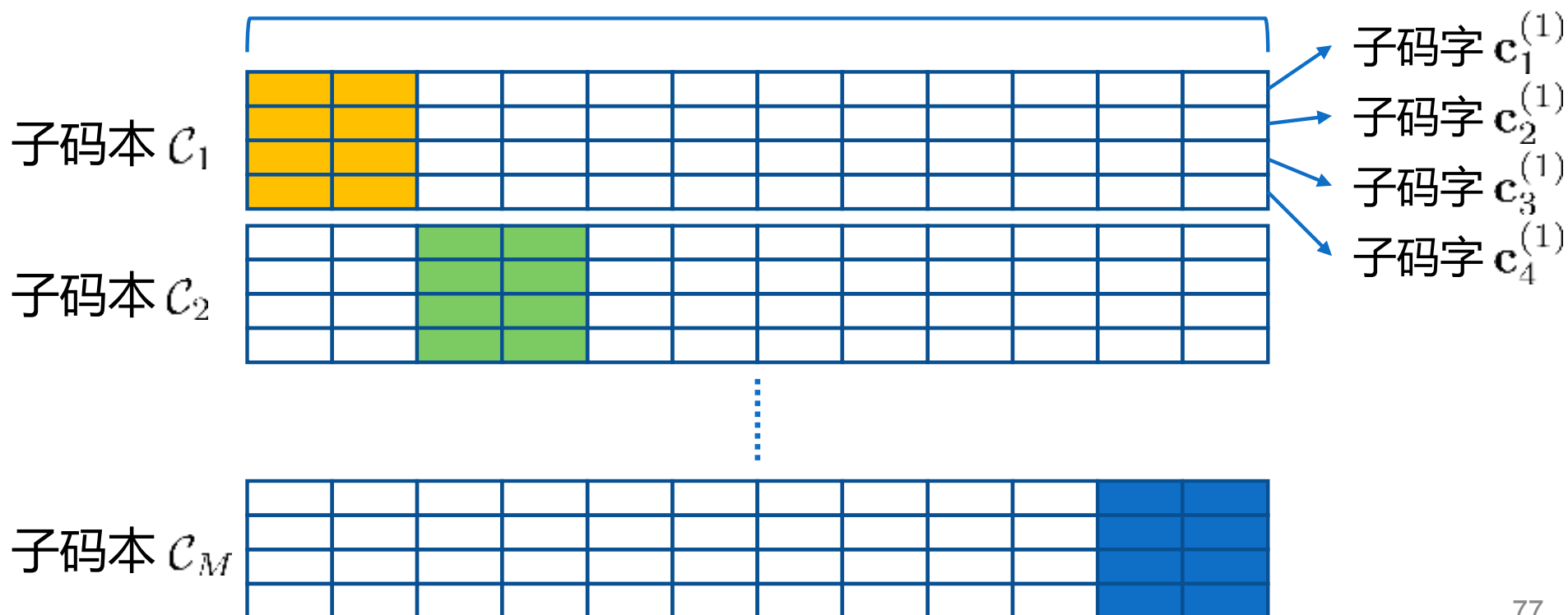
Hash—>PQ

• 乘积量化 (product quantization, PQ)

- 将样本特征量化为 M 个子码字之和的形式

$$q(\mathbf{x}) = \sum_{m=1}^M q_m(\mathbf{x}), \quad q_m(\mathbf{x}) \in \mathcal{C}_m$$

D 维



Hash—>PQ

- 乘积量化

- 优化目标：最小化量化误差

$$\min_{c_1, \dots, c_M} \sum_{\mathbf{x}} \|\mathbf{x} - q(\mathbf{x})\|_2^2$$

- 可等价于分别优化各个子空间下的子码本

$$\min_{c_m} \sum_{\mathbf{x}} \|u_m(\mathbf{x}) - q_m(\mathbf{x})\|_2^2$$

- 求解方法： k 均值聚类
- 时间复杂度： $\mathcal{O}(TkND)$

难以应用于大规模数据集

基于量化的CNN加速与压缩

Fully-connected Layer:

$$T(c_t) = \langle W_{c_t}, S \rangle$$

Convolutional Layer:

$$T_{p_t}(c_t) = \sum_{(p_k, p_s)} \langle W_{c_t, p_k}, S_{p_s} \rangle$$

基于量化的CNN加速与压缩

For the fully-connected layer, we split the weighting vector W_{c_t} and layer input S into M sub-vectors, each of $C'_s = C_s/M$ dimensions:

$$\begin{aligned} T(c_t) &= \langle W_{c_t}, S \rangle \\ &= \sum_{m=1}^M \langle W_{c_t}^{(m)}, S^{(m)} \rangle \end{aligned}$$

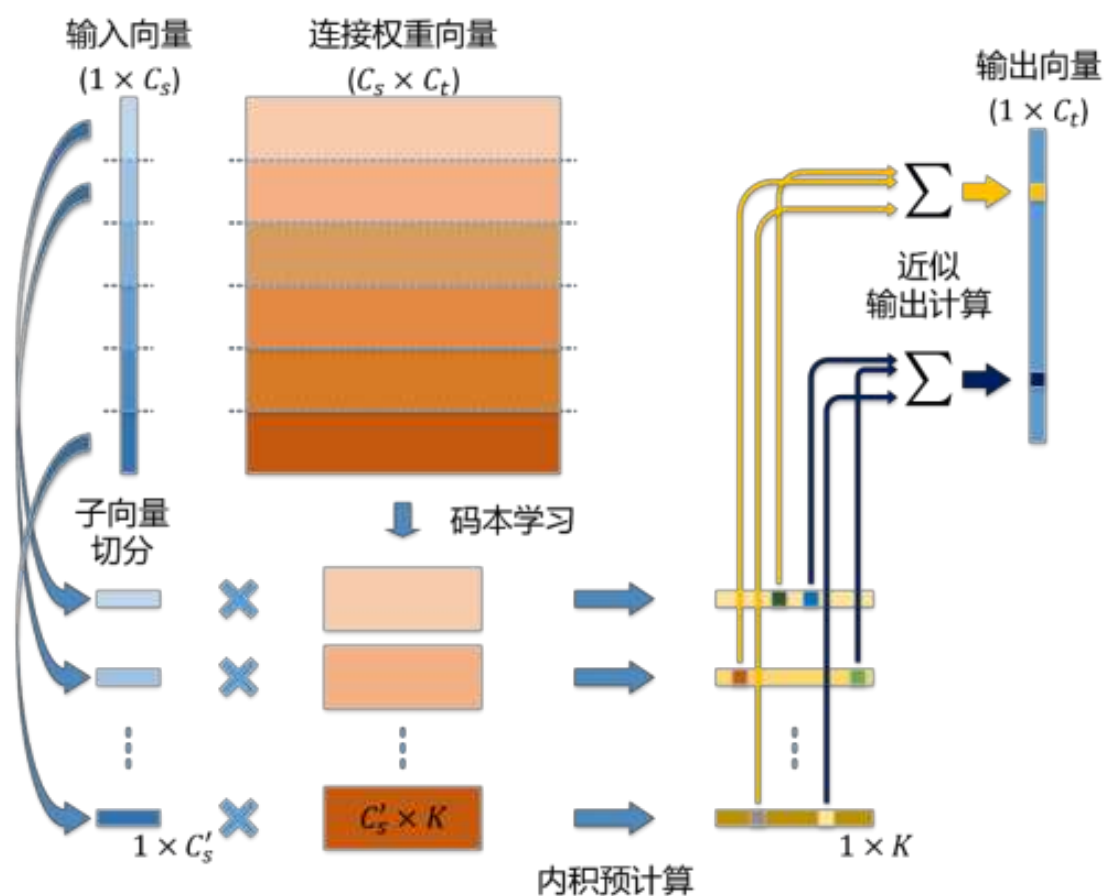
We quantize each sub-vector $W_{c_t}^{(m)}$ with:

$$W_{c_t}^{(m)} \approx D^{(m)} B_{c_t}^{(m)}$$

where $W_{c_t}^{(m)} \in \mathbb{R}^{C'_s \times 1}$, $D^{(m)} \in \mathbb{R}^{C'_s \times K}$, and $B_{c_t}^{(m)} \in \{0, 1\}^{K \times 1}$

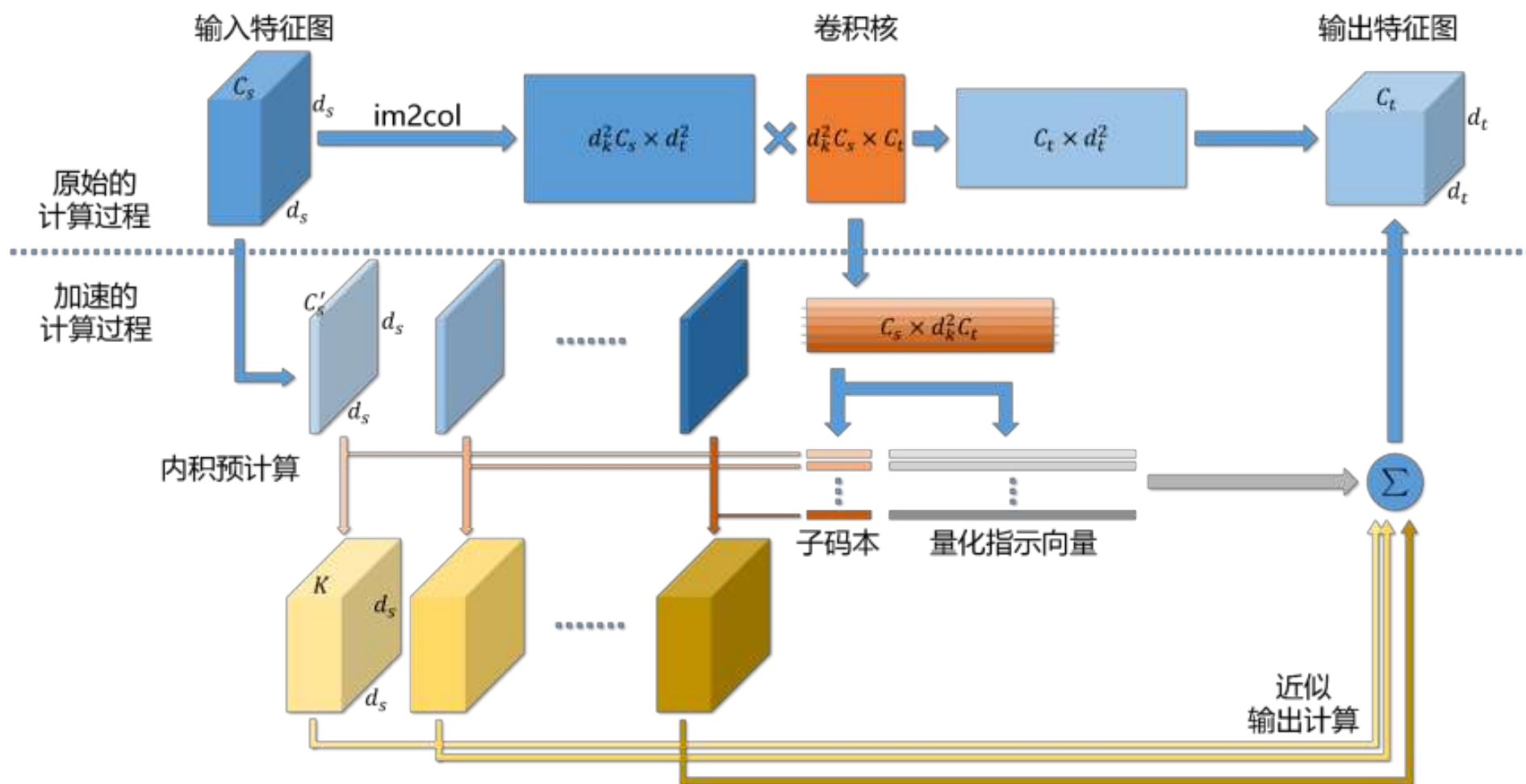
基于量化的CNN加速与压缩

- 全连接层的参数量化与加速计算



基于量化的CNN加速与压缩

• 卷积层的参数量化与加速计算



基于量化的CNN加速与压缩

- 参数量化的求解方法
 - 最小化网络参数的量化误差

$$\begin{aligned} \min_{\mathbf{D}, \mathbf{B}} \quad & \left\| \mathbf{W} - \tilde{\mathbf{W}} \right\|_F^2 \\ \text{s.t.} \quad & \tilde{\mathbf{W}} = g(\mathbf{D}, \mathbf{B}) \end{aligned}$$

- 最小化网络输出的近似误差

$$\begin{aligned} \min_{\mathbf{D}, \mathbf{B}} \quad & \left\| f(\mathbf{X}; \mathbf{W}) - f(\mathbf{X}; \tilde{\mathbf{W}}) \right\|_F^2 \\ \text{s.t.} \quad & \tilde{\mathbf{W}} = g(\mathbf{D}, \mathbf{B}) \end{aligned}$$

基于量化的CNN加速与压缩

- 累积误差问题

- 网络中浅层的参数被量化后，会改变深层的输入

$$\min_{\mathbf{D}, \mathbf{B}} \left\| f(\mathbf{X}; \mathbf{W}) - f(\mathbf{X}; \tilde{\mathbf{W}}) \right\|_F^2$$

$$\text{s.t. } \tilde{\mathbf{W}} = g(\mathbf{D}, \mathbf{B})$$



累积误差修正

$$\min_{\mathbf{D}, \mathbf{B}} \left\| f(\mathbf{X}; \mathbf{W}) - f(\tilde{\mathbf{X}}; \tilde{\mathbf{W}}) \right\|_F^2$$

$$\text{s.t. } \tilde{\mathbf{W}} = g(\mathbf{D}, \mathbf{B})$$

基于量化的CNN加速与压缩

- 分类性能对比
 - 数据集: MNIST

方法	三层网络		五层网络	
	压缩倍数	分类错误率	压缩倍数	分类错误率
原始网络	-	1.35%	-	1.12%
RER [33]	8.0×	2.19%	8.0×	1.24%
LRD [165]	8.0×	1.89%	8.0×	1.77%
DK [50]	8.0×	1.71%	8.0×	1.26%
NN-ES [48]	8.0×	1.69%	8.0×	1.35%
HashNets [48]	8.0×	1.43%	8.0×	1.22%
Q-CNN ¹	10.9×	1.42%	13.0×	1.34%
Q-CNN ²	10.9×	1.39%	13.0×	1.19%

基于量化的CNN加速与压缩

- 分类性能对比
 - 数据集: ILSVRC-12

网络	方法	超参		加速倍数	压缩倍数	Top-1 分类错误率 ↑	Top-5 分类错误率 ↑
		卷积层	全连接层				
AlexNet	BC [43]	-	-	2.00×	32.00×	21.20%	19.20%
	BWN [46]	-	-	2.00×	32.00×	2.80%	3.20%
	DC [38]	-	-	-	9.00/35.00×	0.00%	-0.03%
	Q-CNN ²	8/128	3/32	4.05×	15.10×	1.38%	0.84%
		8/128	4/32	4.15×	18.32×	1.46%	0.97%
CaffeNet	Q-CNN ²	8/128	3/32	4.04×	15.10×	1.43%	0.99%
		8/128	4/32	4.16×	18.32×	1.54%	1.12%
CNN-S	Q-CNN ²	8/128	3/32	5.69×	15.93×	1.48%	0.81%
		8/128	4/32	5.78×	19.57×	1.64%	0.85%
VGG-16	DC [38]	-	-	-	13.00/49.00×	-0.33%	-0.41%
	Q-CNN ²	8/128	3/32	4.92×	16.06×	1.02%	0.38%
		8/128	4/32	4.94×	19.60×	1.13%	0.45%

基于量化的CNN加速与压缩

- 移动设备上的运算效率
 - 数据集: ILSVRC-12

网络	方法	运行时间	硬盘存储	内存占用	Top-5 分类错误率
AlexNet	原始网络	2.93s	232.56MB	264.74MB	19.74%
	Q-CNN ²	0.95s	12.60MB	74.65MB	20.70%
CNN-S	原始网络	10.58s	392.57MB	468.90MB	15.82%
	Q-CNN ²	2.61s	20.13MB	129.49MB	16.68%



大纲

- 背景介绍
- 多媒体内容索引
- 多媒体内容排序
- 近似近邻搜索
- **总结与展望**

总结与展望

- 排序与索引是多媒体搜索的关键问题
- 正在走向成熟，但依然方兴未艾
- 与其它学科交叉融合

Thanks!

Q&A

