



Deep Cross-Modal Hashing

Qing-Yuan Jiang & Wu-Jun Li

LAMDA Group, Department of Computer Science and Technology, Nanjing University, Nanjing, China.

jiangqy@lamda.nju.edu.cn,liwujun@nju.edu.cn



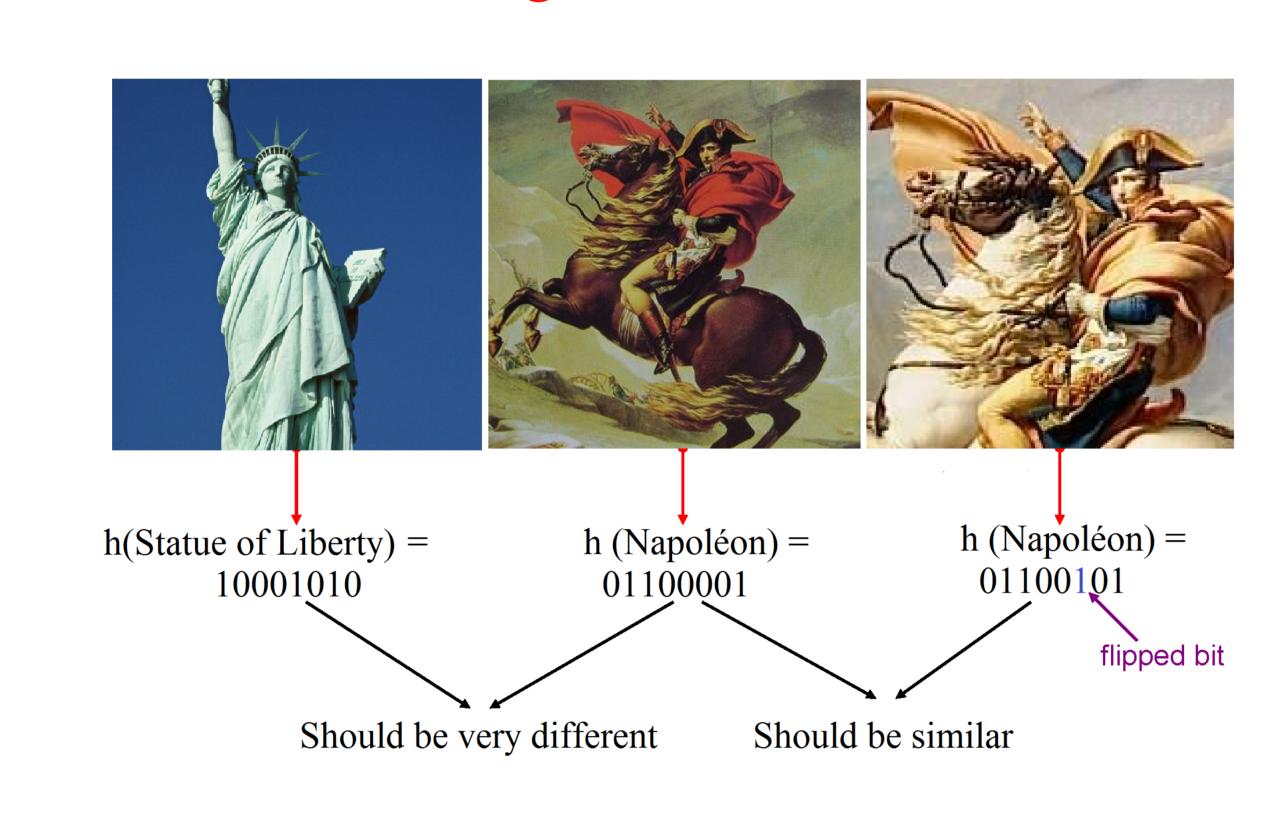


Introduction

Nearest Neighbor Search (NNS)

- ulletGiven a query point q, return the points closest to q in the database (e.g., image retrieval).
- Challenges for NNS in big data applications: curse of dimensionality; storage cost; query speed

- •Similarity preserved hashing is to map the data points from the original space into a Hamming space of binary codes with similarity preserved.
- Hashing can solve the above challenges.



Cross-Modal Hashing (CMH)

- •Cross-modal retrieval: the modality of the query point is different from the modality of the points in database.
- CMH: hashing for cross-modal retrieval. Low storage cost and fast query speed.

- Almost all existing CMH methods are based on hand-crafted features.
- Hand-crafted features might not be compatible for hash-code learning.

Contribution

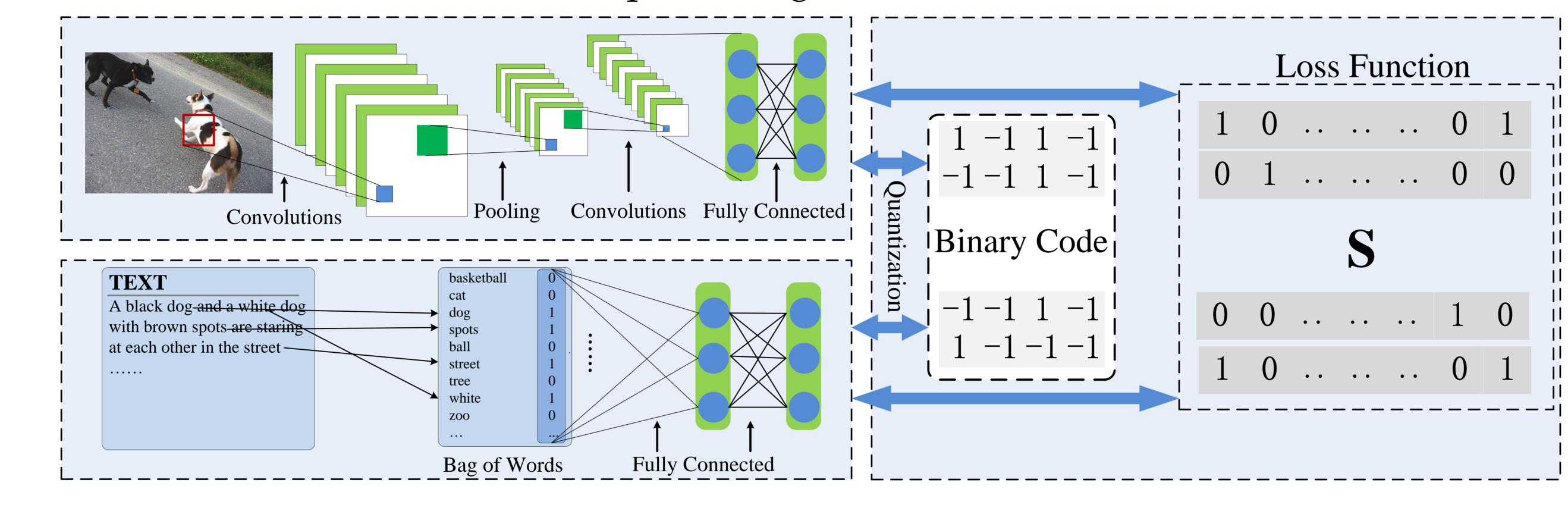
- A novel CMH method, called deep cross-modal hashing (DCMH), for cross-modal retrieval applications.
- •DCMH is an end-to-end learning framework with deep neural networks, one for each modality, to perform feature learning from scratch.
- DCMH achieves the state-of-the-art performance on three datasets.

Notation

- $\mathbf{Y} = \{\mathbf{y}_j\}_{j=1}^n$: n points of text modality.
- •S = $\{S_{ij}\}_{n\times n}$: cross-modal similarities.
- $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$: n points of image modality. $f(\mathbf{x}_i; \theta_x)$: the output of deep neural network for image modality.
 - $g(\mathbf{y}_i; \theta_y)$: the output of deep neural network for text modality.

Model

The end-to-end deep learning framework of DCMH model.



Feature learning part:

This part contains two deep neural networks, one for image modality and the other for text modality. Their configurations are shown in the following tables.

Configuration of the CNN for image modality.

O	
Layer	Configuration
conv1	f. $64 \times 11 \times 11$; st. 4×4 , pad 0, LRN,×2 pool
conv2	f. $265 \times 5 \times 5$; st. 1×1 , pad 2, LRN,×2 pool
conv3	f. $265 \times 3 \times 3$; st. 1×1 , pad 1
conv4	f. $265 \times 3 \times 3$; st. 1×1 , pad 1
conv5	f. $265 \times 3 \times 3$; st. 1×1 , pad $1, \times 2$ pool
full6	4096
full7	4096
f11118	Hash code length c

Configuration of the deep neural network for text modality.

Layer	Configuration
full1	8192
full2	Hash code length \boldsymbol{c}

Deep neural network input:

- -Image deep neural network: raw image.
- -Text deep neural network: Bag-of-words (BOW) feature.
- Hash-code learning part:

$$\min_{\mathbf{B},\theta_{x},\theta_{y}} \mathcal{J} = -\sum_{i,j=1}^{n} (S_{ij}\Theta_{ij} - \log(1 + e^{\Theta_{ij}}))
+ \gamma(\|\mathbf{B} - \mathbf{F}\|_{F}^{2} + \|\mathbf{B} - \mathbf{G}\|_{F}^{2})
+ \eta(\|\mathbf{F}\mathbf{1}\|_{F}^{2} + \|\mathbf{G}\mathbf{1}\|_{F}^{2})
s.t. \mathbf{B} \in \{-1, +1\}^{c \times n}.$$

- $-\mathbf{B} \in \{-1,+1\}^{c \times n}$: binary codes, where c is the code length.
- $-\mathbf{F} \in \mathbb{R}^{c \times n} \text{ with } \mathbf{F}_{*i} = f(\mathbf{x}_i; \theta_x).$
- $-\mathbf{G} \in \mathbb{R}^{c \times n} \text{ with } \mathbf{G}_{*j} = g(\mathbf{y}_j; \theta_y).$
- $-\Theta_{ij} = \frac{1}{2}\mathbf{F}_{*i}^T\mathbf{G}_{*j}$.

Learning

Alternating Learning Algorithm

- Learn θ_x , with θ_y and B fixed. BP for updating θ_x . For each sampled point \mathbf{x}_i , compute the gradient:
 - $\frac{\partial \mathcal{J}}{\partial \mathbf{F}_{*i}} = \frac{1}{2} \sum_{j=1}^{n} (\sigma(\Theta_{ij}) \mathbf{G}_{*j} S_{ij} \mathbf{G}_{*j}) + 2\gamma (\mathbf{F}_{*i} \mathbf{B}_{*i}) + 2\eta \mathbf{F1}.$
- Learn θ_y , with θ_x and B fixed.

BP for updating θ_y . For each sampled point y_j , compute the gradient:

$$\frac{\partial \mathcal{J}}{\partial \mathbf{G}_{*j}} = \frac{1}{2} \sum_{i=1}^{n} (\sigma(\Theta_{ij}) \mathbf{F}_{*i} - S_{ij} \mathbf{F}_{*i}) + 2\gamma (\mathbf{G}_{*j} - \mathbf{B}_{*j}) + 2\eta \mathbf{G} \mathbf{1}.$$

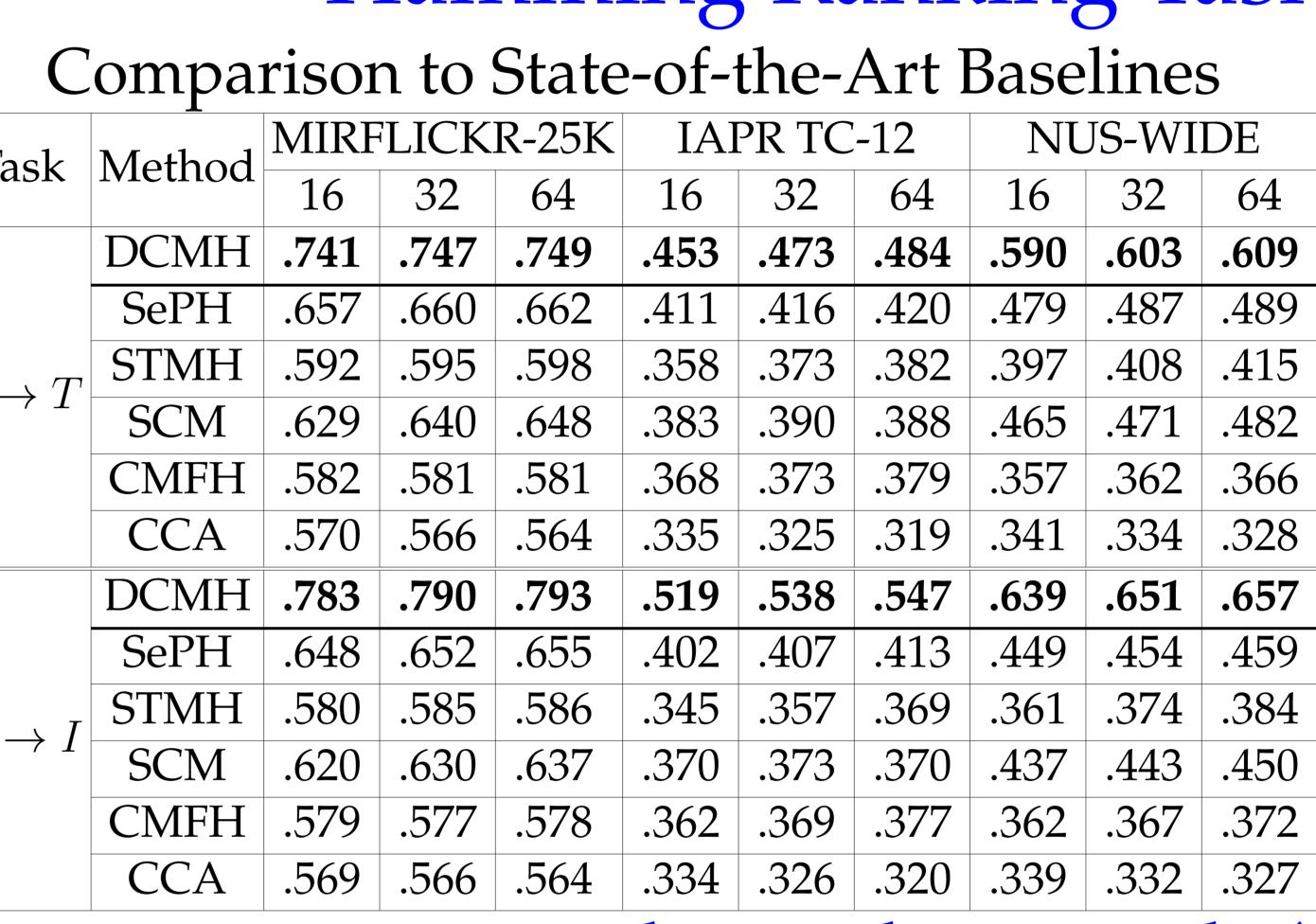
• Learn B, with θ_x and θ_y fixed.

$$\mathbf{B} = \operatorname{sign}(\gamma(\mathbf{F} + \mathbf{G})).$$

Experiment

- •MIRFLICKR-25K: 25,000 image-text pairs which are annotated with one of the 24 unique labels.
- •IAPR TC-12: 20,000 image-text pairs which are annotated using 255 labels.
- •NUS-WIDE: 260,648 image-text pairs. Each point is annotated with one or multiple labels from 81 concept labels. We select 195,834 image-text pairs that belong to the 21 most frequent concepts.
- For MIRFLICKR-25K and IAPR TC-12: 2000/10000 test/training points. For NUS-WIDE: 2100/10500 test/training points.

Hamming Ranking Task (Mean Average Precision)

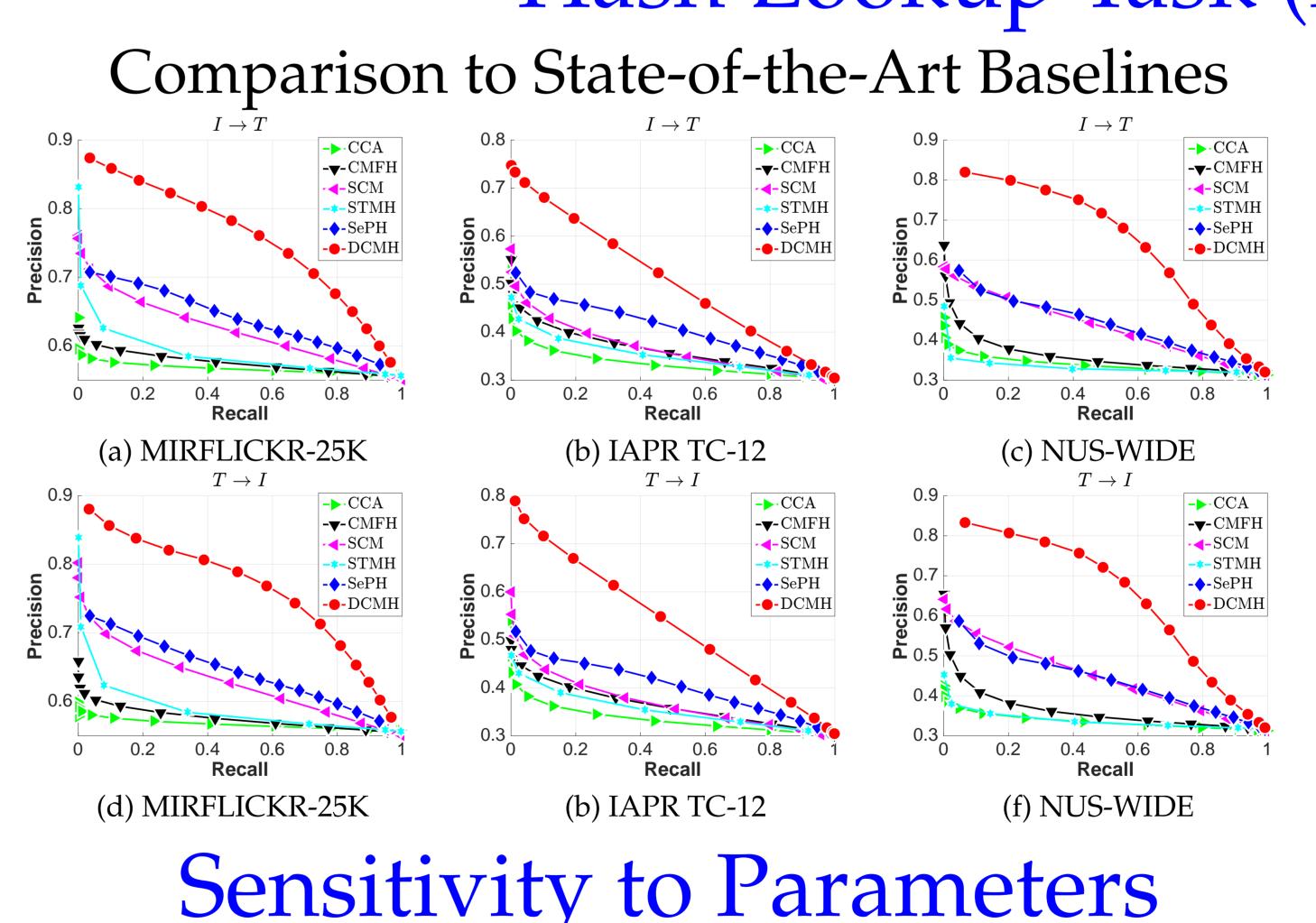


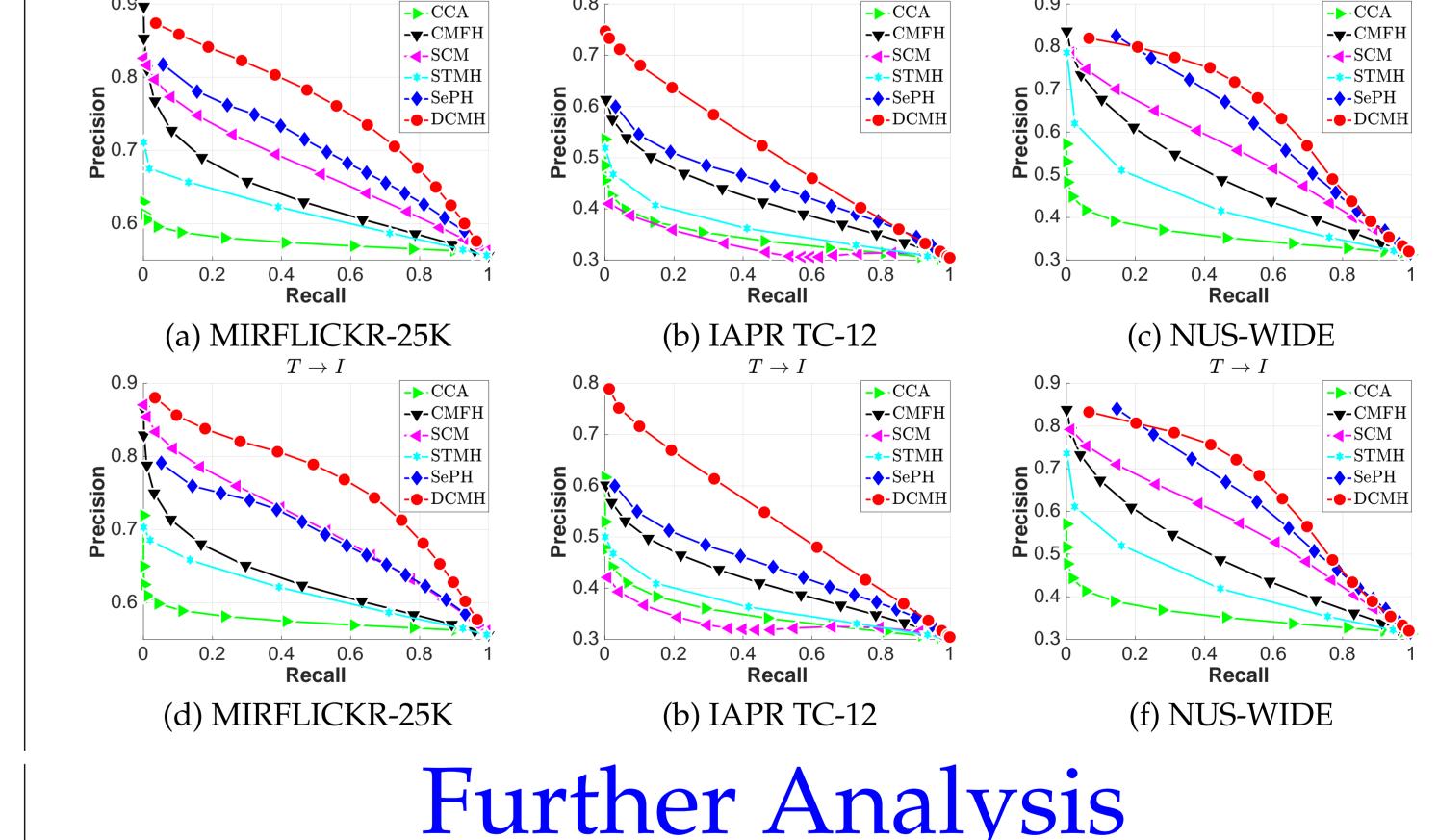
Comparison to Baselines with CNN-F Features SePH .712 .719 .723 .444 .456 .464 .604 .617 .621 SePH .722 .726 .732 .442 .456 .465 .598 .603 .611 CMFH | .637 | .640 | .643 | .417 | .421 | .428 | .503 | .519 | .523

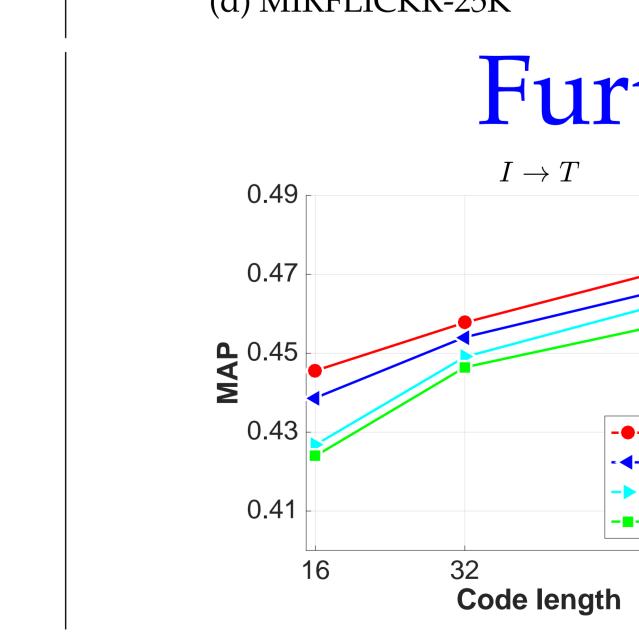
CCA .574 .571 .569 .349 .344 .338 .361 .349 .340

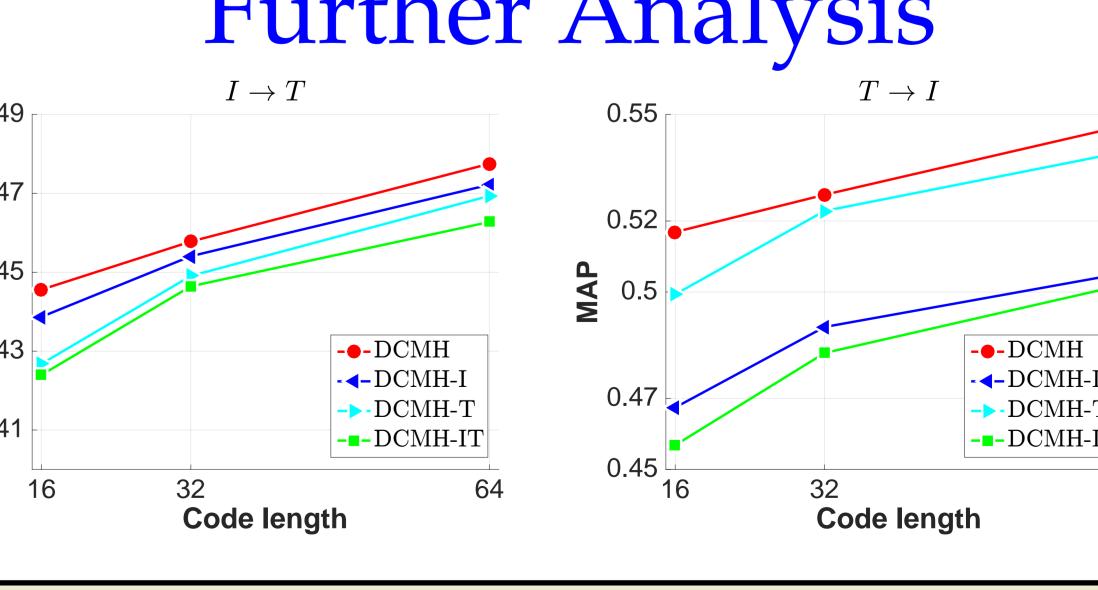
Comparison to Baselines with CNN-F Features

Hash Lookup Task (Precision Recall Curve)









Conclusion

- •DCMH is an end-to-end deep learning framework which can perform simultaneous feature learning and hash-code learning.
- •DCMH can significantly outperform other baselines to achieve the state-of-the-art performance.