



Northeastern University



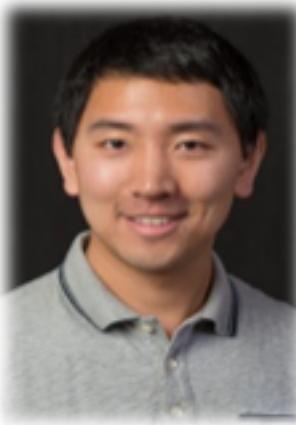
Smile^{lab}
Synergetic Media Learning Lab

Multi-view Visual Data Analytics

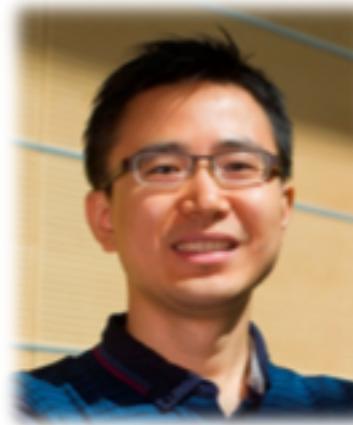
---CVPR-2018 Tutorial



Zhengming Ding



Ming Shao



Yun Fu

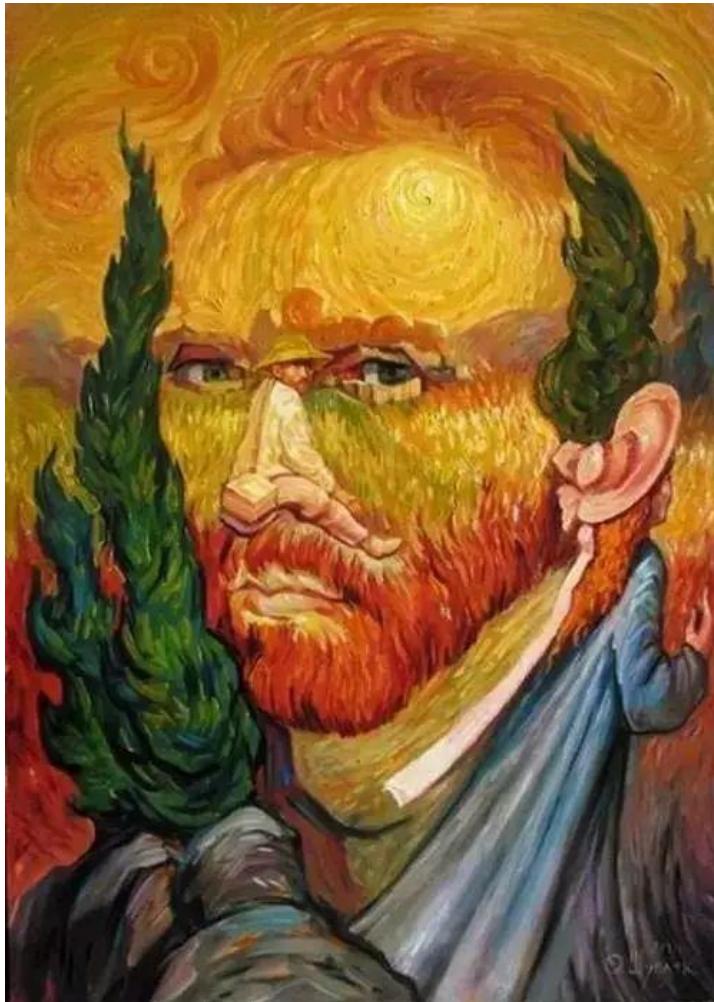


An Example from Life





An Example from Art



Outline

□ Introduction & Background

- Multi-view Visual Data
- Multi-view Learning Problems
- Multi-view Learning Taxonomy

□ Multi-view Learning

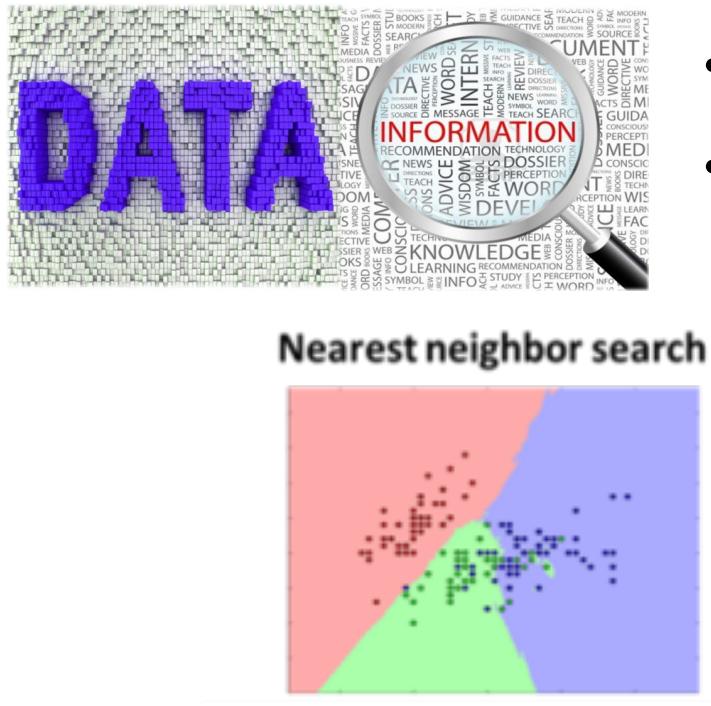
- Projection and Embedding
- Knowledge Fusion
- Multi-view Clustering
- Supervised Multi-view Learning → Zero-shot Learning

□ Domain Adaptation

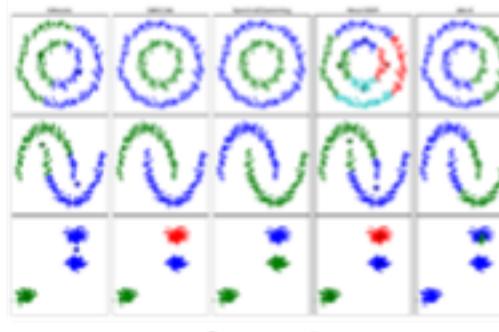
- Transfer Learning → Domain Adaptation
- Multi-Source Domain Adaptation & Domain Generalization

□ Conclusion

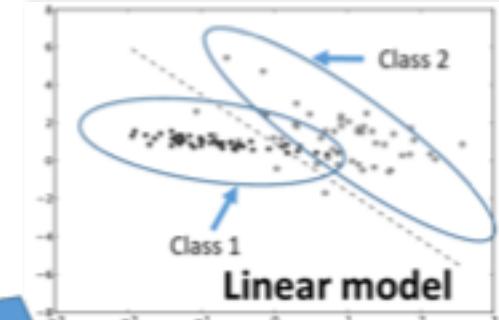
What is Multi-view Learning in General?



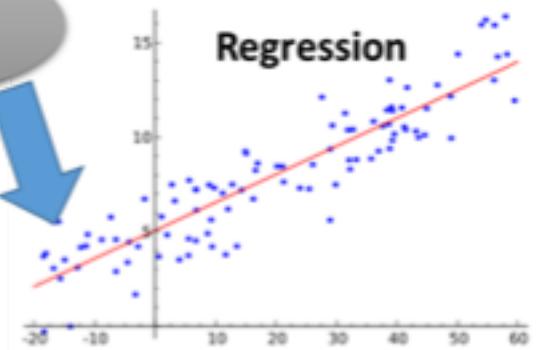
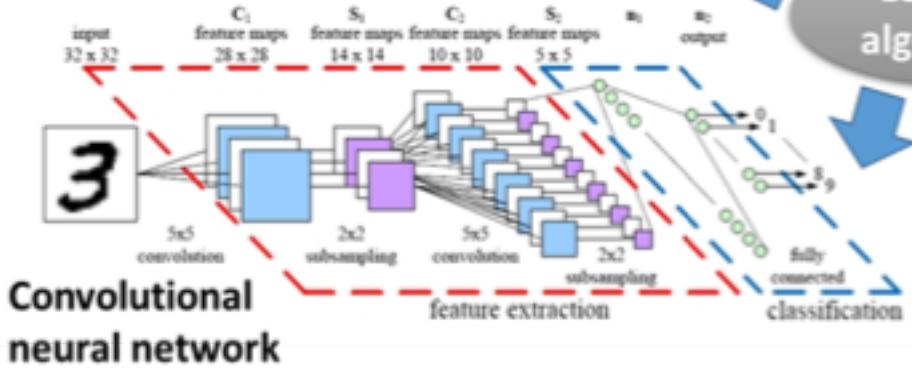
- What **Data** we are interested in??
- What **Learning Problems** we will formulate??



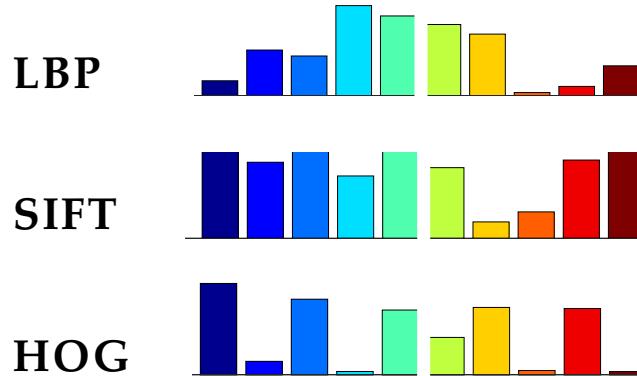
Clustering



Learning
algorithms



Multi-view Visual Data



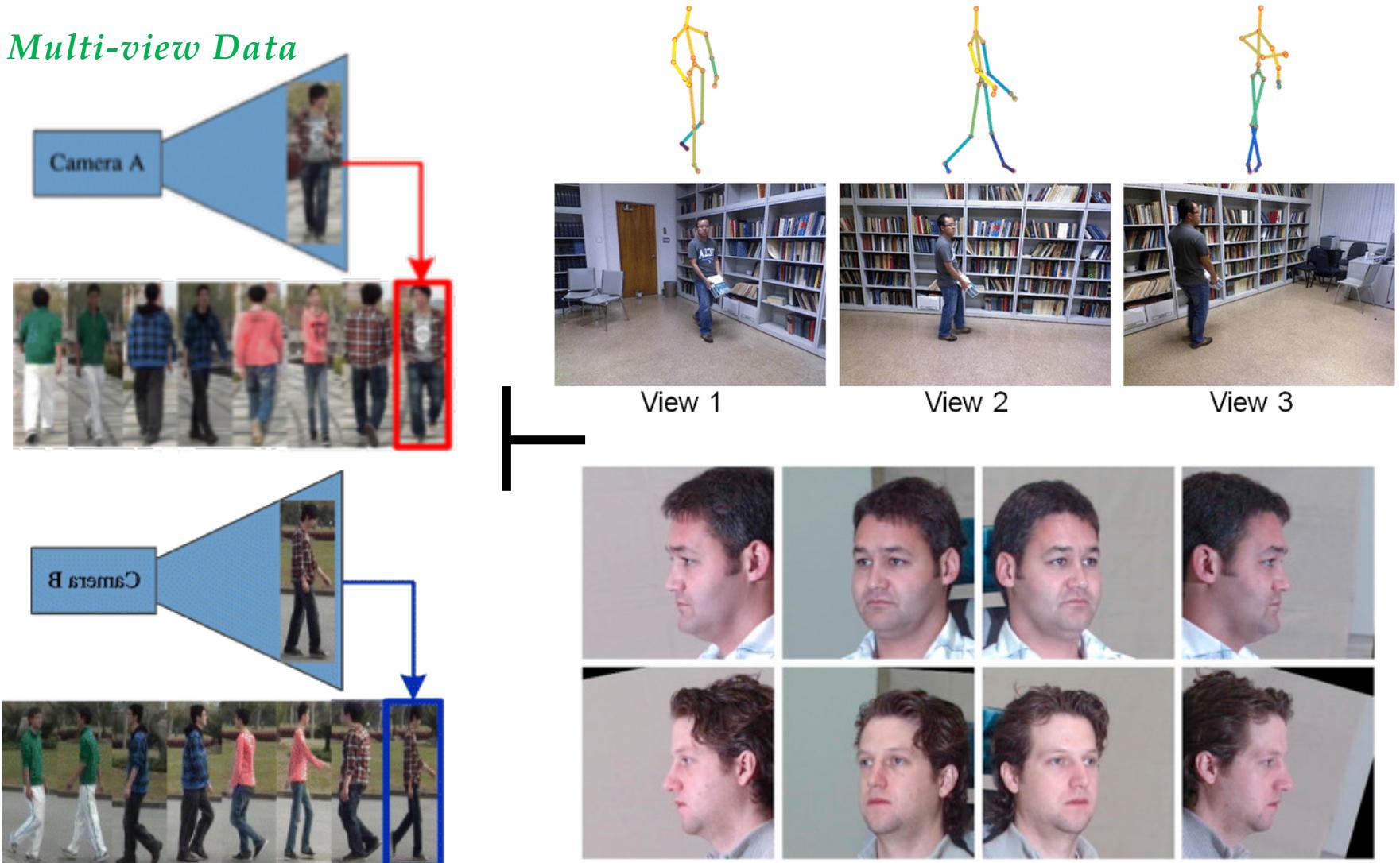
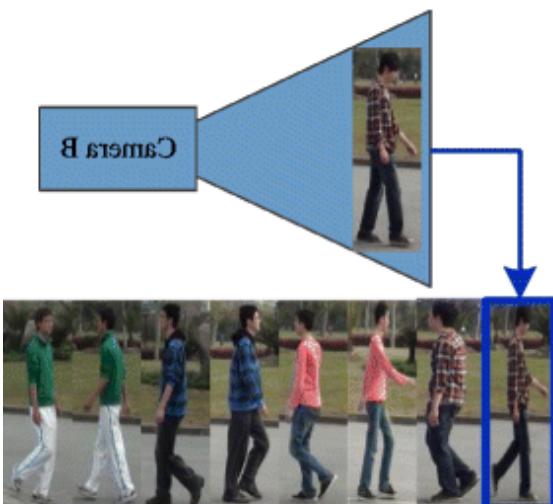
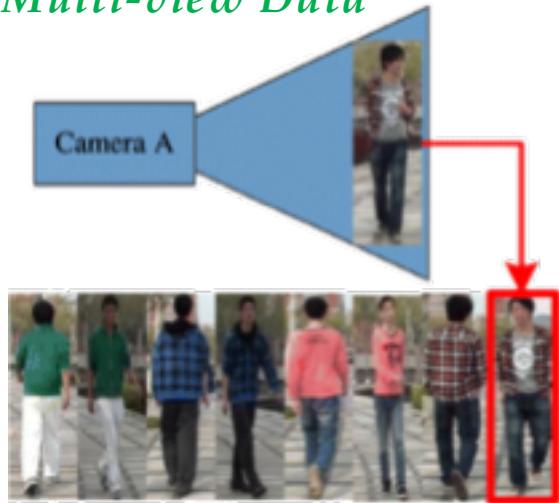
Multiple features



YAHOO!
NEWS

Multi-view Visual Data

Multi-view Data



Multi-view Visual Data

Multi-modal Data



Clip art



Sketches

There is a bed with a striped bedspread. Beside this is a nightstand with a drawer. There is also a tall dresser and a chair with a blue cushion. On the dresser is a jewelry box and a clock.

I am inside a room surrounded by my favorite things. This room is filled with pillows and a comfortable bed. There are stuffed animals everywhere. I have posters on the walls. My jewelry box is on the dresser.

There are brightly colored wooden tables with little chairs. There is a rug in one corner with ABC blocks on it. There is a bookcase with picture books, a larger teacher's desk and a chalkboard.

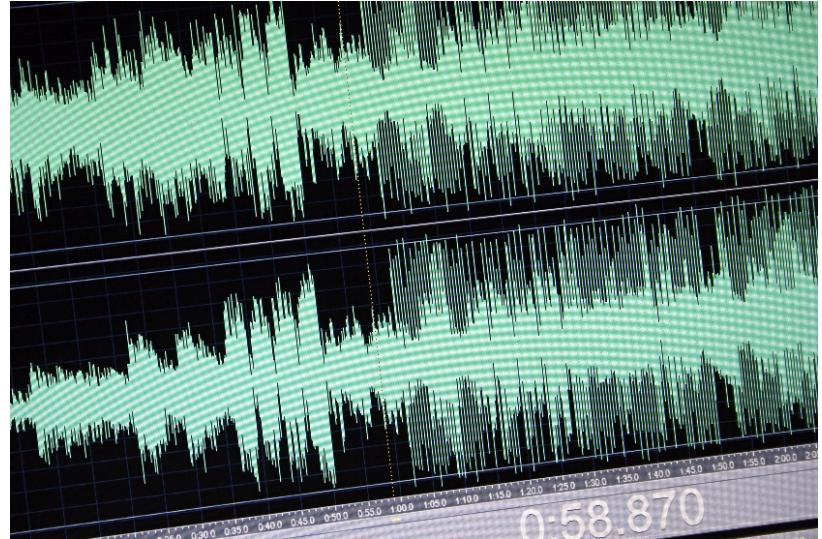
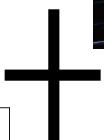
ceiling wall	wall	boat boat building building	boat sky building building	wall railing	wall wall	sky car	sky	ceiling wall wall wall	wall wall wall	ceiling wall	wall wall	building wall table	ceiling wall wall wall floor	
floor wall	boat	water water	wall	wall	road road	road road	floor	floor	floor	person	floor	sky	building	floor

Spatial text

Multi-view Visual Data



Multi-modal Data



automated data mining survey
responses computer transcripts
qualitative root cause
classification insights
ad-hoc analysis product
reviews sentiment analysis
customer dashboards consumer
trends ad-hoc analysis early warning

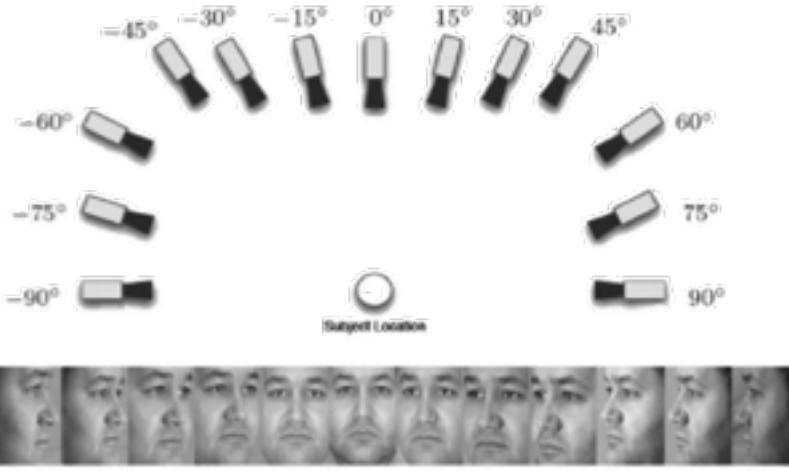
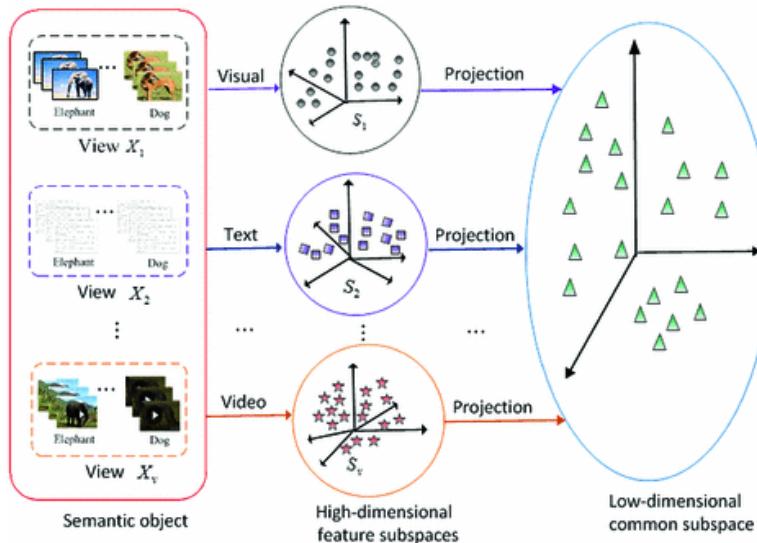


감사합니다 Natick
Danke Ευχαριστίες Dalu
Thank You Köszönöm
Tack Gracias
Спасибо Dank Seé
谢谢 Merci ありがとう



Multi-view Problems

□ Multi-view Projection/Embedding



□ Multi-view Clustering

□ Supervised Multi-view Learning

- Multi-view Classification
- Zero-Shot Learning



This is a cat. Its name is Sam. It is grey. It is fluffy.

It has got a little head and two big yellow eyes. It has got two little ears, a brown nose and a pink mouth. It has got four long legs. It has got a long fluffy tail.

Sam is very funny. It can run, jump and climb. It can't swim and fly.

It likes fish and milk.

fluffy - пухастый, tail - хвостик



Multi-view Problems

□ Transfer Learning

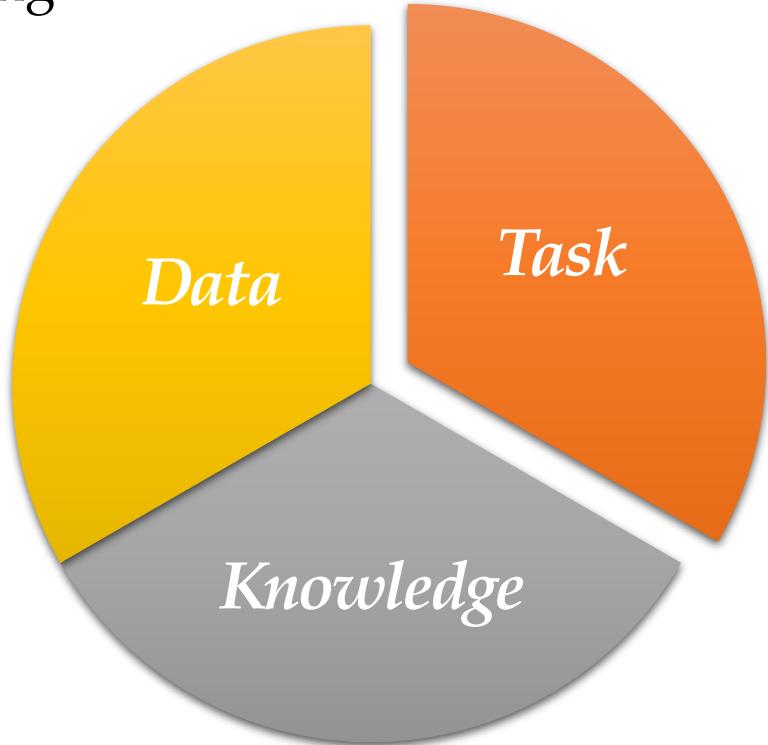
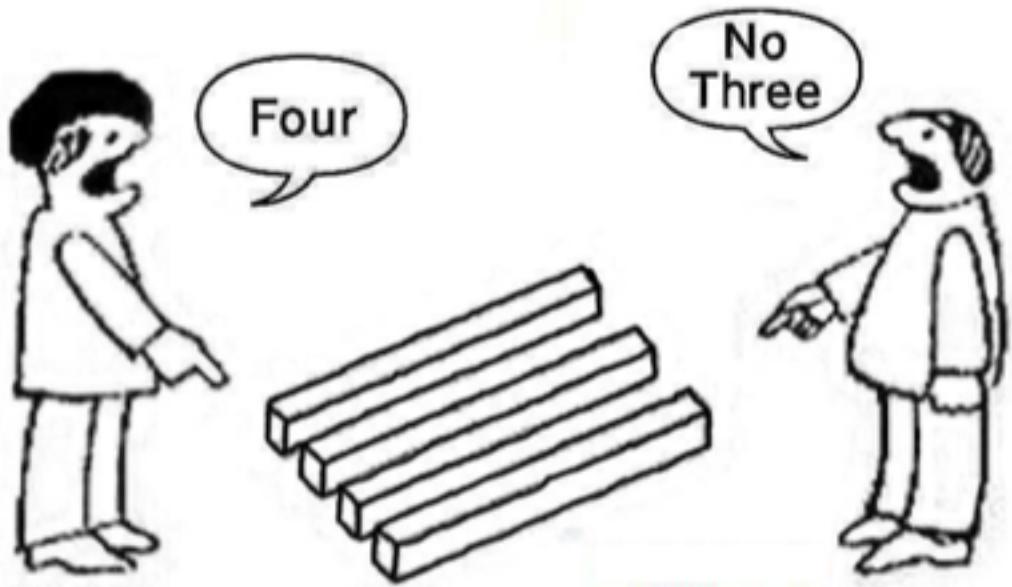


□ Domain Adaptation

- Multi-source Domain Adaptation
- Domain Generalization

Multi-view Learning: How to Categorize??

- Different Data and Different Learning Problems, and possibly more...
- We will offer a better Taxonomy:



Taxonomy [Data View]

□ Category 1 [Sample-Wise Correspondence] (Multi-view Learning)

- **Multiple Features**, e.g., LBP, SIFT, HOG...

Goal: fuse various knowledge from multiple features to boost the final tasks

- **Multi-Modal Visual Data**

Goal: seek a view-invariant space to mitigate the view divergence to facilitate the final task
(adapt knowledge across different views)

□ Category 2 [Class-wise Correspondence] (Transfer learning, Projection/Embedding)

- **Multi-Feature/Multi-Pose/Multi-Modal Visual data**

Goal: transfer knowledge from well-labeled source views to unlabeled target views

Taxonomy [Task View]

□ Unsupervised Learning [Clustering, Projection/Embedding]

Goal: fuse various knowledge from multiple features

Data: multiple features [*unlabeled, sample-wise correspondence*]

□ Supervised Learning [Projection/Embedding, Classification]

Goal: seek a view-invariant space to mitigate the view divergence to facilitate the final task (*adapt knowledge across different views*)

- Sub-Category 1 (Multi-view Learning) [sample-wise correspondence]

Training Stage: multiple labeled view data

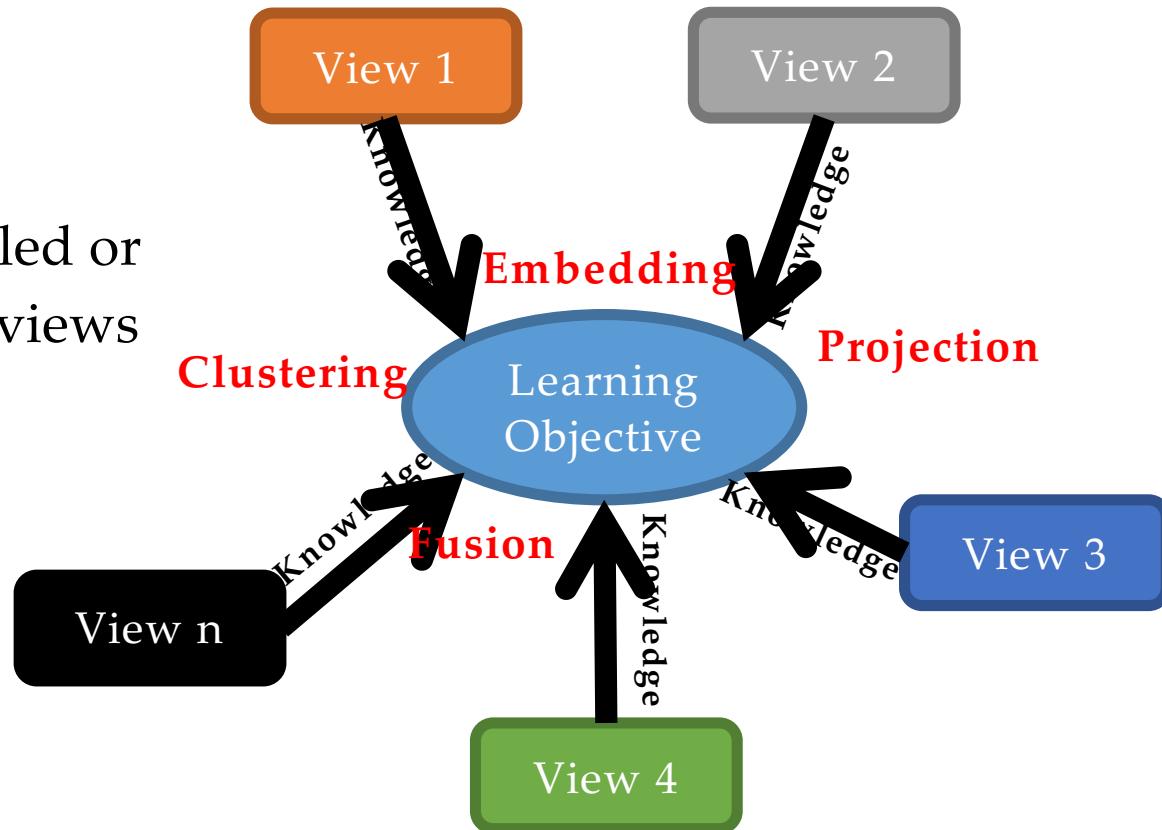
Test Stage: some labeled views, some unlabeled views

- Sub-Category 2 (Transfer learning) [class-wise correspondence]

Training Stage: some source labeled views & some target unlabeled views

Taxonomy [Knowledge View]

- Knowledge integration from different views
- **Goal:** Find a way to better integrate knowledge (labeled or unlabeled) from different views for a common aim:
 - Multi-view Clustering
 - Multi-view Projection and Embedding
 - Knowledge Fusion

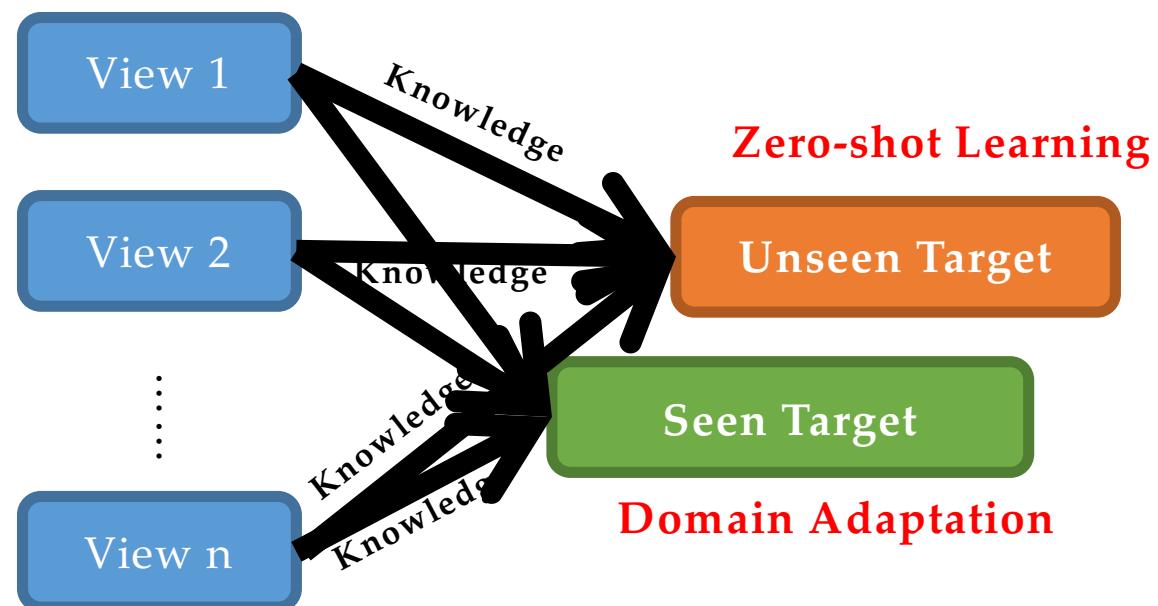


Taxonomy [Knowledge View]

- Knowledge transfer from one view(s) to another view(s)

- **Goal:** Reuse knowledge from well-established data in a new problem (seen or unseen)

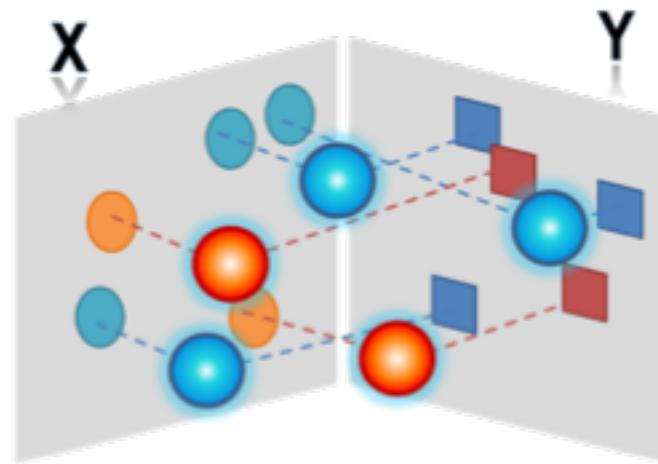
- Transfer Learning
- Domain Adaptation
- Zero-shot Learning



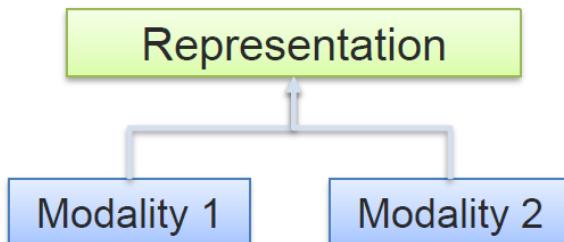
Multi-view Algorithms

❖ Representation learning

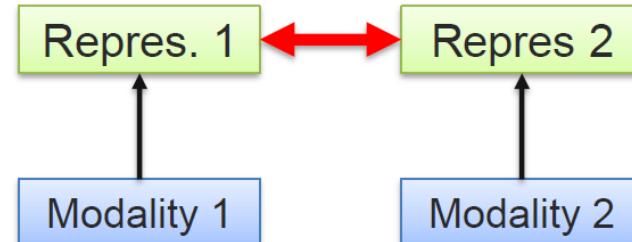
- Joint representations
- Coordinated representations
- Projection and embedding
- Similarity metrics, dictionary learning
- Multi-view AutoEncoder, CNN, GAN



A Joint representations:



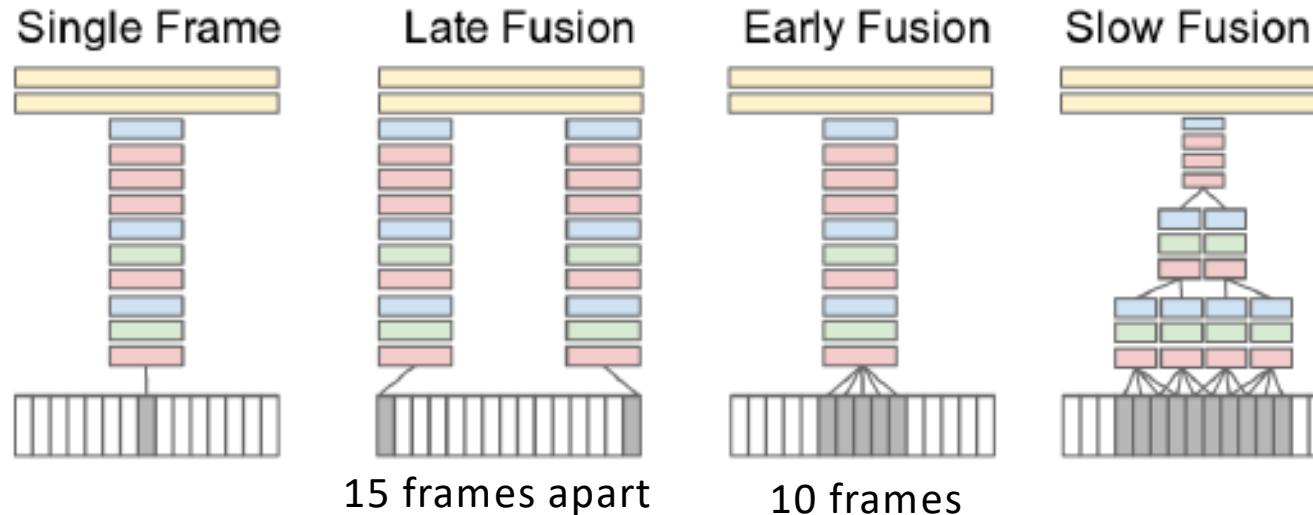
B Coordinated representations:



Multi-view Algorithms

❖ Fusion

To join information from two or more views to perform a prediction task



[1] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, Fei-Fei Li: Large-Scale Video Classification with Convolutional Neural Networks. CVPR 2014: 1725-1732

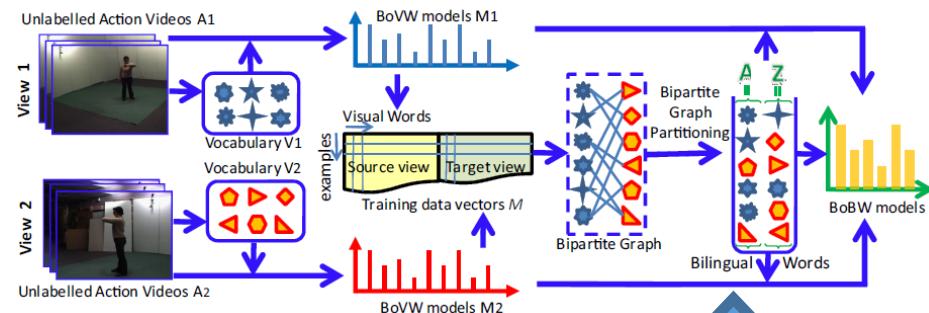
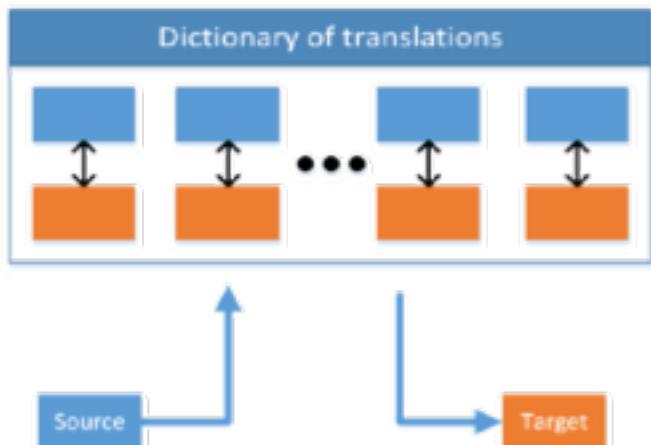
[2] Vrigkas, Michalis, Christophoros Nikou, and Ioannis A. Kakadiaris. "A review of human activity recognition methods." *Frontiers in Robotics and AI* 2 (2015): 28.

Multi-view Algorithms

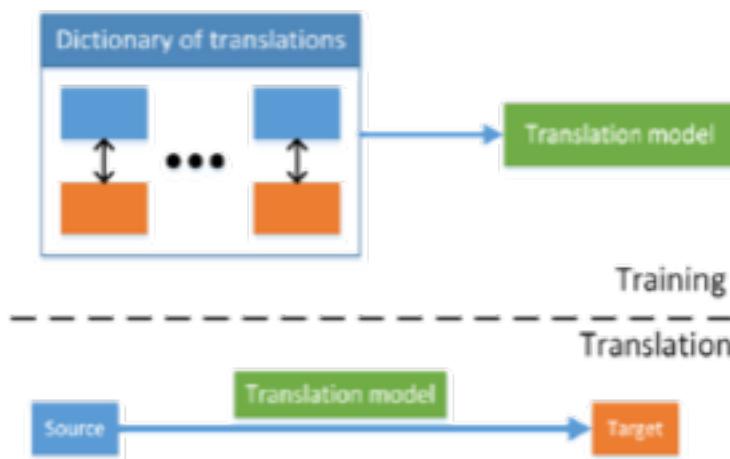
❖ Translation

Process of changing data from one modality to another, where the translation relationship can often be open-ended or subjective.

A Example-based



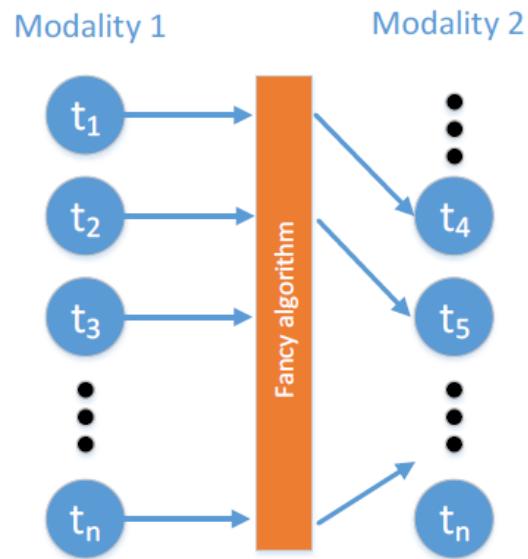
B Model-driven



Multi-view Algorithms

❖ Alignment

Identify the direct relations between (sub)elements from two or more different modalities.



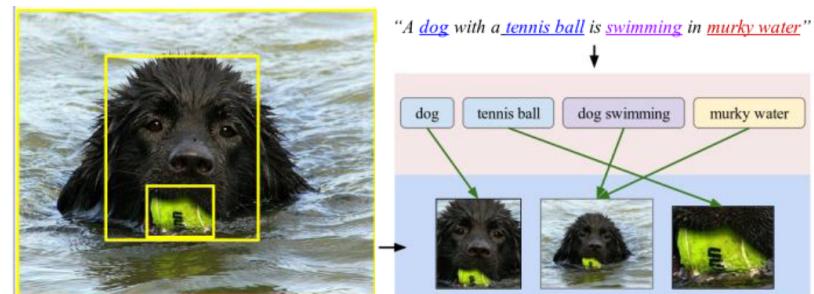
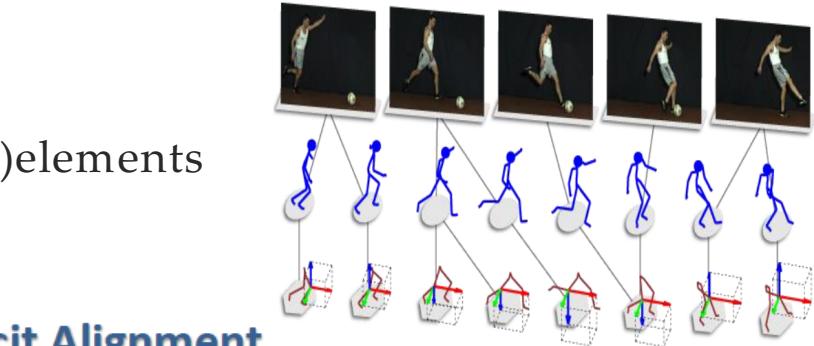
A Explicit Alignment

The goal is to directly find correspondences between elements of different modalities

B Implicit Alignment

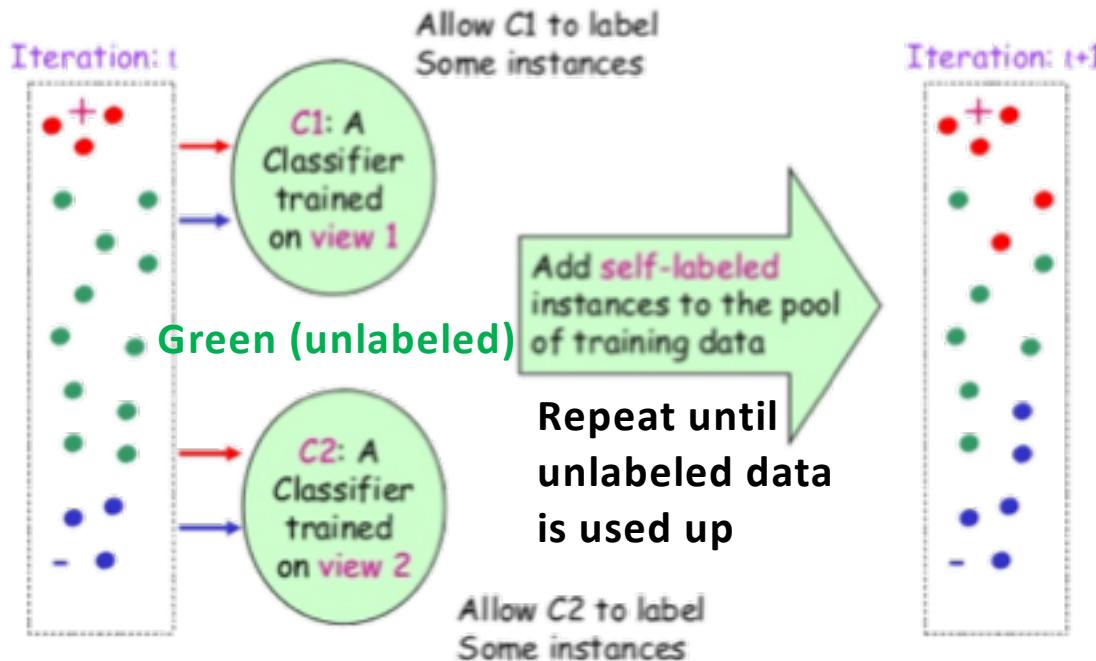
Uses internally latent alignment of modalities in order to better solve a different problem

Karpathy et al., Deep Fragment Embeddings for Bidirectional Image Sentence Mapping,
<https://arxiv.org/pdf/1406.5679.pdf>

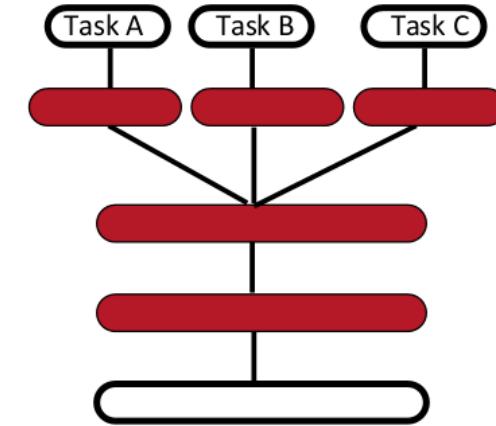


Multi-view Algorithms

Co-Learning & Multi-task Learning



Two views are conditionally independent



Learning tasks in parallel while using a shared representation; what is learned for each task can help other tasks be learned better



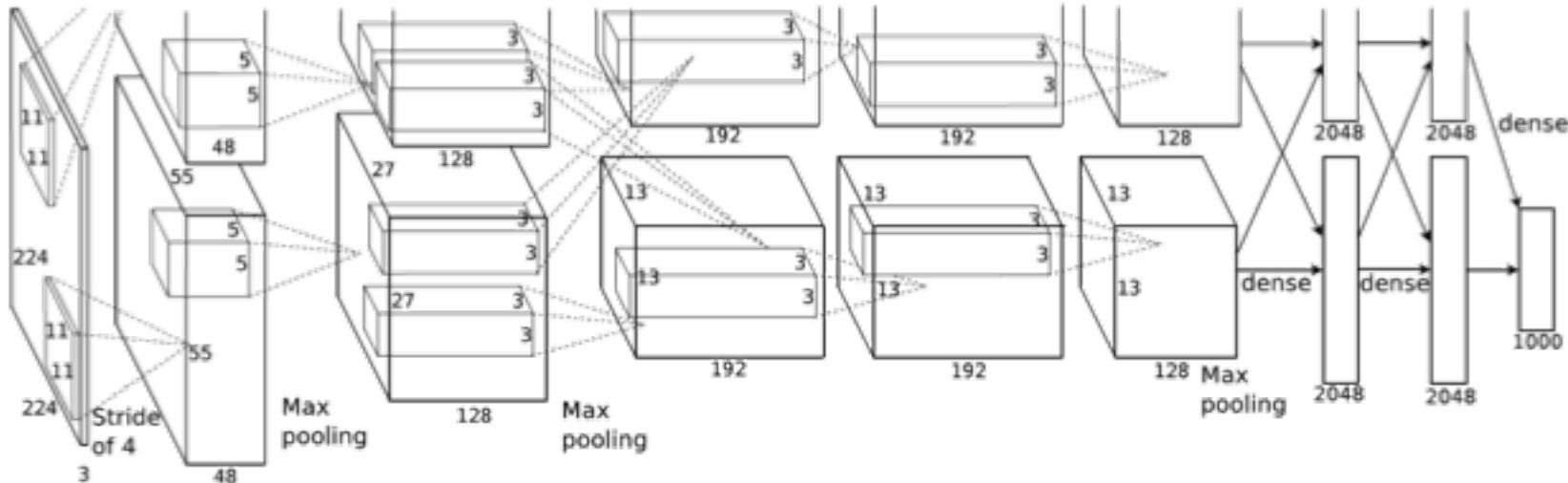
Unified Model

- The Tutorial will focus on the
Representation Learning and Knowledge Fusion/Alignment
- Assume we have more than one data domains: X_1, X_2, \dots, X_i . In general, the discussed methods are organized in three lines:
 - ❖ Modeling Features Representation for each source: $f_i(X_i)$
 - ❖ Modeling Coherence/Alignment between different sources: $A(X_i, X_j)$
 - ❖ Regularization term regarding the label information and data underlying distribution, and semantics: $R(X_i)$

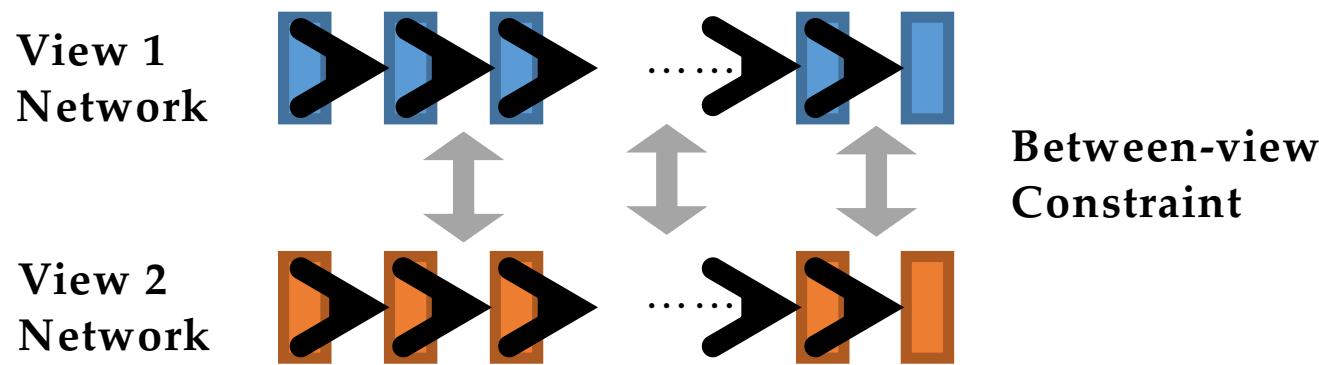
$$\min_{f_1(\cdot), \dots, f_v(\cdot)} \sum_{i=1, i < j}^v \mathcal{A}(f_i(X_i), f_j(X_j)) + \lambda \sum_{k=1}^v \mathcal{R}(f_k(X_k))$$

Unified Model

- We barely consider the deep features in the learning process...



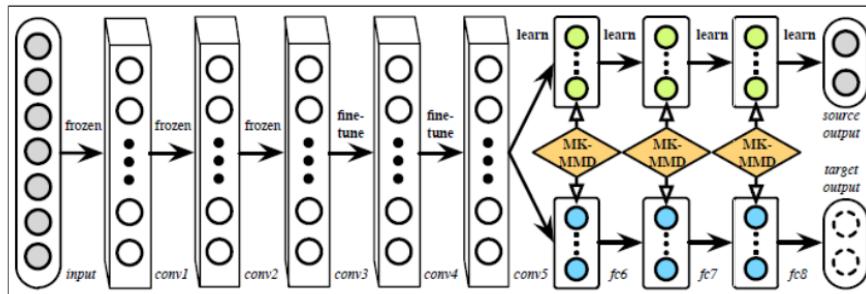
- When consider multi-view learning, what that would be??



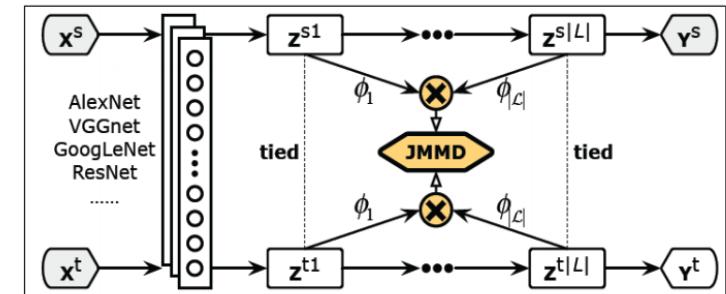
Unified Model

□ From Shallow to Deep Learning

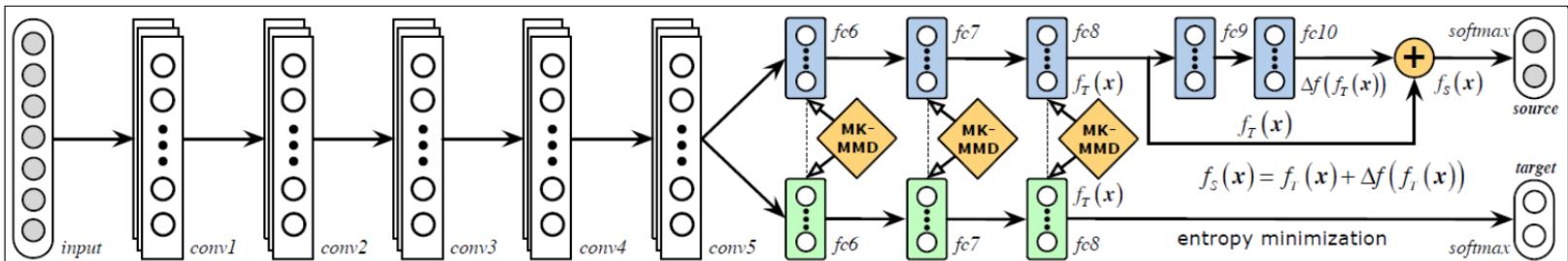
- Redistribute the constraints above into the deep structure, then which layer??
- Usually modeled on top layers, but may be from the first layer
- Multiple networks to model different aspects, e.g., one for different domains, the other for different classes



(a) The Deep Adaptation Network (DAN) architecture



(b) The Joint Adaptation Network (JAN) architecture



(c) The Residual Transfer Network (RTN) architecture

Unified Model

Representation Learning + Multi-view Alignment (Fusion)

$$\min_{f_1(\cdot), \dots, f_v(\cdot)} \sum_{i=1, i < j}^v \mathcal{A}(f_i(X_i) | f_j(X_j)) + \lambda \sum_{k=1}^v \mathcal{R}(f_k(X_k))$$

Linear Mapping → Kernel → Tensor

Dictionary Learning (Sparse/Low-Rank Coding)

Auto-Encoder & Neural Networks

Convolutional Neural Networks

- Category 1 [Sample-Wise Correspondence] (Multi-view Learning)
- Category 2 [Class-wise Correspondence] (Transfer learning)

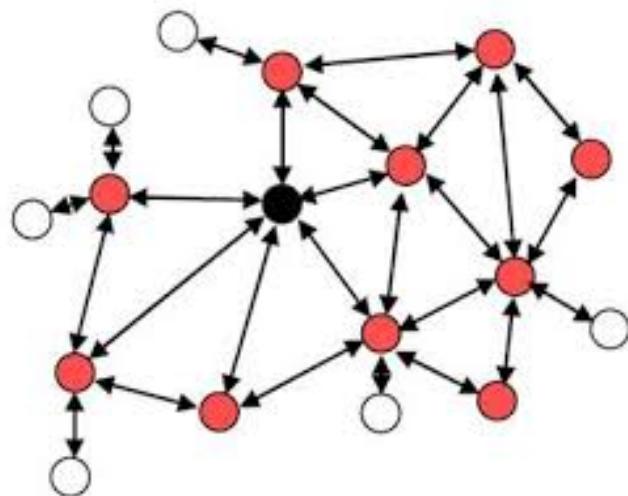
- Maximum Mean Discrepancy
- Reconstruction-based Alignment
- Adversarial loss [0/1]

Unified Model

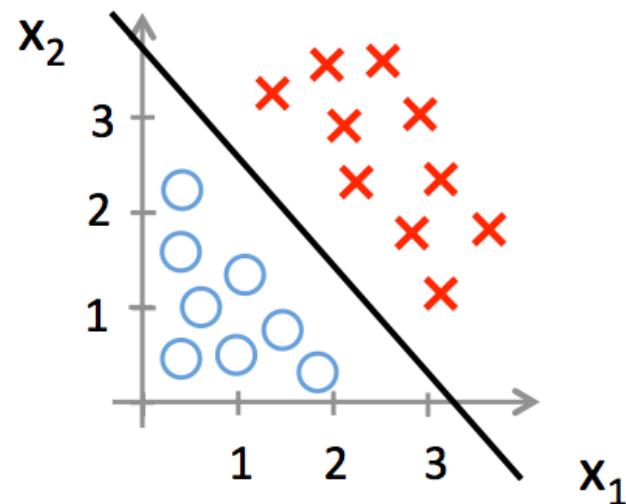
Representation Learning + Multi-view Alignment (Fusion)

$$\min_{f_1(\cdot), \dots, f_v(\cdot)} \sum_{i=1, i < j}^v \mathcal{A}(f_i(X_i), f_j(X_j)) + \lambda \sum_{k=1}^v \boxed{\mathcal{R}(f_k(X_k))}$$

Graph Regularizer



Regression model



Outline

Northeastern University



Smile
Synergetic Media Learning Lab

□ Introduction & Background

- Multi-view Visual Data
- Multi-view Learning Problems
- Multi-view Learning Taxonomy

□ Multi-view Learning

- Projection and Embedding
- Knowledge Fusion
- Multi-view Clustering
- Supervised Multi-view Learning → Zero-shot Learning

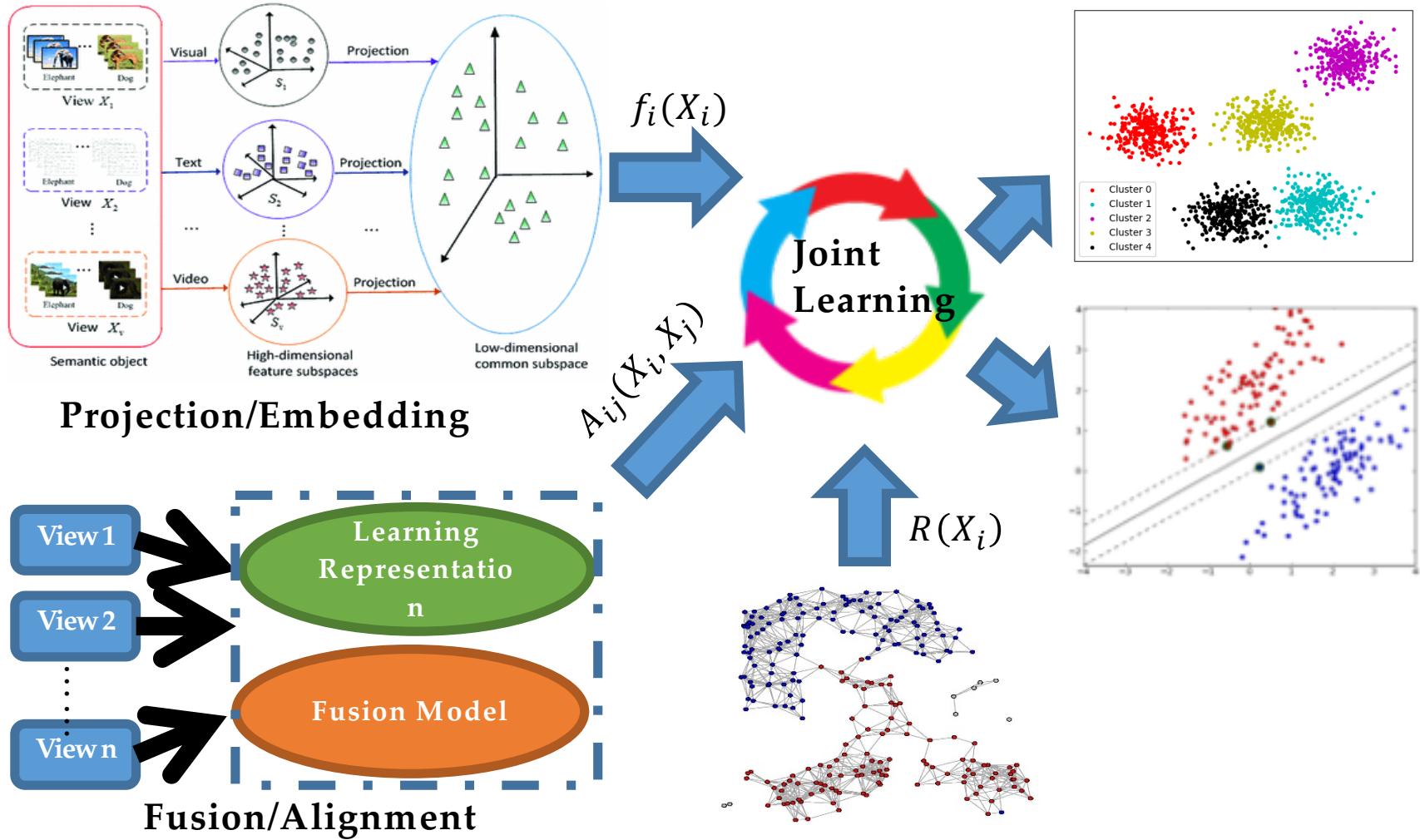
□ Domain Adaptation

- Transfer Learning → Domain Adaptation
- Domain Generalization → Zero-shot Learning

□ Conclusion



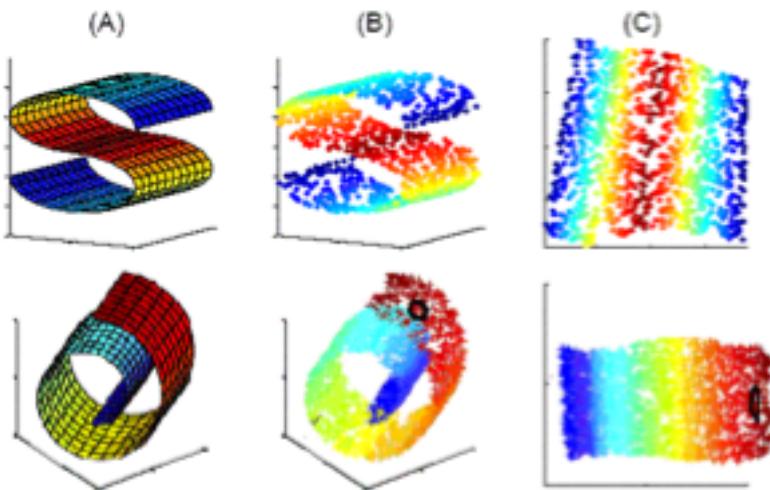
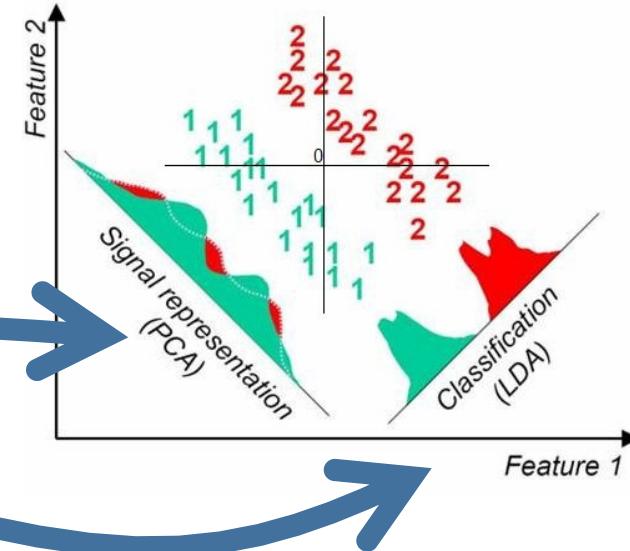
Multi-view Learning RoadMap



Background (Dimensionality Reduction)

• Projection:

- Given the input data $X \in \mathbb{R}^{n \times D}$, find a Linear Mapping: $f(X) = Y$ that will map X into a lower dimensional space $Y \in \mathbb{R}^{n \times d}$, and f is usually a projection matrix.
- Typical methods: PCA, LDA, LPP



• Nonlinear Embedding:

- Given the input data $X \in \mathbb{R}^{n \times D}$, find a function: $f(X) = Y$ that will map X into a lower dimensional space $Y \in \mathbb{R}^{n \times d}$, and f is usually a mapping with implicit formulation.
- Typical methods: LLE, ISOMAP

Multi-view Projection: Canonical Correlation Analysis (CCA)

Canonical Correlation Analysis (CCA)

Formally, for two views $X \in R^{d \times n}$ and $Y \in R^{k \times n}$, CCA computes two projection vectors, $w_x \in R^d$ and $w_y \in R^k$, such that the following correlation coefficient is maximized:

$$\rho = \frac{w_x^T X Y^T w_y}{\sqrt{(w_x^T X X^T w_x)(w_y^T Y Y^T w_y)}}$$

Since ρ is invariant to the scaling of w_x and w_y , CCA can be formulated equivalently as

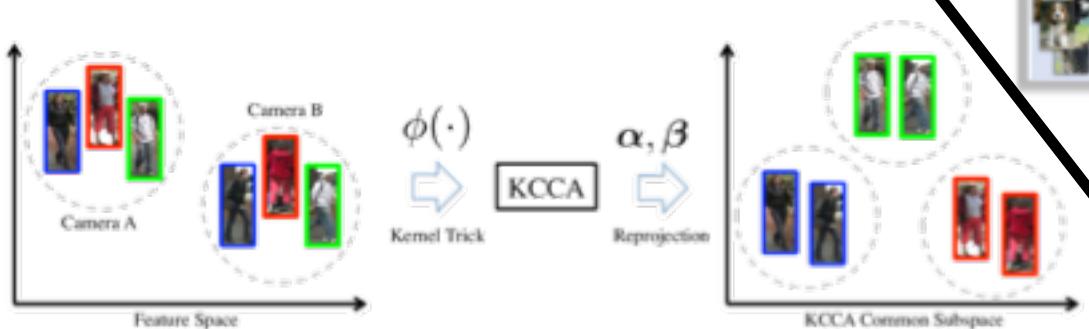
$$\begin{aligned} \max_{w_x, w_y} \quad & w_x^T X Y^T w_y \\ \text{s.t.} \quad & w_x^T X X^T w_x = 1, \quad w_y^T Y Y^T w_y = 1 \end{aligned}$$

✓ It has a sense of consensus.

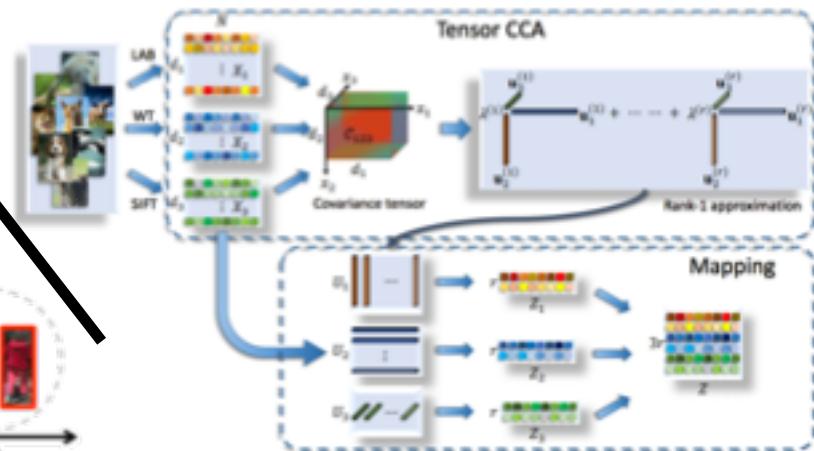
Multi-view Projection: Canonical Correlation Analysis (CCA)

Extensions of CCA

Kernel-based



Tensor-based

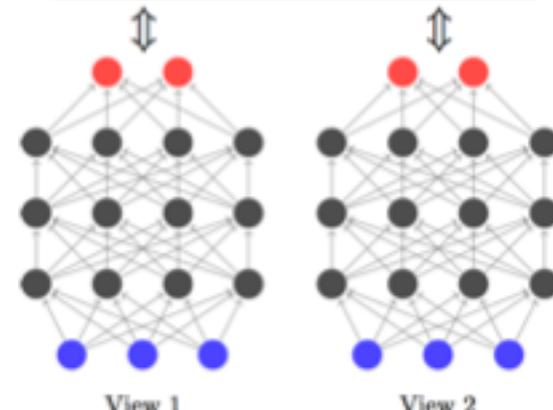


[Kernel-based] David R. Hardoon, Sándor Székely, John Shawe-Taylor: Canonical Correlation Analysis: An Overview with Application to Learning Methods. *Neural Computation* 16(12): 2639-2664 (2004)

[Tensor-based] Tensor Canonical Correlation Analysis for Multi-view Dimension Reduction, Yong Luo, Dacheng Tao, Kotagiri Ramamohanarao, Chao Xu, and Yonggang Wen, *IEEE Transactions on Knowledge and Data Engineering (T-KDE)*, vol. 27, no. 11, pp. 3111-3124, 2015.

[DeepNN-based] Galen Andrew, Raman Arora, Jeff A. Bilmes, Karen Livescu: Deep Canonical Correlation Analysis. *ICML* (3) 2013: 1247-1255

Canonical Correlation Analysis



DNN-based

Multi-view Projection: Canonical Correlation Analysis (CCA)

Extensions of CCA

Multiple Sets

Assume we will apply CCA on multiple views (more than two)

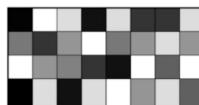


Even more... 

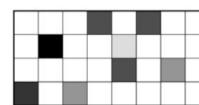
We will conduct this in a pairwise fashion: X and Y , Y and Z , Z and Y

$$\rho = \frac{w_x^T C_{xy} w_y}{\sqrt{(w_x^T C_{xx} w_x)(w_y^T C_{yy} w_y)}} + \frac{w_y^T C_{yz} w_z}{\sqrt{(w_y^T C_{yy} w_y)(w_z^T C_{zz} w_z)}} + \frac{w_z^T C_{xz} w_x}{\sqrt{(w_z^T C_{zz} w_z)(w_x^T C_{xx} w_x)}}$$

Sparsity



Dense Projection



Sparse Projection

Motivated by good feature selection and stability of the features, we pursue sparsity in CCA:

$$\rho = \max_{w_x, w_y} \frac{w_x^T C_{xy} w_y}{\sqrt{w_x^T C_{xx} w_x w_y^T C_{yy} w_y}} \\ \text{s.t. } \|w_x\|_0 \leq s_x, \quad \|w_y\|_0 \leq s_y.$$

[**Multiple Sets**] Kettenring, J.R. Canonical analysis of several sets of variables. *Biometrika* (1971)

[**Sparsity**] A. Wiesel, M. Kliger, and A. O.

Hero, III, "A greedy approach to sparse canonical correlation analysis," ArXiv e-prints, 2008.



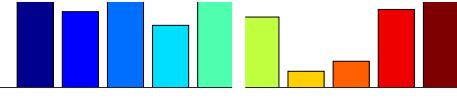
Multi-view Embedding

- Embedding offers a way to find the low-dimensional representation through **an implicit mapping**
- For multi-view data, how to do embedding in a joint manner?

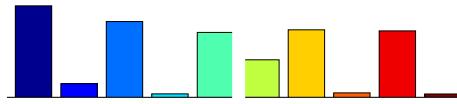
LBP



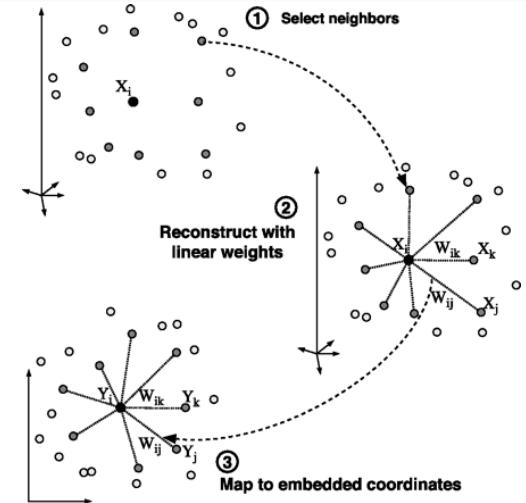
SIFT



HOG



From nD
to 2D



$$\arg \min_{Y, c} \sum_{v=1}^V c_v^T \text{tr}(Y \mathcal{L}^v Y^T)$$

Solved in an
alternative fashion

$$\text{s.t. } YY^T = I; \sum_{v=1}^V c_v = 1, c_v \geq 0.$$

[1] Tian Xia, Dacheng Tao, Tao Mei, Yongdong Zhang: Multiview Spectral Embedding. IEEE Trans. Systems, Man, and Cybernetics, Part B 40(6): 1438-1446 (2010)

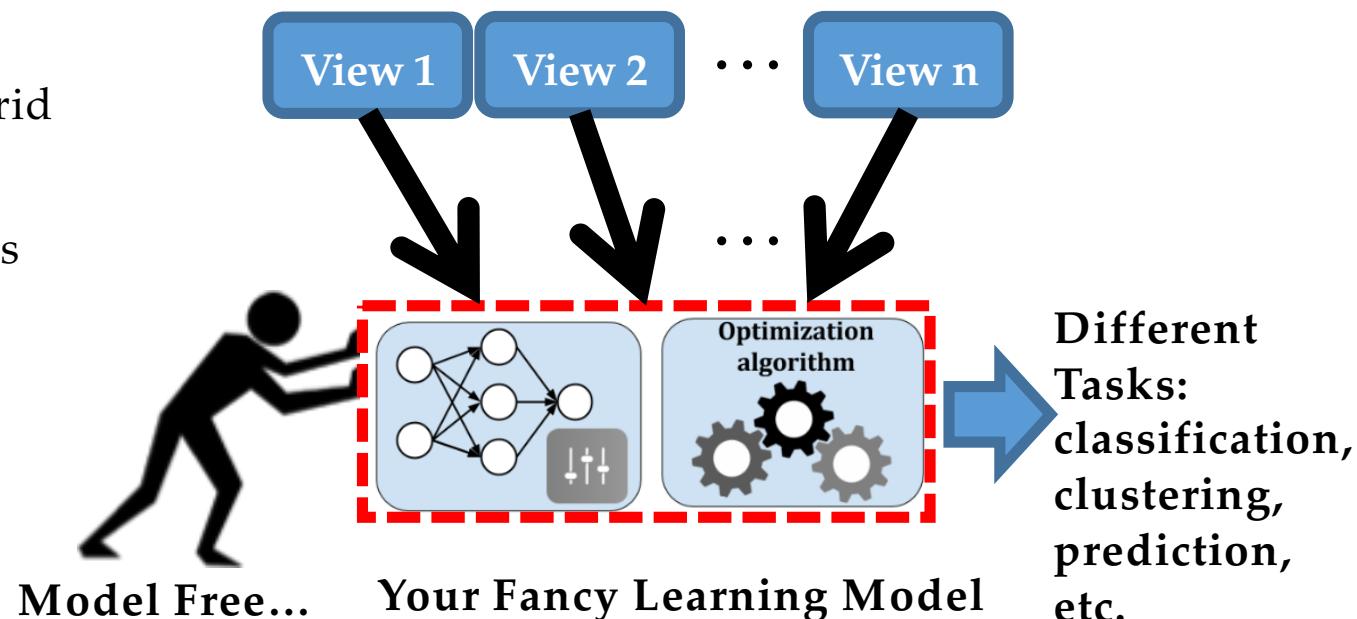
[2] Shen H, Tao D, Ma D (2013) Multiview Locally Linear Embedding for Effective Medical Image Retrieval. PLoS ONE 8(12): e82409.



Knowledge Fusion

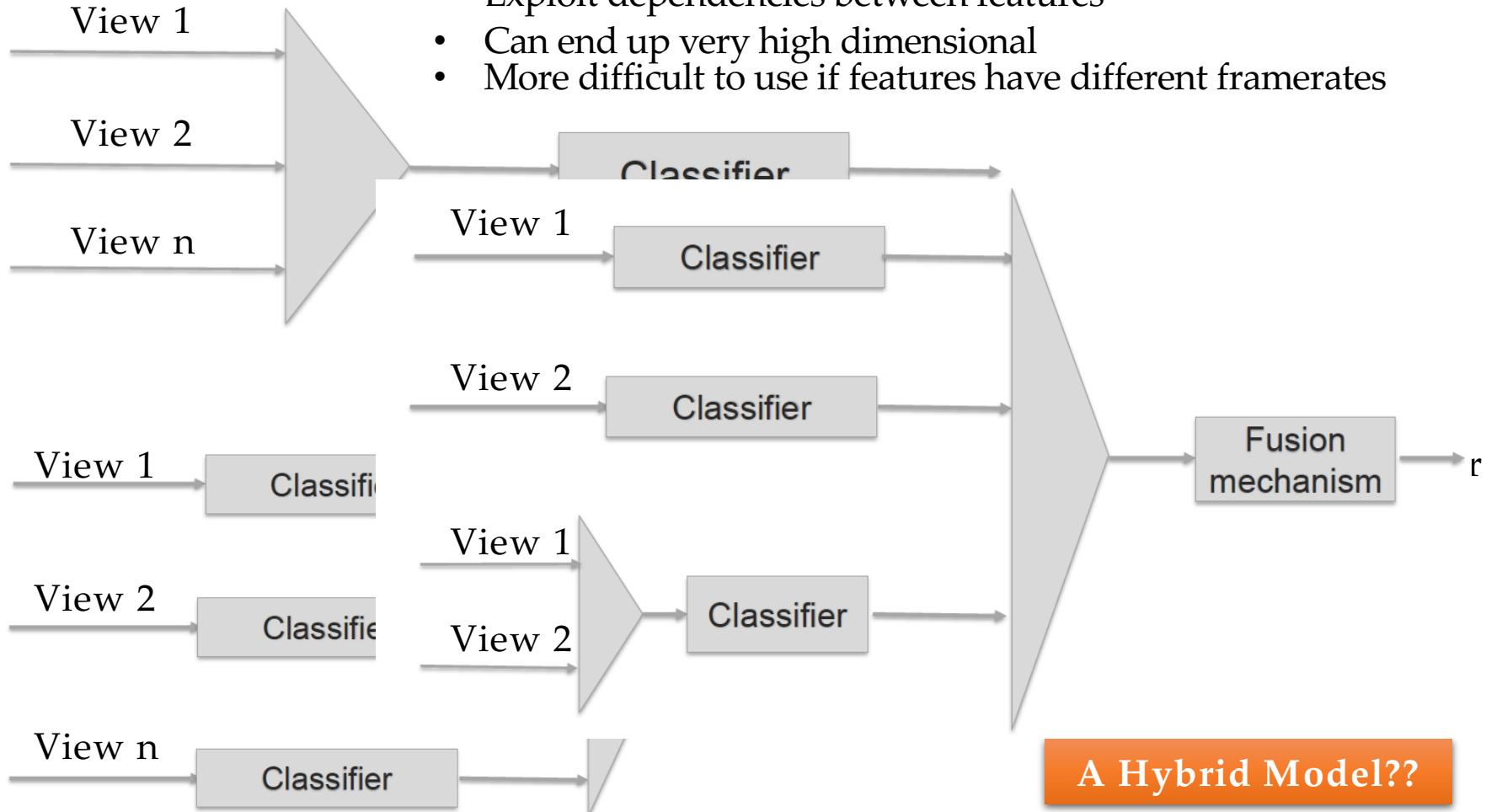
- Process of joining information from two or more modalities to perform a prediction
- One of the earlier and more established problems, e.g. audio-visual speech recognition, multimedia event detection, multimodal emotion recognition

- Model Free
 - Early, late, hybrid
- Model Based
 - Kernel Methods
 - Graph



Model Free Approaches – Early vs. Late Fusion

- Easy to implement – just concatenate the features
- Exploit dependencies between features
- Can end up very high dimensional
- More difficult to use if features have different framerates





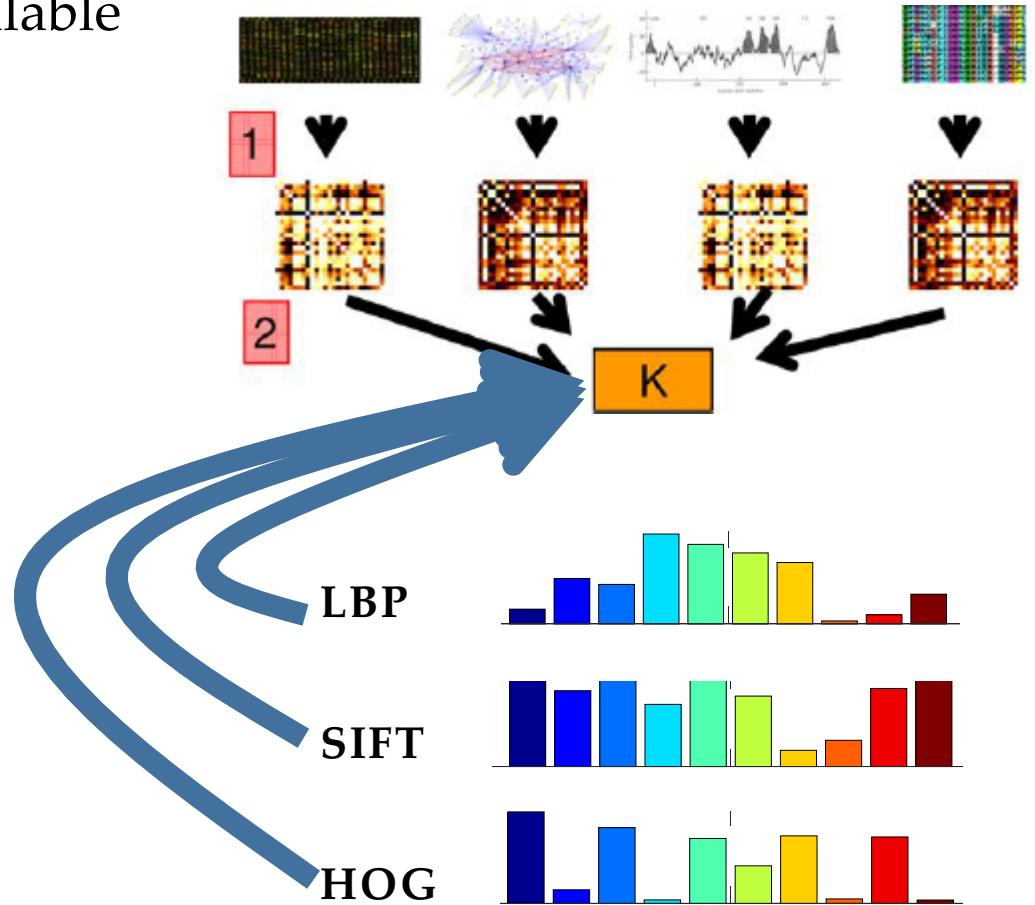
Multi-Kernel Learning

- The overall MKL framework:
- Extract features from all available sources
- Construct kernel matrices
 - Different features
 - Different kernel types
 - Different kernel parameters
- Find the optimal kernel combination and the kernel classifier, e.g., SVM

$$\mathbf{K}(\boldsymbol{\beta}) = \sum_{j=1}^s \beta_j \mathbf{K}_j$$

The new kernel is a linear combination of different kernels

Create individual kernels for each source (string kernel, diffusion kernel)

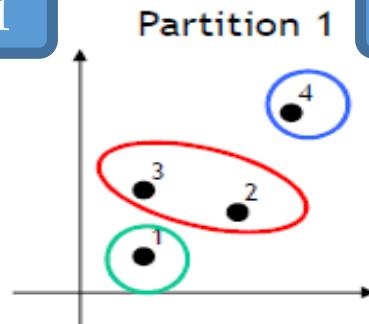




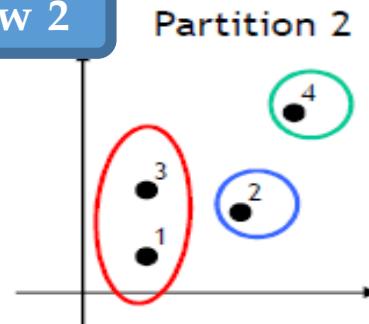
Graph based Fusion

- We may have more than one partitions from different views or features, and fusing them using graph (co-association matrix) is straightforward

View 1



View 2



	1	2	3	4
1	1	0	0	0
2	0	1	1	0
3	0	1	1	0
4	0	0	0	1

Co-association matrices

	1	2	3	4
1	1	0	1	0
2	0	1	0	0
3	1	0	1	0
4	0	0	0	1

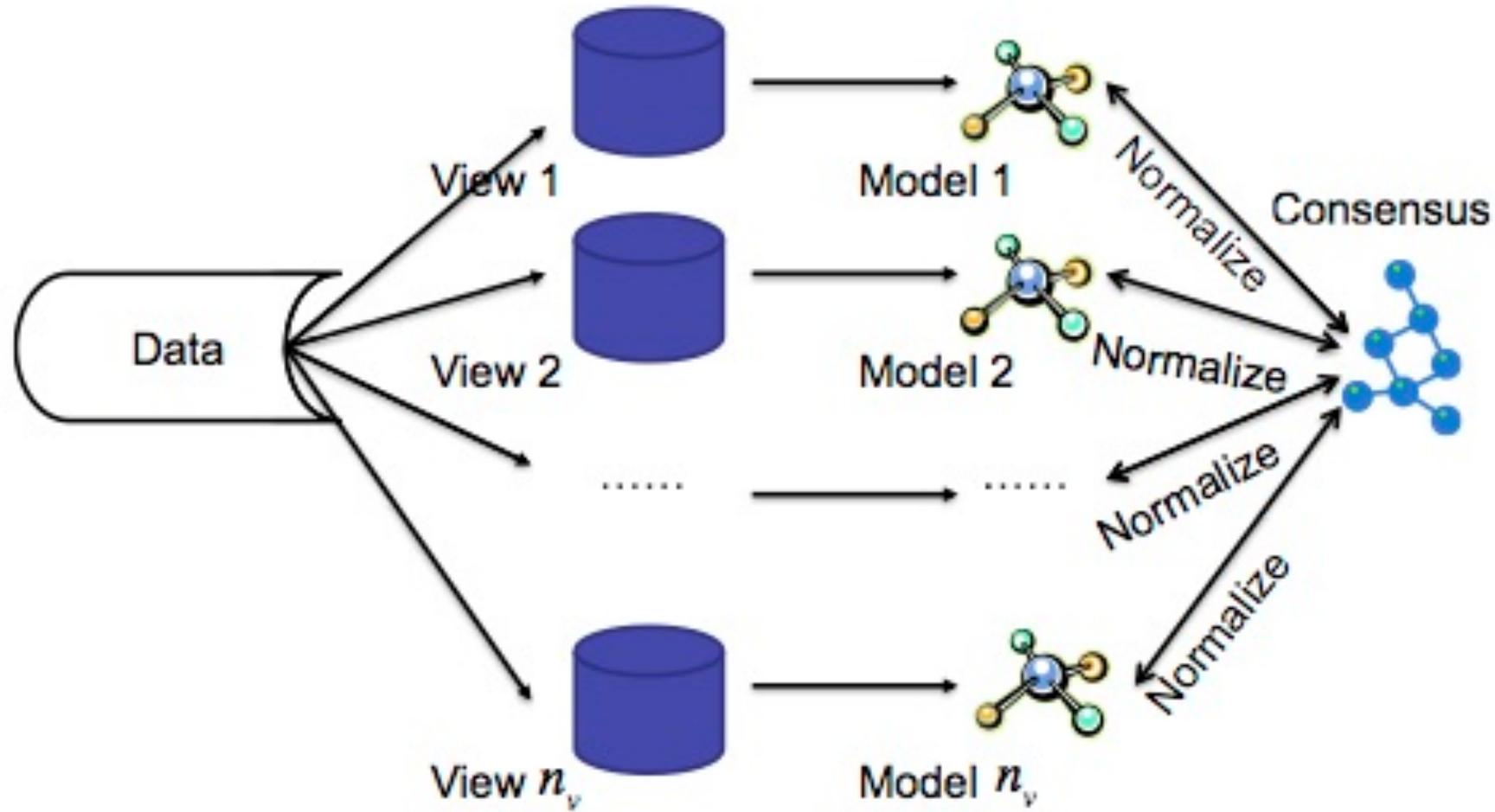


Resultant partition

	1	2	3	4
1	1	0	½	0
2	0	1	½	0
3	½	½	1	0
4	0	0	0	1



Multi-view Clustering

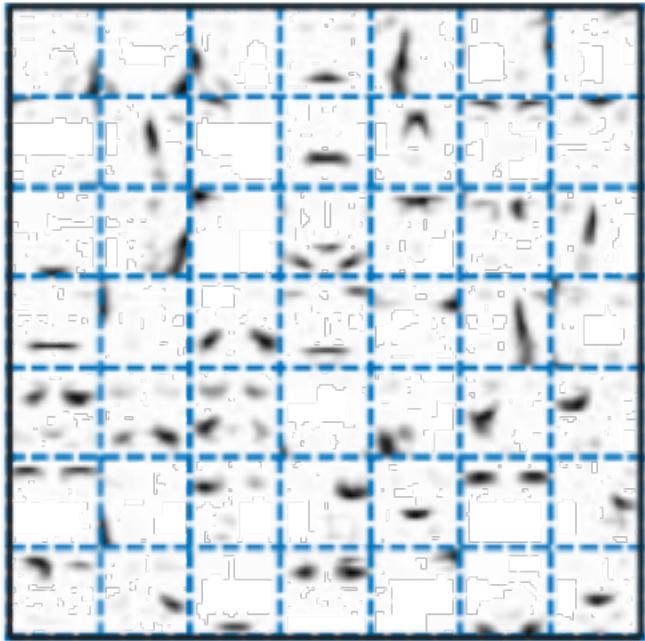


Methodology

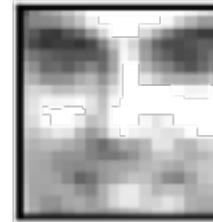
- Multi-view Clustering

Multi-view Clustering via Joint Nonnegative Matrix Factorization

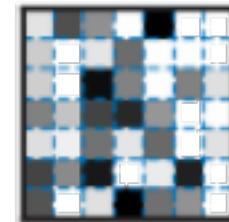
NMF



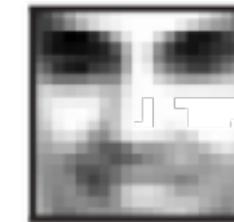
Original



\times



$=$



Basis Coefficients

Objective function:

$$\min_{U,V} \| \boxed{X} - UV^T \|_F^2, \text{ s.t. } \boxed{U} \geq 0, \boxed{V} \geq 0$$

$M \times N$ $M \times K$ $N \times K$

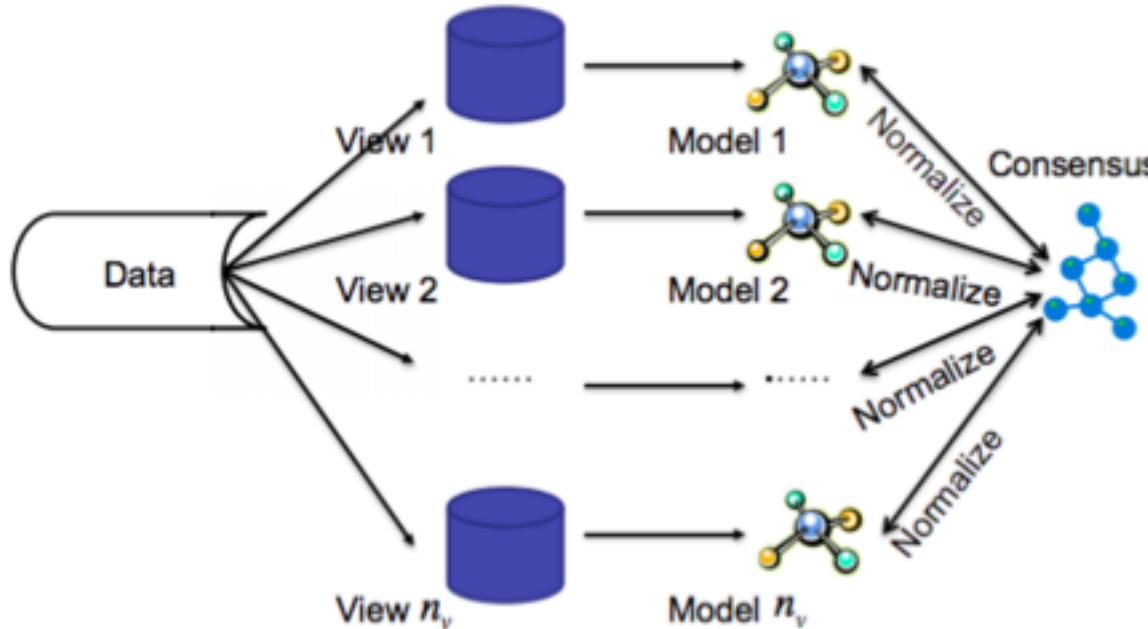
Learning the parts of objects by non-negative matrix factorization – Nature'99

Daniel D Lee, H Sebastian Seung

Methodology

- Multi-view Clustering

Multi-view Clustering via Joint Nonnegative Matrix Factorization



NMF has a good interpretability, and it is reported to achieve competitive performance compared with most of the state-of-the-art unsupervised algorithms.

The latent representations $V^{(v)}$ in different views are forced to be close to the consensus one V^* .

Objective function:

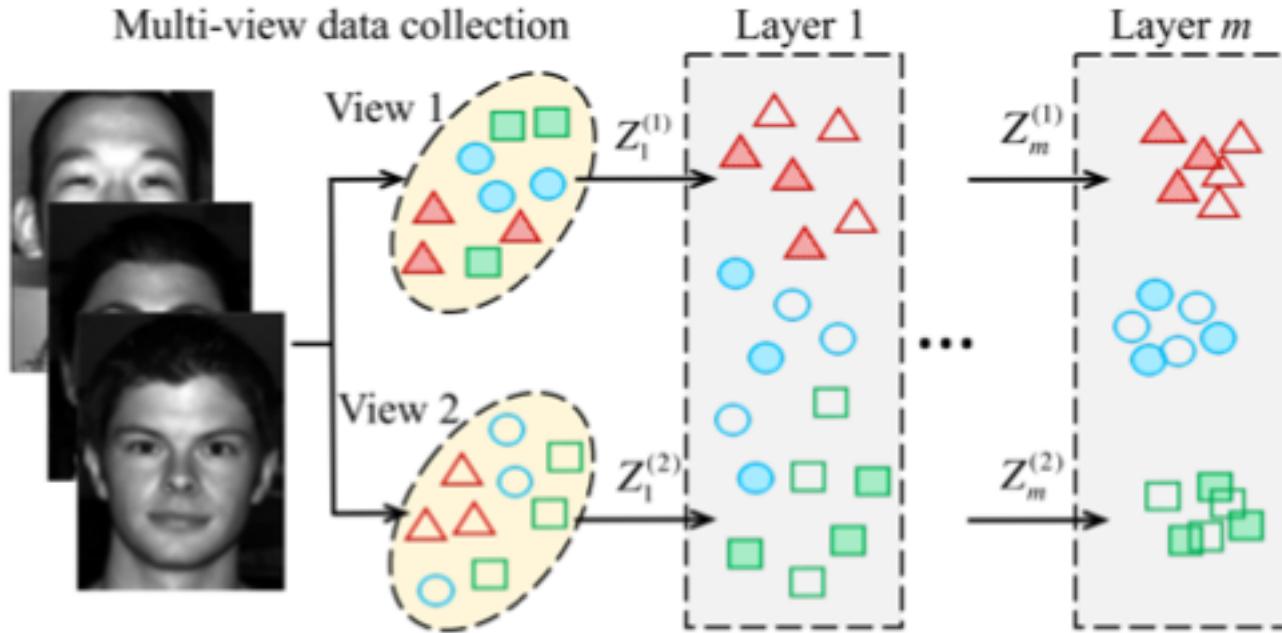
$$\sum_{v=1}^{n_v} \|X^{(v)} - U^{(v)}(V^{(v)})^T\|_F^2 + \sum_{v=1}^{n_v} \lambda_v \|V^{(v)} - V^*\|_F^2$$

s.t. $\forall 1 \leq k \leq K, \|U_{:,k}^{(v)}\|_1 = 1, U^{(v)}, V^{(v)}, V^* \geq 0$

Methodology

- Multi-view Clustering

Multi-View Clustering via Deep Matrix Factorization



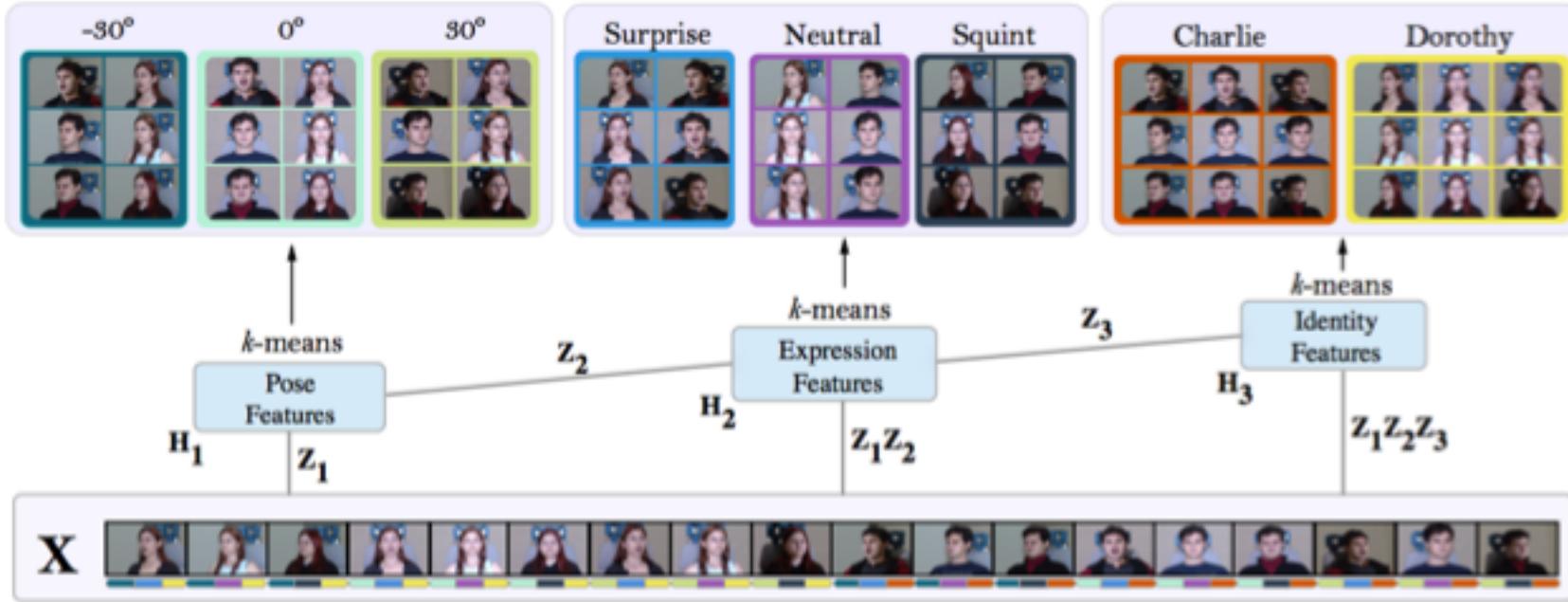
Motivation:

To learn the hierarchical semantics of multi-view data in a layer-wise fashion, semi-nonnegative matrix factorization is adopted.

Methodology

- Multi-view Clustering

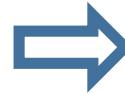
Single-View Clustering via Deep Semi-NMF



Layer-wise formulation: $\mathbf{X}^\pm \approx \mathbf{Z}_1^\pm \mathbf{H}_1^+$

$$\mathbf{X}^\pm \approx \mathbf{Z}_1^\pm \mathbf{Z}_2^\pm \mathbf{H}_2^+$$

$$\mathbf{X}^\pm \approx \mathbf{Z}_1^\pm \mathbf{Z}_2^\pm \mathbf{Z}_3^\pm \mathbf{H}_3^+$$



$$\mathbf{X}^\pm \approx \mathbf{Z}_1^\pm \mathbf{Z}_2^\pm \cdots \mathbf{Z}_m^\pm \mathbf{H}_m^+$$

Methodology

- Multi-view Clustering

Multi-View Clustering via Deep Matrix Factorization

Objective function:

$$\min_{\substack{Z_i^{(v)}, H_i^{(v)}, \\ H_m, \alpha^{(v)}}} \sum_{v=1}^V (\alpha^{(v)})^\gamma \left(\|X^{(v)} - Z_1^{(v)} Z_2^{(v)} \dots Z_m^{(v)} H_m\|_F^2 + \beta \text{tr}(H_m L^{(v)} H_m^T) \right)$$

Decomposition on all views, where the representations on the last layer $H_m^{(v)}$ are forced to be same H_m .

The hidden representation H are non-negative, with good interpretability.

$L^{(v)}$ is the graph Laplacian of the graph for view v , where each graph is constructed in k-nearest neighbor fashion.

$$\text{s.t. } H_i^{(v)} \geq 0, H_m \geq 0 \quad \sum_{v=1}^V \alpha^{(v)} = 1, \alpha^{(v)} \geq 0$$

Application

- Multi-view Clustering

Extension: Incomplete Scenario

When the data from one modality/more modalities are inaccessible because of sensor failure or other reasons, most traditional MVC methods would inevitably degenerate or even fail.



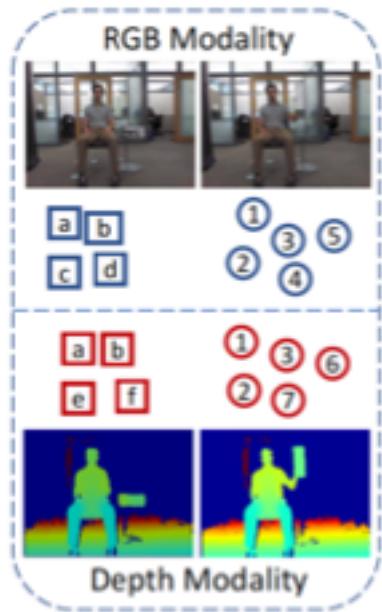
Application

- Multi-view Clustering

Extension: Incomplete Scenario

Incomplete Multi-Modal Visual Data Grouping

Motivation:



Objective function:

$$\begin{aligned} & \min_{P_c, \hat{P}^{(1)}, \hat{P}^{(2)}, U^{(1)}, U^{(2)}, A} \left\| \begin{bmatrix} X_c^{(1)} \\ \hat{X}^{(1)} \end{bmatrix} - \begin{bmatrix} P_c \\ \hat{P}^{(1)} \end{bmatrix} U^{(1)} \right\|_F^2 + \\ & \quad \left\| \begin{bmatrix} X_c^{(2)} \\ \hat{X}^{(2)} \end{bmatrix} - \begin{bmatrix} P_c \\ \hat{P}^{(2)} \end{bmatrix} U^{(2)} \right\|_F^2 + \mathcal{G}(P, A) + \mathcal{R}(U, A). \\ & \text{s.t. } \forall i A_i^T \mathbf{1} = 1, A_i \succeq 0. \end{aligned}$$

$$\mathcal{G}(P, A) = \beta \text{tr}(P^T L_A P),$$

$$\mathcal{R}(U, A) = \lambda(\|U^{(1)}\|_F^2 + \|U^{(2)}\|_F^2) + \gamma \|A\|_F^2$$

Application

- Multi-view Clustering

Extension: Incomplete Scenario

Incomplete Multi-Modal Visual Data Grouping

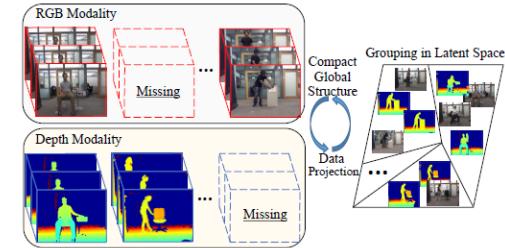
Objective function:

$$\min_{\substack{P_c, P^{(1)}, P^{(2)} \\ U^{(1)}, U^{(2)}, A}} \left\| \begin{bmatrix} X_c^{(1)} \\ X^{(1)} \end{bmatrix} - \begin{bmatrix} P_c \\ P^{(1)} \end{bmatrix} U^{(1)} \right\|_F^2 + \left\| \begin{bmatrix} X_c^{(2)} \\ X^{(2)} \end{bmatrix} - \begin{bmatrix} P_c \\ P^{(2)} \end{bmatrix} U^{(2)} \right\|_F^2 \\ + \mathcal{G}(P, A) + \boxed{\mathcal{R}(U, A)}$$

$$\text{s.t. } \forall i \ A_i^T \mathbf{1} = 1, \ A_i \succeq 0.$$

$$\mathcal{G}(P, A) = \beta \text{tr}(P^T L_A P)$$

$$\mathcal{R}(U, A) = \lambda(\|U^{(1)}\|_F^2 + \|U^{(2)}\|_F^2) + \gamma \|A\|_F^2$$



- Use the shared data $X_c^{(v)}$ in each view to learn the common representation P_c in low-dimensional subspace space.
- Regularizers to prevent trivial solution. Here we choose a simple Frobenius norm. Its alternatives include ℓ_1 and others.

Application

- Multi-view Clustering

Extension: Incomplete Scenario

Incomplete Multi-Modal Visual Data Grouping

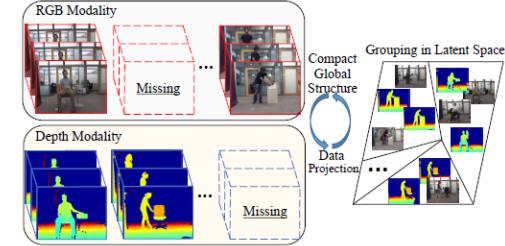
Objective function:

$$\min_{\substack{P_c, \hat{P}^{(1)}, \hat{P}^{(2)} \\ U^{(1)}, U^{(2)}, A}} \left\| \begin{bmatrix} X_c^{(1)} \\ \hat{X}^{(1)} \end{bmatrix} - \begin{bmatrix} P_c \\ \hat{P}^{(1)} \end{bmatrix} U^{(1)} \right\|_F^2 + \left\| \begin{bmatrix} X_c^{(2)} \\ \hat{X}^{(2)} \end{bmatrix} - \begin{bmatrix} P_c \\ \hat{P}^{(2)} \end{bmatrix} U^{(2)} \right\|_F^2 \\ + \boxed{\mathcal{G}(P, A)} + \mathcal{R}(U, A)$$

s.t. $\forall i A_i^T \mathbf{1} = 1, A_i \succeq 0.$

$$\boxed{\mathcal{G}(P, A) = \beta \text{tr}(P^T L_A P)}$$

$$\mathcal{R}(U, A) = \lambda(\|U^{(1)}\|_F^2 + \|U^{(2)}\|_F^2) + \gamma \|A\|_F^2$$



- $P = [P_c; \hat{P}^{(1)}; \hat{P}^{(2)}]$
- Graph Laplacian term to preserve locality information, where L_A is the Laplacian matrix of similarity matrix A , which is learned on the latent representation P .

Outline



□ Introduction & Background

- Multi-view Visual Data
- Multi-view Learning Problems
- Multi-view Learning Taxonomy

□ Multi-view Learning

- Projection and Embedding
- Knowledge Fusion
- Multi-view Clustering
- Supervised Multi-view Learning → Zero-shot Learning

□ Domain Adaptation

- Transfer Learning → Domain Adaptation
- Multi-Source Domain Adaptation & Domain Generalization

□ Conclusion

Supervised Multi-View Learning & Domain Adaptation

- Supervised Multi-view Learning
[sample-wise correspondence]

Training Stage: multiple labeled view data

Test Stage:

[Setting 1]: labeled views → unlabeled views

[Setting 2]: multiple unlabeled view data

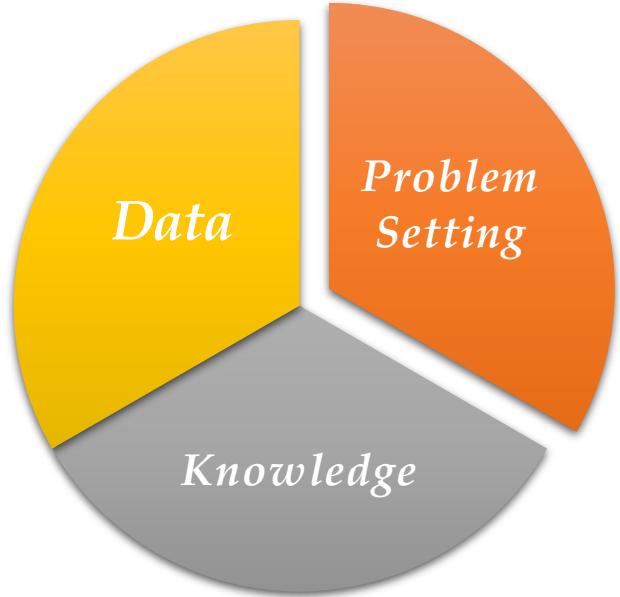
Goal: fuse various knowledge from multiple views

Goal: adapt knowledge across different views

- Domain Adaptation **[class-wise correspondence]**

Training Stage: some source labeled domains & some target **[un]labeled** domains

Test Stage: source to predict target data



Supervised Multi-View Learning

Training Stage: $X_s = \{X_s^1, \dots, X_s^v | y\}$

Test Stage: one view to recognize others $\{y_*\}$
or $\{X_s | y\}$ to recognize test data $\{X_t | y\}$



- Cross-pose Face Recognition
- Multi-modal Recognition
- **Sample-wise Correspondence**
- **Labels Information**

$$\min_{f_1(\cdot), \dots, f_v(\cdot)} \sum_{i=1, i < j}^v \mathcal{A}(f_i(X_i), f_j(X_j)) + \lambda \sum_{k=1}^v \mathcal{R}(f_k(X_k))$$

- Feature Learning
→ Subspace Learning, Deep Learning
- ◆ Alignment & fusion
→ joint & coordinated representation

- ***Label information***
 - Supervised Graph
 - Regression loss

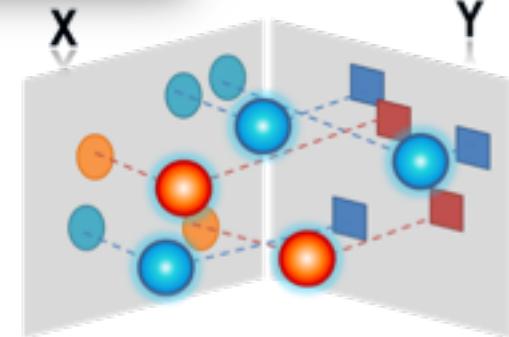
CCA → Supervised CCA

CCA

$$\begin{aligned} \max_{\mathbf{P}, \mathbf{Q}} \quad & \text{trace}(\mathbf{P}^\top \mathbf{X}_T \mathbf{X}_S^\top \mathbf{Q}^\top), \\ \text{s.t.} \quad & \mathbf{P}^\top \mathbf{X}_T \mathbf{X}_T^\top \mathbf{P} = \mathbf{I}, \quad \mathbf{Q}^\top \mathbf{X}_S \mathbf{X}_S^\top \mathbf{Q} = \mathbf{I} \end{aligned}$$



$$||\mathbf{X}_T^\top \mathbf{P} - \mathbf{X}_S^\top \mathbf{Q}||_F^2$$



Supervised CCA

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{P}, \mathbf{Q}} \quad & C \sum_{i=1}^{n_L} \xi_{L,i} + \frac{1}{2} \|\mathbf{w}\|^2 + \frac{\gamma}{2} ||\mathbf{X}_T^\top \mathbf{P} - \mathbf{X}_S^\top \mathbf{Q}||_F^2 \\ \text{s.t.} \quad & \mathbf{P}^\top \mathbf{X}_T \mathbf{X}_T^\top \mathbf{P} = \mathbf{I}, \quad \mathbf{Q}^\top \mathbf{X}_S \mathbf{X}_S^\top \mathbf{Q} = \mathbf{I}, \end{aligned}$$

cross-entropy loss function

$$\xi_{L,i} = \log \left(\sum_{k=1}^K \exp(\mathbf{w}_k^\top \tilde{\mathbf{x}}_{L,i} - \mathbf{w}_{y_{L,i}}^\top \tilde{\mathbf{x}}_{L,i}) \right)$$

pair-wise constraints

Multi-view Discriminant Analysis

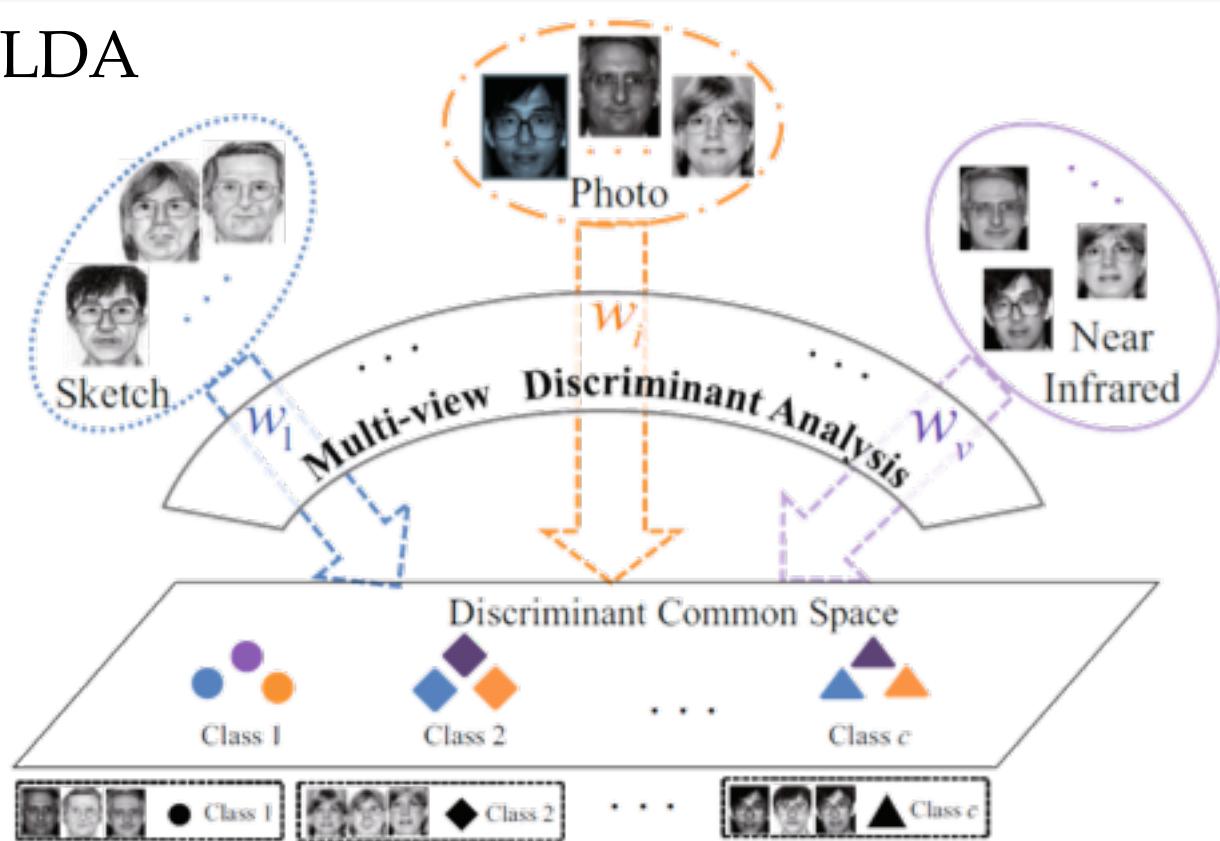
LDA → Multi-view LDA



multiple projections

+ *Fisher loss*

$$\max_{\mathbf{w}_1, \dots, \mathbf{w}_v} \frac{\text{Tr}(\mathbf{S}_B^y)}{\text{Tr}(\mathbf{S}_W^y)}$$



\mathbf{S}_B^y : between-class scatter matrix

\mathbf{S}_W^y : within-class scatter matrix

Multi-view Discriminant Analysis – ECCV 2012 & IEEE TPAMI 2016

Meina Kan, Shiguang Shan, Haihong Zhang, Shihong Lao, and Xilin Chen

Supervised Multi-View Learning

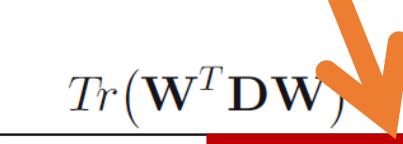
Extension

$$(\mathbf{w}_1^*, \mathbf{w}_2^*, \dots, \mathbf{w}_v^*)$$

$$= \arg \max_{\mathbf{w}_1, \dots, \mathbf{w}_v} \frac{\text{Tr}(\mathbf{W}^T \mathbf{D} \mathbf{W})}{\text{Tr}(\mathbf{W}^T \mathbf{S} \mathbf{W}) + \lambda \sum_{i,j=1}^v \|\beta_i - \beta_j\|_2^2}$$

$$= \arg \max_{\mathbf{w}_1, \dots, \mathbf{w}_v} \frac{\text{Tr}(\mathbf{W}^T \mathbf{D} \mathbf{W})}{\text{Tr}(\mathbf{W}^T (\mathbf{S} + \lambda \mathbf{M}) \mathbf{W})}$$

view-consistent regularizer



Projection & Data

$$\mathbf{w}_i = \mathbf{X}_i \boldsymbol{\beta}_i$$

$$\begin{aligned} \mathbf{X}_1 &= \mathbf{R} \mathbf{X}_2 \\ \mathbf{w}_1 &= \mathbf{R} \mathbf{w}_2 \end{aligned}$$

$$\mathbf{X}_1 \boldsymbol{\beta}_1 = \mathbf{R} \mathbf{X}_2 \boldsymbol{\beta}_2 = \mathbf{X}_1 \boldsymbol{\beta}_2$$

Results on CUFSF and HFB Datasets

		CCA[4]**	CCA[4]+LDA	CDFE[19]	CSR[21]	PLS[6]	U-LDA[31]	GMA[28]	MvDA	MvDA-VC
CUFSF	Photo-Sketch	45.5%	45.0%	45.6%	50.2%	48.6%	46.8%	-	53.4%	56.3%
	Sketch-Photo	47.5%	50.6%	47.6%	49.0%	51.0%	53.4%	-	55.5%	61.5%
HFB	NIR-VIS	36.7%	40.0%	40.8%	26.7%	38.3%	39.1%	47.5%	53.3%	59.2%
	VIS-NIR	30.0%	40.0%	36.7%	32.5%	40.8%	40.0%	45.0%	50.0%	59.2%

Multi-view Discriminant Analysis. –ECCV 2012 & IEEE TPAMI 2016

Meina Kan, Shiguang Shan, Haihong Zhang, Shihong Lao, and Xilin Chen



Supervised Multi-View Learning

Deep Multi-View Learning

$$\min_{\mathbf{g}_c, \mathbf{f}_1, \dots, \mathbf{f}_v} Tr \left(\frac{\mathbf{S}_W^y}{\mathbf{S}_B^y} \right)$$

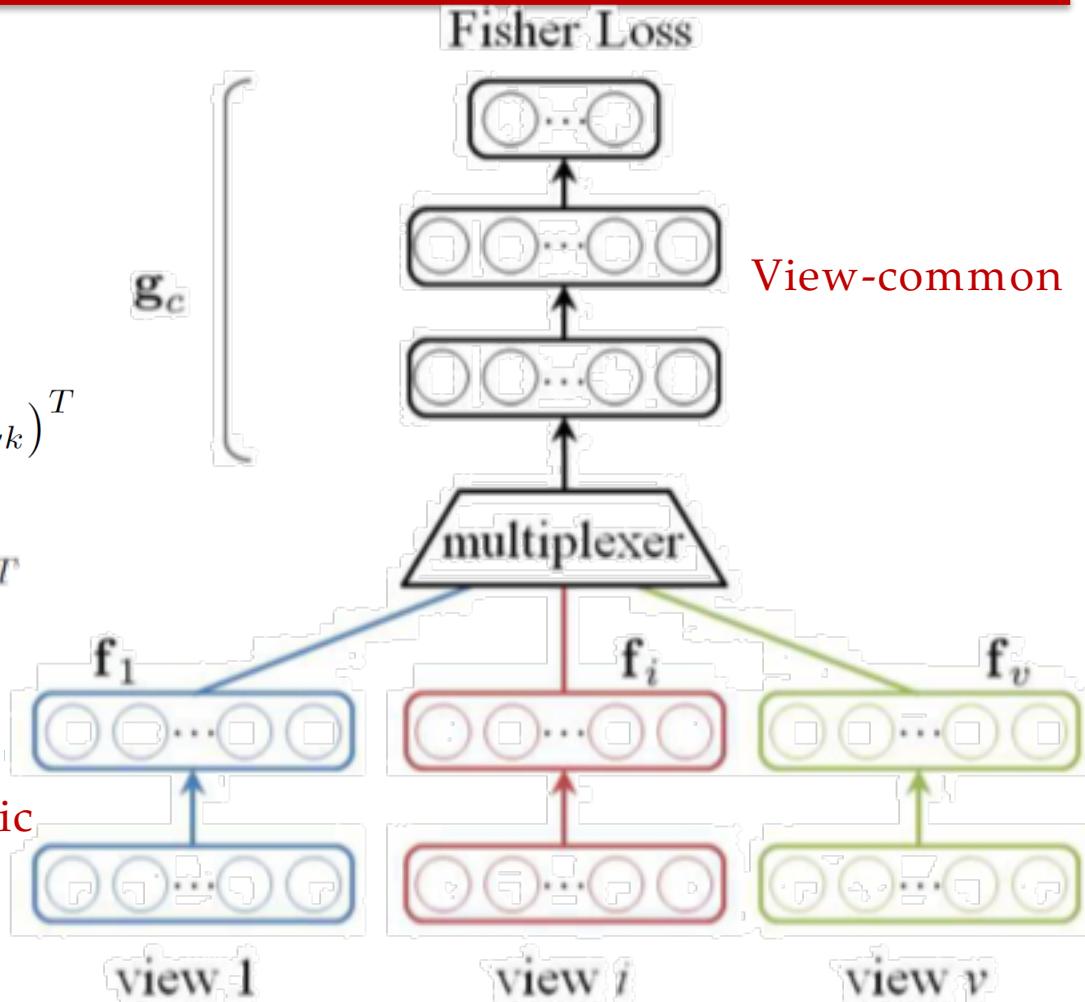
$$\mathbf{S}_W^y = \sum_{k=1}^c \sum_{i=1}^v \sum_{j=1}^{n_{ki}} (\mathbf{y}_{jk}^i - \boldsymbol{\mu}_k) (\mathbf{y}_{jk}^i - \boldsymbol{\mu}_k)^T$$

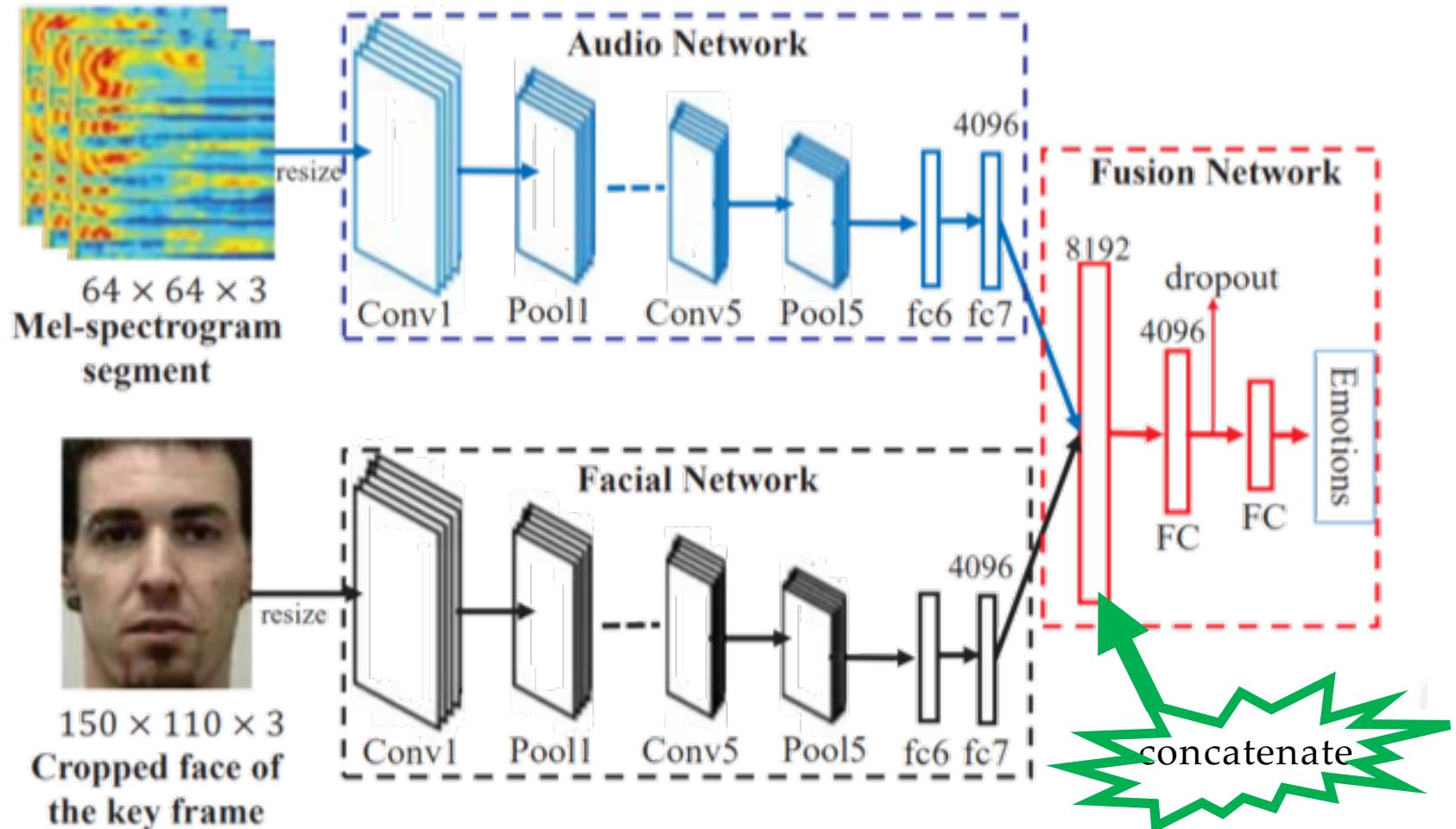
$$\mathbf{S}_B^y = \sum_{k=1}^c n_k (\boldsymbol{\mu}_k - \boldsymbol{\mu}) (\boldsymbol{\mu}_k - \boldsymbol{\mu})^T$$

Fisher-like Loss

$$\mathbf{y}_j^i = \mathbf{g}_c (\mathbf{f}_i (\mathbf{x}_j^i))$$

View-specific







Special Case: Zero-Shot Learning

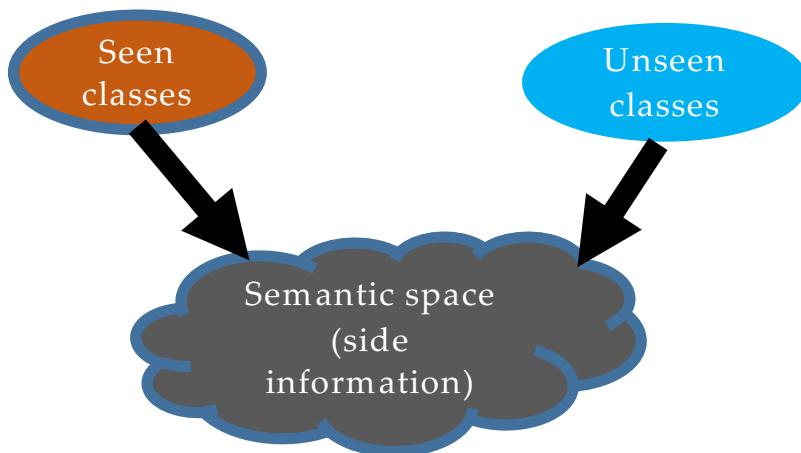
Supervised Multi-View Learning [Setting 1]



Zero-Shot Learning

- **Seen classes:** labeled data available, training classes, 'side' information available
- **Unseen classes:** no training data (zero-shot), test classes
 - Zero-Shot Recognition
 - Zero-shot Annotation
 - Zero-shot Retrieval

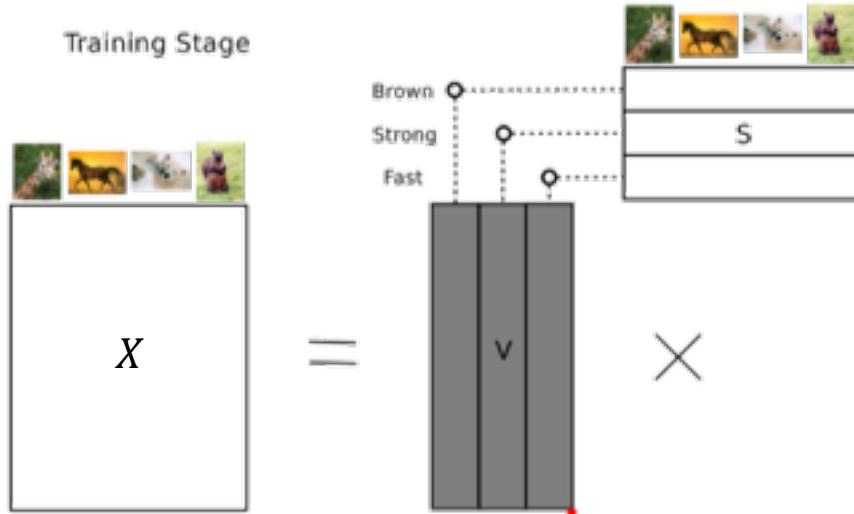
find relationship between visual and semantic views



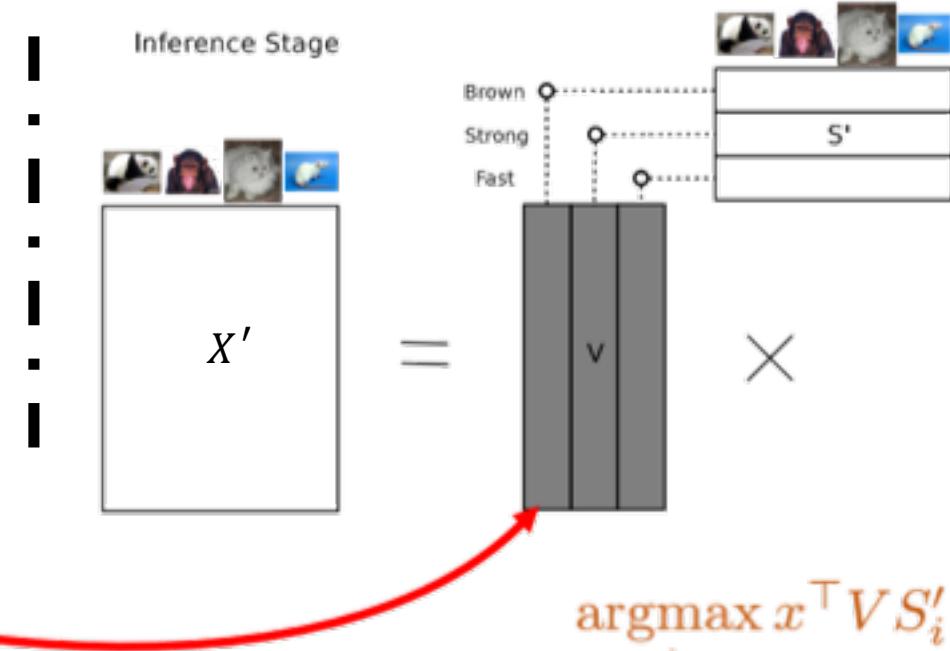


Zero-Shot Learning

Training Stage



Inference Stage



Sample-wise correlation ship

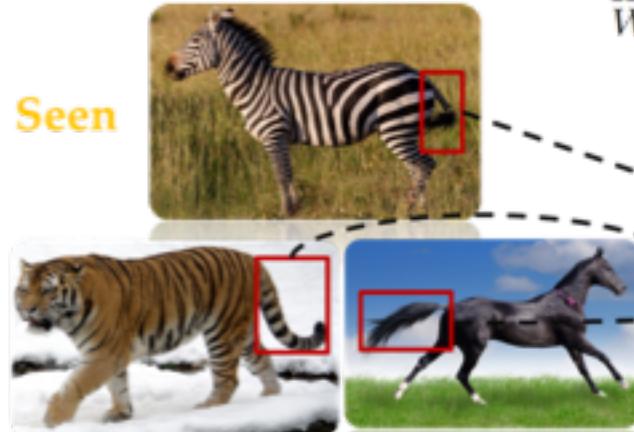
$$Y = \{-1, 1\}$$

$$\underset{V \in \mathbb{R}^{d \times a}}{\text{minimise}} L(X^T V S, Y) + \Omega(V; S, X)$$

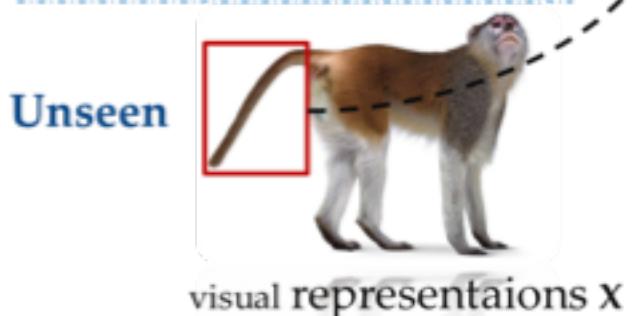
$$\Omega(V; S, X) = \gamma \|VS\|_{\text{Fro}}^2 + \lambda \|X^T V\|_{\text{Fro}}^2 + \beta \|V\|_{\text{Fro}}^2$$



Zero-Shot Learning

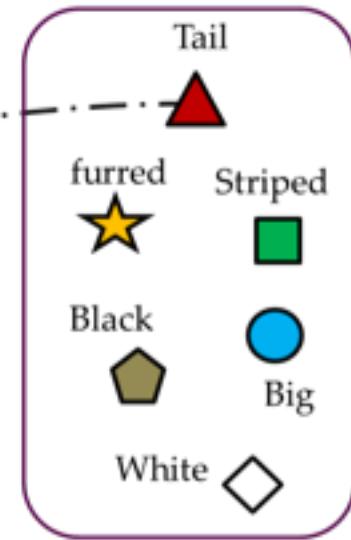
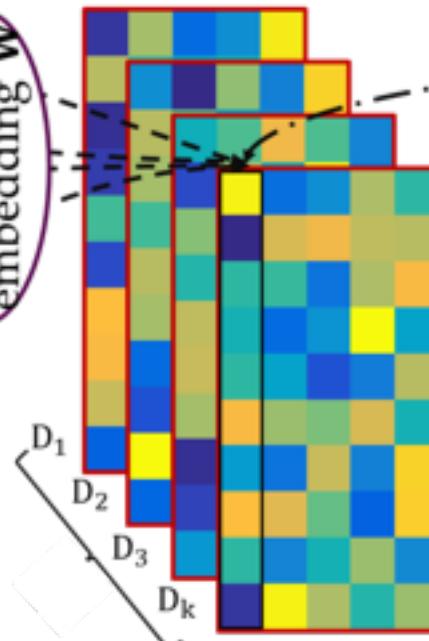


The same attribute **has a tail** across different categories.



$$\begin{aligned} & \min_{W,D} \|WX - DA\|_F^2 + \alpha \text{rank}(W) \\ & \text{s.t. } \|d_j\|_2^2 \leq 1, \forall j, \end{aligned}$$

low-rank embedding W



semantic representations A



Zero-Shot Learning

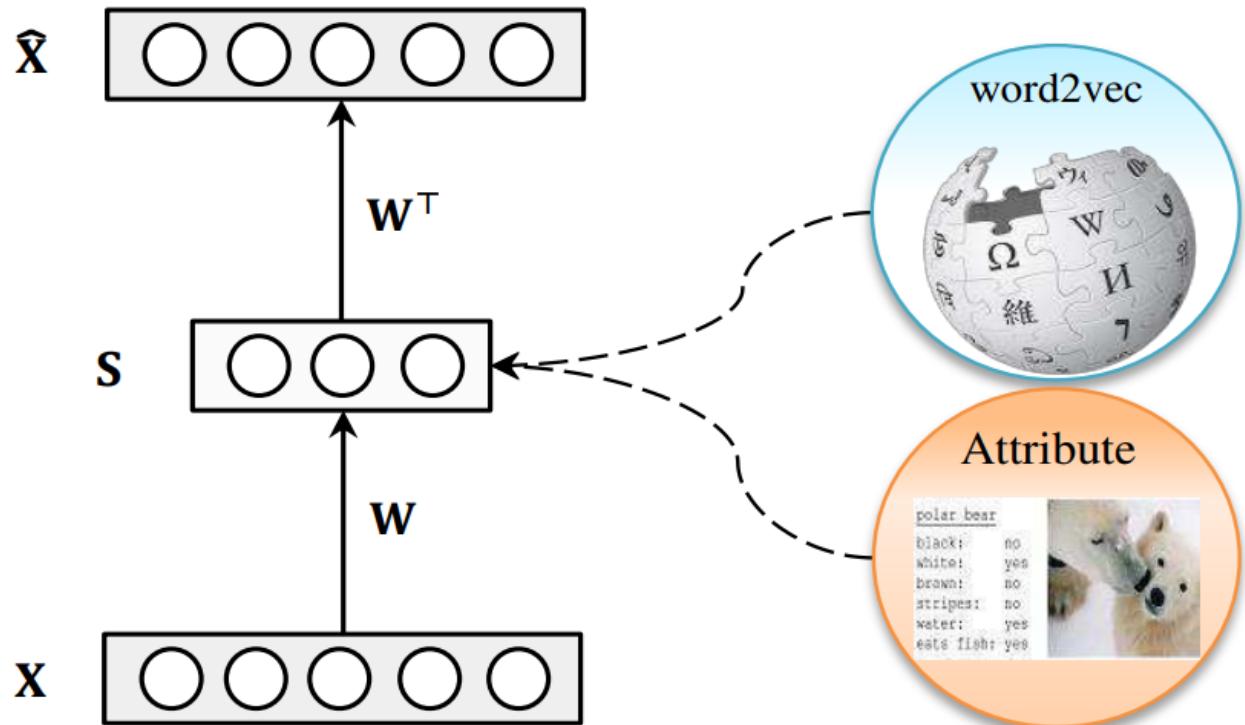
The proposed semantic auto-encoder leverages the semantic side information such as attributes and word vector, while learning an encoder and a decoder

$$\min_{\mathbf{W}} \|\mathbf{X} - \mathbf{W}^\top \mathbf{S}\|_F^2$$

$$s.t. \quad \mathbf{W}\mathbf{X} = \mathbf{S}$$

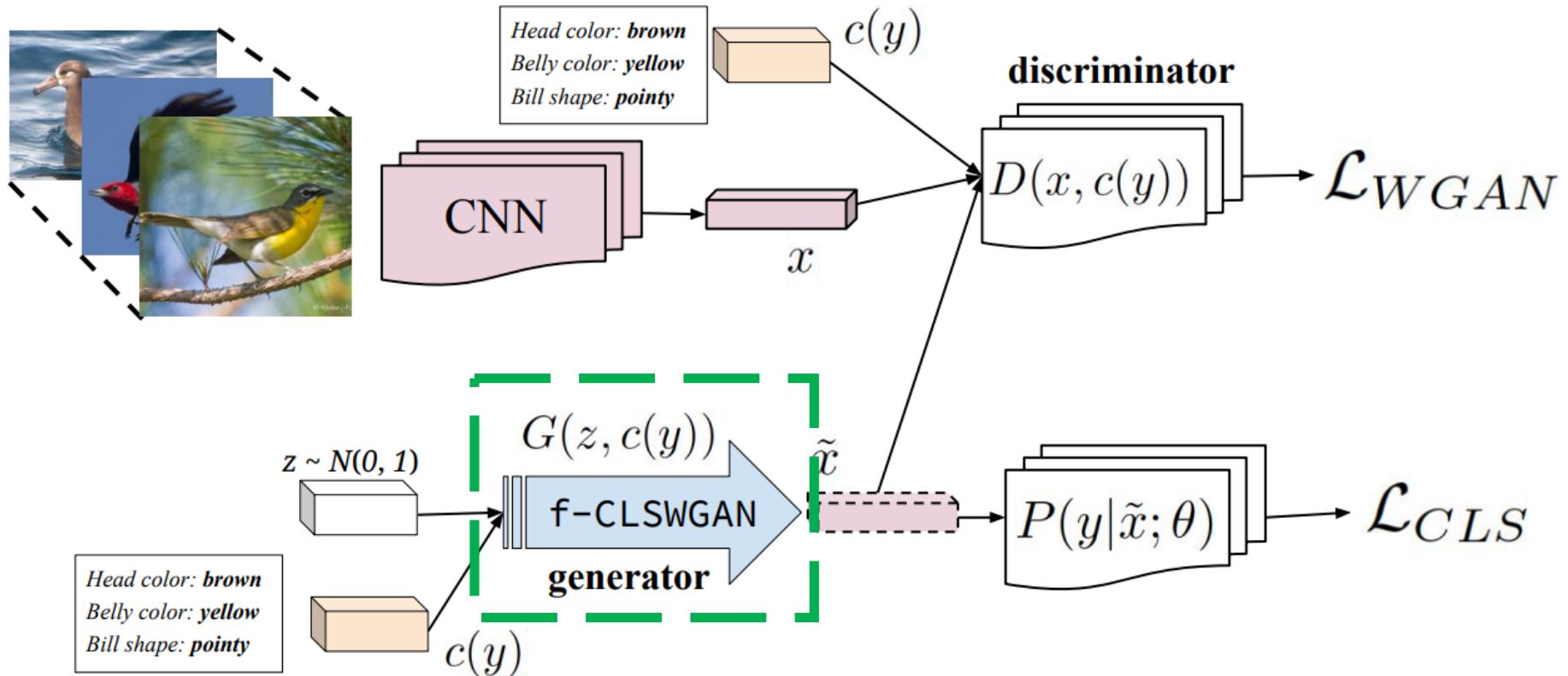


$$\min_{\mathbf{W}} \|\mathbf{X} - \mathbf{W}^\top \mathbf{S}\|_F^2 + \lambda \|\mathbf{W}\mathbf{X} - \mathbf{S}\|_F^2$$





Zero-Shot Learning



minimize the classification loss over the generated features and the Wasserstein distance with gradient penalty

Outline

□ Introduction & Background

- Multi-view Visual Data
- Multi-view Learning Problems
- Multi-view Learning Taxonomy

□ Multi-view Learning

- Projection and Embedding
- Knowledge Fusion
- Multi-view Clustering
- Supervised Multi-view Learning → Zero-shot Learning

□ Domain Adaptation

- Transfer Learning → Domain Adaptation
- Multi-Source Domain Adaptation & Domain Generalization

□ Conclusion

Examples

amazon.com[®]



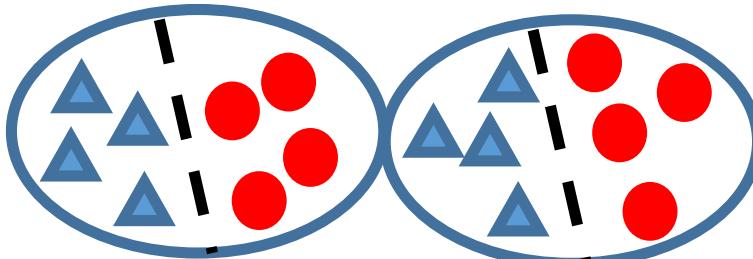
The Office



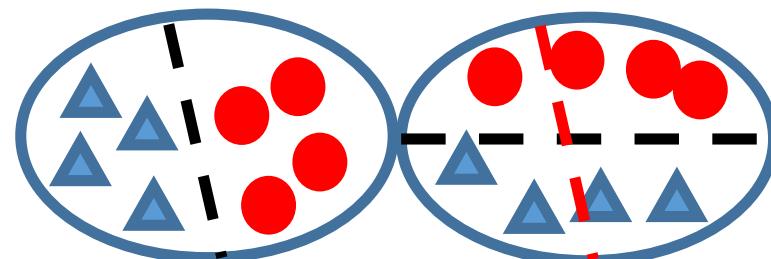
TRAIN



TEST



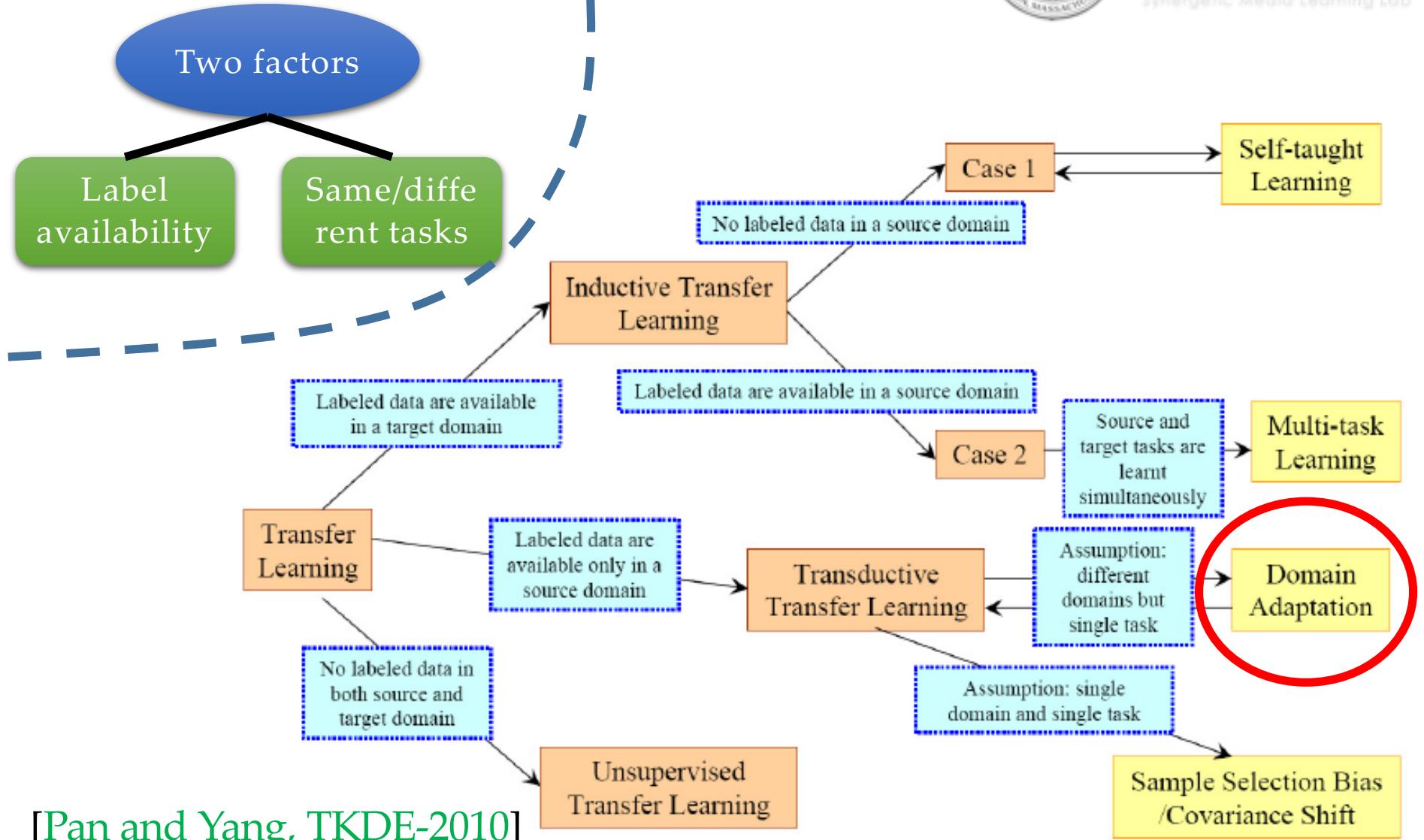
Training and test are
from the **same** domain



Training and test are
from **different** domains



Taxonomy



Domain Adaptation

- Source Views (labeled)

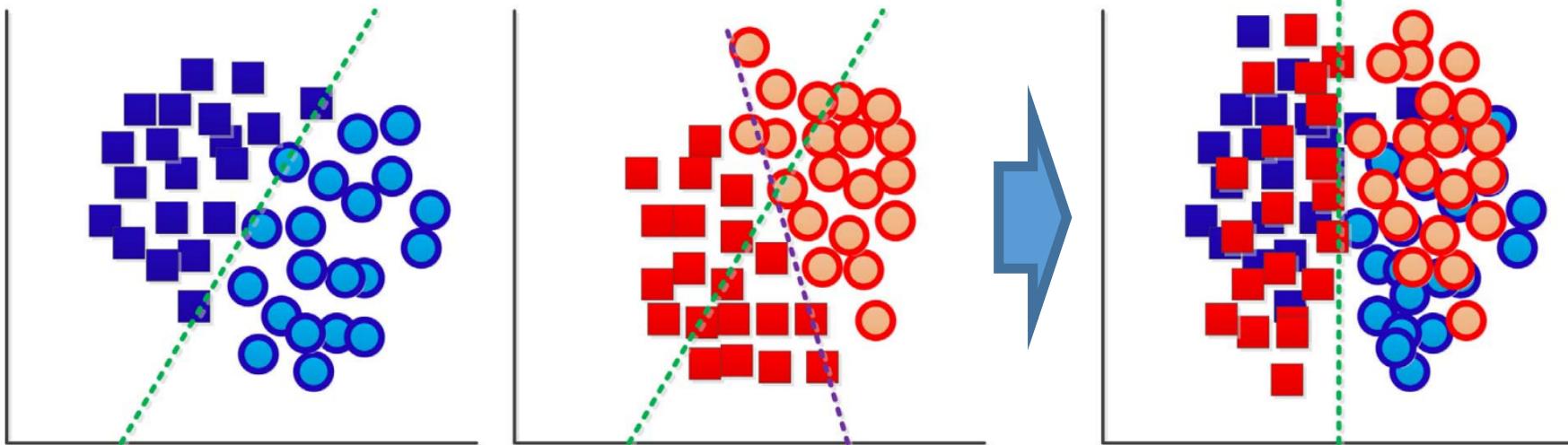
$$D_{S,M} = \{(x_{i,M}, y_{i,M}), i = 1, 2 \dots, N, P_S(X, Y)\}$$

- Target Views ([un]labeled)

$$D_T = \{(x_i, ?), i = 1, 2 \dots, N, P_T(X, Y)\}$$

mismatch

Performance degrades significantly!



Domain Adaptation Approaches

Instance based approach

Parameter based approach

Relational knowledge

Feature based approach

- ❖ **Instance:** partial source data are **reusable**
 - Instance selection approach: **TrAdaBoost**
- ❖ **Parameter:** individual models for related tasks should share some **parameters or priors**
 - **Multi-task learning**
- ❖ **Relational:** **no i.i.d.** assumptions, transfer relationship among data between domains
- ❖ **Feature:** good representations to mitigate **domain divergence**
 - **Subspace learning → Linear Projection**
 - **Dictionary learning → New Representation**
 - **Deep learning → Convolutional Neural Network**

We focus on Representation

Domain Alignment

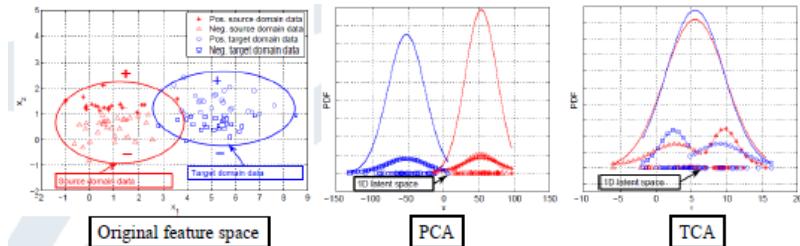


Representation Learning

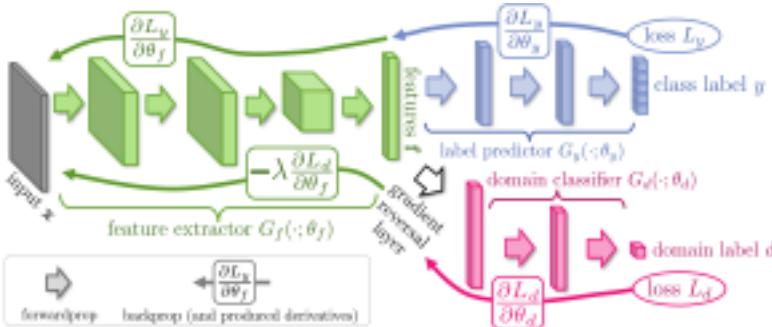
- Source and target have some overlapping features
- Have support in either source or target domain
- Find the transform φ

Projection learning

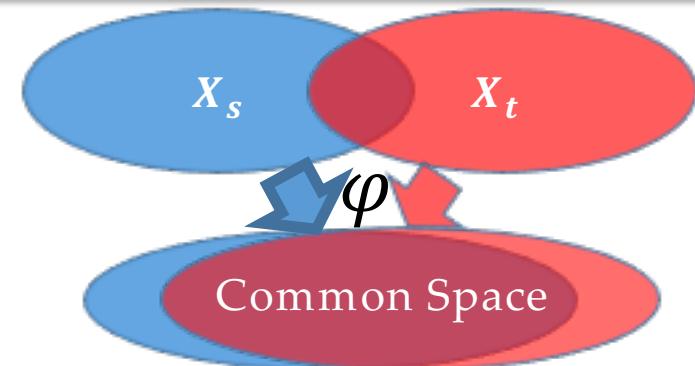
[Pan et al AAAI-08,
Saenko et al ECCV-2010,
Ding et al AAAI-14]



[Ganin et al ICML-15, Long et al NIPS-16]

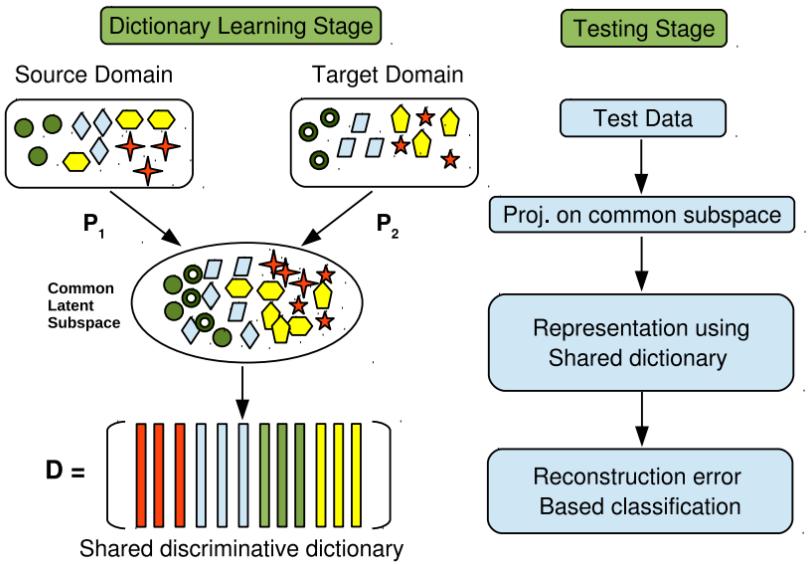


Deep learning



Dictionary learning

[Shekhar et al CVPR-13,
Nguyen et al TIP-15]



Domain Alignment

➤ Marginal distribution

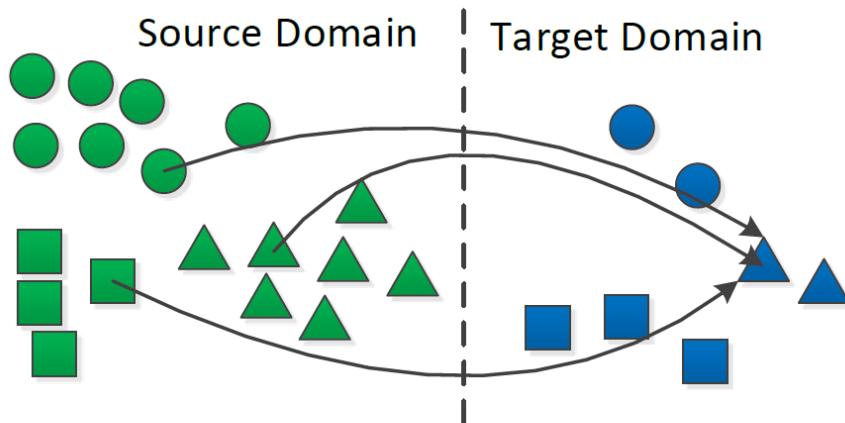
$$P_s(X_s), P_t(X_t)$$

➤ Conditional distribution

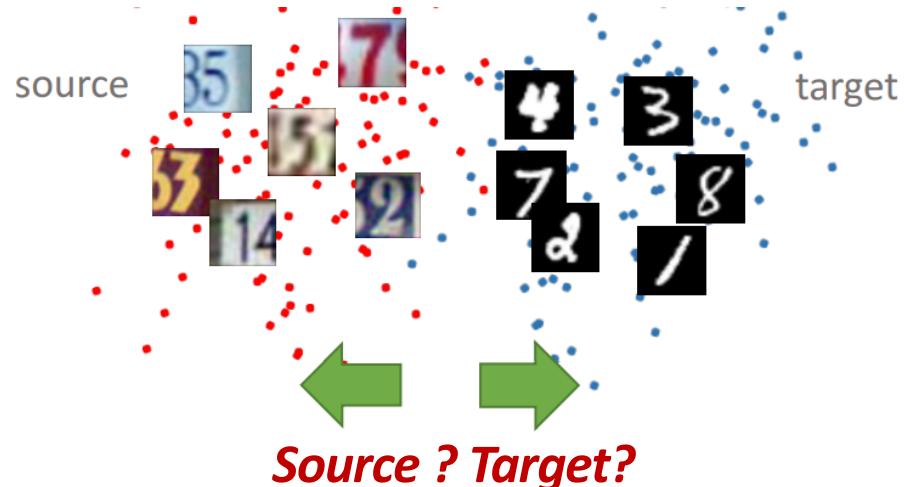
$$Q_s(y_s|X_s), Q_t(y_t|X_t)$$

- MMD loss $MMD^2(s, t) = \sup_{\|\phi\|_{\mathcal{H}} \leq 1} \|E_{x^s \sim s}[\phi(x^s)] - E_{x^t \sim t}[\phi(x^t)]\|_{\mathcal{H}}^2$

- Reconstruction loss

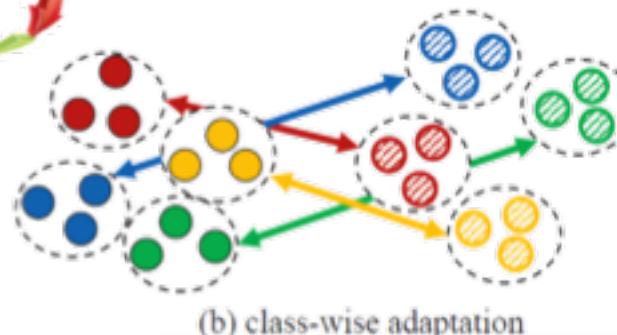
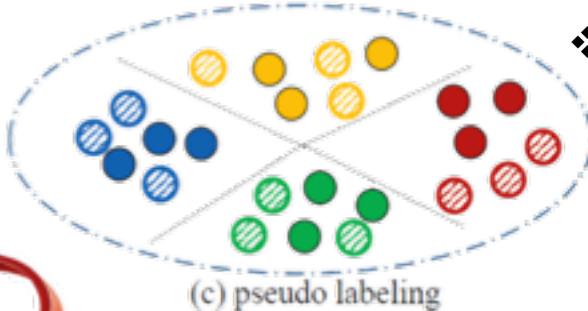
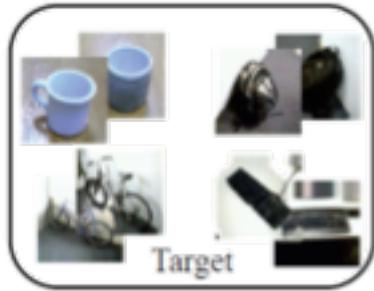
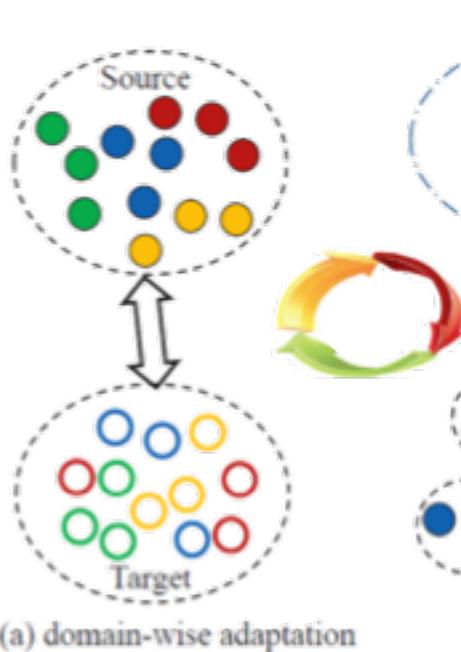


- Adversarial loss



- Parameters Sharing, e.g., weights/dictionary

Conditional MMD loss +Projection learning



❖ Class-wise MMD

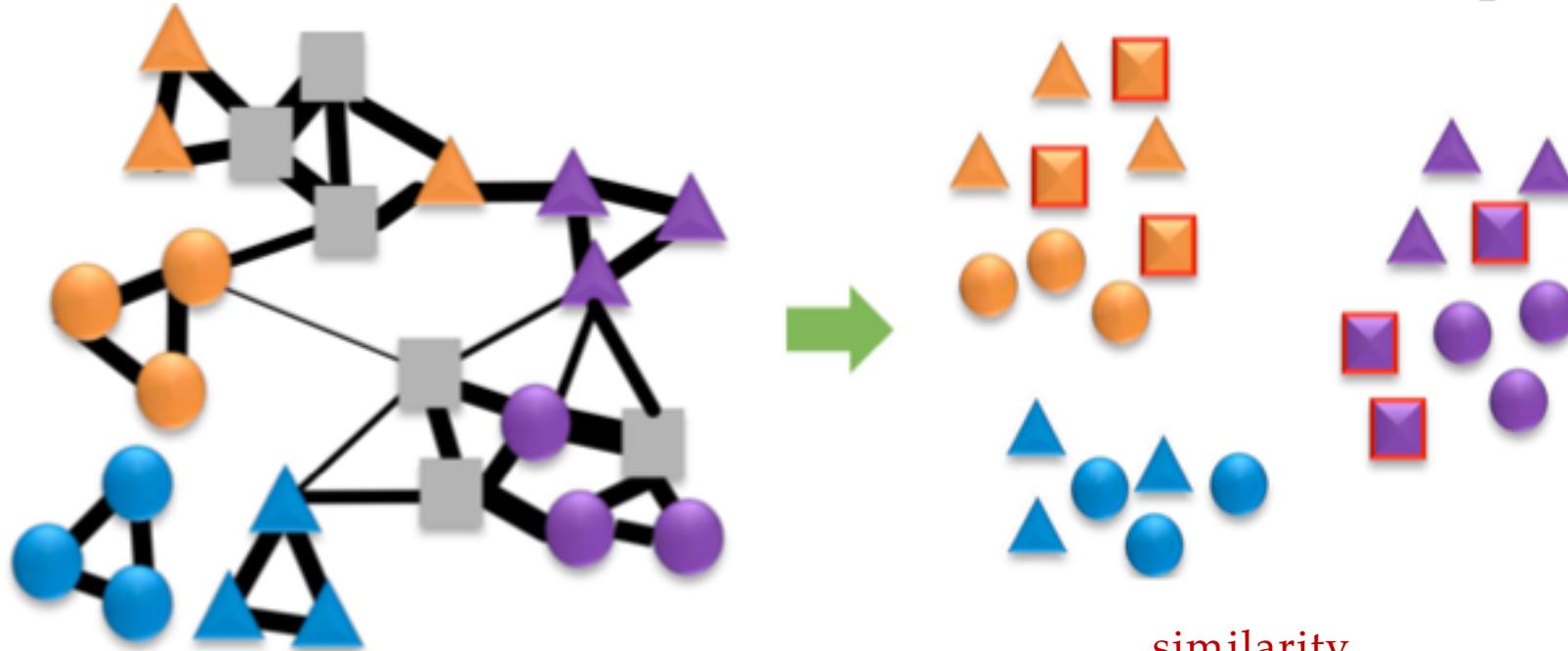
- Conditional distribution alignment
- Iteratively optimize target labels

❖ Domain-wise MMD *only aligns marginal distribution*

$$\left\| \frac{1}{n_s^{(c)}} \sum_{\mathbf{x}_i \in \mathcal{D}_s^{(c)}} \mathbf{A}^T \mathbf{x}_i - \frac{1}{n_t^{(c)}} \sum_{\mathbf{x}_j \in \mathcal{D}_t^{(c)}} \mathbf{A}^T \mathbf{x}_j \right\|^2$$

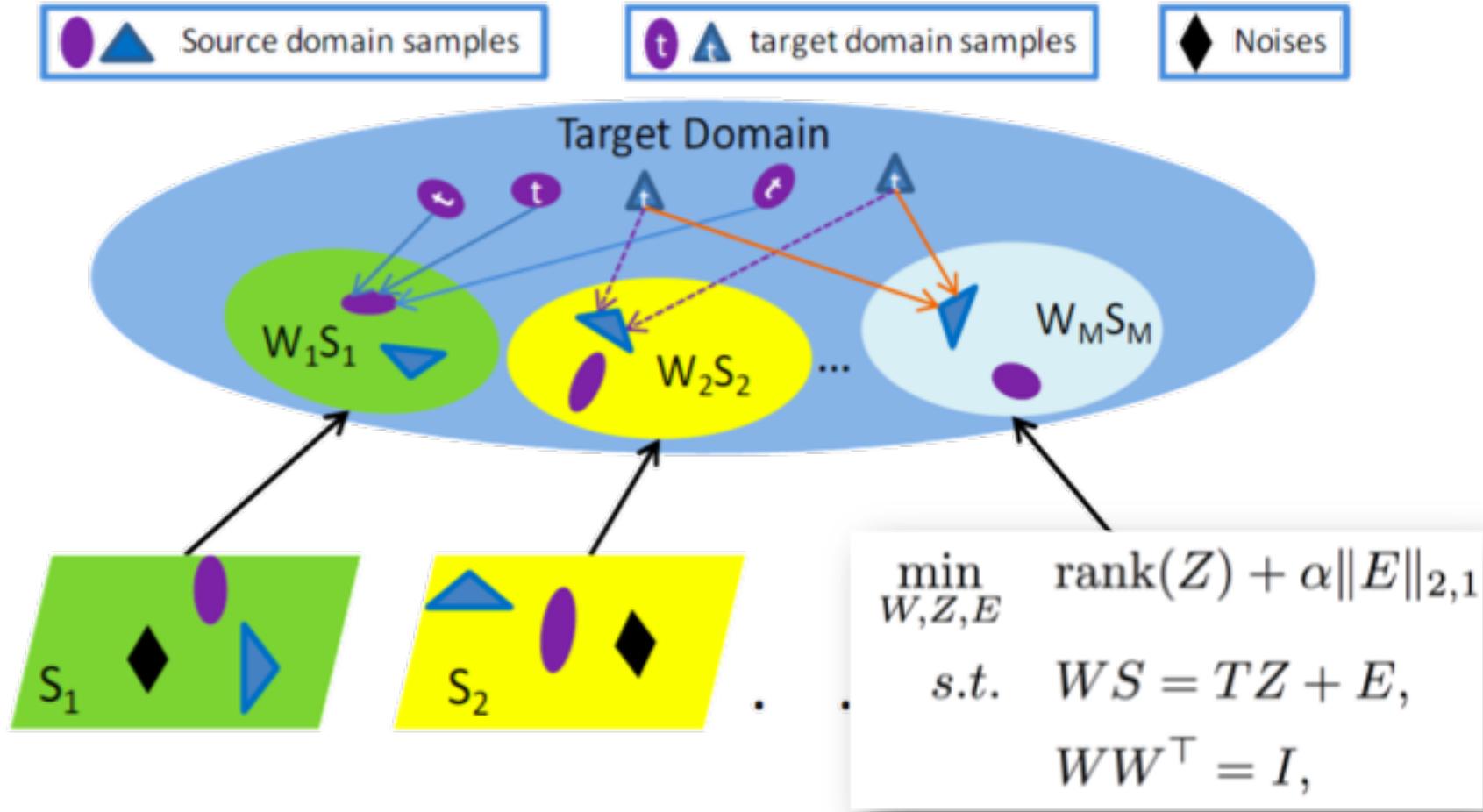
Modified MMD loss +Projection learning

- Exploiting label and latent-domain information within and across domains

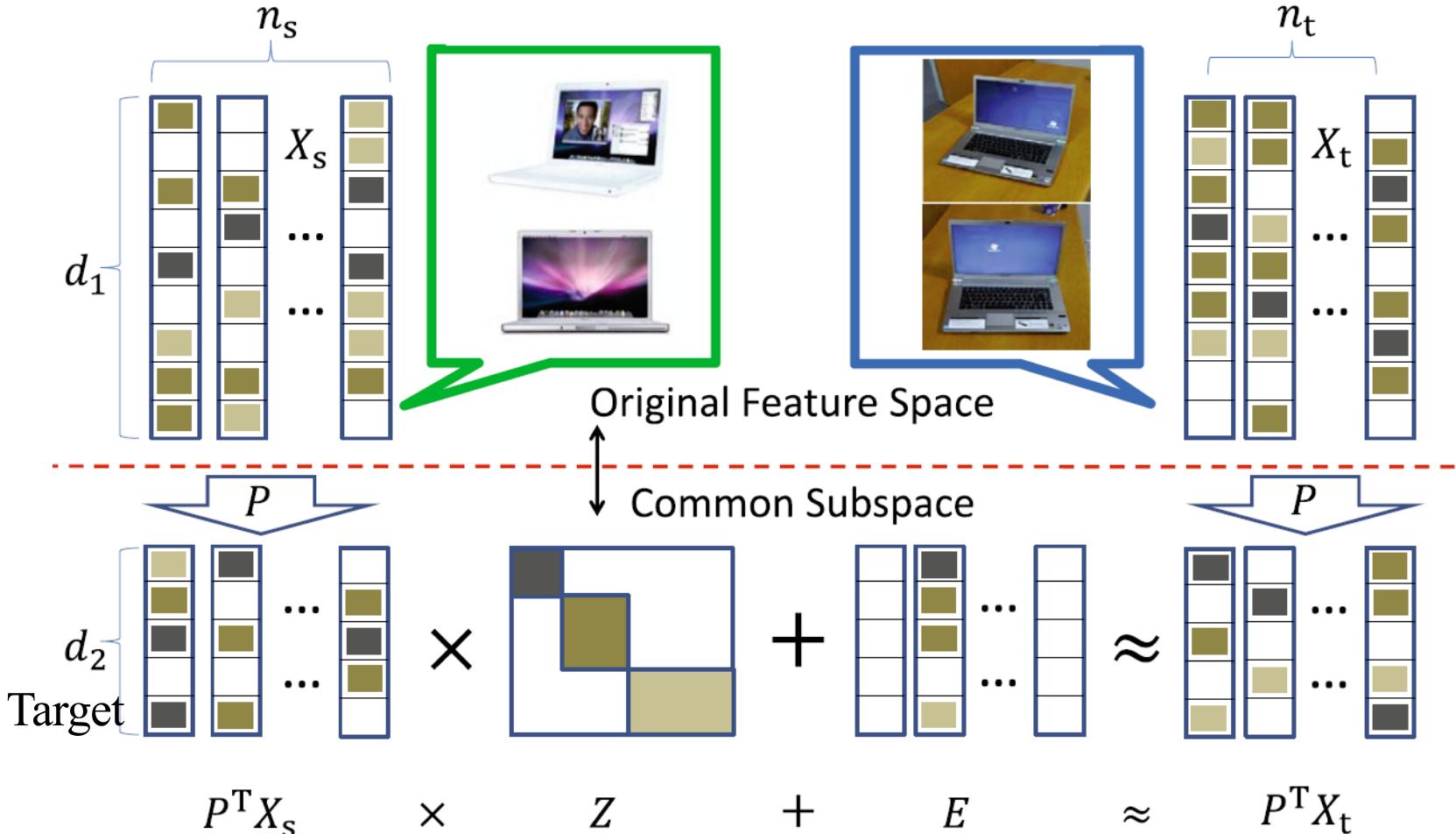


$$\mathcal{M}_{\phi,d}(\mathcal{P}_S(\mathbf{X}_S|\mathbf{y}_S), \mathcal{P}_T(\mathbf{X}_T|\mathbf{y}_T)) = \sum_{i,j} \frac{m_{ij}^{ST}}{\sum_k m_{ki}^{SS} \sum_l m_{lj}^{TT}} \left\| \hat{\phi}(\mathbf{x}_i^S) - \hat{\phi}(\mathbf{x}_j^T) \right\|^2$$

Reconstruction loss +Projection learning



Reconstruction loss +Projection learning



Dictionary Learning + MMD loss

□ Common Dictionary + Graph Regularizer [1]

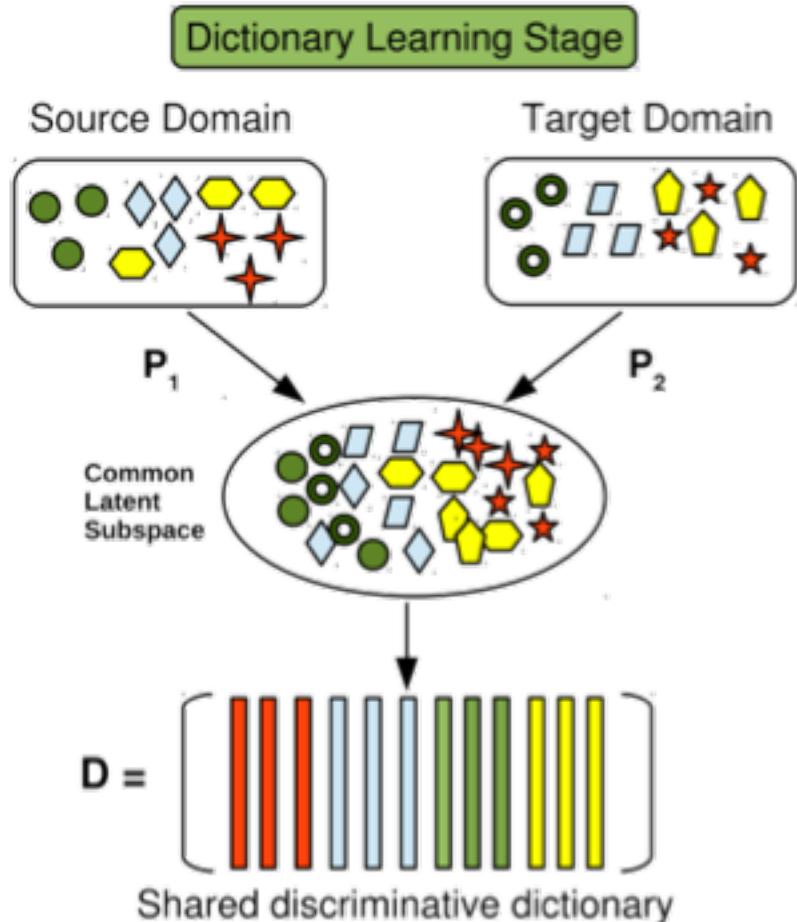
$$\begin{aligned} \min_{\mathbf{B}, \mathbf{S}} & \|\mathbf{X} - \mathbf{BS}\|_F^2 + \gamma \text{tr}(\mathbf{SLS}^T) + \lambda \sum_{i=1}^n |\mathbf{s}_i| \\ \text{s.t. } & \|\mathbf{b}_i\|^2 \leq c, i = 1, \dots, k \end{aligned}$$

□ MMD loss

$$\left\| \frac{1}{n_l} \sum_{i=1}^{n_l} \mathbf{s}_i - \frac{1}{n_u} \sum_{j=n_l+1}^{n_l+n_u} \mathbf{s}_j \right\|^2 = \sum_{i,j=1}^n \mathbf{s}_i^T \mathbf{s}_j M_{ij} = \text{tr}(\mathbf{SMS}^T)$$

-
- [1]. Gao et al. Local features are not lonely – Laplacian sparse coding for image classification, CVPR, 2010
 - [2]. Long et al. Transfer sparse coding for robust image representation. CVPR. 2013.

Common Dictionary +Projection learning



$$\{D^*, \tilde{P}^*, \tilde{X}^*\} = \arg \min_{D, \tilde{P}, \tilde{X}} C_1(D, \tilde{P}, \tilde{X}) + \lambda C_2(\tilde{P})$$

$$\text{s.t. } P_i P_i^T = I, i = 1, 2 \text{ and } \|\tilde{x}_j\|_0 \leq T_0, \forall j$$

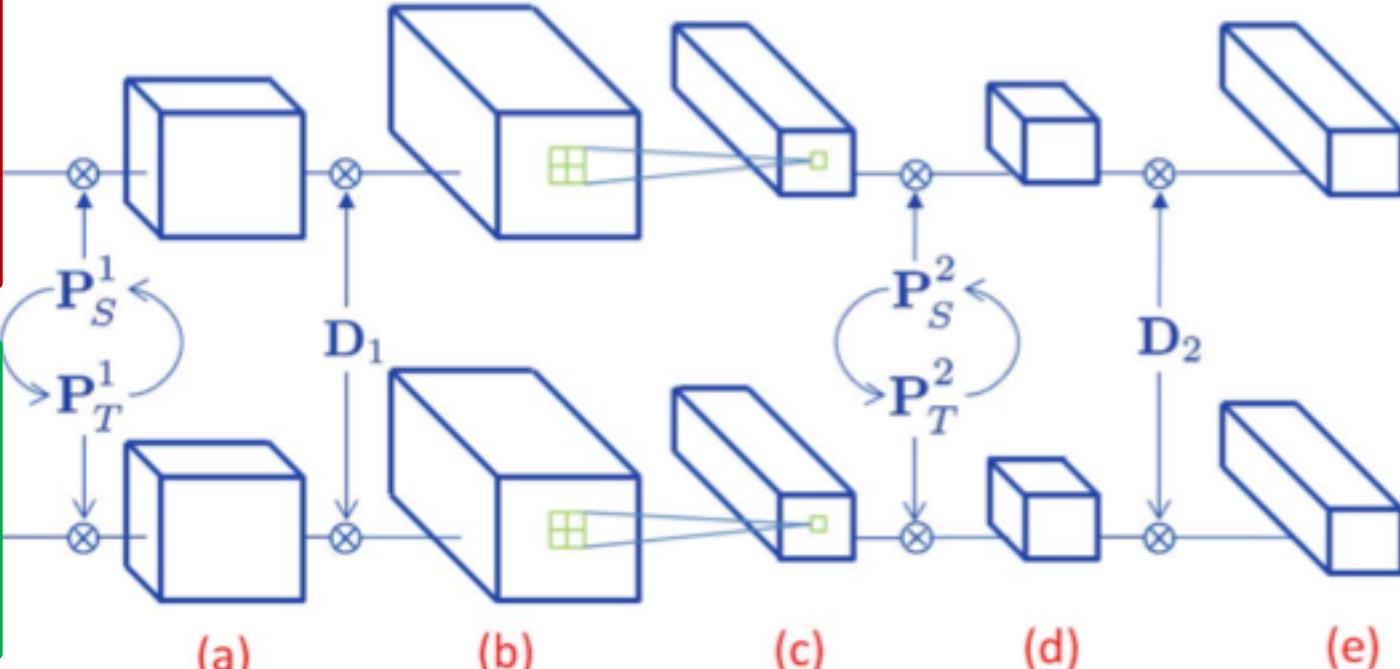
□ Common dictionary

$$C_1(D, \tilde{P}, \tilde{X}) = \|P_1 Y_1 - D X_1\|_F^2 + \\ \|P_2 Y_2 - D X_2\|_F^2$$

□ PCA Reconstruction

$$C_2(\tilde{P}) = \|Y_1 - P_1^T P_1 Y_1\|_F^2 + \\ \|Y_2 - P_2^T P_2 Y_2\|_F^2.$$

Common Dictionaries +Projection learning



$$\mathcal{L}(\mathbf{Y}_S, \mathbf{P}_S, \mathbf{D}, \mathbf{X}_S, \alpha, \beta) + \lambda \mathcal{L}(\mathbf{Y}_T, \mathbf{P}_T, \mathbf{D}, \mathbf{X}_T, \alpha, \beta)$$

$$\text{s.t. } \mathbf{P}_S \mathbf{P}_S^T = \mathbf{P}_T \mathbf{P}_T^T = \mathbf{I}, \quad \|\mathbf{d}_i\|_2 = 1, \quad \forall i \in [1, K],$$

$$\mathcal{L}(\mathbf{Y}, \mathbf{P}, \mathbf{D}, \mathbf{X}, \alpha, \beta)$$

$$= \|\mathbf{PY} - \mathbf{DX}\|_F^2 + \alpha \|\mathbf{Y} - \mathbf{P}^T \mathbf{PY}\|_F^2 + \beta \|\mathbf{X}\|_1$$

[1] Nguyen et al. DASH-N: Joint hierarchical domain adaptation and feature learning. *IEEE TIP* 2015

[2]. Ding et al. Deep Low-rank Coding for Transfer Learning, *IJCAI* 2015

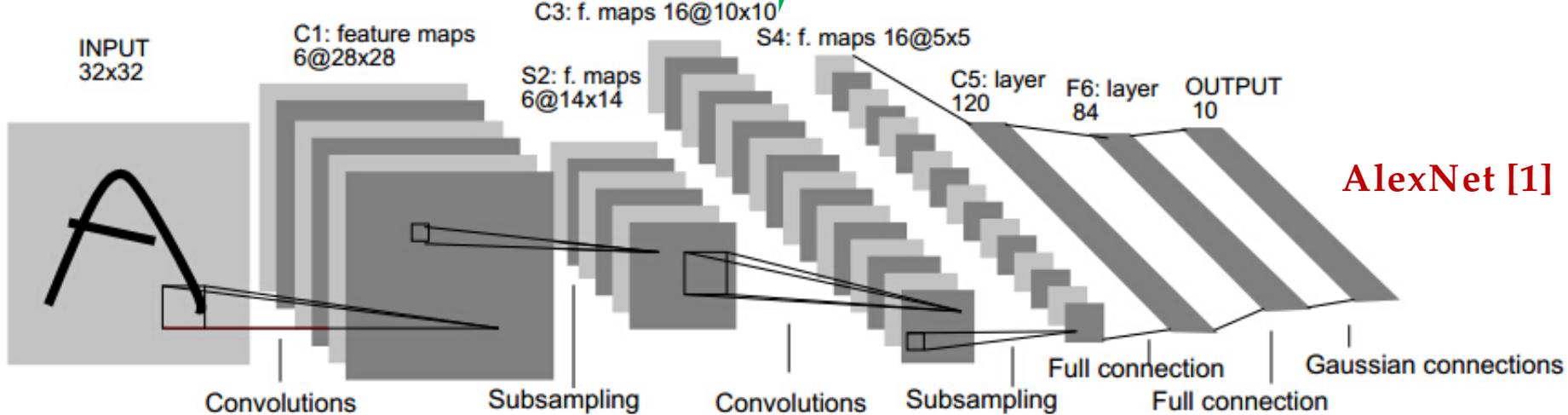


Deep Learning

Learner: $f : x \rightarrow y$



Conduct knowledge transfer



Pre-trained model on large-scale datasets, e.g., ImageNet



Top layers will be more task-specific [2]



Align different domains

[1]. Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. *NIPS*. 2012.

[2]. Yosinski, J., Clune, J., Bengio, Y., and Lipson, H. How transferable are features in deep neural networks? In *NIPS*, 2014

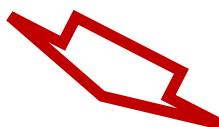
Unified Deep Domain Adaptation

❖ Two-Step Domain Adaptation

Deep Features + Shallow DA methods

$$\min_{f_1(\cdot), \dots, f_v(\cdot)} \sum_{i=1, i < j}^v \mathcal{A}(f_i(X_i), f_j(X_j)) + \lambda \sum_{k=1}^v \mathcal{R}(f_k(X_k))$$

❖ Deep Domain Adaptation



Feature Learning +
Domain Alignment

- MMD loss
- Adversarial loss

- Parameters Sharing



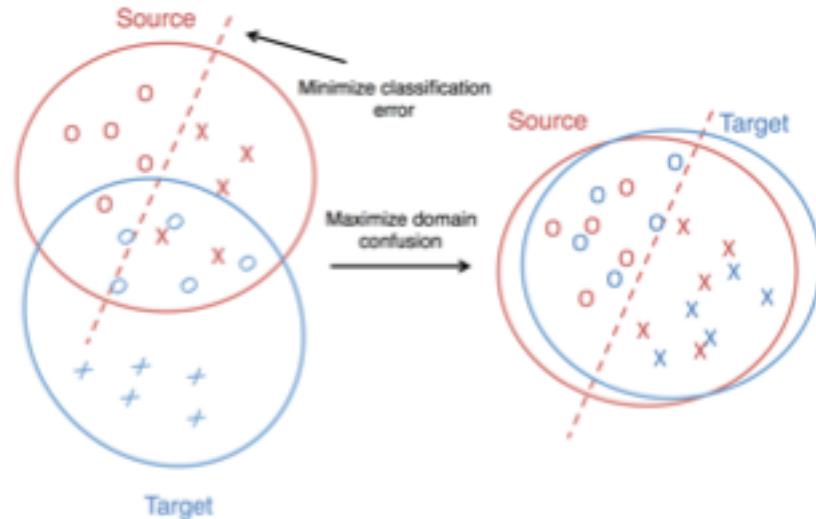
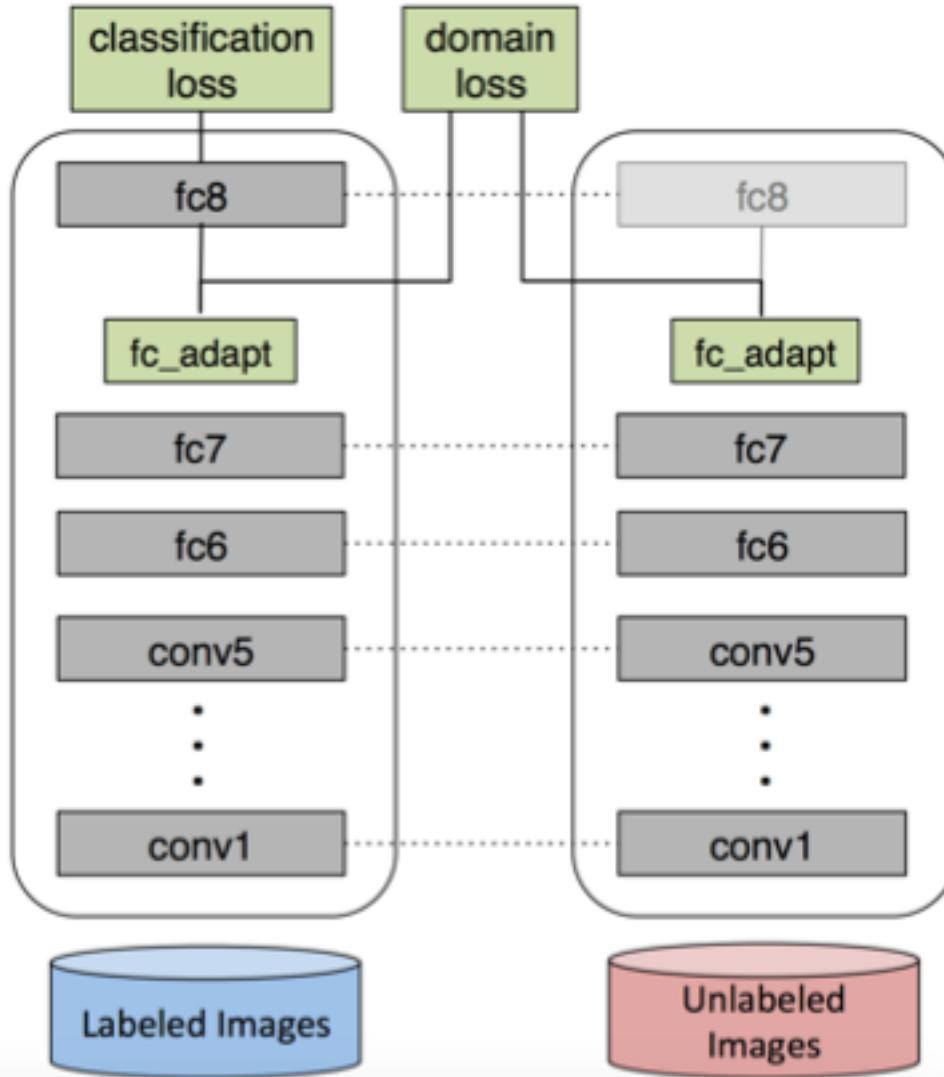
Labeled source domain

$L_c(X_s, y_s)$



Soft labels of target domain

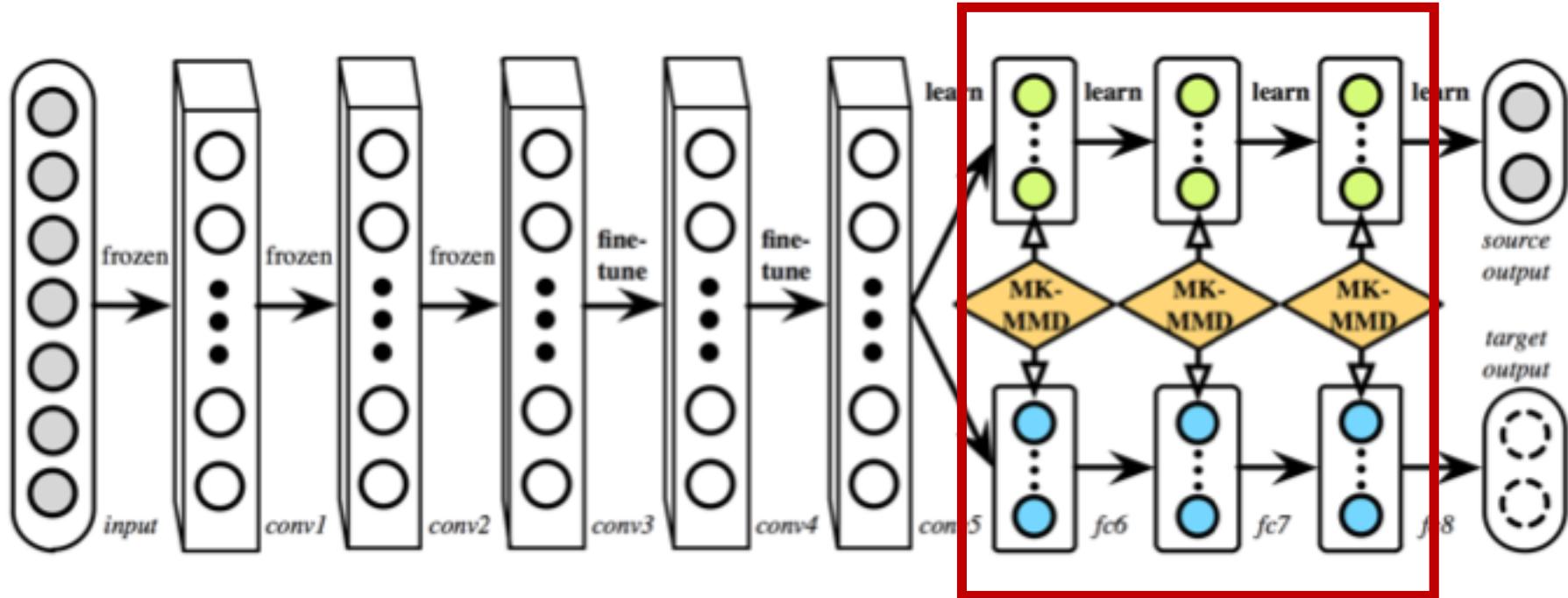
Deep Learning+ MMD Loss



$$\begin{aligned}\mathcal{L} = & \mathcal{L}_C(X_L, y) \\ & + \lambda \text{MMD}^2(X_S, X_T)\end{aligned}$$

Tzeng, Eric, et al. "Deep domain confusion: Maximizing for domain invariance." *arXiv preprint arXiv:1412.3474* (2014).

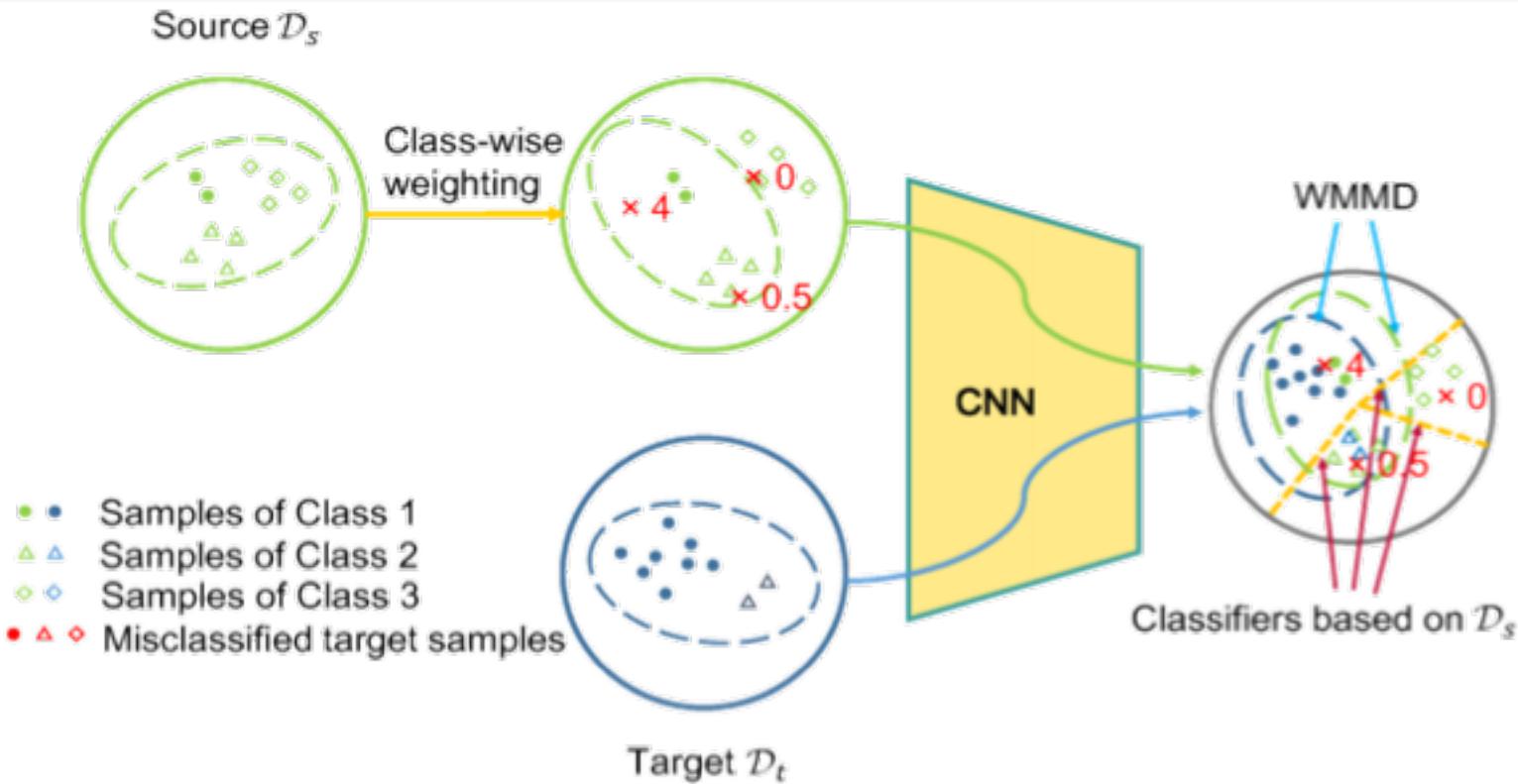
Deep Learning+ MK-MMD Loss



$$\min_{\theta \in \Theta} \max_{k \in \mathcal{K}} \frac{1}{n_a} \sum_{i=1}^{n_a} J(\theta(\mathbf{x}_i^a), y_i^a) + \lambda \sum_{\ell=l_1}^{l_2} d_k^2(\mathcal{D}_s^\ell, \mathcal{D}_t^\ell) \sigma_k^{-2}$$

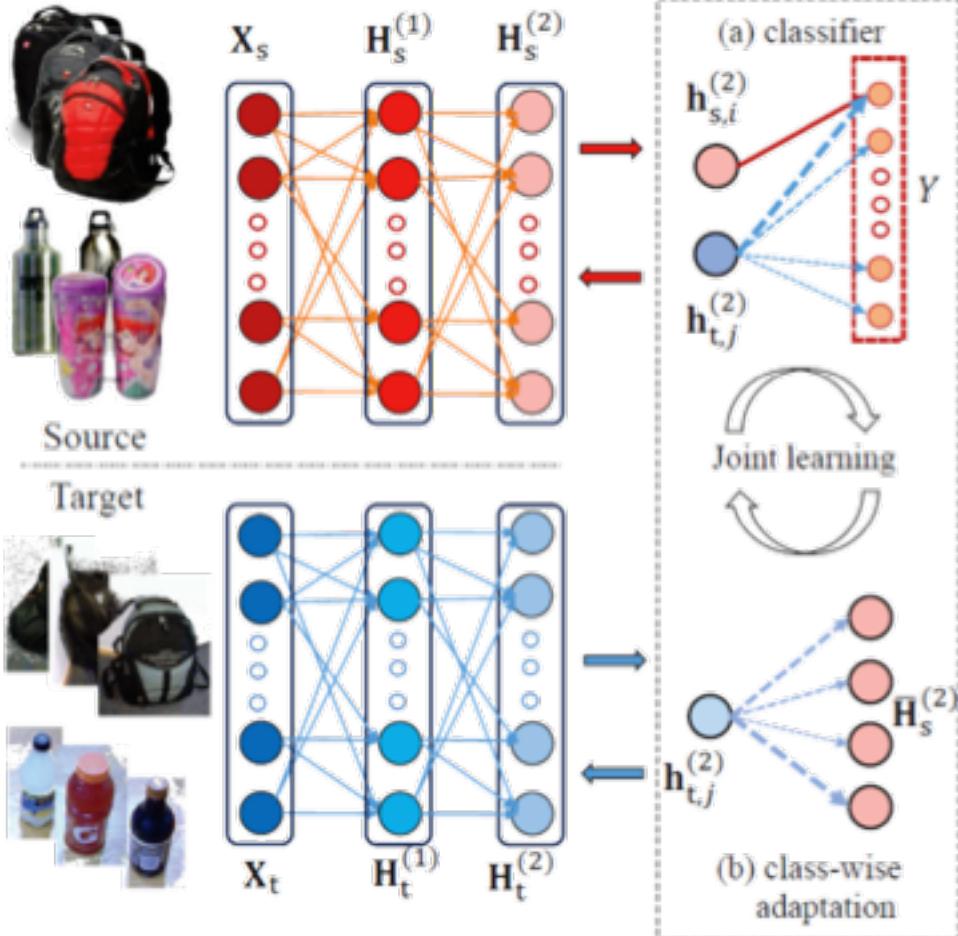
multi-layer adaptation

Deep Learning+ Weighted MMD loss



$$\text{MMD}^2(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^N \phi(\mathbf{x}_j^t) \right\|_{\mathcal{H}}^2 \quad \Rightarrow \quad \text{MMD}_w^2(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{\sum_{i=1}^M \alpha_{y_i^s}} \sum_{i=1}^M \alpha_{y_i^s} \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^N \phi(\mathbf{x}_j^t) \right\|_{\mathcal{H}}^2$$

Deep Learning+ Probabilistic MMD loss



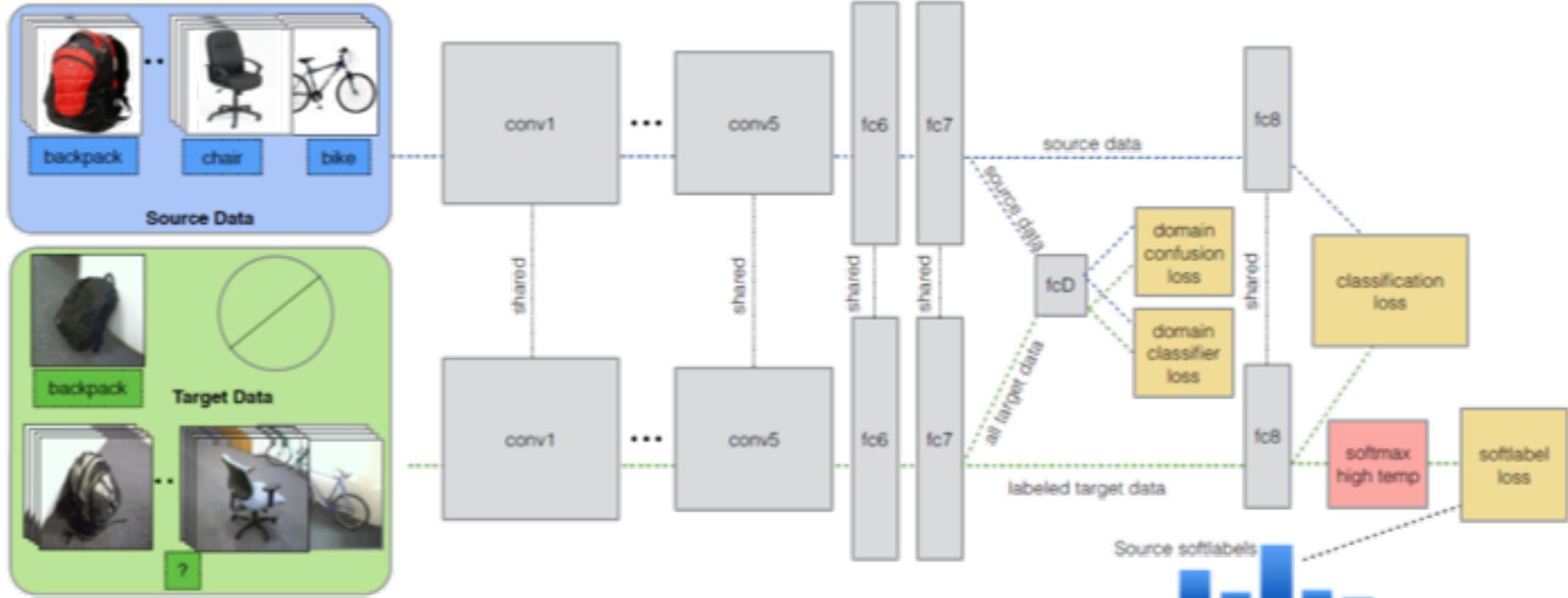
$$\text{MMD}^2(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^N \phi(\mathbf{x}_j^t) \right\|_{\mathcal{H}}^2$$



$$\frac{1}{n_t} \sum_{j=1}^{n_t} \|\mathbf{h}_{t,j}^{(L)} - \sum_{c=1}^C p_{t,c}^j \bar{\mathbf{h}}_{s,c}^{(L)}\|_2^2$$

- the center of the c -th class source data;
- the probability of the j -th target point to be assigned to the label of the c -th class (*Softmax output of target data*)

Deep Learning+ Adversarial loss



$$\mathcal{L}(x_S, y_S, x_T, y_T, \theta_D; \theta_{\text{repr}}, \theta_C) = \mathcal{L}_C(x_S, y_S, x_T, y_T; \theta_{\text{repr}}, \theta_C)$$

Classifier loss

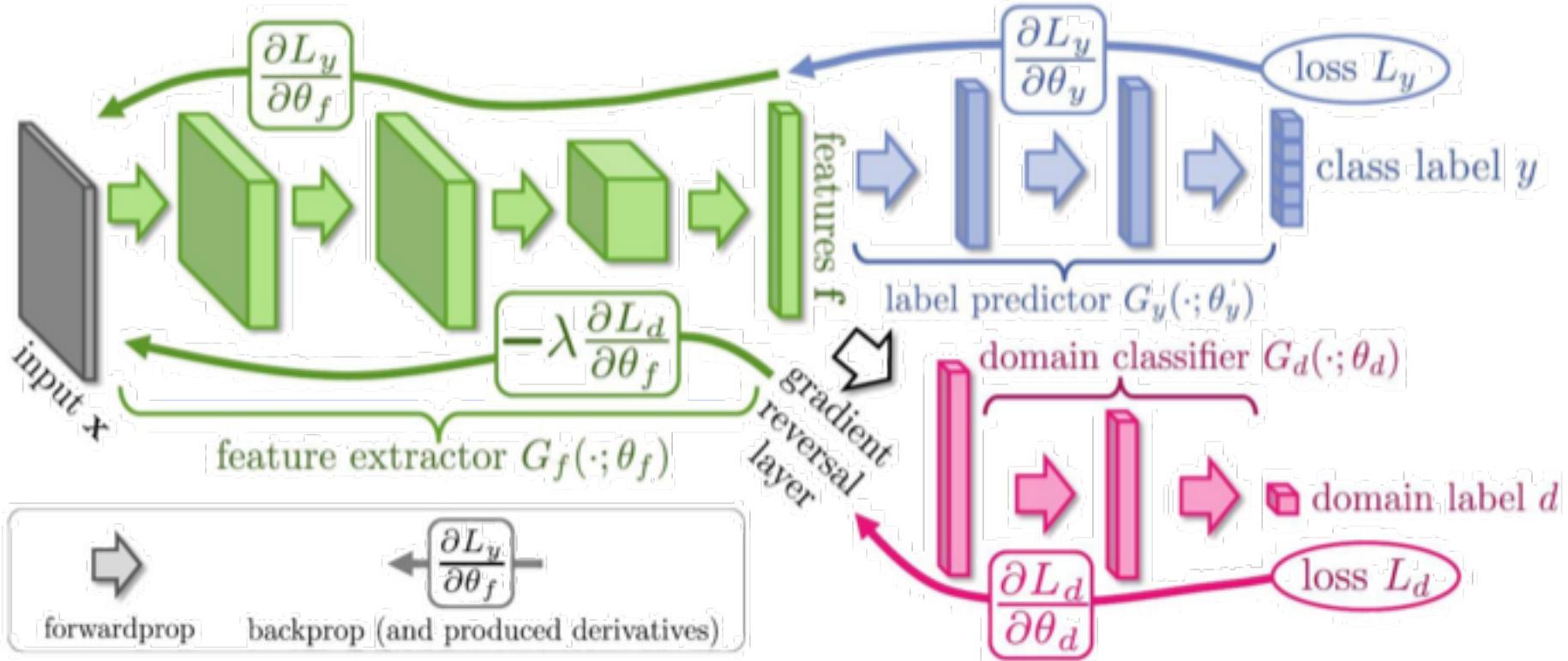
$$+ \lambda \mathcal{L}_{\text{conf}}(x_S, x_T, \theta_D; \theta_{\text{repr}})$$

$$+ \nu \mathcal{L}_{\text{soft}}(x_T, y_T; \theta_{\text{repr}}, \theta_C).$$

Domain confusion loss

Soft label loss

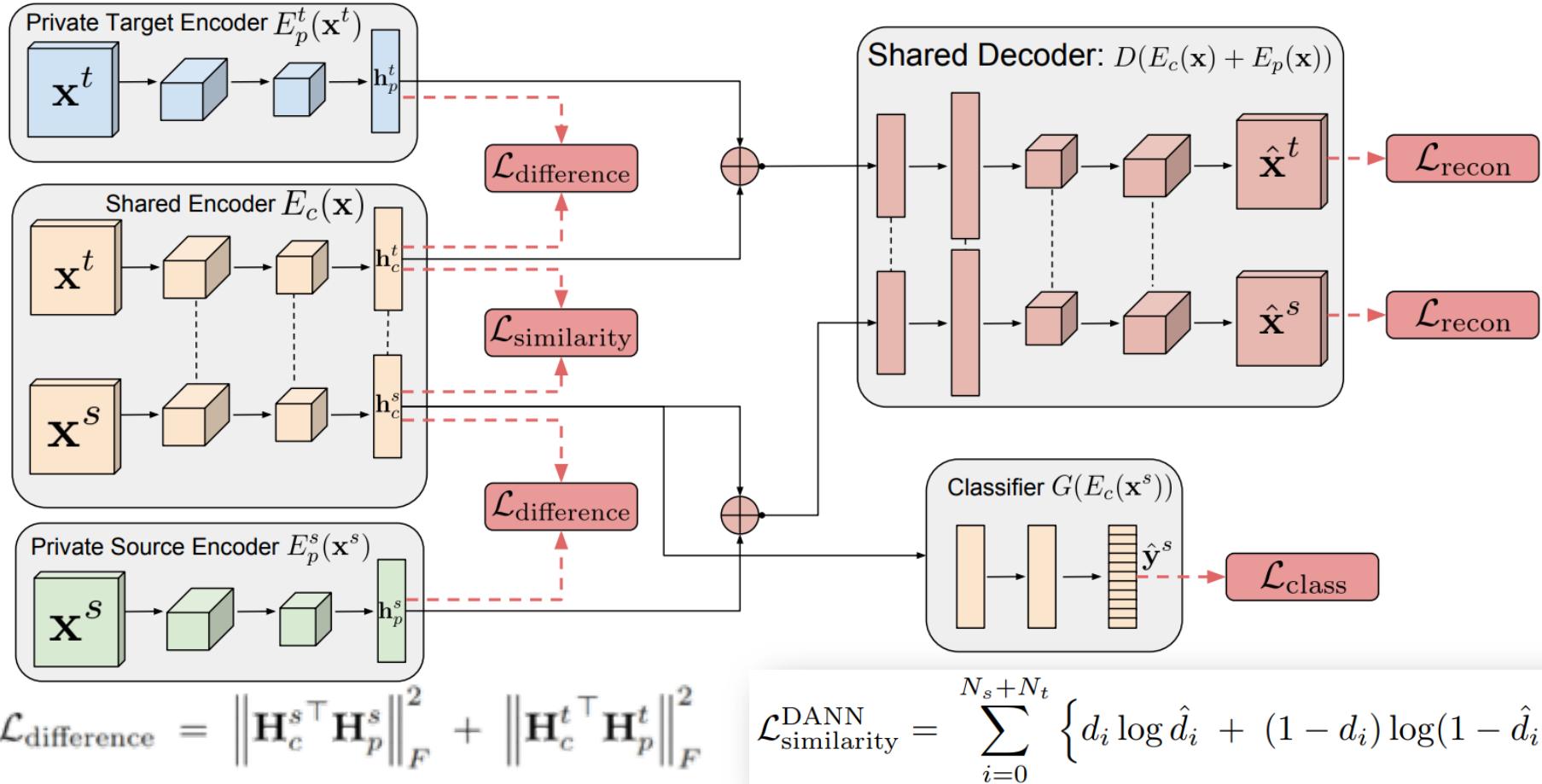
Deep Learning+ Adversarial loss



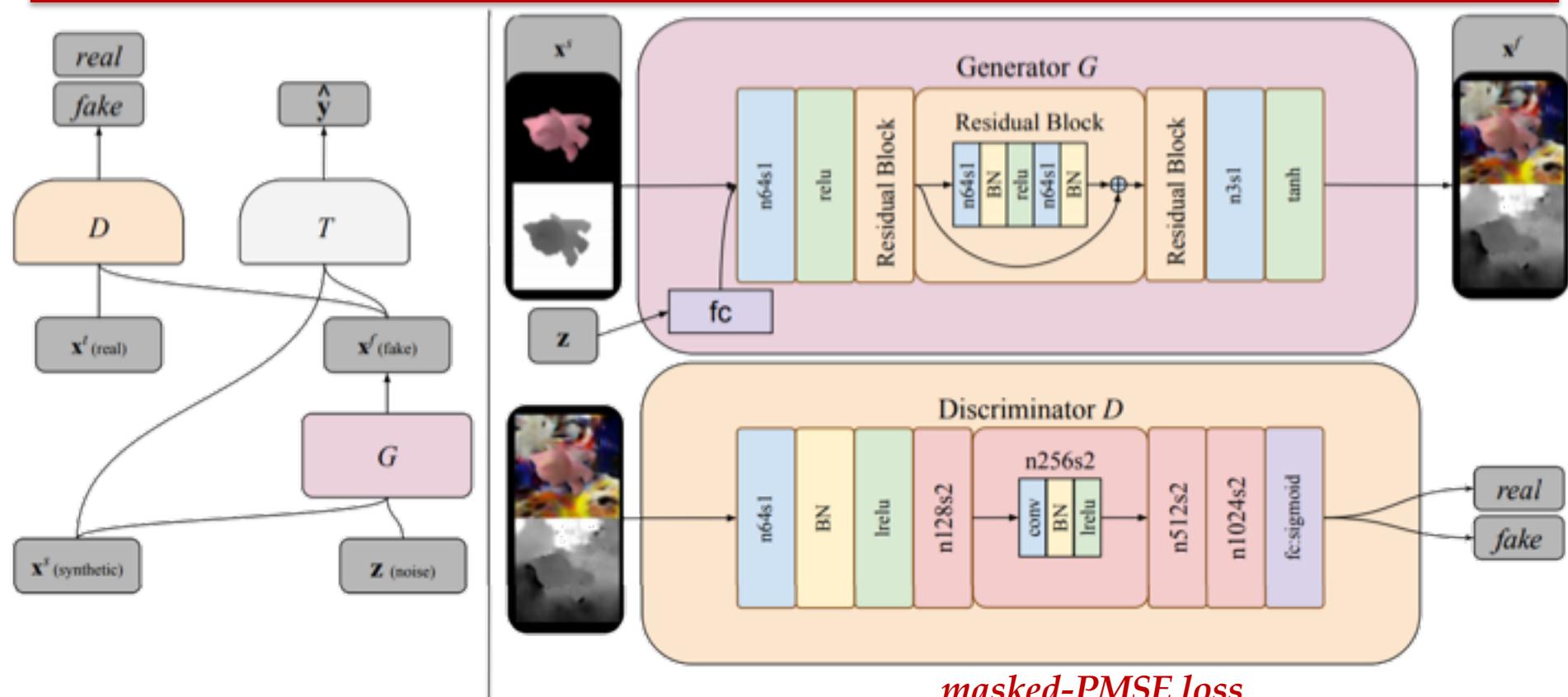
$$E(\theta_f, \theta_y, \theta_d) = \sum_{\substack{i=1..N \\ d_i=0}} L_y(G_y(G_f(\mathbf{x}_i; \theta_f); \theta_y), y_i) - \lambda \sum_{i=1..N} L_d(G_d(G_f(\mathbf{x}_i; \theta_f); \theta_d), d_i)$$

Ganin, Yaroslav, and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. *ICML*, 2015

Deep Learning+ Adversarial loss



Deep Learning+ Adversarial loss



$$\min_{\theta_G, \theta_T} \max_{\theta_D} \alpha \mathcal{L}_d(D, G) + \beta \mathcal{L}_t(T, G) + \gamma \boxed{\mathcal{L}_c(G)}$$

masked-PMSE loss

$$\begin{aligned} \mathcal{L}_c(G) = & \mathbb{E}_{\mathbf{x}^s, \mathbf{z}} \left[\frac{1}{k} \|(\mathbf{x}^s - G(\mathbf{x}^s, \mathbf{z}; \theta_G)) \circ \mathbf{m}\|_2^2 \right. \\ & \left. - \frac{1}{k^2} ((\mathbf{x}^s - G(\mathbf{x}^s, \mathbf{z}; \theta_G))^T \mathbf{m})^2 \right] \end{aligned}$$

Outline

□ Introduction & Background

- Multi-view Visual Data
- Multi-view Learning Problems
- Multi-view Learning Taxonomy

□ Multi-view Learning

- Projection and Embedding
- Knowledge Fusion
- Multi-view Clustering
- Supervised Multi-view Learning → Zero-shot Learning

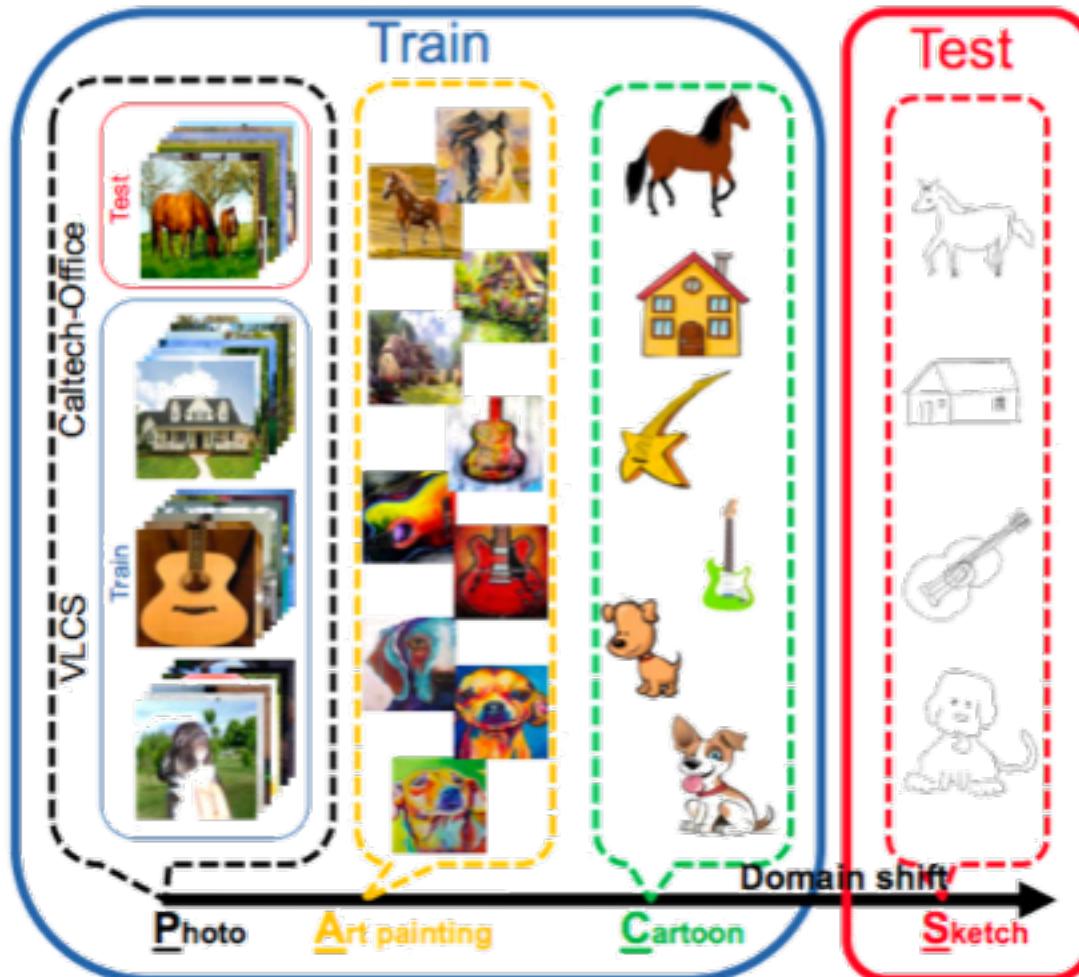
□ Domain Adaptation

- Transfer Learning → Domain Adaptation
- Multi-Source Domain Adaptation & Domain Generalization

□ Conclusion



Multi-Source Domains



Domain Adaptation

- Labeled Source
- Unlabeled or Limited Labeled Target



Multi-Source Domain Adaptation

- Multiple Labeled Sources
- Unlabeled or Limited Labeled Target

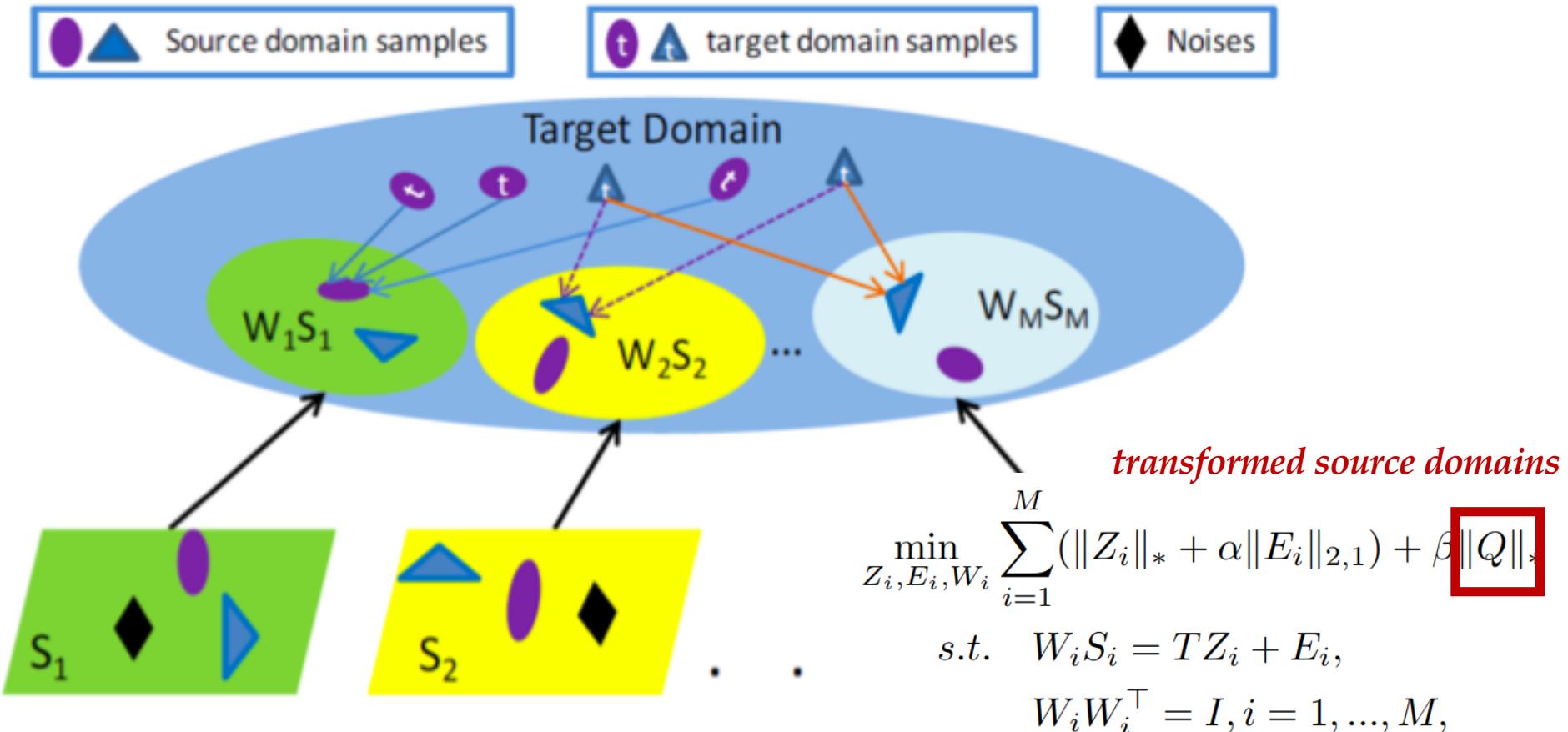


Domain Generalization

- Multiple Labeled Sources
- Missing Targets

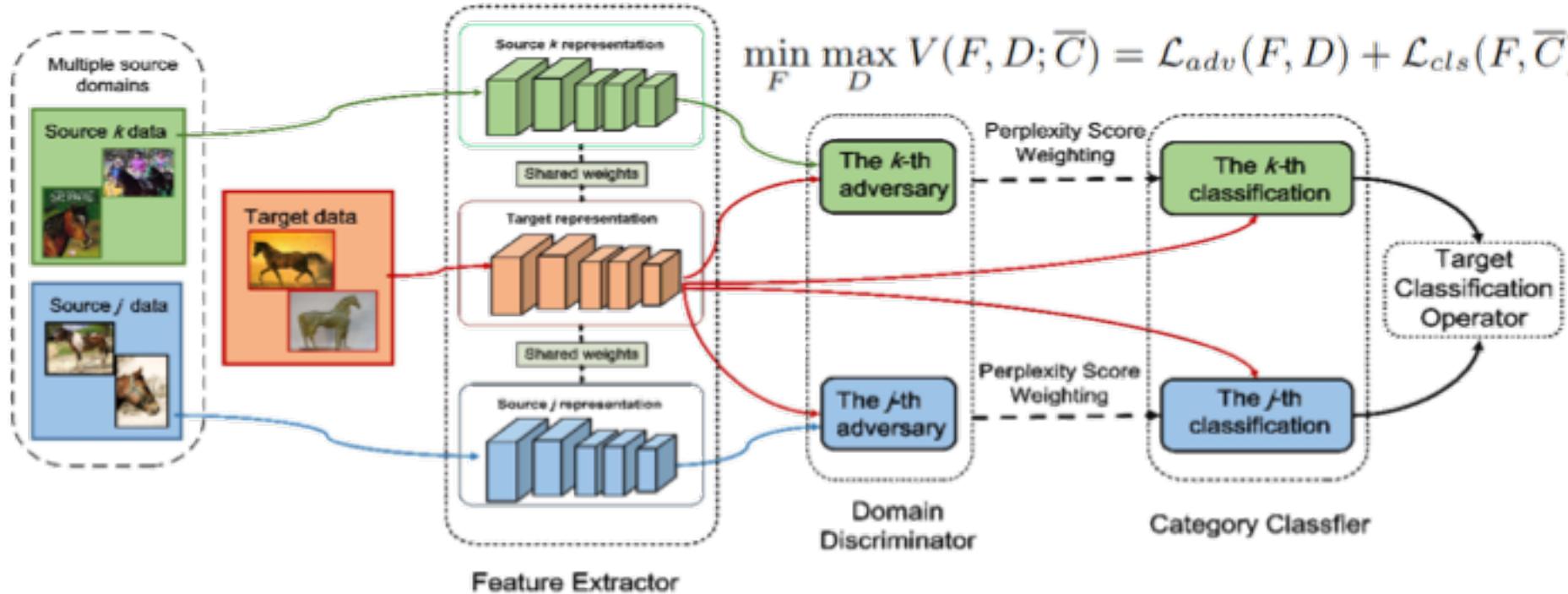


Reconstruction loss





Multi-Source Adaptation



$$\mathcal{L}_{adv}(F, D) = \frac{1}{N} \sum_j^N \mathbb{E}_{x \sim X_{s_j}} [\log D_{s_j}(F(x))] + \mathbb{E}_{x^t \sim X_t} [\log(1 - D_{s_j}(F(x^t)))]$$

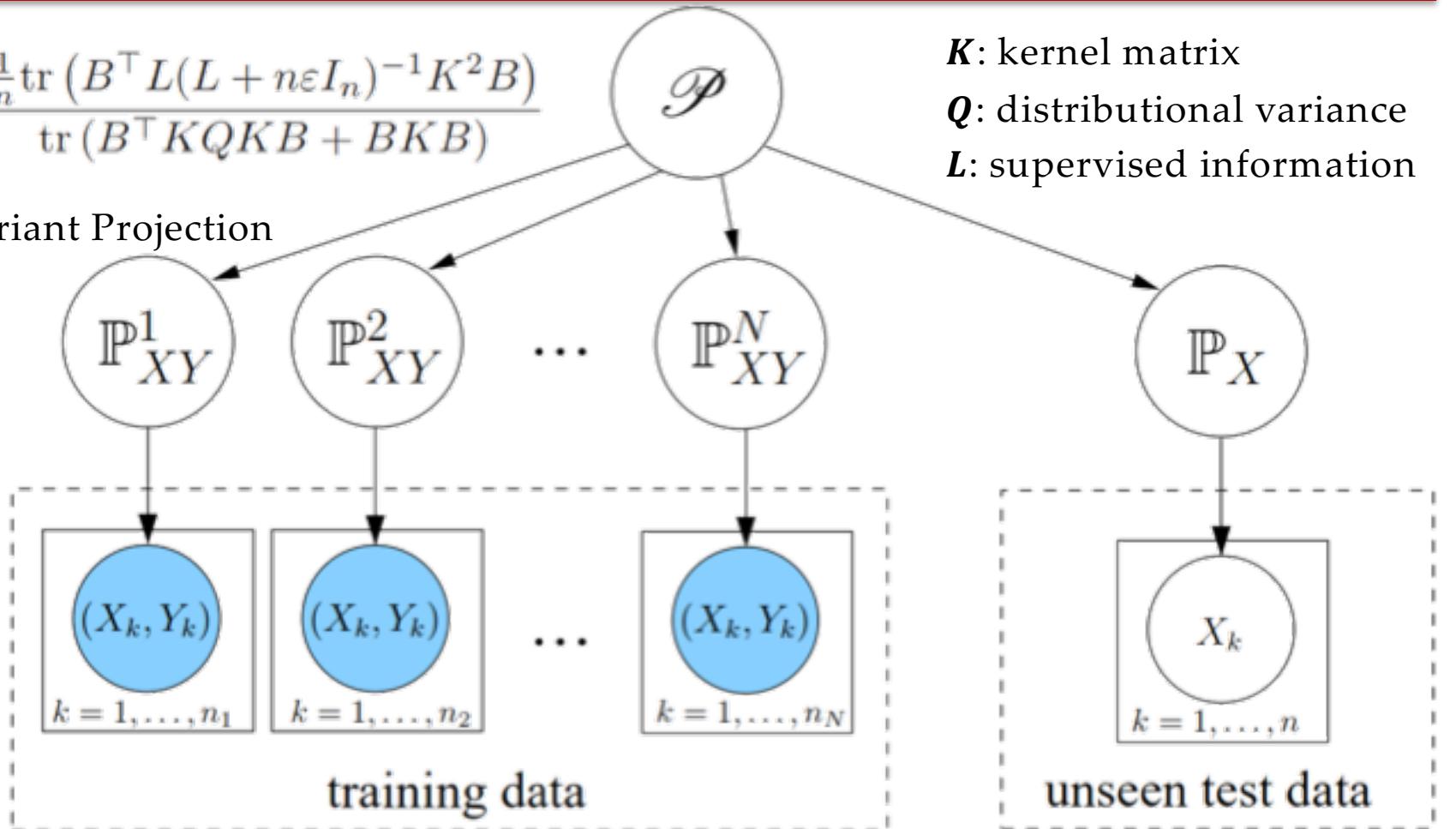
$$\min_{F, C} \mathcal{L}_{cls}(F, C) = \sum_j^N \mathbb{E}_{(x, y) \sim (X_{s_j}, Y_{s_j})} [\mathcal{L}(C_{s_j}(F(x)), y)] + \mathbb{E}_{(x^t, \hat{y}) \sim (X_t^p, Y_t^p)} [\sum_{\hat{y} \in \mathcal{C}_{\hat{s}}} \mathcal{L}(C_{\hat{s}}(F(x^t)), \hat{y})]$$



Domain Generalization

$$\max_{B \in \mathbb{R}^{n \times m}} \frac{\frac{1}{n} \text{tr} (B^\top L(L + n\varepsilon I_n)^{-1} K^2 B)}{\text{tr} (B^\top K Q K B + BKB)}$$

B : Invariant Projection



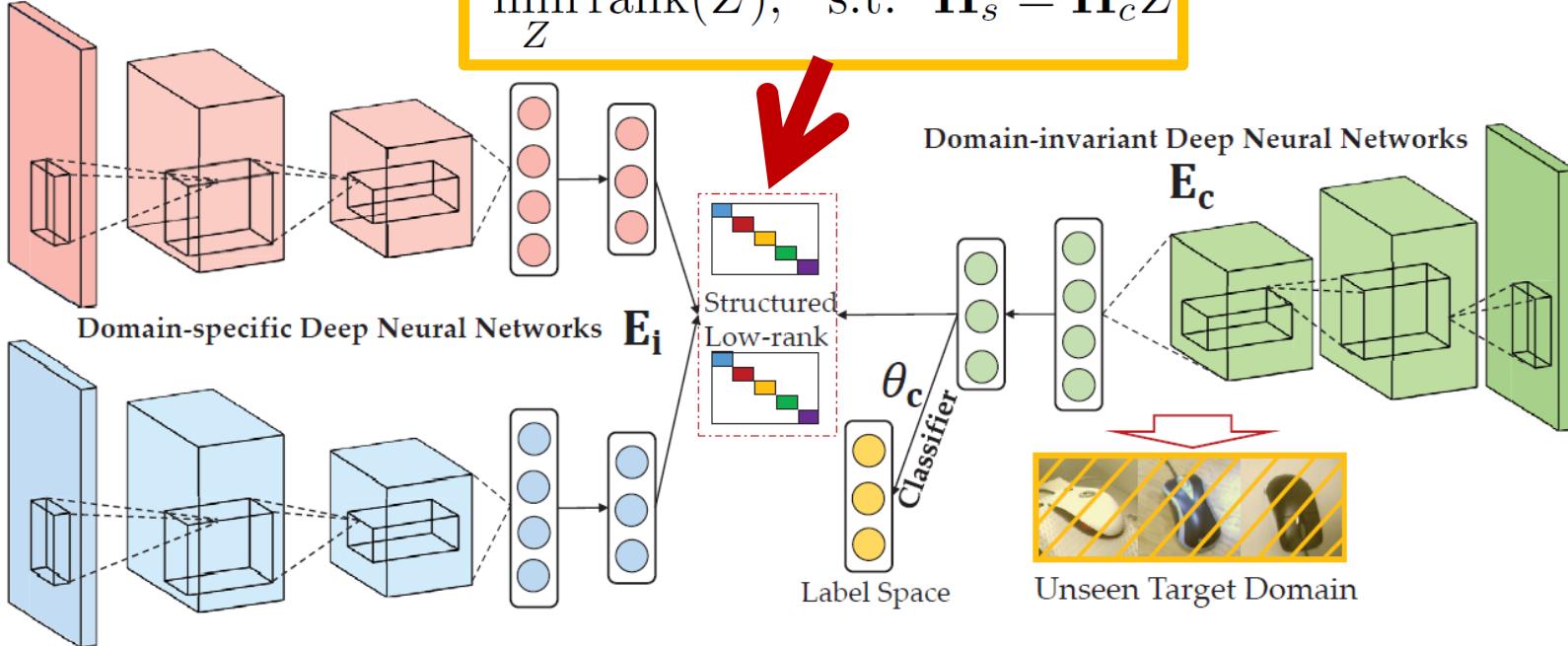
Muaned, Krikamol, David Balduzzi, and Bernhard Schölkopf. "Domain generalization via invariant feature representation." *International Conference on Machine Learning*. 2013.



Domain Generalization



$$\min_{Z} \text{rank}(Z), \quad \text{s.t. } \mathbf{H}_s = \mathbf{H}_c Z$$

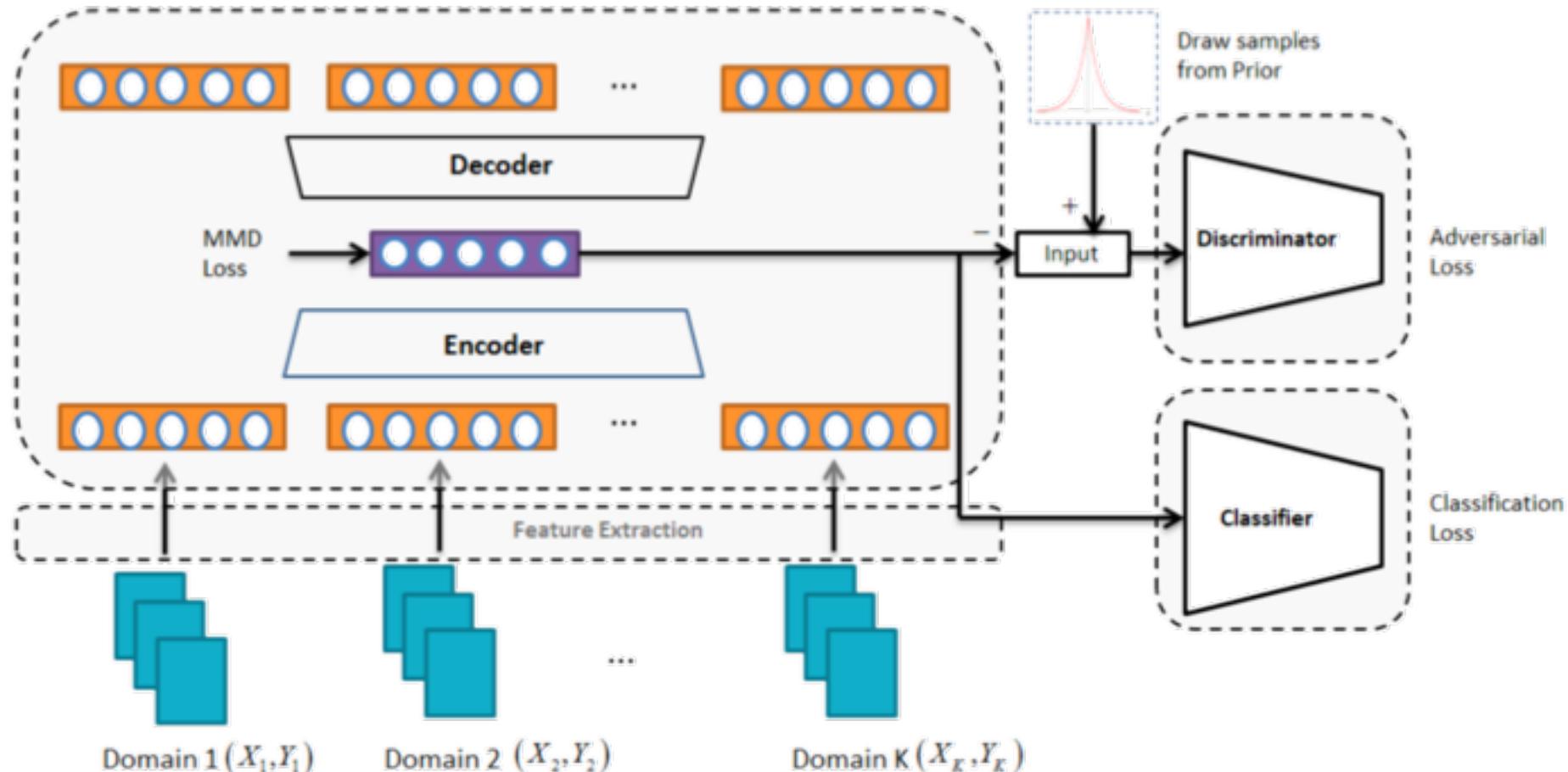


- (a) Multiple domain-specific deep structures tend to be learned to capture the rich information from each source.
- (b) A domain-invariant deep structure is built for all the domains, and further generalize to the unseen domain in the test stage with learned classifier .
- (c) Low-rank reconstruction is adopted to align two types of networks in structured low-rank fashion.

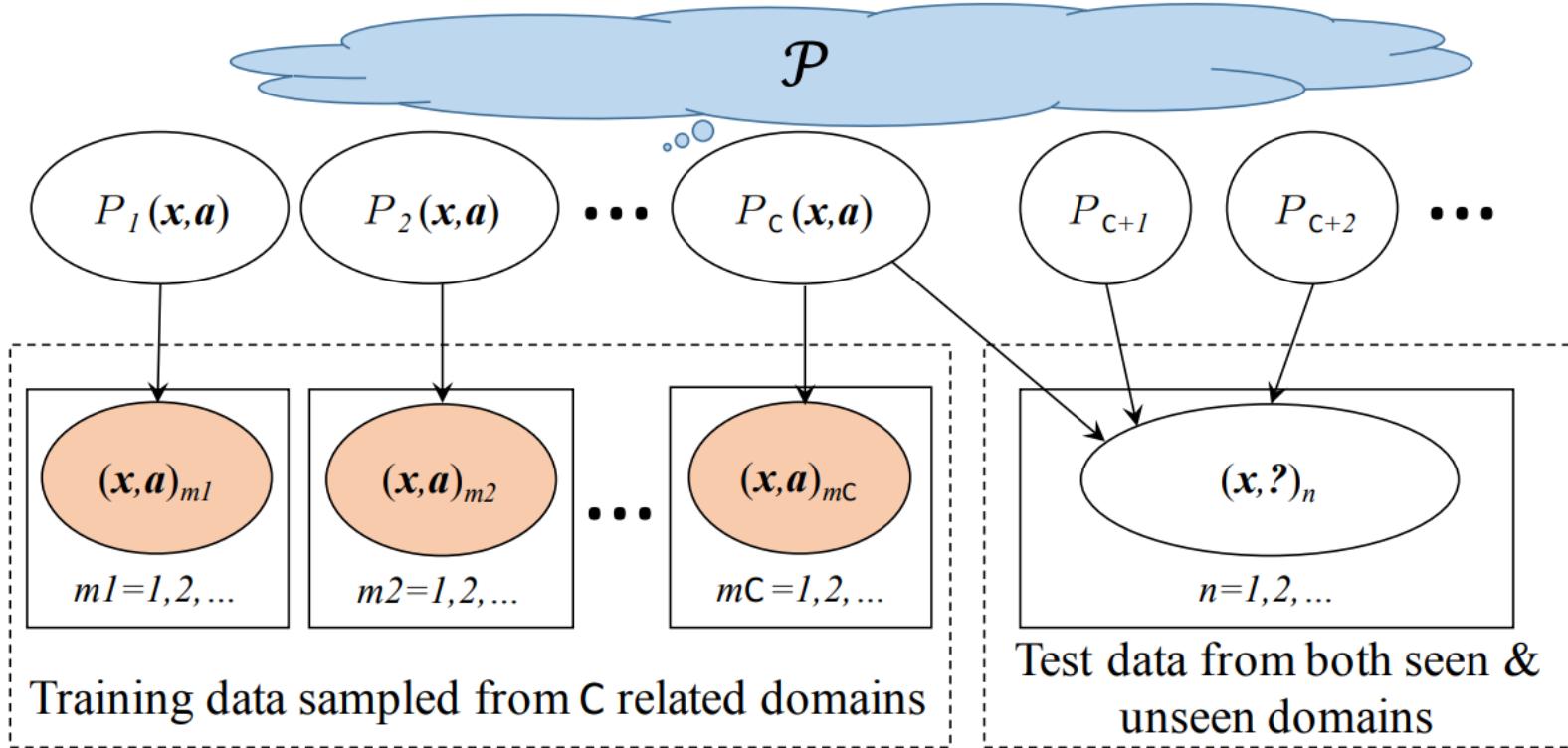


Domain Generalization

$$\min_{C, Q, P} \max_D \mathcal{L}_{\text{err}} + \lambda_0 \mathcal{L}_{\text{ae}} + \lambda_1 \mathcal{R}_{\text{mmd}} + \lambda_2 \mathcal{J}_{\text{gan}}$$



Domain Generalization & Zero-Shot Learning



$$\max_B \quad \frac{\text{tr}(\gamma B^T K^2 B / M + (1 - \gamma) B^T K L K B)}{\text{tr}(B^T K Q K B + B^T K B)}$$

Outline

□ Introduction & Background

- Multi-view Visual Data
- Multi-view Learning Problems
- Multi-view Learning Taxonomy

□ Multi-view Learning

- Projection and Embedding
- Knowledge Fusion
- Multi-view Clustering
- Supervised Multi-view Learning → Zero-shot Learning

□ Domain Adaptation

- Transfer Learning → Domain Adaptation
- Multi-Source Domain Adaptation & Domain Generalization

□ Conclusion

Conclusion

□ Taxonomy

■ Multi-view Data

sample-wise
correspondence

class-wise
correspondence

■ Multi-view Problems

projection/embedding; clustering; classification
zero-shot learning, domain adaptation/generalization

■ Multi-view Knowledge

Knowledge Integration; Knowledge Transfer

□ Unified model

From Shallow To Deep

Feature Learning + View Alignment



$$\min_{f_1(\cdot), \dots, f_v(\cdot)} \sum_{i=1, i < j}^v \mathcal{A}(f_i(X_i), f_j(X_j)) + \lambda \sum_{k=1}^v \mathcal{R}(f_k(X_k))$$

Northeastern University



Smile^{lab}
Synergetic Media Learning Lab

Thank you!

Q& A

