

《神经网络理论与应用》第七讲

Neural Network Theory and Applications

主讲教师：郑伟龙

助教：尹昊龙、史涵雯

上海交通大学计算机科学与工程系

weilong@sjtu.edu.cn

<http://bcmi.sjtu.edu.cn>

Outline of Lecture Seven

- Generative adversarial networks (GAN)
- Self-supervised Learning
- Transfer Learning

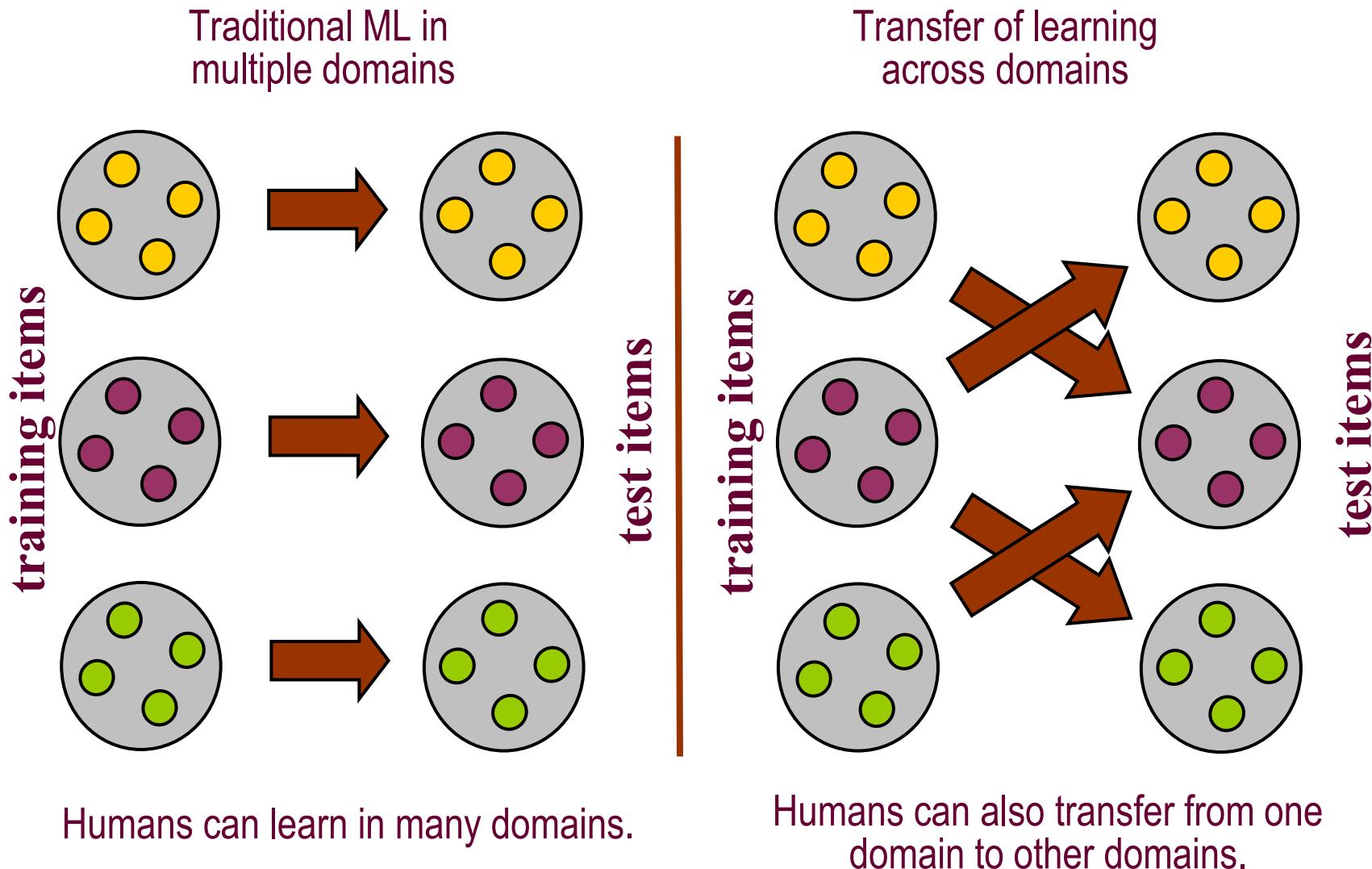
训练集和测试集的独立同分布假设

- 传统机器学习理论基于训练数据集与测试数据集具有独立分布的假设
- 上述假设过于理想，在解决实际问题时，通常难以成立。例如，脑电信号存在严重的个体差异或场景差异问题。
- 迁移学习方法可以克服上述问题

Transfer Learning



Traditional ML vs. TL



Notations

Domain:

It consists of two components: A feature space \mathcal{X} , a marginal distribution

$\mathcal{P}(X)$, where $X = \{x_1, x_2, \dots, x_n\} \in \mathcal{X}$

In general, if two domains are different, then they may have different feature spaces or different marginal distributions.

Task:

Given a specific domain and label space \mathcal{Y} . for each x_i in the domain, to predict its corresponding label y_i , where $y_i \in \mathcal{Y}$

In general, if two tasks are different, then they may have different label spaces or different conditional distributions

$\mathcal{P}(Y|X)$, where $Y = \{y_1, \dots, y_n\}$ and $y_i \in \mathcal{Y}$

Notations (2)

For simplicity, we only consider at most two domains and two tasks.

Source domain:

$\mathcal{P}(X_S)$, where $X_S = \{x_{S_1}, x_{S_2}, \dots, x_{S_{n_S}}\} \in \mathcal{X}_S$

Task in the source domain:

$\mathcal{P}(Y_S|X_S)$, where $Y_S = \{y_{S_1}, y_{S_2}, \dots, y_{S_{n_S}}\}$ and $y_{S_i} \in \mathcal{Y}_S$

Target domain:

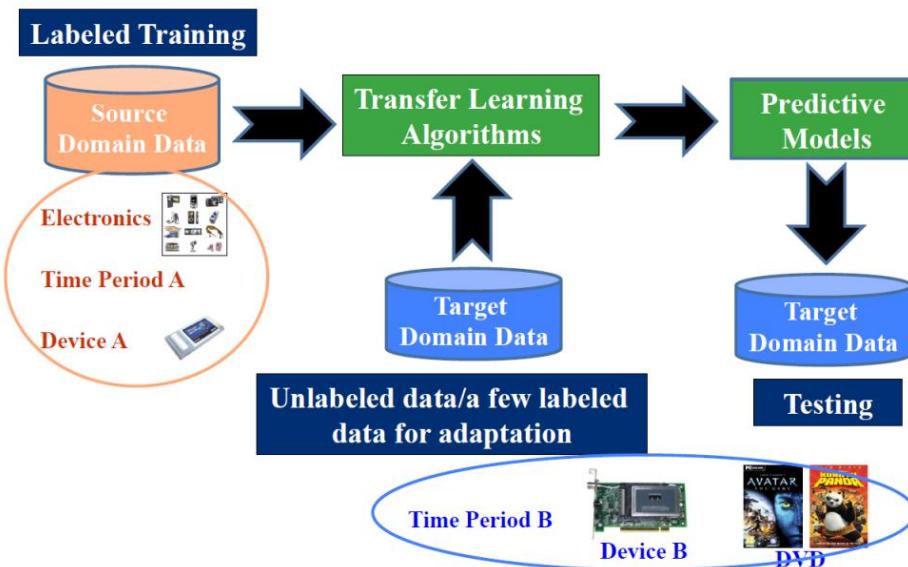
$\mathcal{P}(X_T)$, where $X_T = \{x_{T_1}, x_{T_2}, \dots, x_{T_{n_T}}\} \in \mathcal{X}_T$

Task in the target domain

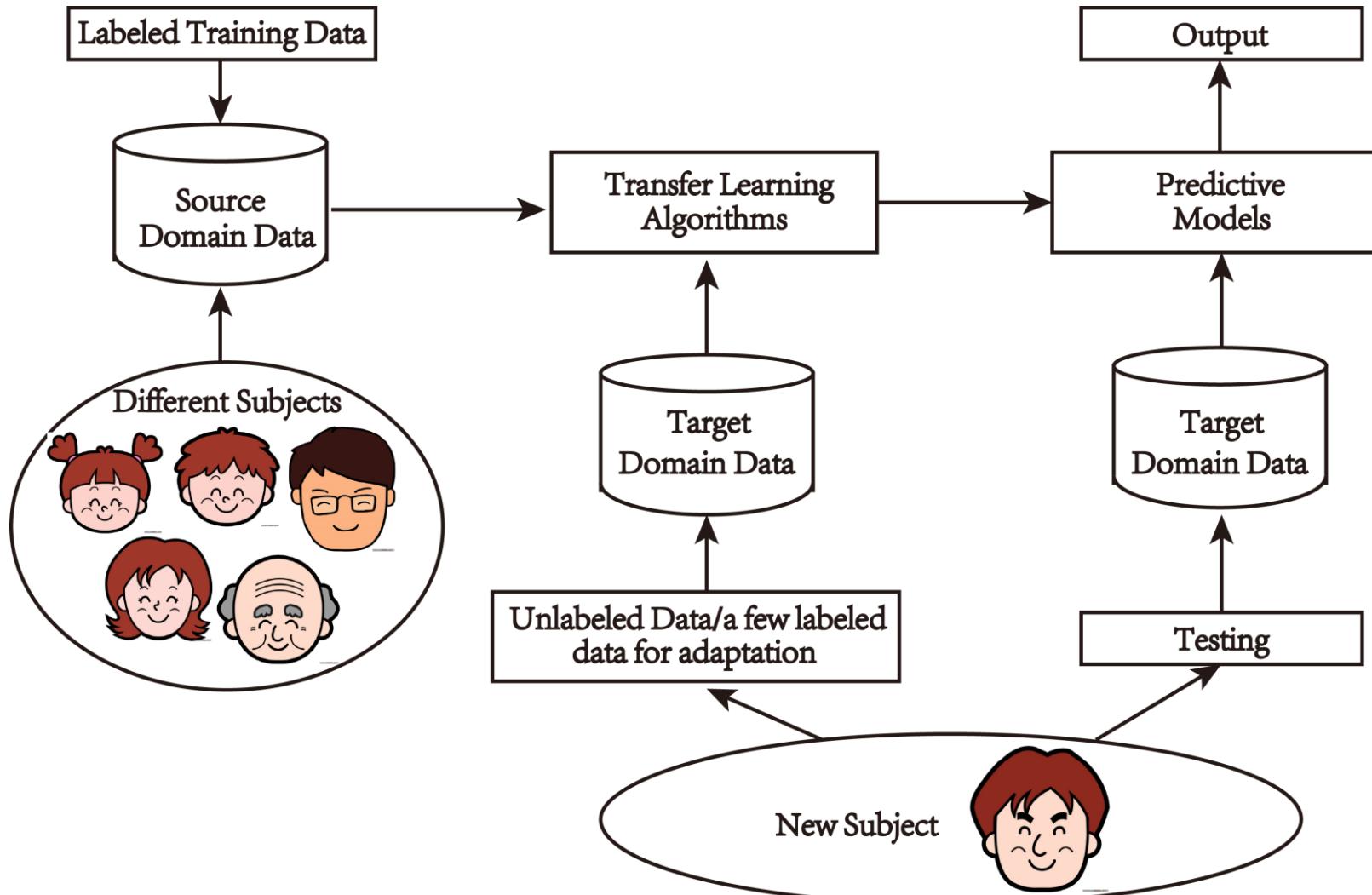
$\mathcal{P}(Y_T|X_T)$, where $Y_T = \{y_{T_1}, y_{T_2}, \dots, y_{T_{n_T}}\}$ and $y_{T_i} \in \mathcal{Y}_T$

Why Transfer Learning?

- In some domains, labeled data are in short supply.
- In some domains, the calibration effort is very expensive.
- In some domains, the learning process is time consuming.
- How to extract knowledge learnt from related domains to help learning in a target domain with a few labeled data?
- How to extract knowledge learnt from related domains to speed up learning in a target domain?



Transfer Learning Framework



Settings of Transfer Learning

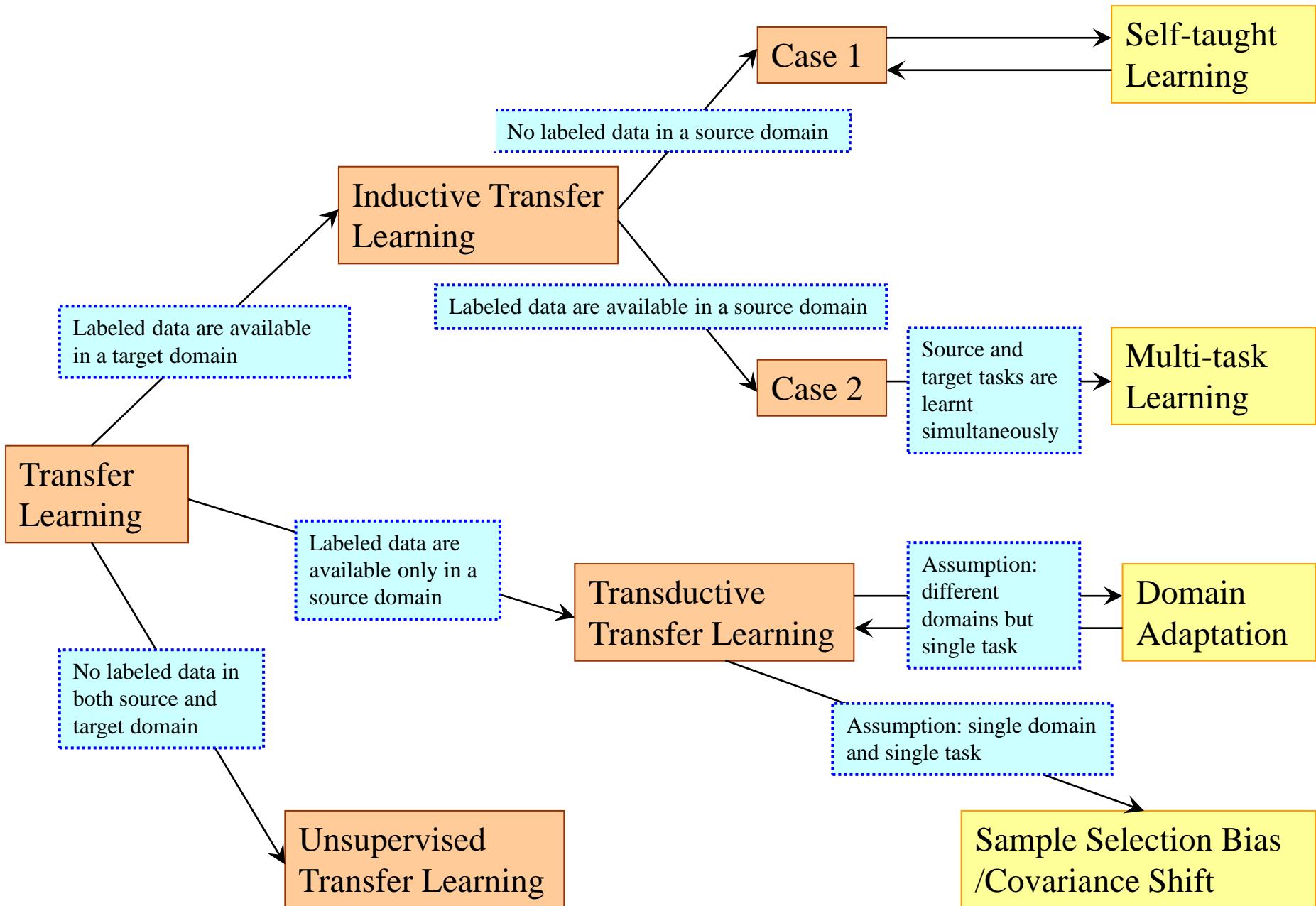
Relationship between Traditional Machine Learning and Various Transfer Learning Settings

Learning Settings		Source and Target Domains	Source and Target Tasks
Traditional Machine Learning		the same	the same
Transfer Learning	<i>Inductive Transfer Learning / Unsupervised Transfer Learning</i>	the same	different but related
	<i>Transductive Transfer Learning</i>	different but related	different but related
		different but related	the same

Different Settings of Transfer Learning

Transfer Learning Settings	Related Areas	Source Domain Labels	Target Domain Labels	Tasks
<i>Inductive Transfer Learning</i>	Multi-task Learning	Available	Available	Regression, Classification
	Self-taught Learning	Unavailable	Available	Regression, Classification
<i>Transductive Transfer Learning</i>	Domain Adaptation, Sample Selection Bias, Co-variate Shift	Available	Unavailable	Regression, Classification
<i>Unsupervised Transfer Learning</i>		Unavailable	Unavailable	Clustering, Dimensionality Reduction

Pan S J, Yang Q. A survey on transfer learning , IEEE Transactions on knowledge and data engineering, 2010, 22(10): 1345-1359.



Three Research Issues

- "What to transfer?"

- Which part of knowledge can be transferred across domain? e.g. feature representation, parameter settings, latent feature distribution, etc.

- "How to transfer?"

- Specific learning algorithms to transfer the knowledge.

- "When to transfer?"

- Asks in which situations, transferring skills should be done. Likewise, we are interested in knowing in which situations, knowledge should not be transferred.

Approaches to Transfer Learning

Transfer learning approaches	Description
<i>Instance-transfer</i>	<i>To re-weight some labeled data in a source domain for use in the target domain</i>
<i>Feature-representation-transfer</i>	Find a “good” feature representation that reduces difference between a source and a target domain or minimizes error of models
<i>Model-transfer</i>	Discover shared parameters or priors of models between a source domain and a target domain
<i>Relational-knowledge-transfer</i>	Build mapping of relational knowledge between a source domain and a target domain.

Inductive Transfer Learning

Instance-transfer Approaches

- **Assumption:** the source domain and target domain data use exactly the same features and labels.
- **Motivation:** Although the source domain data can not be reused directly, there are some parts of the data that can still be reused by re-weighting.
- **Main Idea:** Discriminatively adjust weights of data in the source domain for use in the target domain.

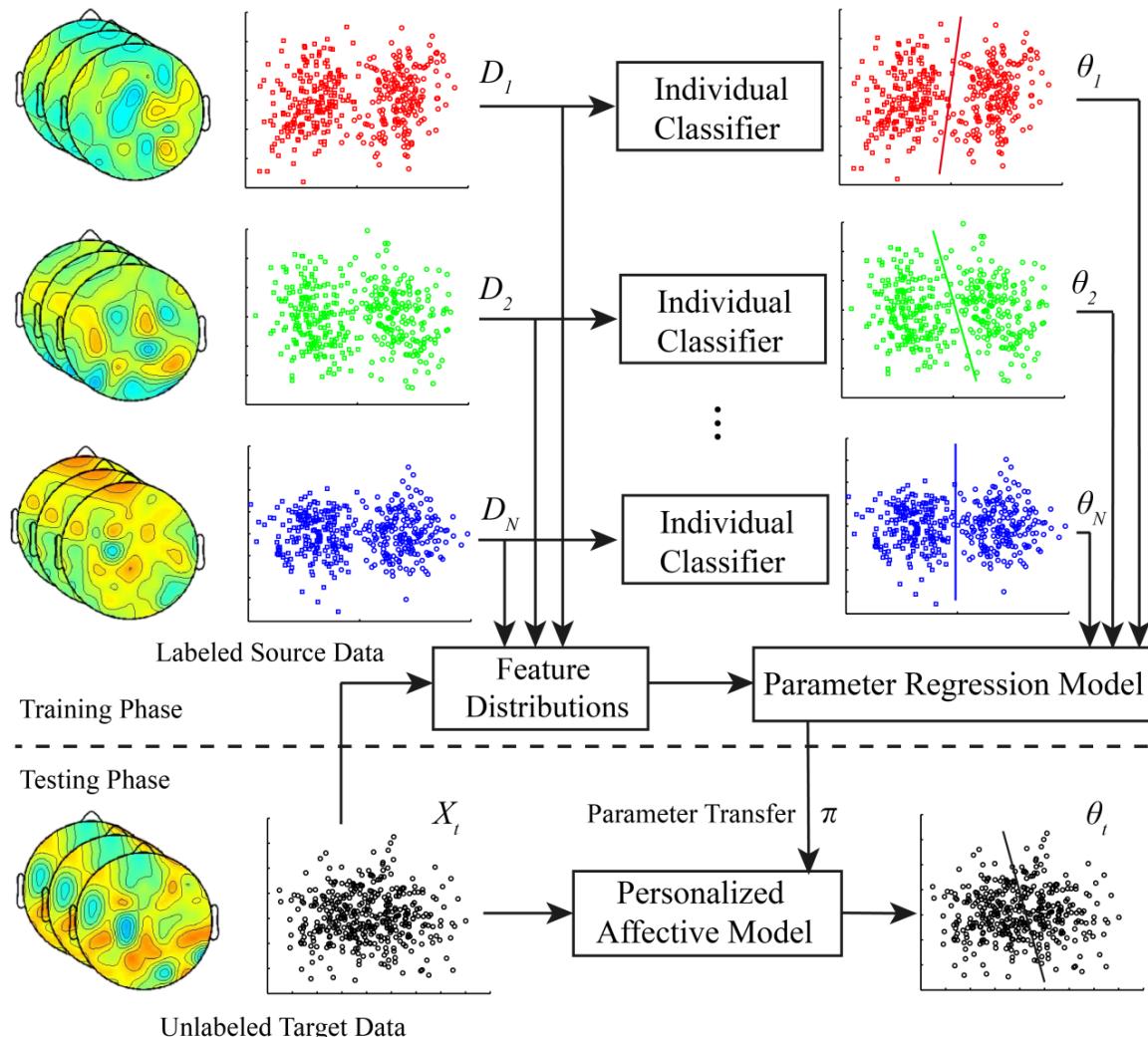
Feature Representation Transfer Approach

- **Assumption:** If t tasks are related to each other, then they may share some common features which can benefit for all tasks.
- **Input:** t tasks, each of them has its own training data.
- **Output:** Common features learnt across t tasks and t models for t tasks, respectively.

Negative Transfer

- Most approaches to transfer learning assume transferring knowledge across domains be always positive.
- However, in some cases, when two tasks are too dissimilar, brute-force transfer may even hurt the performance of the target task, which is called **negative transfer**.
- Some researchers have studied how to measure relatedness among tasks.
- How to design a mechanism to avoid negative transfer needs to be studied theoretically.

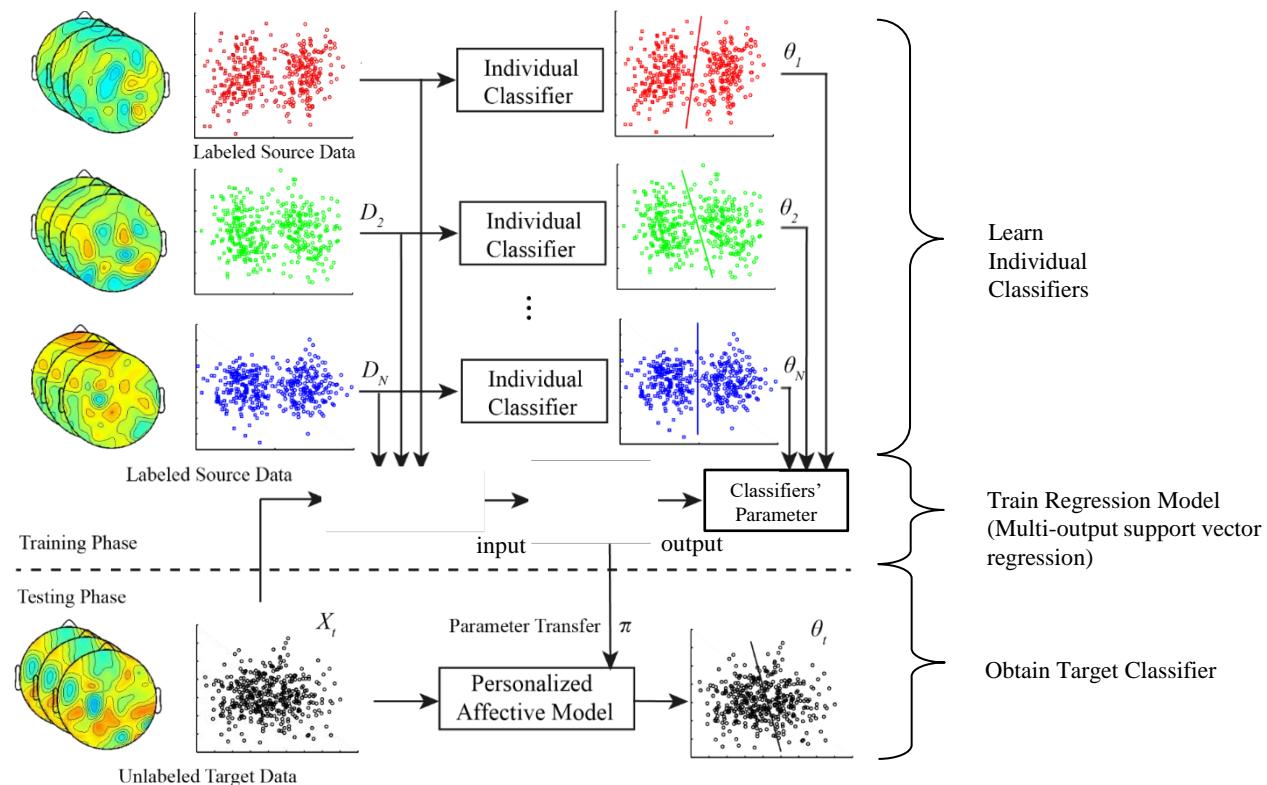
Transductive Parameter Transfer



直推参数迁移

Enver Sangineto, Gloria Zen, Elisa Ricci, and Nicu Sebe. We are not all equal: Personalizing models for facial expression analysis with transductive parameter transfer. In ACM International Conference on Multimedia, pages 357–366. ACM, 2014.

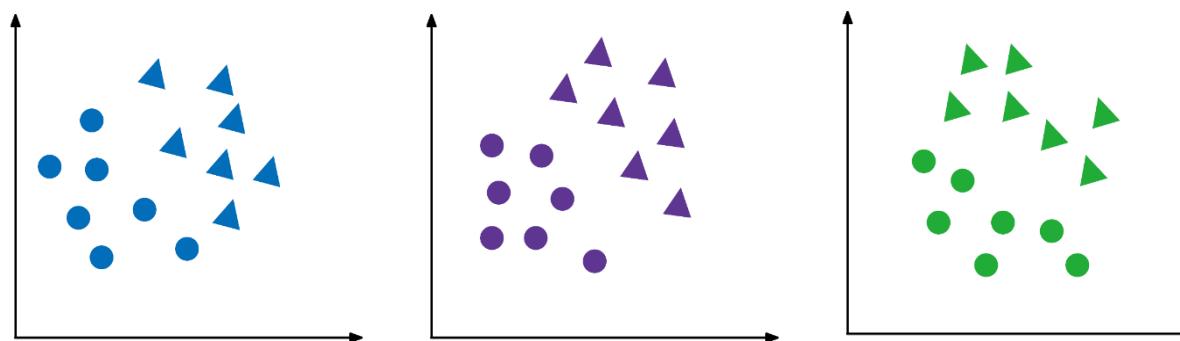
Transductive Parameter Transfer



Transductive Parameter Transfer

An introductory example

Suppose we have labeled data sets from three domains (source domains) and an unlabeled target domain data set.



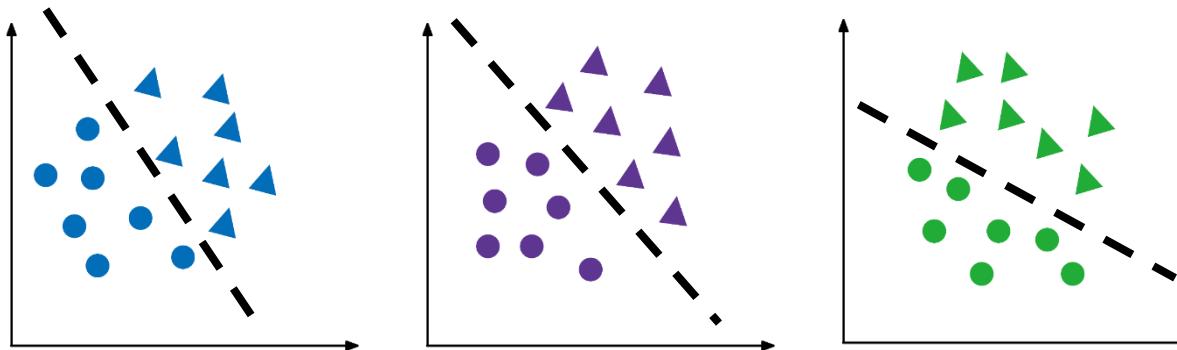
E. Sangineto, G. Zen, E. Ricci, and N. Sebe. We are not all equal: Personalizing models for facial expression analysis with transductive parameter transfer. In ACM Multimedia, 2014.

Transductive Parameter Transfer

An introductory example

Suppose we have labeled data sets from three domains (source domains) and an unlabeled target domain data set.

For the source domains, we can easily train SVMs on the labeled data.



E. Sangineto, G. Zen, E. Ricci, and N. Sebe. We are not all equal: Personalizing models for facial expression analysis with transductive parameter transfer. In ACM Multimedia, 2014.

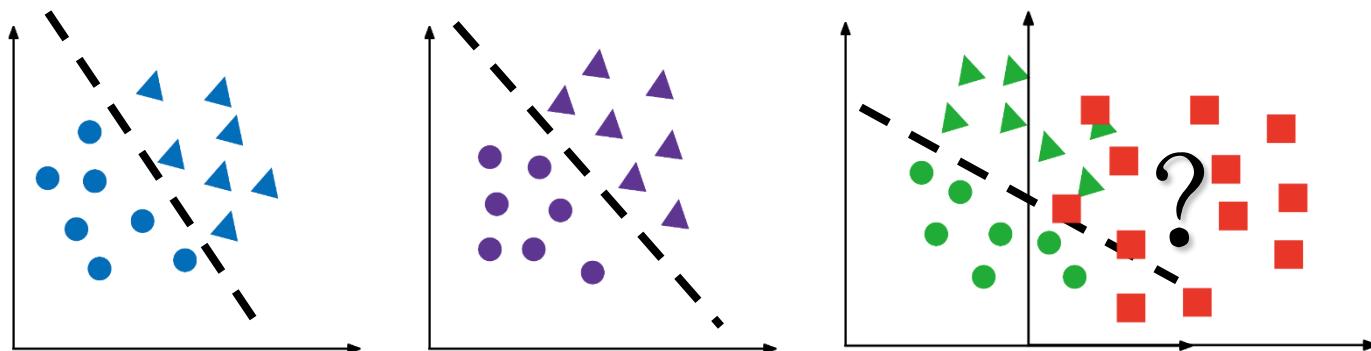
Transductive Parameter Transfer

An introductory example

Suppose we have labeled data sets from three domains (source domains) and an unlabeled target domain data set.

For the source domains, we can easily train SVMs on the labeled data.

Now the problem is how can we transfer the parameters in the trained SVMs to the target domain data set.



E. Sangineto, G. Zen, E. Ricci, and N. Sebe. We are not all equal: Personalizing models for facial expression analysis with transductive parameter transfer. In ACM Multimedia, 2014.

Transductive Parameter Transfer

Demonstration of the objective

We denote the source domain data sets as

$$D_1^s, \dots, D_N^s, D_i^s = \{\mathbf{x}_j^s, y_j^s\}_{j=1}^{n_i^s}$$

Where D_i^s indicates the i th data set. Also let us use X_i^s to indicate the data with label. Then we assume that X_i^s is drawn from the marginal distribution $P_{X_i^s} : X_i^s \sim P_{X_i^s}$.

For each data set, we can easily train an SVM parameterized with θ_i^s .

The question is: for target domain data $X^t \sim P_{X^t}$, how can we obtain the corresponding SVM parameter θ^t .

Transductive Parameter Transfer

Analysis

By observing the problem setting, we assume there is a function from marginal distribution to the SVM parameter vector:

$$f : \mathcal{P} \rightarrow \Theta$$

where \mathcal{P} and Θ indicate the space of the marginal distribution and the SVM parameter vectors respectively.

So we have

$$\theta_i^s = f(P_{X_i^s})$$

If the function can be obtained, the target domain SVM parameter vector is simply

$$\theta^t = f(P_{X^s}).$$

Transductive Parameter Transfer

Analysis (Cont.)

However, the marginal distributions $P_{X_i^s}$ can not be easily depicted.

Because we know $X_i^s \sim P_{X_i^s}$, we can use X_i^s instead of $P_{X_i^s}$.

We then assume another function:

$$\hat{f} : 2^{\mathcal{X}} \rightarrow \Theta.$$

\hat{f} is a function from the feature set space $2^{\mathcal{X}}$ to the parameter space.

Similarly, we have $\theta_i^s = \hat{f}(X_i^s)$ and $\theta^t = \hat{f}(X^t)$.

Now the problem becomes how to learn the function \hat{f} .

Transductive Parameter Transfer

Learning problem of \hat{f} with M-SVR

The training set is $\mathcal{T} = \{X_i^s, \theta_i^s\}_{i=1}^N$.

It is a simple multivariate regression problem. Various methods can be used to solve it. In the original paper, M-SVR (Tuia et al.) was applied.

D. Tuia, J. Verrelst, L. Alonso, F. Perez-Cruz, and G. Camps-Valls. Multioutput support vector regression for remote sensing biophysical parameter estimation. *IEEE Geoscience and Remote Sensing Letters*, 8(4):804–808, July 2011.

Tranductive Parameter Transfer

Learning problem of \hat{f} with M-SVR Cont.

In the M-SVR framework, \hat{f} can be defined by

$$\hat{f}(X) = \phi(X)'B + c'$$

where B is a matrix and c is a vector. $\phi(X)$ is a nonlinear mapping from the matrix space to a vector space.

The corresponding loss is defined as

$$L(B, c; \mathcal{T}) = \frac{1}{2} \|B\|_F^2 + \lambda_E \sum_{i=1}^N E(\|\theta_i^{s'} - \hat{f}(X_i^s)\|)$$

where $\|\cdot\|_F$ indicates the Frobenius norm: $\|A\|_F = \sqrt{\text{trace}(A'A)}$, the λ_E is a tradeoff hyperparameter, and $E(\cdot)$ is a loss function parameterized by ϵ :

$$E(u) = \begin{cases} 0 & u < \epsilon \\ u^2 - 2u\epsilon + \epsilon^2 & u \geq \epsilon \end{cases}$$

Tranductive Parameter Transfer

Learning problem of \hat{f} with M-SVR Cont.

By applying the kernel trick, the prediction function can be rewritten as

$$\hat{f}(X) = [\kappa(X_1^s, X), \dots, \kappa(X_N^s, X)] \cdot V + \mathbf{c}'$$

where V is obtained by solving the dual problem.

To optimize V and \mathbf{c} , we adopt iterative procedure. In iteration k , suppose we have the temporal values of B and \mathbf{c} being B^k and \mathbf{c}^k . Let us denote the error $e = \hat{f}(X)' - \theta$, then the loss function can be expand with Taylor expansion:

$$L^*(B, \mathbf{c}) = \frac{1}{2} \|B\|_F^2 + \lambda_E \left(\sum_{i=1}^N E(e_i^k) + \frac{dE(u)}{du} \Big|_{u_i^k} \frac{(e_i^k)'}{u_i^k} [e_i - e_i^k] \right)$$

where $u = \|e\|$. The superscripts ‘ k ’s indicate the values at the k th iteration.

Transductive Parameter Transfer

Learning problem of \hat{f} with M-SVR Cont.

Then we can further approximate the loss function by using a quadratic approximation:

$$L^{**}(B, \mathbf{c}) = \frac{1}{2} \|B\|_F^2 + \lambda_E \frac{1}{2} \sum_{i=1}^N a_i u_i^2 + \tau$$

where

$$a_i = \left. \frac{\lambda_E}{u_i^k} \frac{dE(u)}{du} \right|_{u_i^k}$$

and τ is a value that does not depend on B and \mathbf{c} .

It can be easily observed that the $L^{**}(\cdot)$ is simply a quadratic function of B and \mathbf{c} , so the solution B^{sol} and \mathbf{c}^{sol} can be obtained.

Transductive Parameter Transfer

Learning problem of \hat{f} with M-SVR Cont.

The update rule of B and c is

$$\begin{bmatrix} B^{k+1} \\ (\mathbf{c}^{k+1})' \end{bmatrix} = \begin{bmatrix} B^k \\ (\mathbf{c}^k)' \end{bmatrix} + \eta_k \begin{bmatrix} B^{sol} - B^k \\ (\mathbf{c}^{sol})' - (\mathbf{c}^k)' \end{bmatrix}$$

where the step size η_k is decided by back tracking algorithm. Specifically, we initially set it to 1, and check whether the loss function drops. If not, the step size is multiplied by a value less than 1 and retry.

Transductive Parameter Transfer

Learning problem of \hat{f} with M-SVR Cont.

Solving B is not convenient because map function $\phi(X)$ is usually not analytical, so we substitute V into the previous equations by utilizing $B_{\cdot,j} = KV_{\cdot,j}$ where K is the kernel matrix defined by

$$K_{ij} = \kappa(X_i, X_j)$$

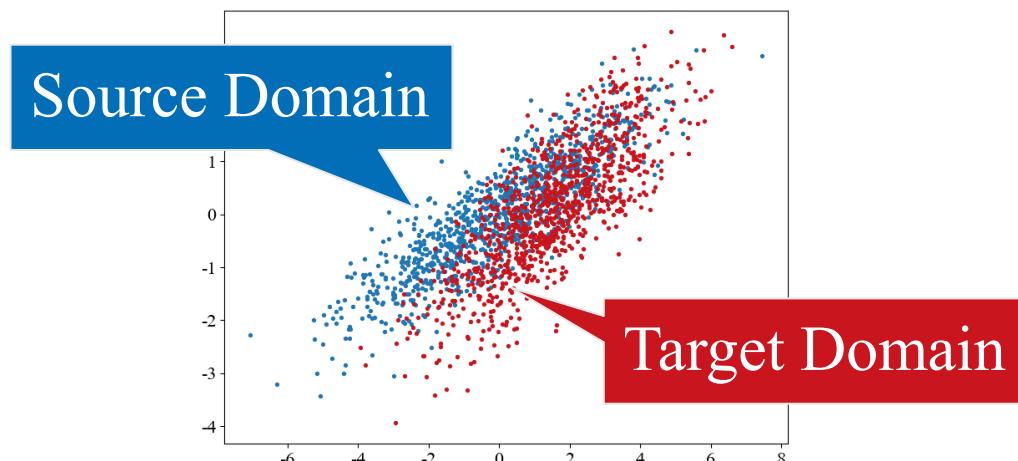
Tansductive Parameter Transfer

Some points

1. Other types of parameterized models (e.g., LDA) should also work well with TPT.
2. When learning, the features are matrices. So some kind of mapping should be applied to project the features to vector spaces. For M-SVR, the mapping corresponds to specially designed kernels.
3. The regression method used in TPT is also free to choose. But the authors found M-SVR performed better on their tasks.

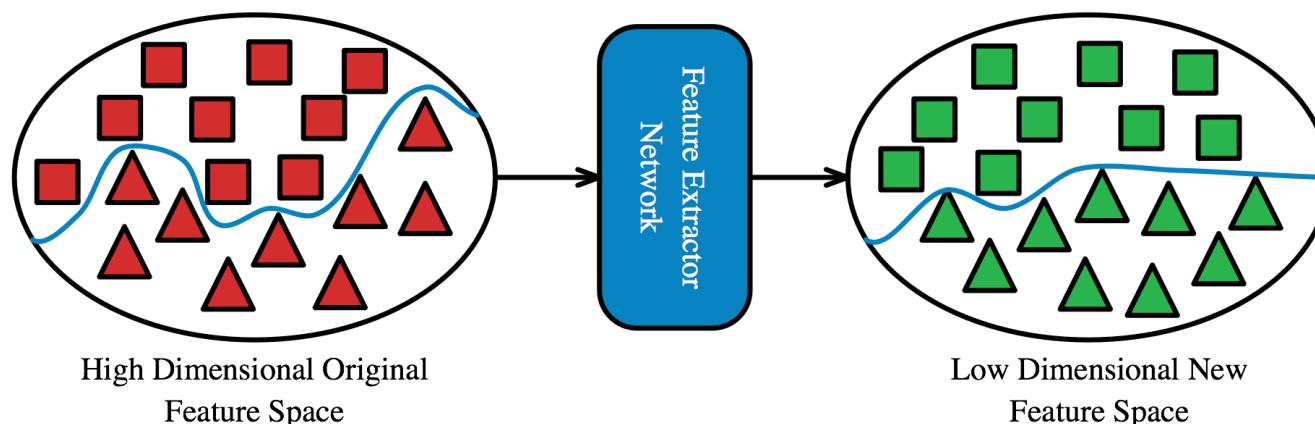
Domain Shift and Domain Adaptation

- Training data are drawn from source domain $\{X_s, \varphi_s\}$, and test data are drawn from target domain $\{X_t, \varphi_t\}$.
- Both domains have the same set of features ($X_s, X_t \in R^m$).
- Domain shift ($P(X_s) \neq P(X_t)$) makes ordinary learning methods degenerate.
- We can use domain adaptation methods to eliminate or reduce the domain shift.



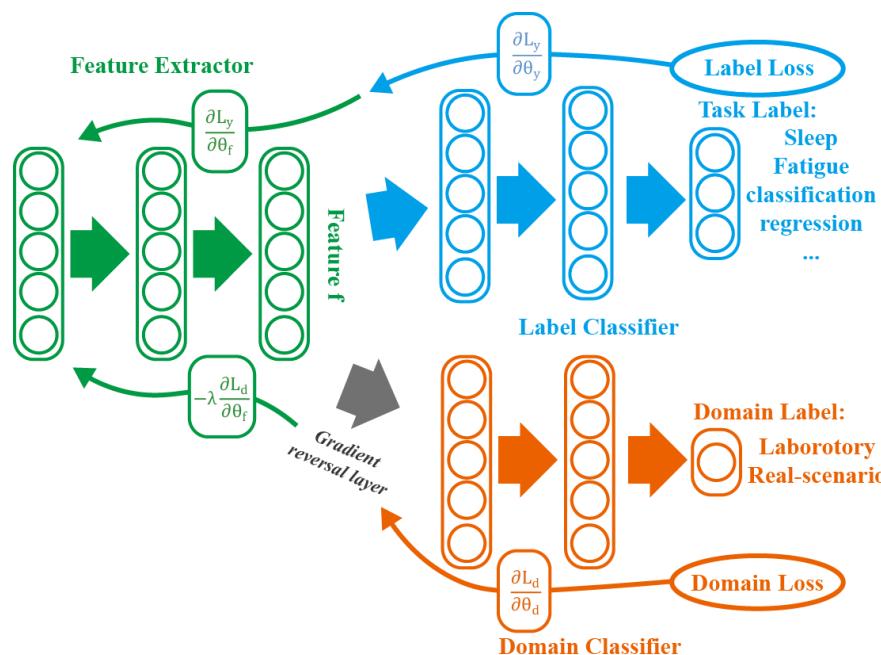
Domain Adaptation with Adversarial Networks

- To reduce the domain shift, neural networks can be used to transform the original feature into a new feature space where two conditions must be satisfied:
 - The new features keep essential information for label prediction.
 - Domain shift is reduced (features from both domains share similar distributions in the new feature space).



Domain Adversarial Neural Network

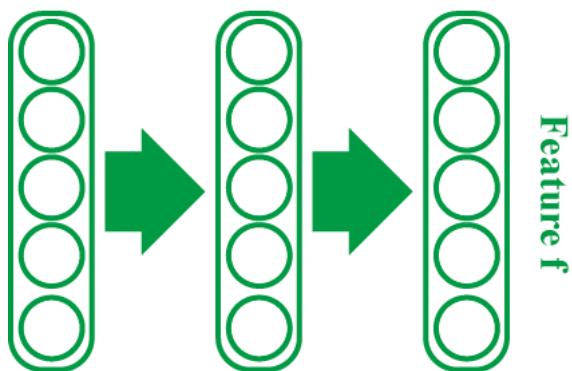
DANN is a domain adaptation approach with deep architectures that can be trained with labeled source domain data and unlabeled target domain data. Its adaptation behavior is achieved by augmenting a normal feed-forward model with an adversarial manner.



Paper: Y. Ganin and V. Lempitsky, Unsupervised domain adaptation by backpropagation, ICML2015, pp. 1180–1189.

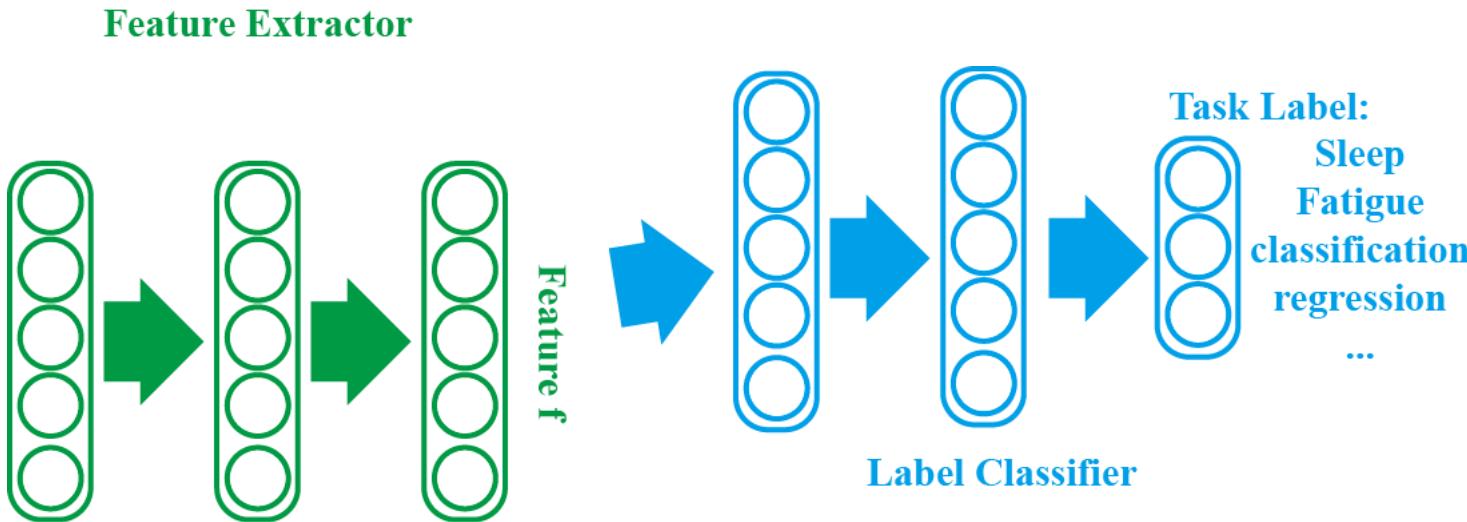
Domain Adversarial Neural Network

Feature Extractor



The first layers works as a feature extractor, it maps the low level features to a non-linear transformed high level feature f

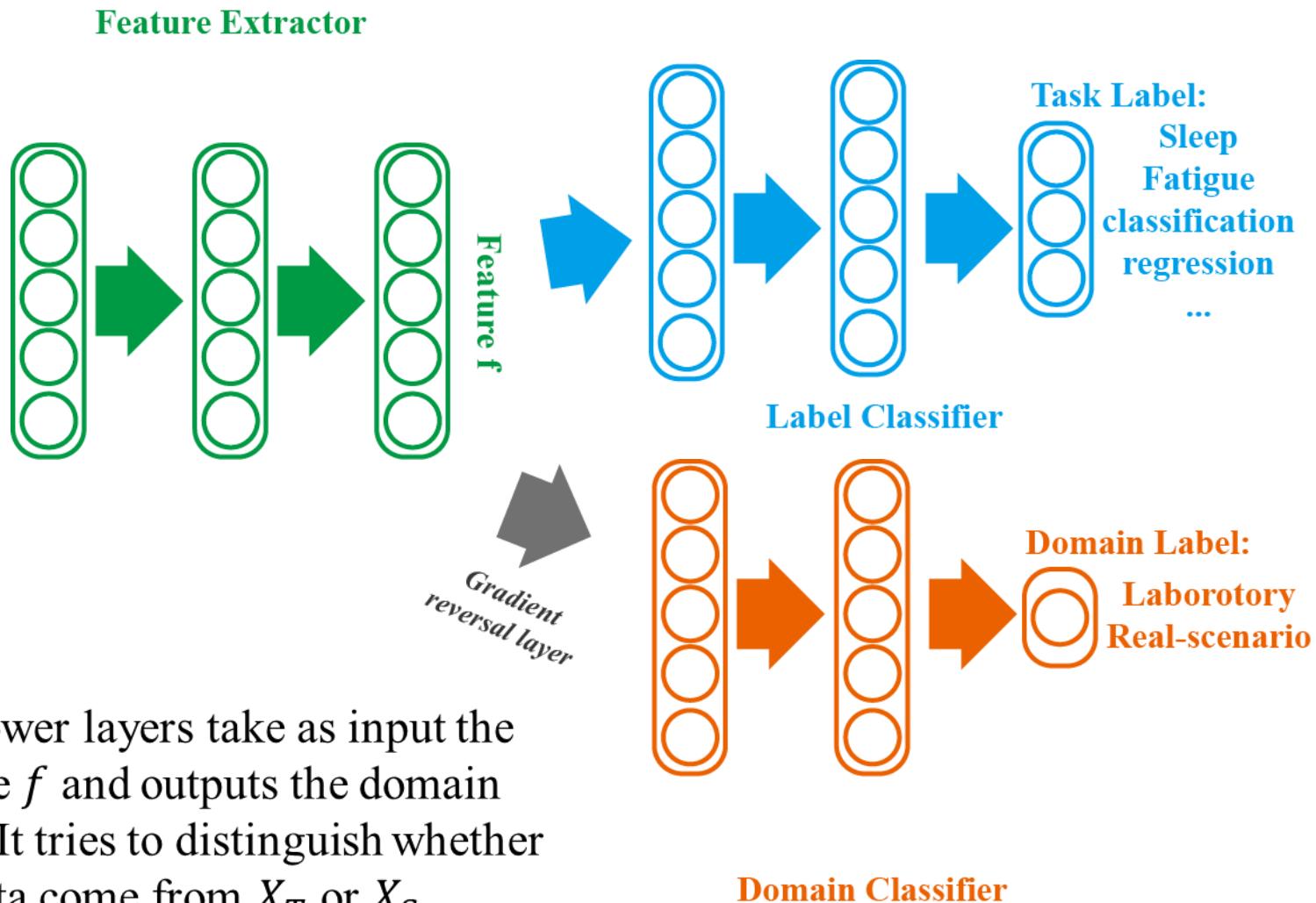
Domain Adversarial Neural Network



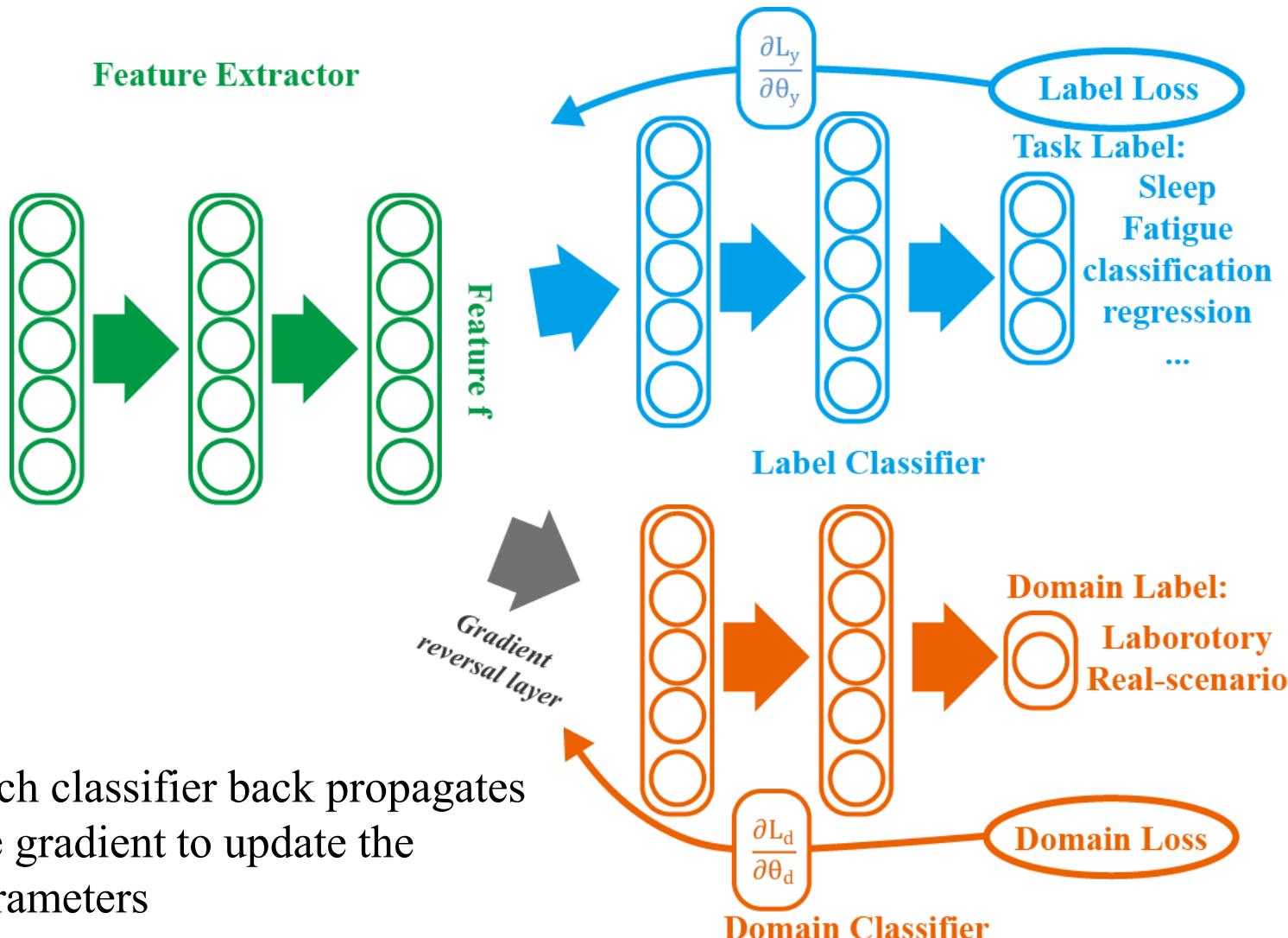
The upper layers take as input the feature f and outputs the label prediction for the required task

These two parts work like a normal multilayer network.

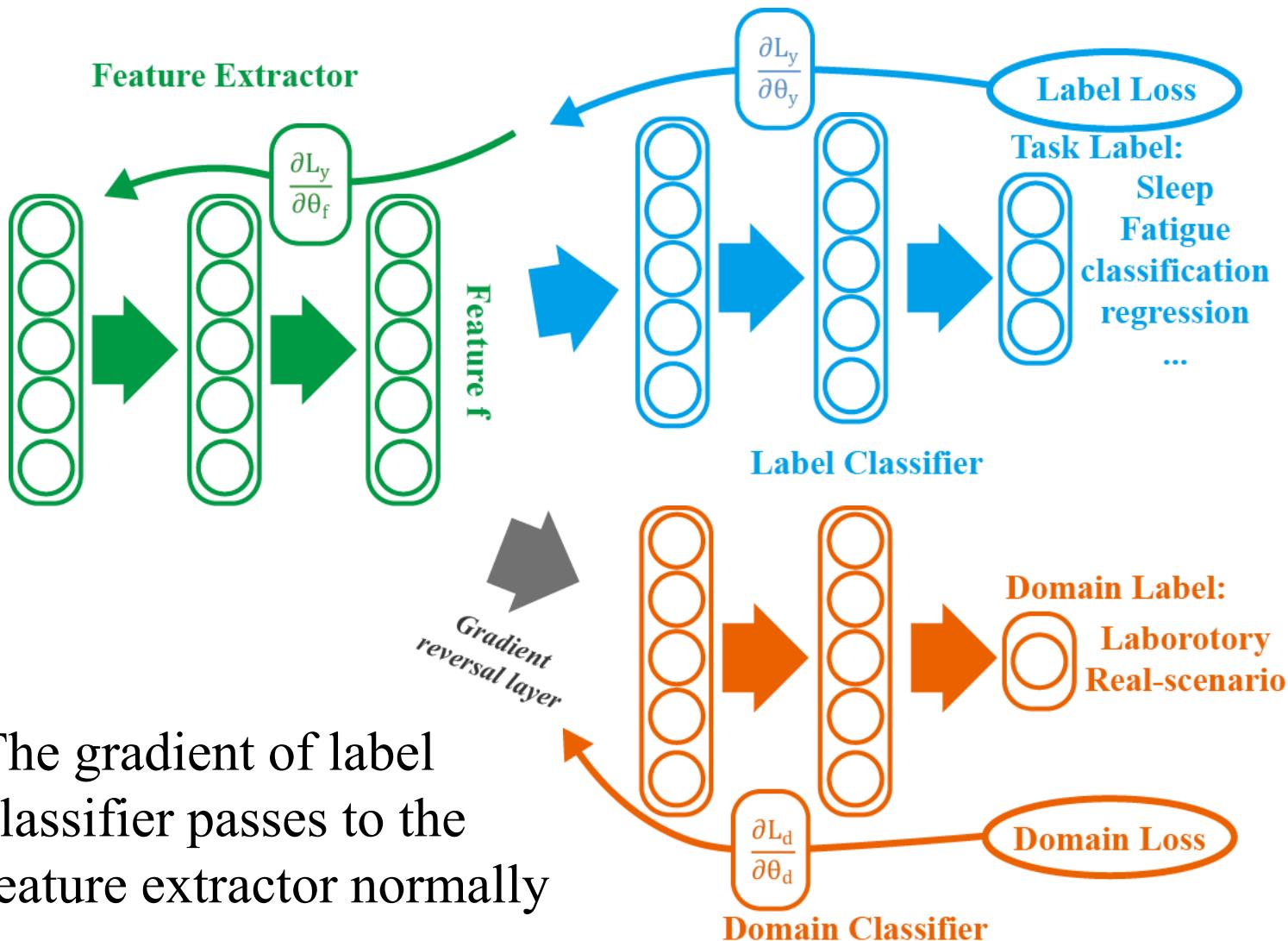
Domain Adversarial Neural Network



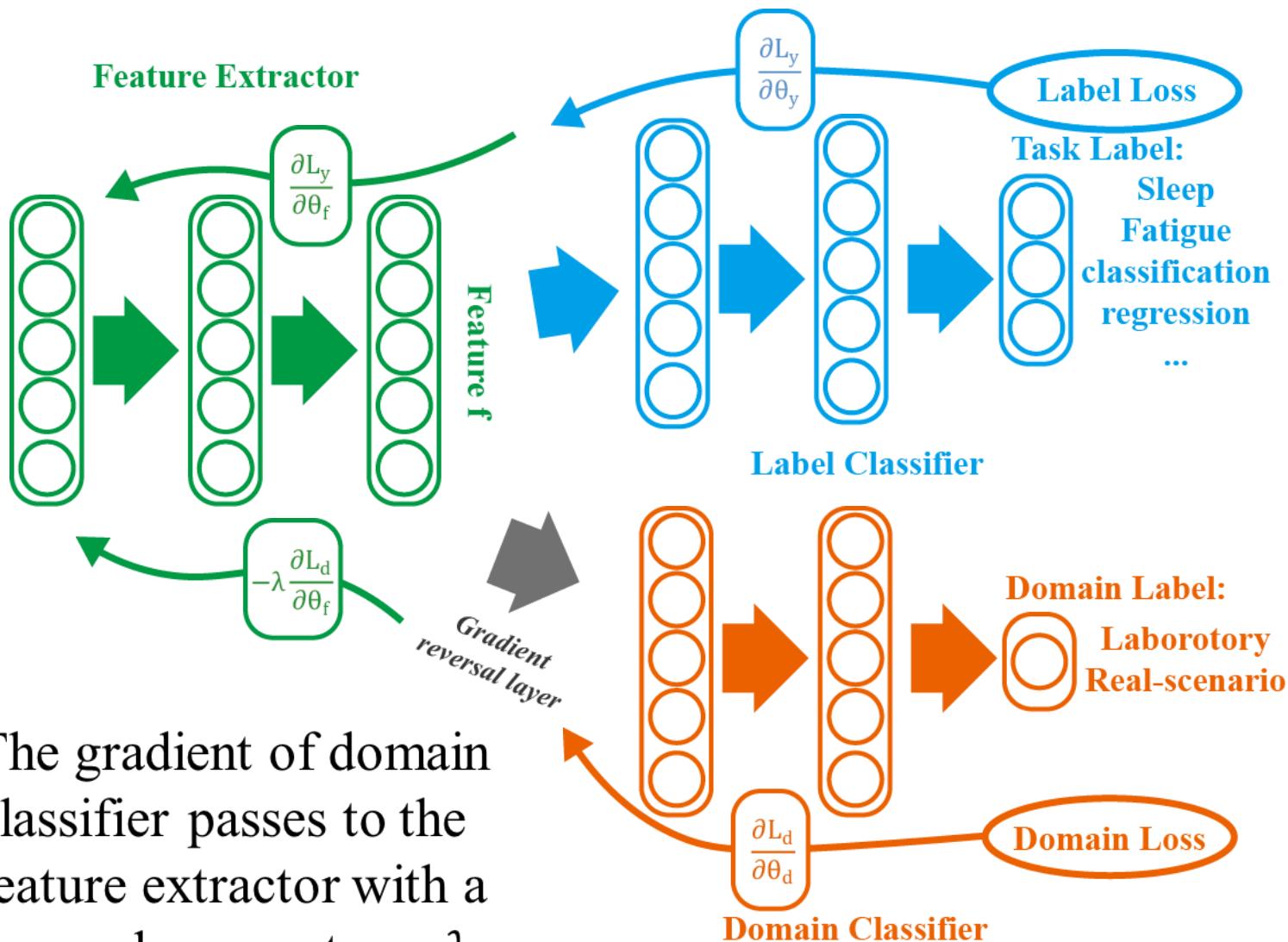
Domain Adversarial Neural Network



Domain Adversarial Neural Network

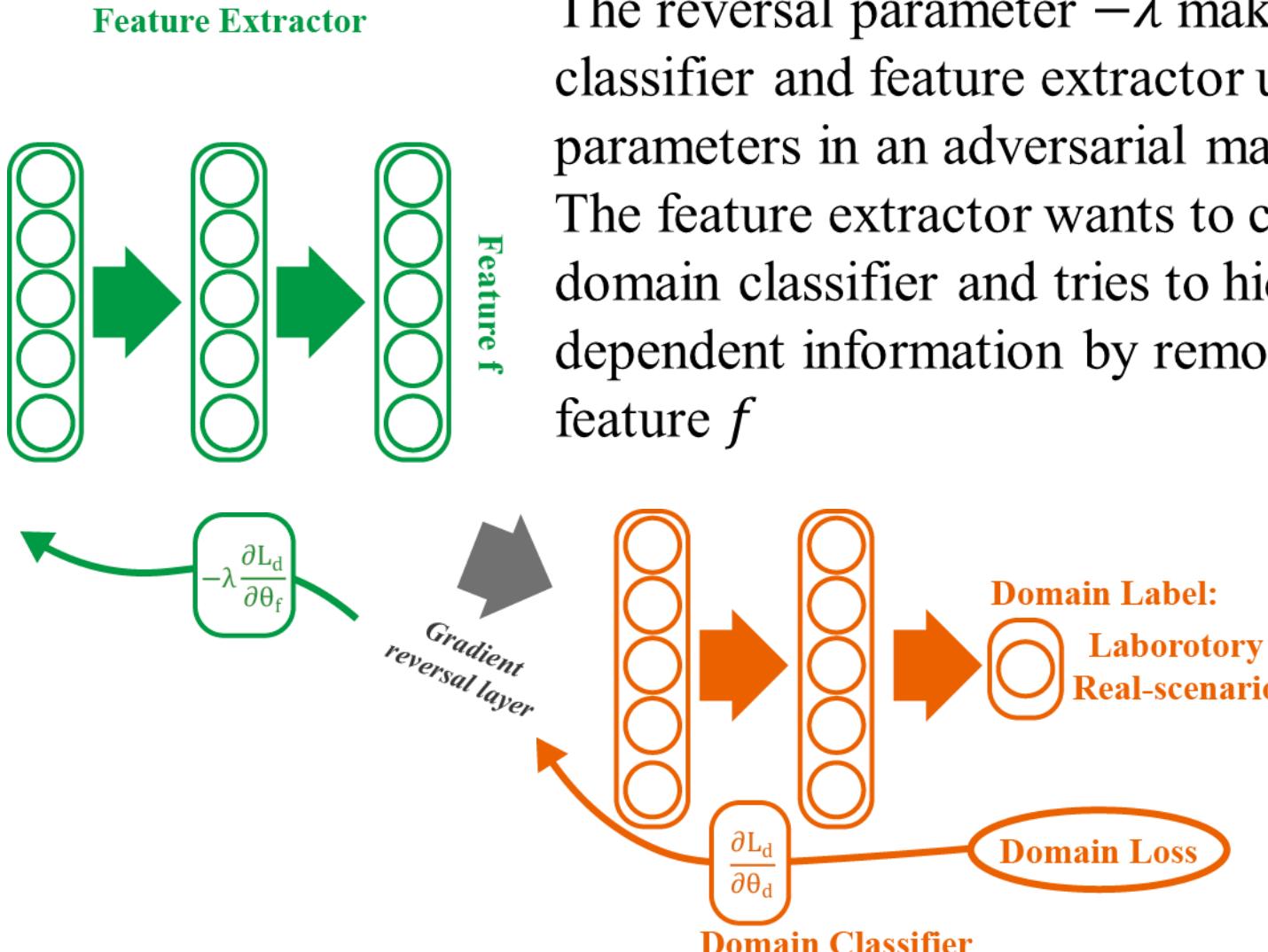


Domain Adversarial Neural Network



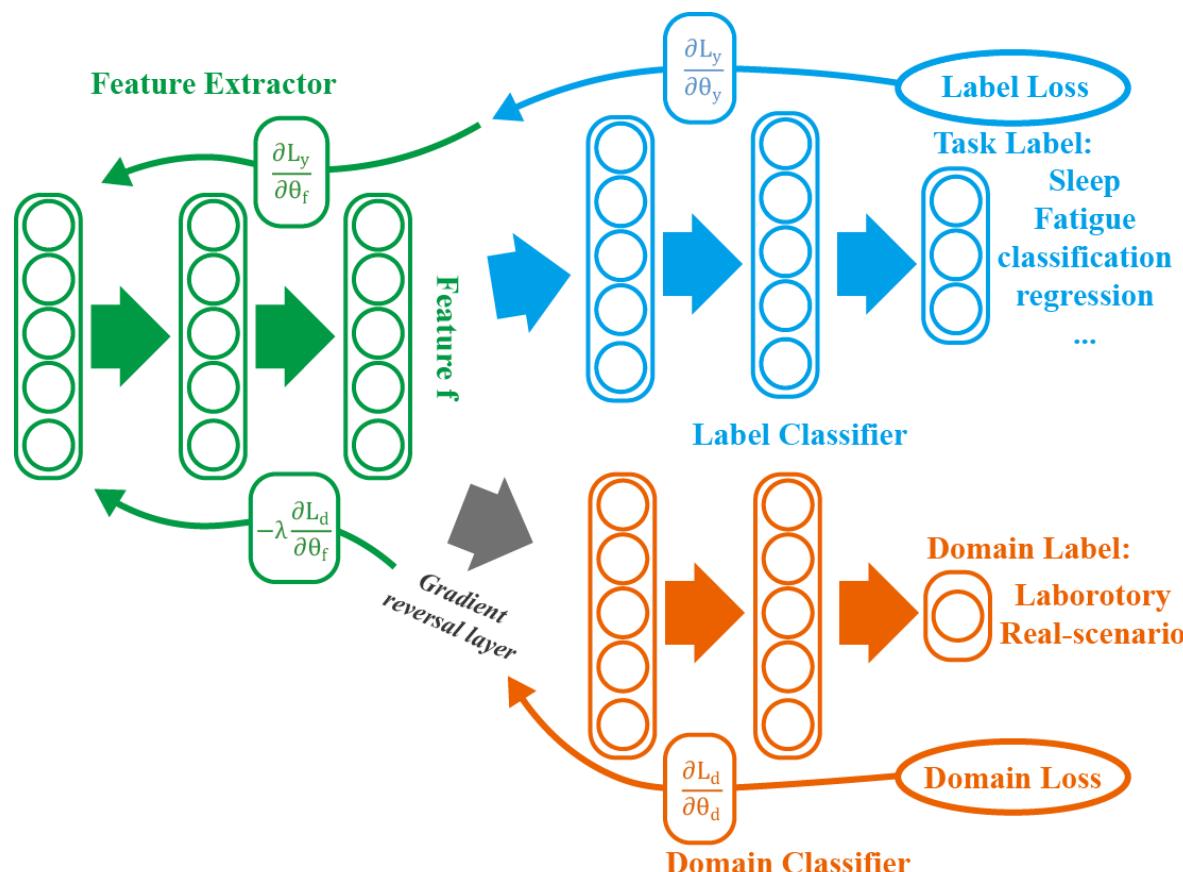
The gradient of domain classifier passes to the feature extractor with a reversal parameter $-\lambda$

Domain Adversarial Neural Network



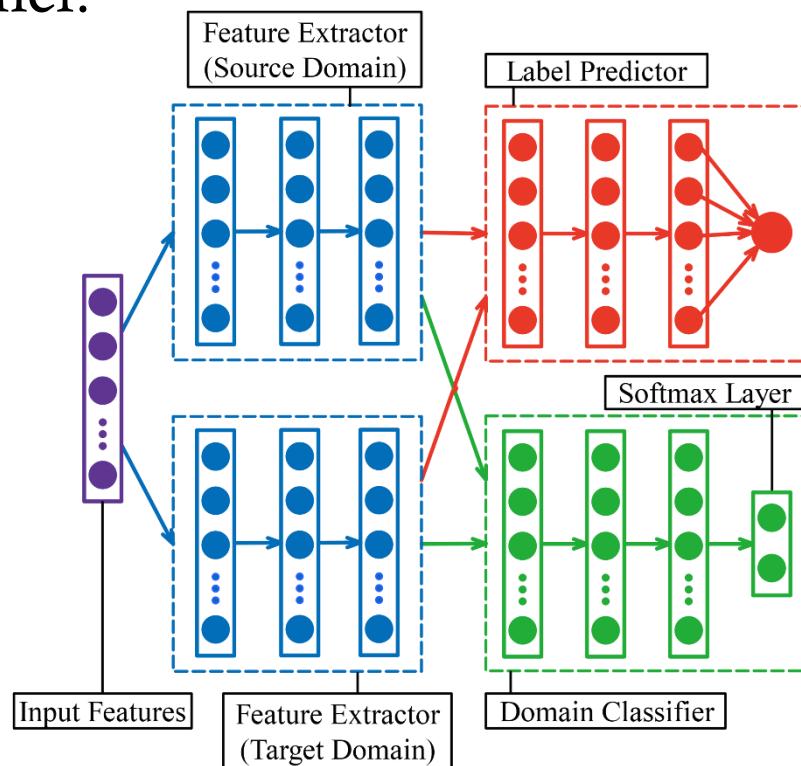
Domain Adversarial Neural Network

When the model converges, the feature f would be containing no domain related information. Thus it works as function of domain adaptation.



Adversarial Discriminative Domain Adaptation

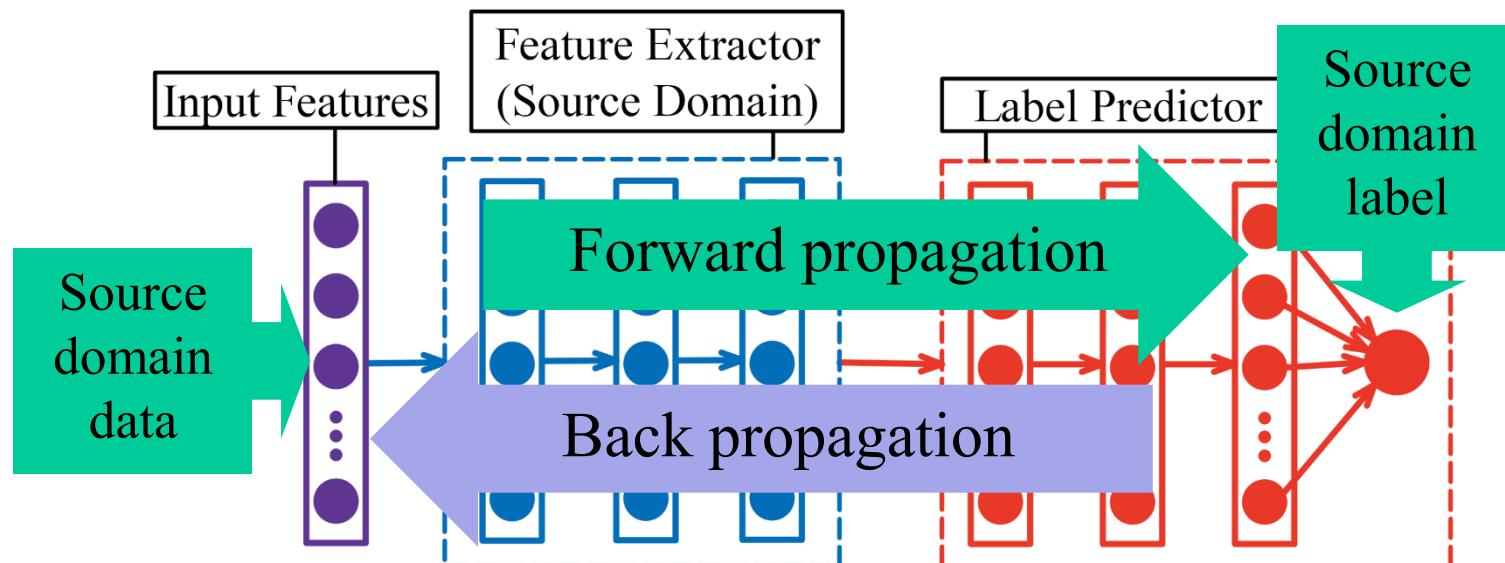
- ADDA is based on deep adversarial network. It aligns the distributions of source domain data and target domain data by training a feature extractor and a domain classifier in an adversarial manner.



E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, “Adversarial discriminative domain adaptation,” in IEEE Conference on Computer Vision and Pattern Recognition, July 2017, pp. 2962–2971.

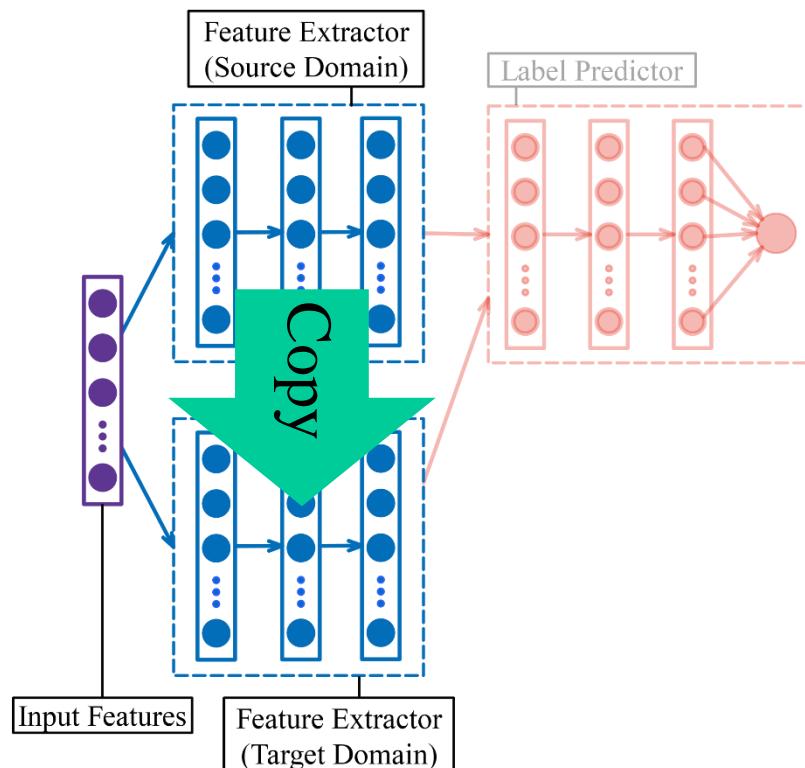
Training of ADDA (step 1)

- First, we train a simple multilayer perceptron on the source domain data with source domain data as input. The network is divided into 3 parts: input layer, feature extractor and label predictor. We can follow the general training procedure for this step (e.g., BP algorithm).



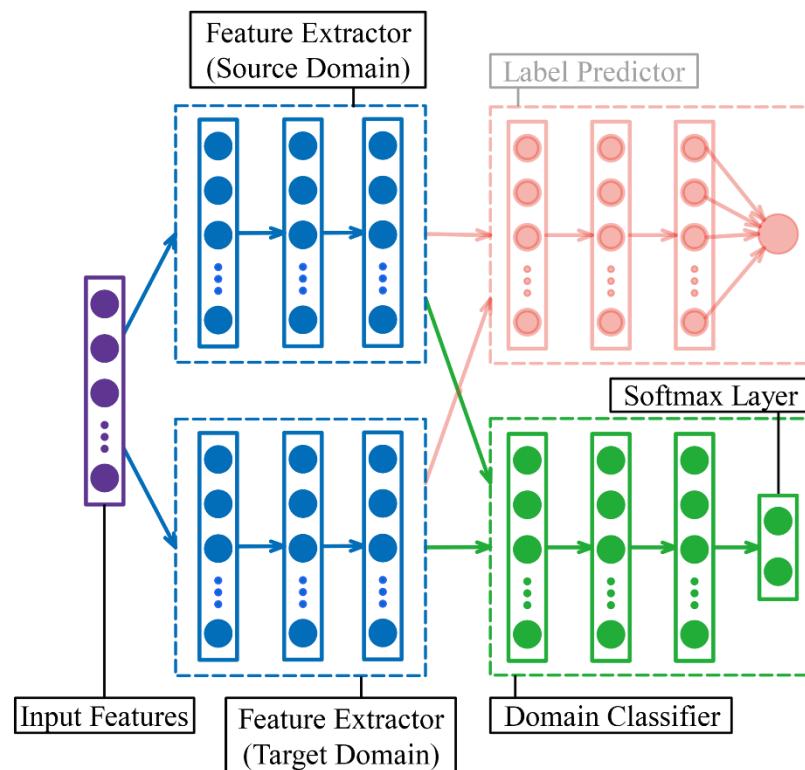
Training of ADDA (step 2)

- After the network converge, another feature extractor is constructed. The new feature extractor only accepts input from the target domain. Its parameters are initialized by copying the ones from the original feature extractor.



Training of ADDA (step 3)

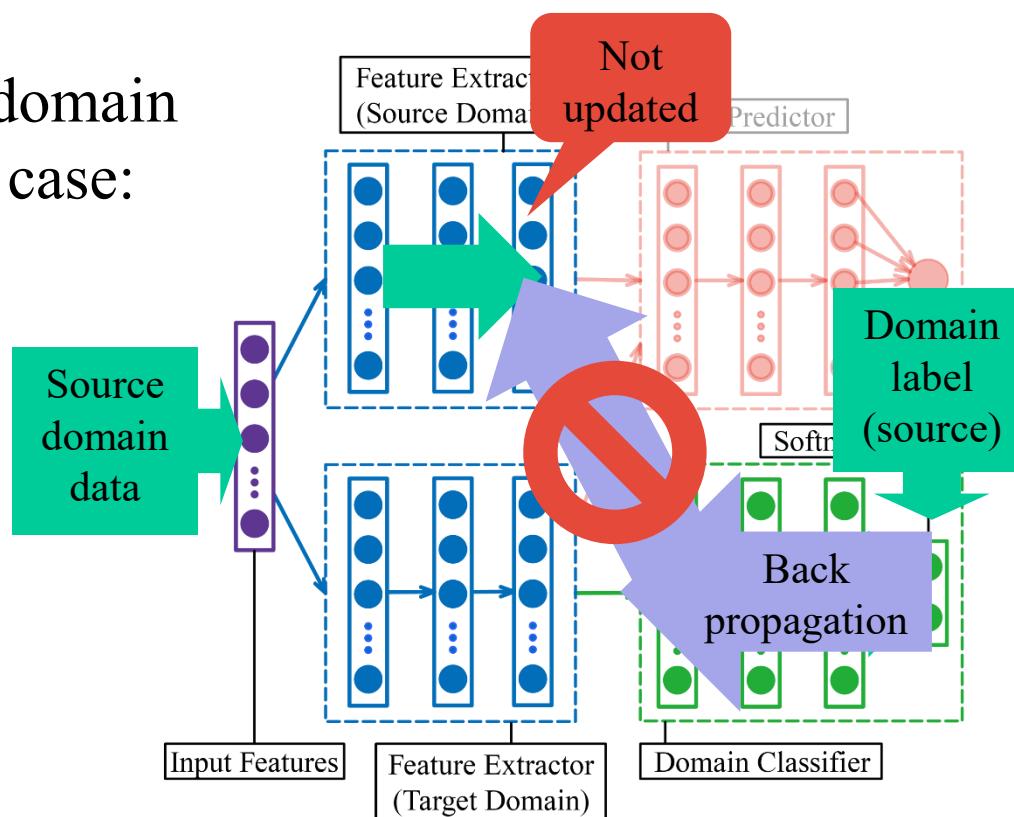
- In this step, a domain classifier sub-network is trained when the label predictor and the source domain feature extractor are not updated. The domain classifier is trained to discriminate which domain (source or target?) the input is from.



Training of ADDA (step 3 Cont.)

- However, the target domain feature extractor is trained to deceive the domain classifier. The domain classifier and the target domain feature extractor compete with each other (i.e., adversarial training).

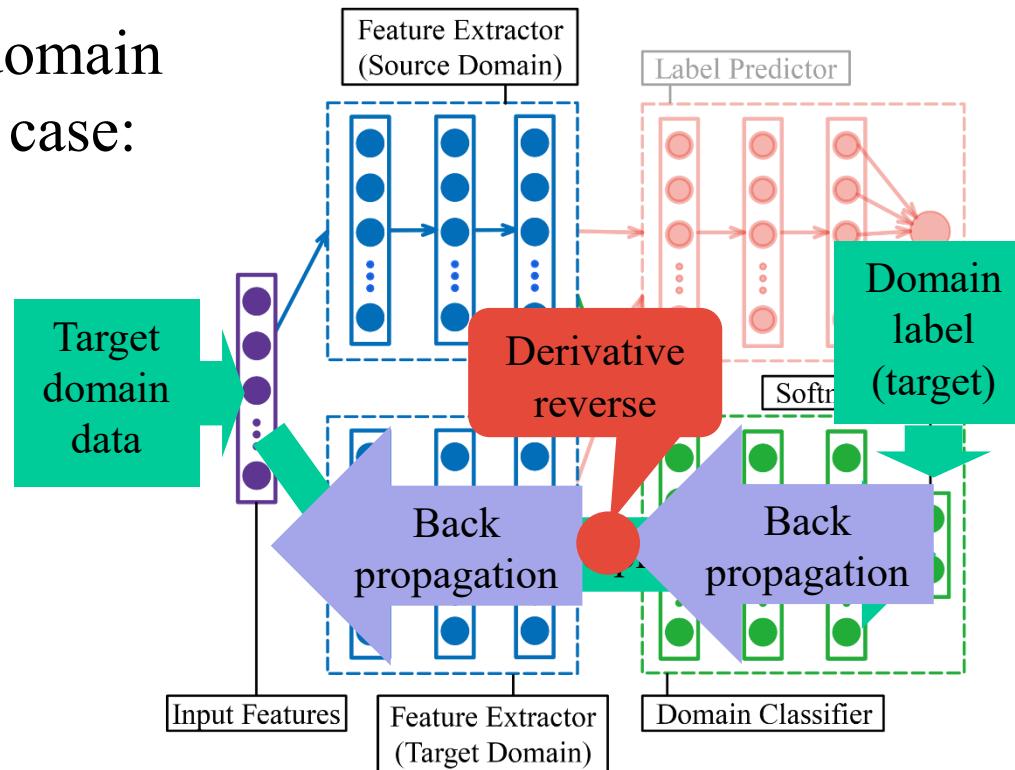
Source domain
training case:



Training of ADDA (step 3 Cont.)

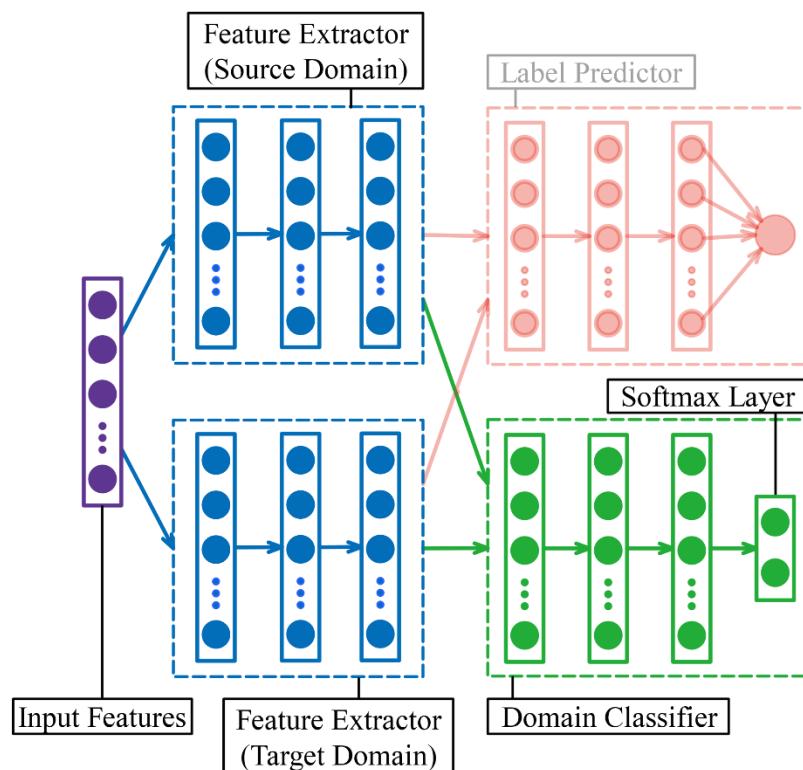
- However, the target domain feature extractor is trained to deceive the domain classifier. The domain classifier and the target domain feature extractor compete with each other (i.e., adversarial training).

Target domain
training case:



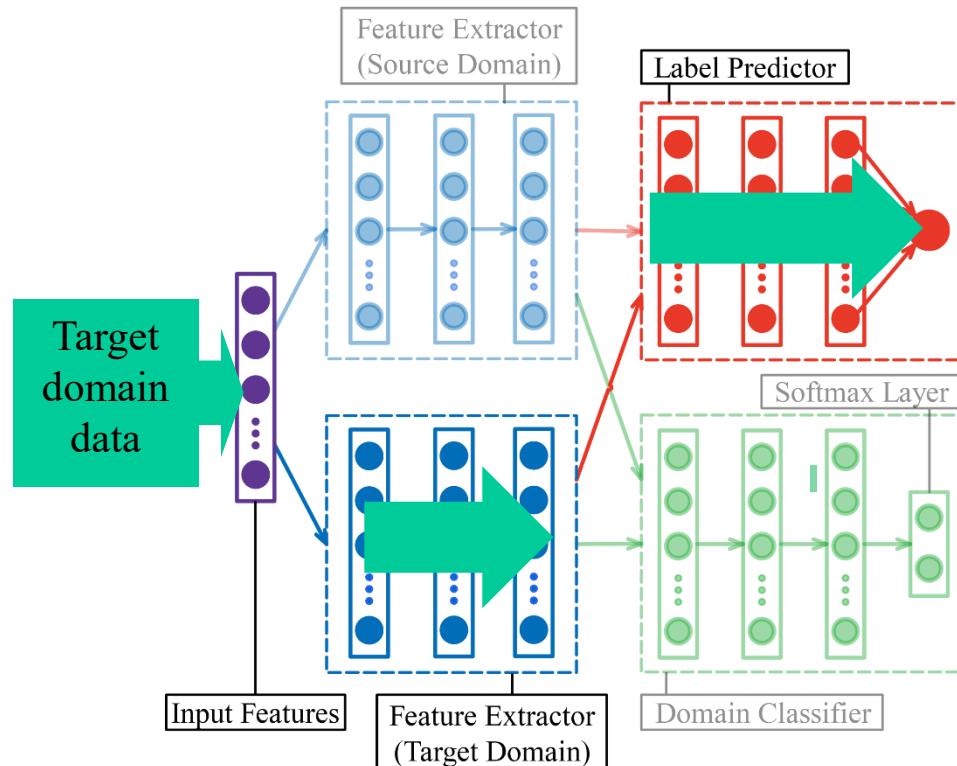
Training of ADDA (step 3 Cont.)

- As the training procedure progresses, the target feature extractor becomes better and better in deceiving the domain classifier. In another word, the target feature extractor extracts features that are similar to the source ones. In this way, the domain discrepancy is reduced.



Training of ADDA (step 4)

- After the adversarial training converges, we can predict the labels of the target domain data by using the pass as shown in the figure.



Personalizing EEG-based Affective Models with Transfer Learning

Wei-Long Zheng and Bao-Liang Lu, Personalizing EEG-based Affective Models with Transfer Learning,
Proc. of the 25th International Joint Conference on Artificial Intelligence (IJCAI-16).

Feature Reduction based Subject Transfer

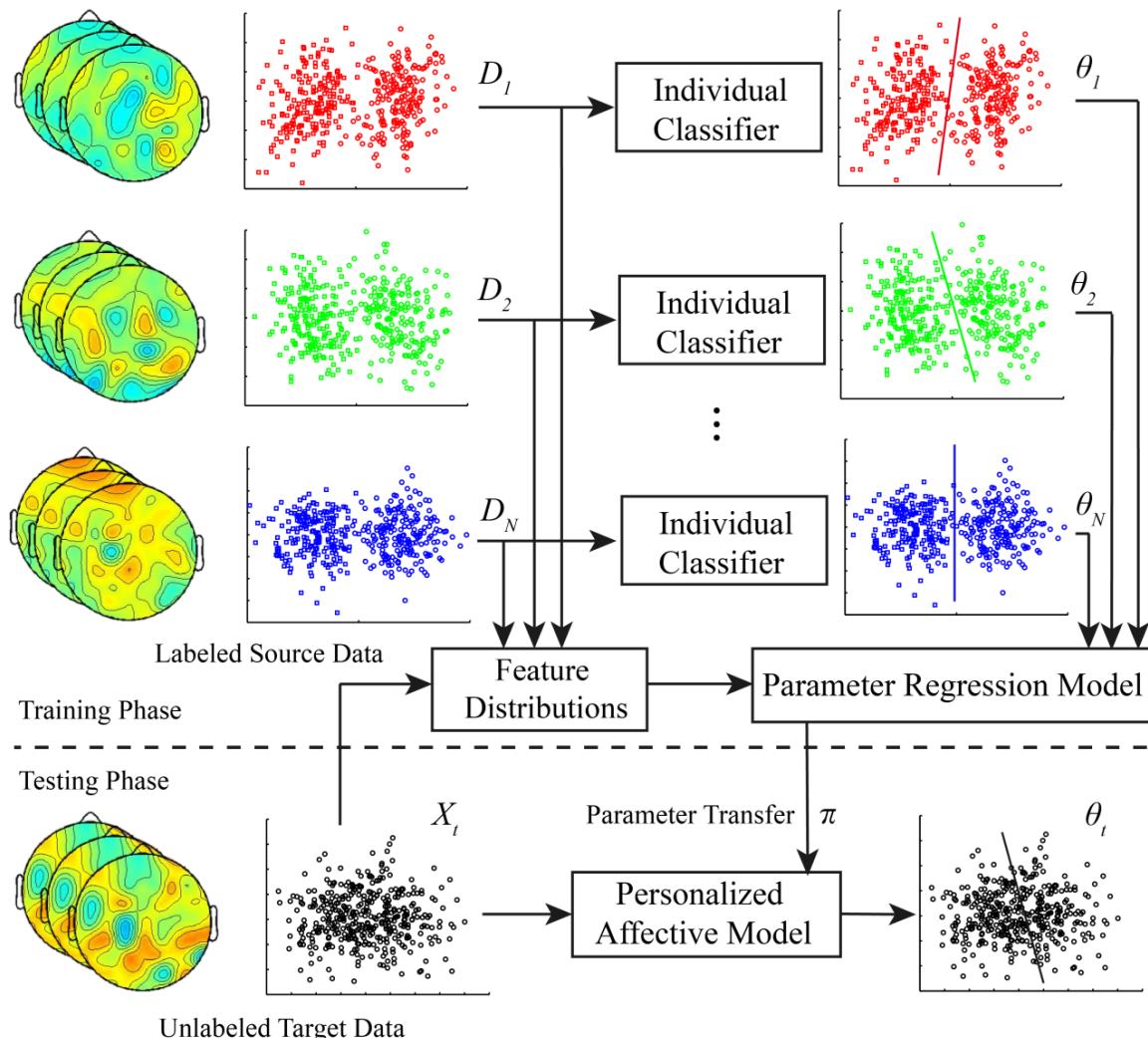
Although the distributions of source domain and target domain in high dimensional space are different, we may find a low dimensional manifold space where the distributions of both domains are similar.

TCA and KPCA try to learn a set of common transfer components underlying both the source domain and the target domain. When projected to this subspace, the difference of feature distributions of both domains can be reduced.

$$P(\phi(X_S)) \approx P(\phi(X_T))$$

$$P(Y_S | \phi(X_S)) \approx P(Y_T | \phi(X_T))$$

Transductive Parameter Transfer

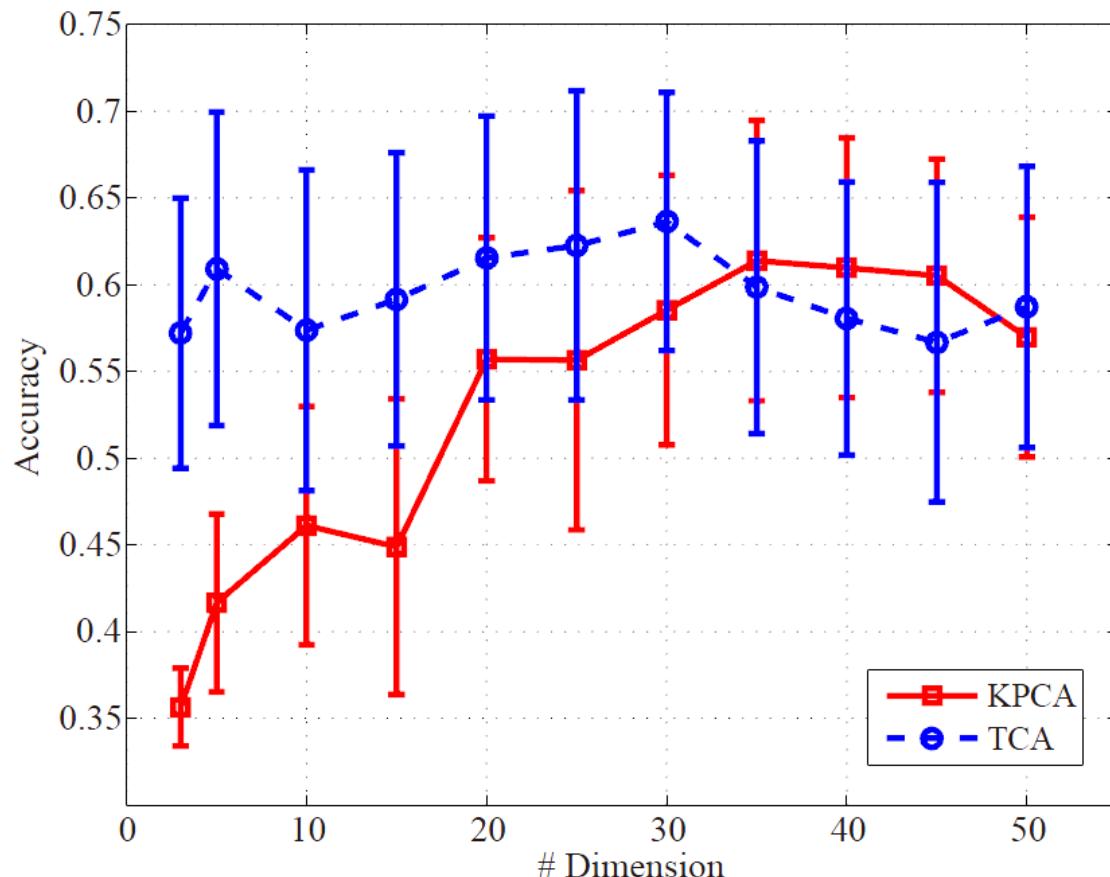


Enver Sangineto, Gloria Zen, Elisa Ricci, and Nicu Sebe. We are not all equal: Personalizing models for facial expression analysis with transductive parameter transfer. In ACM International Conference on Multimedia, pages 357–366. ACM, 2014.

Evaluation Details

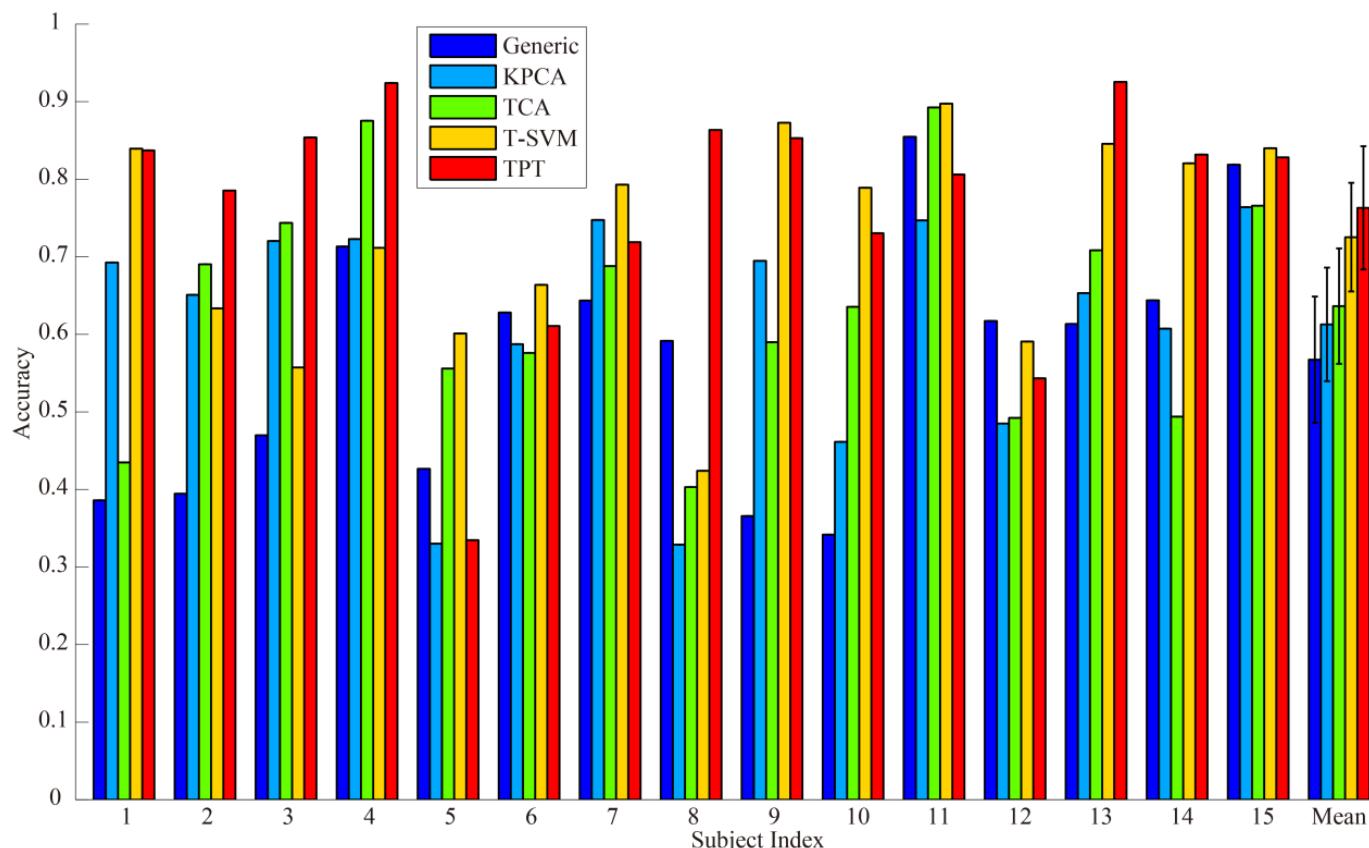
- **Leave-one-subject-out cross validation method**
- **Baseline method**
 - We concatenate data from all available subjects as training data and train a generic classifier with linear SVM.
- **Transductive SVM (T-SVM)**
 - It is developed to learn a decision boundary and maximize the margin with unlabeled data.
- **One vs. one strategy for multi-class classification**
- **All the algorithms are implemented in MATLAB.**

Experimental Results



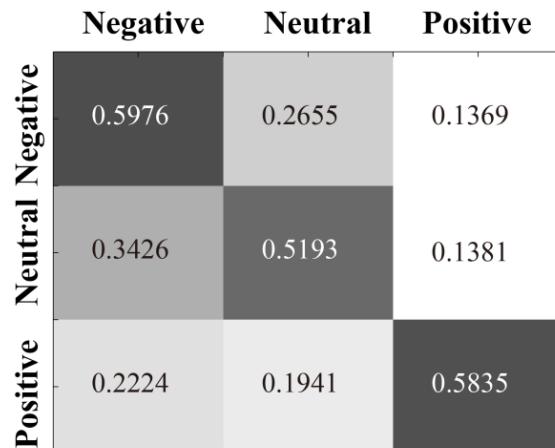
Comparison of KPCA and TCA approaches for different dimensionality of the subspace.

Experimental Results (2)

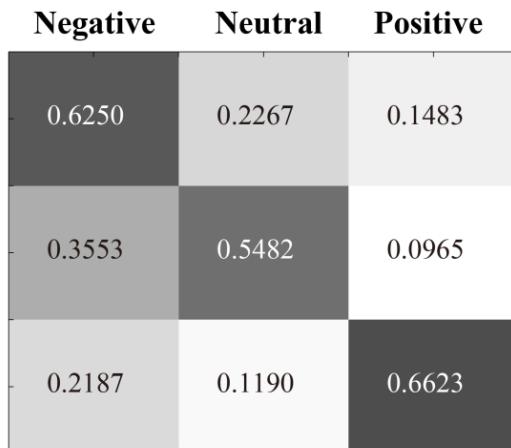


Stats.	Generic	KPCA	TCA	T-SVM	TPT
Mean	0.5673	0.6128	0.6364	0.7253	0.7631
Std.	0.1629	0.1462	0.1488	0.1400	0.1589

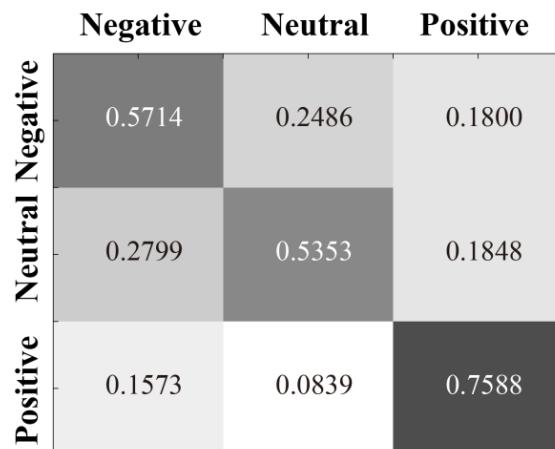
Generalization Performance



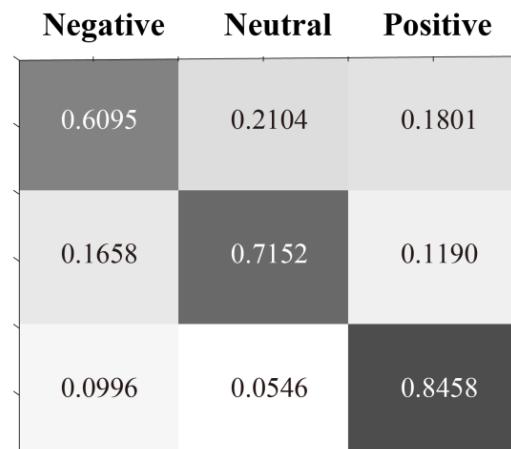
(a) Generic



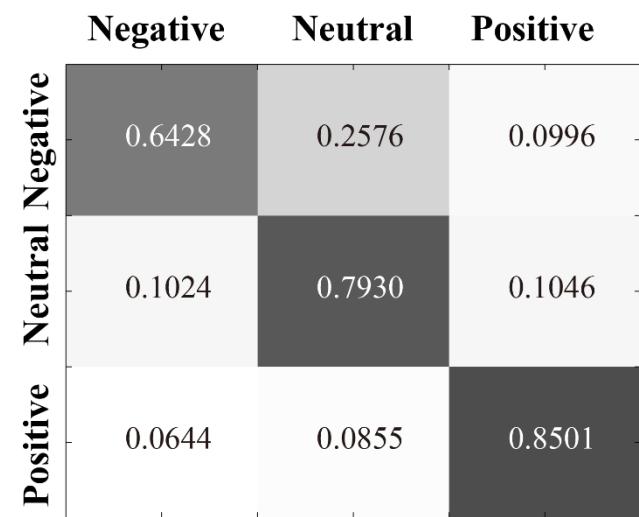
(b) KPCA



(c) TCA



(d) T-SVM



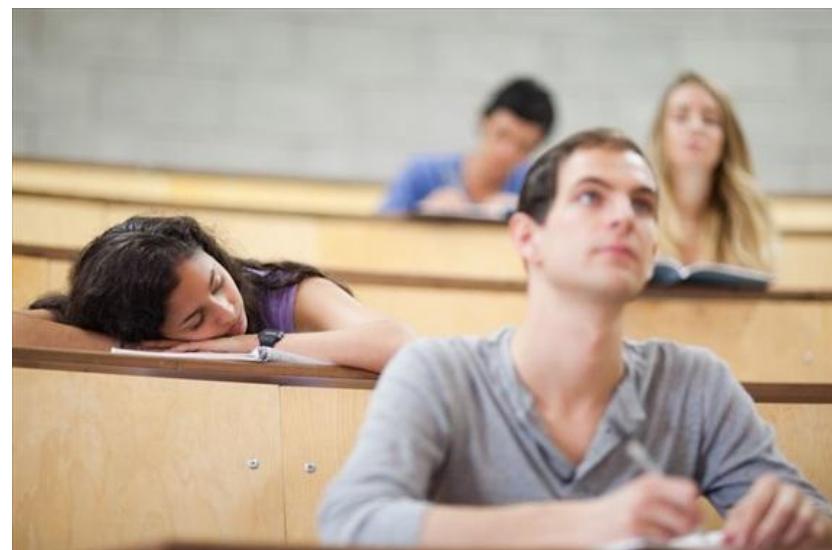
(e) TPT

Sleep Quality Estimation: From Laboratory to Real Scenario

Jia-Jun Tong, Yun Luo, Bo-Qun Ma, Wei-Long Zheng, Bao-Liang Lu, Xiao-Qi Song and Shi-Wei Ma, Proc. IJCNN2018

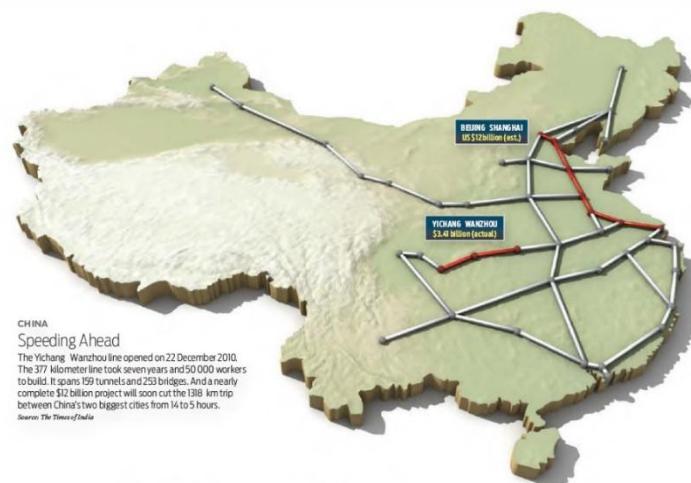
Introduction

- Sleep quality evaluation has remarkable value.
- An objective and effective measurement of sleep quality is quite valuable in transportation, medicine, health care, and neuroscience.

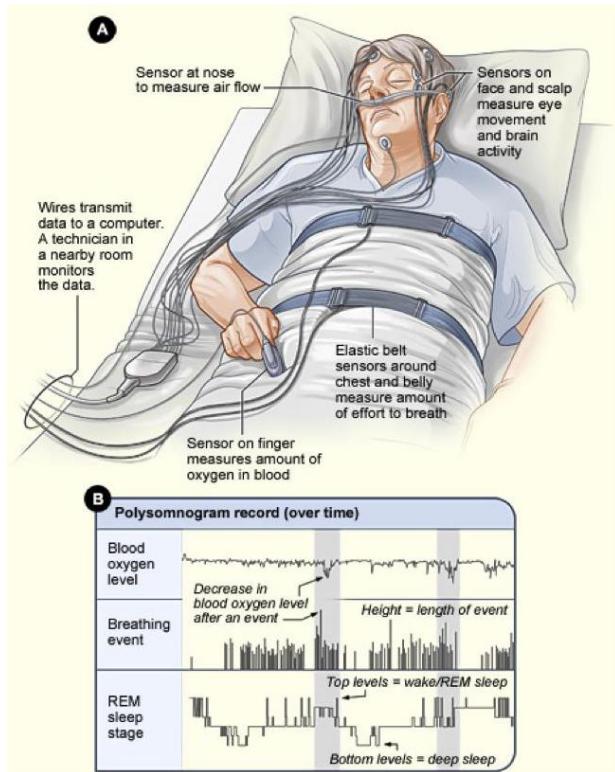


Background: High-speed railways in China

- 20380 kms in use. It is about half the world's supply
- 16775 kms are developing
- 30000 kms will be finished by 2020
- Active safety and mental state monitoring technologies become very important
 - Check sleep quality before work
 - Predict drowsy driving during work
 - Monitor mental state before work, during work, and after work



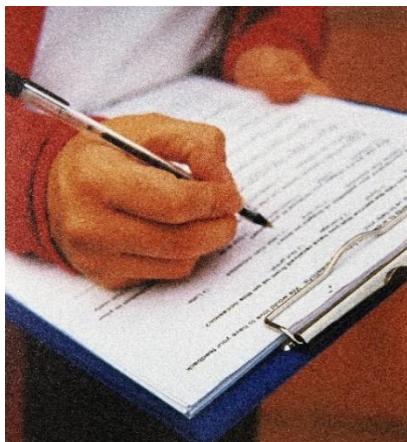
Monitoring Whole Sleep Procedure



- **Polysomnography (PSG) and smart bands require the subjects to wear equipments such as EEG cap, eye sensors, nose sensors, elastic belt sensors during the whole sleep procedure**

Our Research Goal

- At present, sleep quality assessment is performed by using questionnaire at CRH (China Railway High-speed)
 - Subjective
 - inaccurate
- An objective and convenient method as alcometer is perfect



Our previous work

We proposed an approach which is characterized by

- EEG-based Evaluation
- Data acquisition in 30 minutes, instead of whole-process physiological signal acquisition
- Assessment is data driven
- Subject-independent (transfer learning)

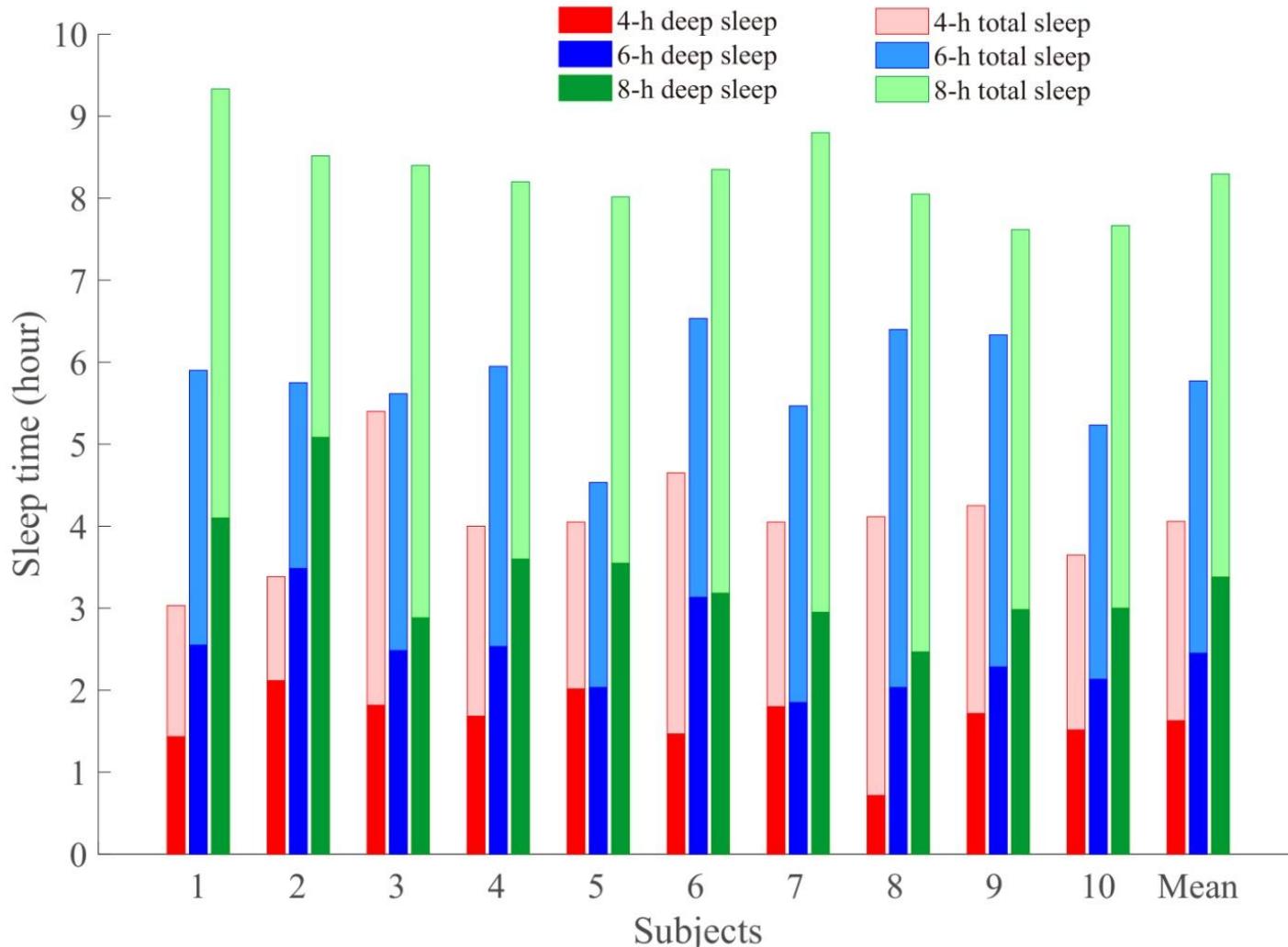
Xing-Zan Zhang, Wei-Long Zheng, and Bao-Liang Lu, EEG-based sleep quality evaluation with deep transfer learning, **Proc. ICONIP 2017**.

Li-Li Wang, Wei-Long Zheng, Hai-Wei Ma and Bao-Liang Lu, Measuring Sleep Quality from EEG with Machine Learning Approaches, **Proc. IJCNN 2016**

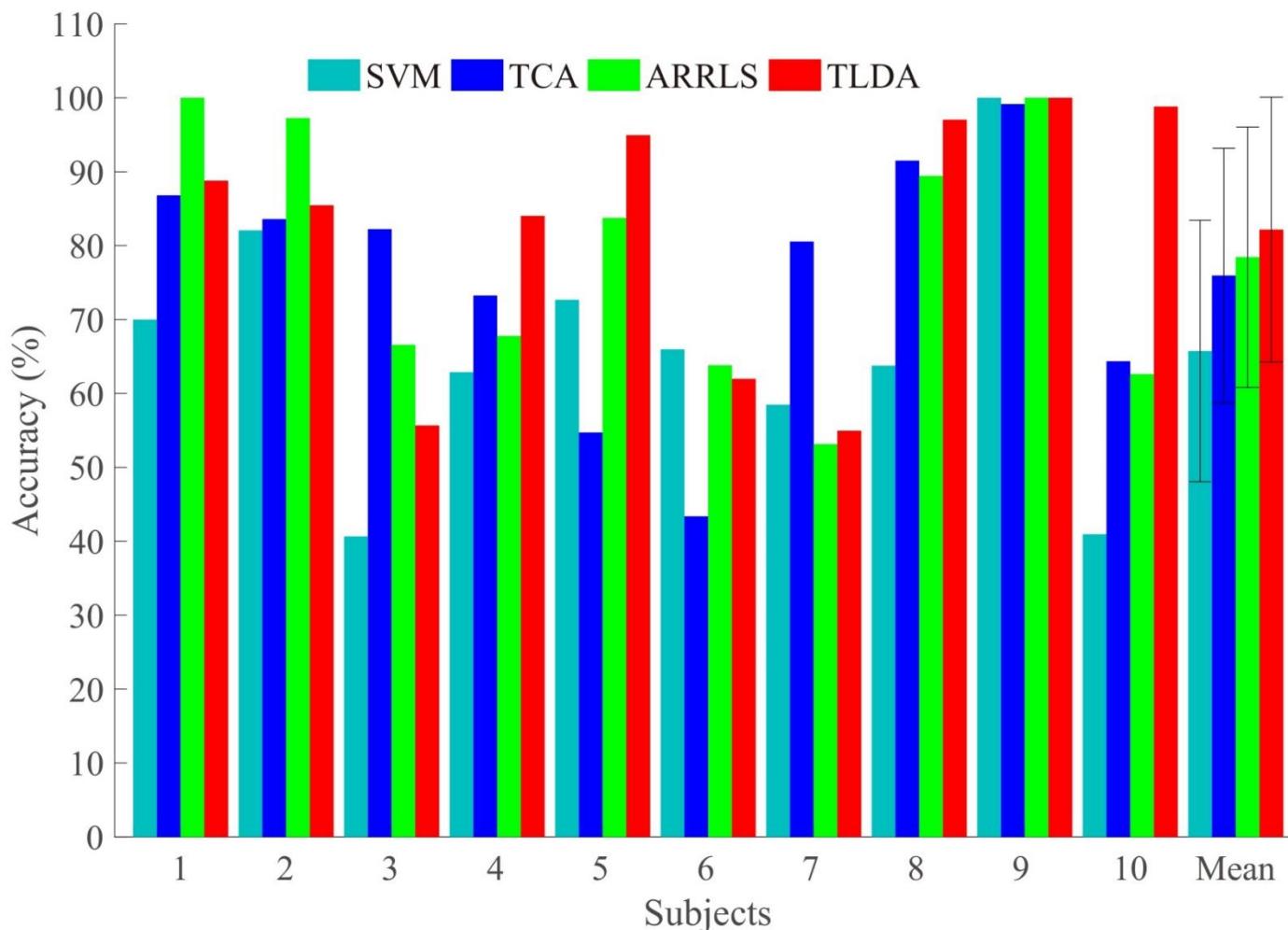
Experiment Settings

- Deep sleep is the main factor that counts for sleep quality according to National Sleep Foundation (NSF)
- 4-h (poor), 6-h (normal), 8-h (good)
- 3:00-7:00, 1:00-7:00, 23:00-7:00,
- Subjects: six males and four females
- Age range: 21-26, mean: 23.57, std: 1.62

Total and Deep Sleep Time

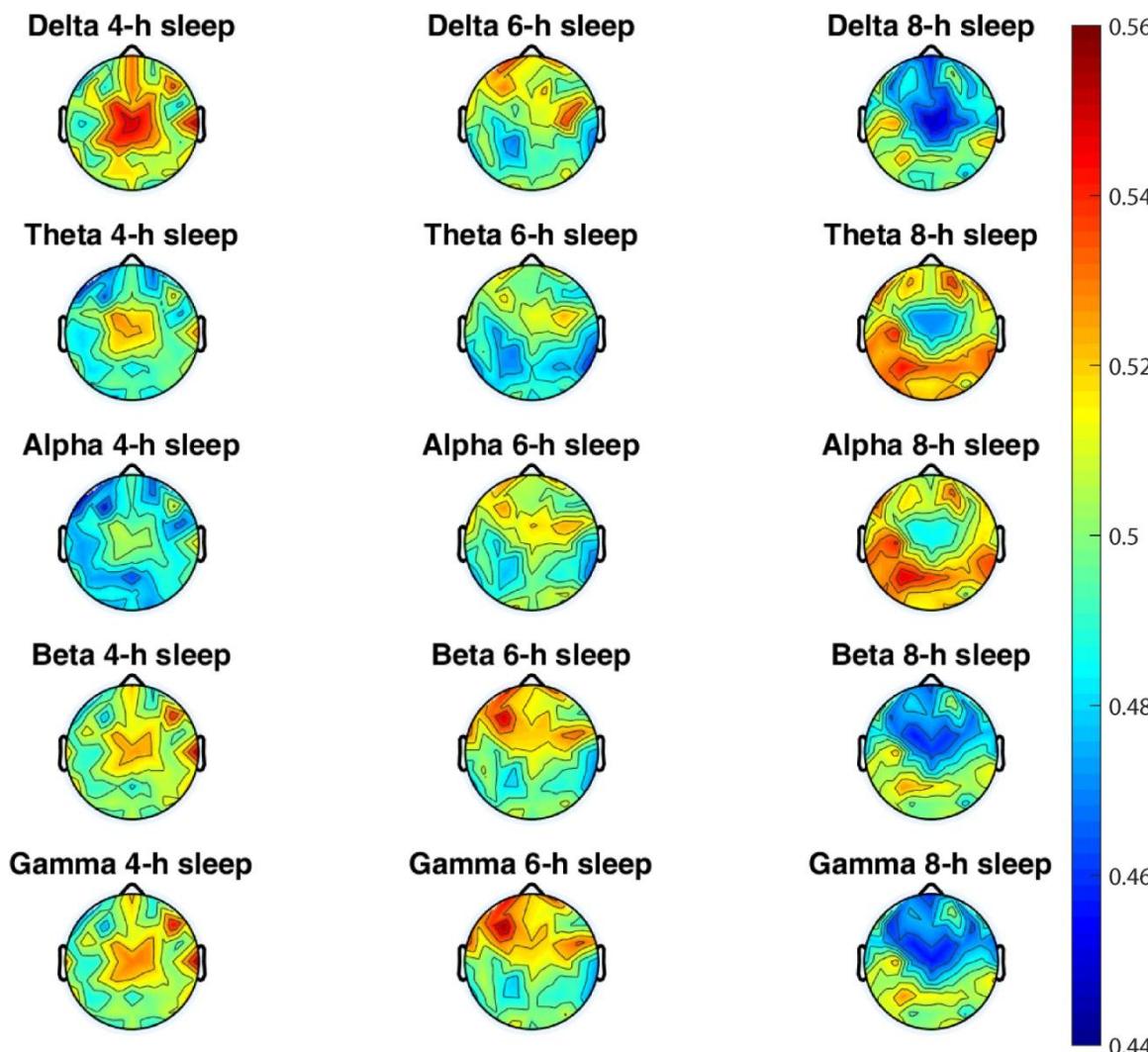


Results and discussion



Accuracy comparison of SVM, TCA, ARRLS and TLDA for each subject

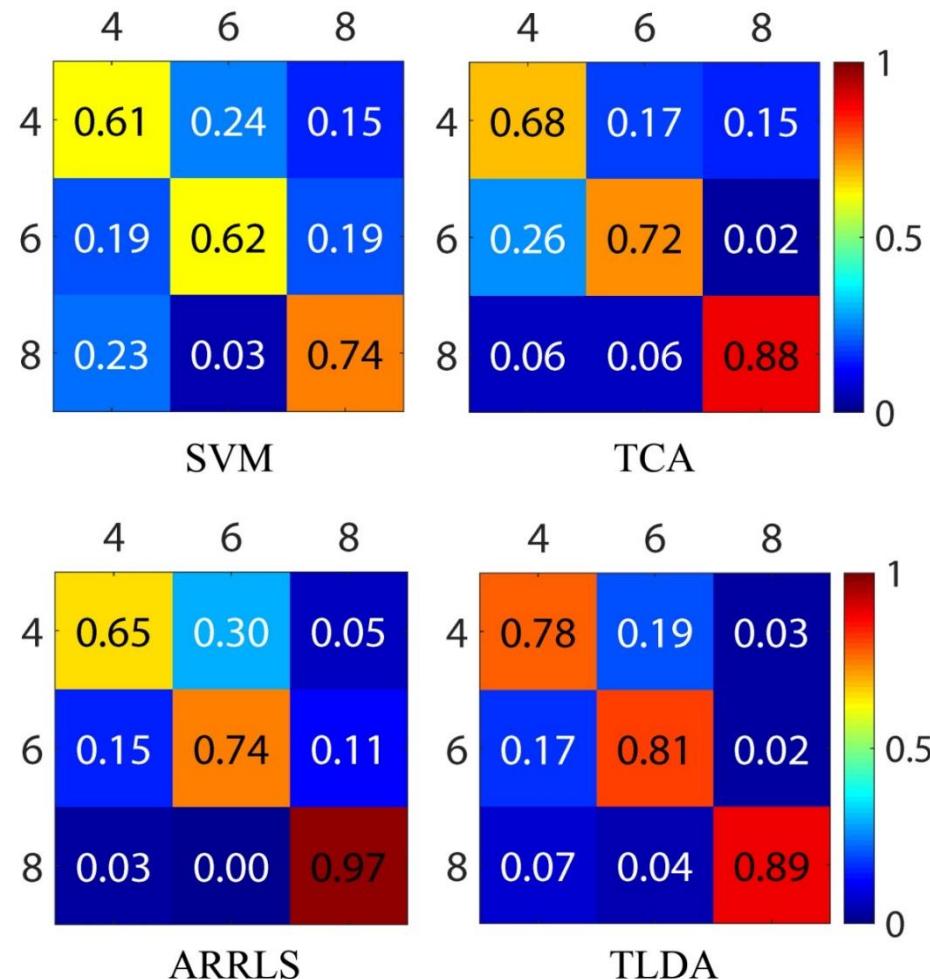
Results and discussion-DE pattern



Neural patterns of three kinds of sleep quality

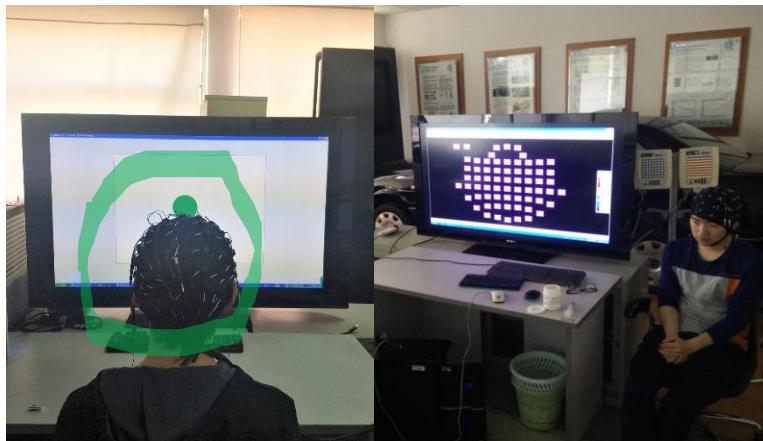
Results and discussion-confusion matrices

- 8-hour sleep can be more effectively recognized.
- 4-hour and 6-hour sleep are more likely to be misclassified with each other.
- DE patterns of 4-hour sleep and 6-hour sleep are much closer than the patterns of 8-hour sleep.



Motivation

- The data collection and annotation in real scenarios are usually costly and difficult
- Transferring knowledge from laboratory to real scenario
 - Laboratory: Students with controlled sleep time
 - Real Scenario: CRH railway drivers before work
- Cross-scenario & Cross-subject tasks



Modeling
Predicting



In the Laboratory

- We take 4-hour sleep, 6-hour sleep and 8-hour sleep as poor, normal and good in terms of the sleep quality in our study
- The sleep time and wake up time for three experiments are 3:00-7:00, 1:00-7:00 and 23:00-7:00, respectively.
- Total sleep time and deep sleep time of all subjects are recorded by smart bands
- High quality EEG acquisition: 62-channel wet electrode cap
- Different state for same subjects: 3 times, 30 min each experiments



4h

6h

8h

In the Real Scenario (CRH)

- To keep the drivers' daily routine, the experiment settings have to be adjusted to meet the need of real-scenario application.
 - No control on sleep time: range of 2~10 h
 - Easy to set up devices: 18-channel dry electrode cap (DSI-24)
 - One experiment for each subjects: 1 time, 6 min each



Sleep:
2~10h



(CRH: China Railway High-speed)

Summary of Domain Differences

Categories	Specifications	Laboratory	Real-scenario
Subjects	Age	21-26	25-49
	Sex-distribution	4 males, 6females	70 males
	Occupation	Students	High-speed train drivers
	Experiments per subject	3	1
Controlled conditions	Sleep time	4,6,8 hours	Not controlled
	Monitoring sleep	Yes	No
	Head Cleaning	Yes	No
	Experiment Starts	≤1 hour after wake up	1-10 hour after wake up
Experiment settings	Experiment duration	30 min	6 min
	Task	Eyes open	Eyes open and close
	Environment	Quite and isolated	Noisy
Devices	Electrods type	Wet electrodes	Dry electrodes
	Resolution	62 channels	18 channels
	Reference	REF(Between CPZ and PZ)	A1 and A2 (On the ears)
	Common Mode Follower	0	1 (CMF)
	Sample rate	1000 Hz	300 Hz

Eight Domain Adaptation Methods

- Baseline methods
 - TCA: Transfer Component Analysis
 - ITL: Information-Theoretical Learning
 - GFK: Geodesic Flow Kernel
 - JDA: Joint Distribution Adaptation
 - SA: Subspace Alignment
 - TJM: Transfer Joint Matching
 - MIDA: Maximum Independence Domain Adaptation
- Adversarial Networks
 - Domain Adversarial Neural Network

Cross-subject and Cross-scenario Tasks

- After the preprocessing that unifies the format of data, the two domains still possess a lot of discrepancies
- View the Laboratory data as source domain, we can then use domain adaptation methods to transfer knowledge to the real scenario
- Cross-subject task: training and testing on real-world data
- Cross-scenario task: training on laboratory data but testing on real-world data

Results: Cross-subject Task

ACCURACY (%) OF CROSS-SUBJECT TASK

Method	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean±Std
SVM	59.57	48.41	57.42	55.71	49.74	54.17±4.87
TCA	58.37	72.98	63.77	69.48	53.85	63.69±7.82
ITL	66.07	66.79	57.14	76.43	56.71	64.63±8.13
GFK	59.41	66.79	53.73	58.45	58.49	59.37±4.70
JDA	59.41	59.52	76.11	60.44	55.91	62.28±7.92
SA	77.42	67.38	60.36	70.68	46.47	64.46±11.7
TJM	65.91	74.72	58.73	62.22	59.37	64.19±6.53
MIDA	71.35	69.76	66.91	69.21	56.51	66.75±5.94
DANN	67.50	74.64	71.23	76.98	78.25	73.72±4.38

The model performance is evaluated in a 5-fold cross validation manner. The real-world data is split into 5 parts, each containing 14 subjects and no subject's data should appear in multiple validation parts.

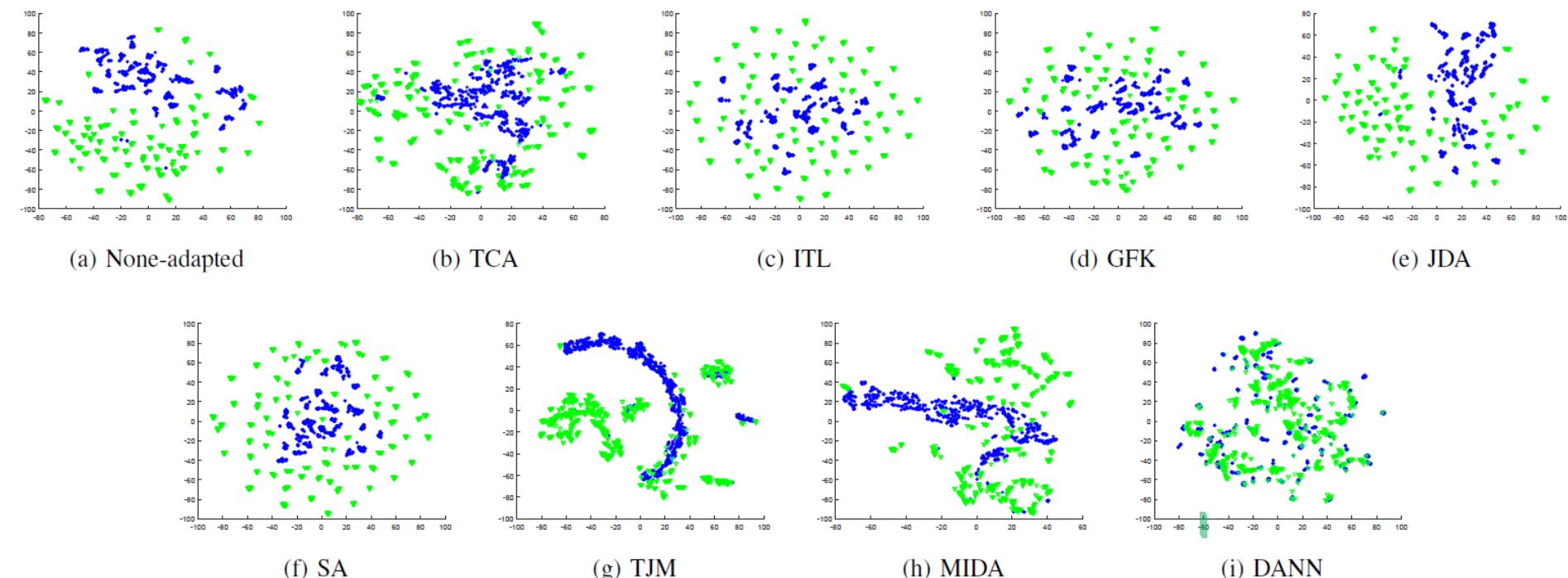
Results: Cross-scenario Task

ACCURACY (%) OF CROSS-SCENARIO TASK

Method	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean±Std
SVM	43.86	40.35	42.18	38.98	43.59	41.79±2.10
TCA	54.60	61.19	38.89	67.34	56.87	55.60±10.6
ITL	53.69	63.10	24.68	60.71	49.41	50.32±13.7
GFK	49.64	62.46	40.99	70.95	49.25	54.66±11.9
JDA	47.70	61.43	35.71	66.11	66.94	55.58±13.5
SA	42.74	62.54	37.50	61.23	54.17	51.64±11.1
TJM	41.67	60.28	38.97	59.88	60.60	52.28±10.9
MIDA	56.47	66.35	38.69	68.93	56.87	57.46±11.9
DANN	69.10	64.68	60.79	62.6296	69.25	65.29±3.80

Experimental Results (cont.)

The effects of adaptation on cross-scenario task



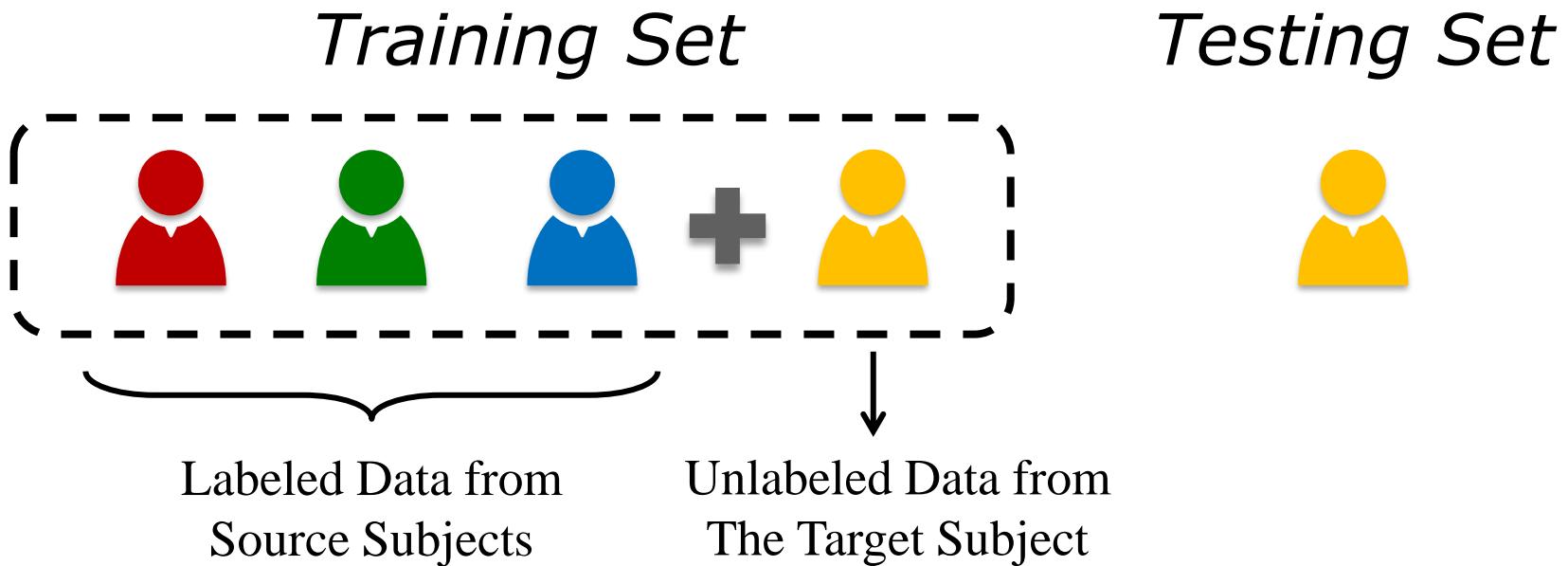
Conclusions

- To estimate the sleep quality of high-speed train drivers, we have evaluated two tasks: the cross-subject and the cross-scenario.
- Eight domain adaption methods have been introduced to both tasks to reduce the domain discrepancy between source and target domains.
- The experimental results have shown that Domain Adversarial Neural Network (DANN) approach outperforms other seven DA models significantly on both tasks in terms of classification accuracy.
- The DANN approach is also stable and robust on the cross-scenario task.

Reducing EEG Subject Variability by Domain Generalization with Deep Adversarial Network

Reducing the Subject Variability

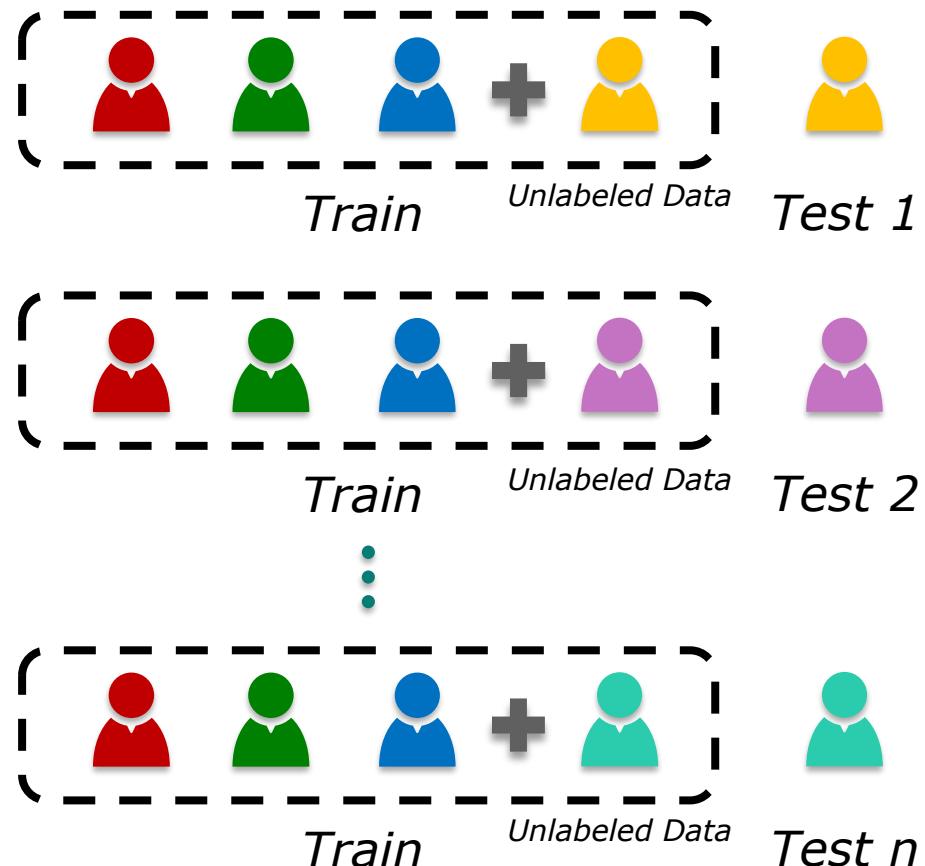
- Subject-dependent (Domain Adaptation)
 - Utilize a small amount of individualized training data
 - Require an entire training session for each test subject



Wei C S, Lin Y P, Wang Y T, et al. A subject-transfer framework for obviating inter-and intra-subject variability in EEG-based drowsiness detection. NeuroImage, 2018, 174: 407-419.

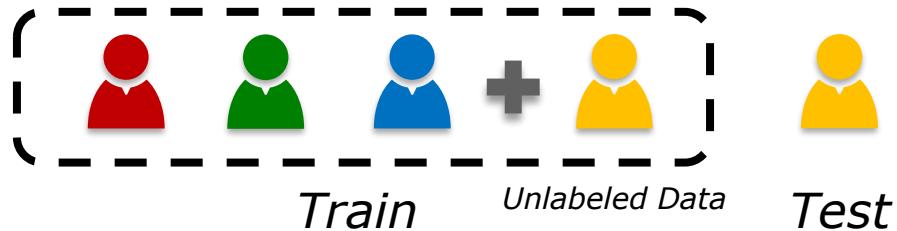
Limitations

- ❑ More new subjects means more models to train
- ❑ Subject domain information loss in traditional transfer learning problem
- ❑ Previous training provides little intuition for the new domain feature extractor

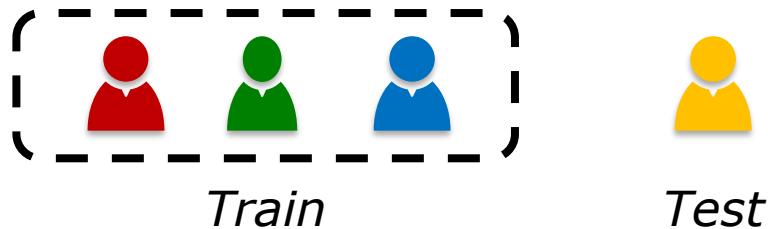


Reducing the Subject Variability

- Subject-dependent (Domain Adaptation)
 - Utilize a small amount of individualized training data
 - Require an entire training session for each test subject



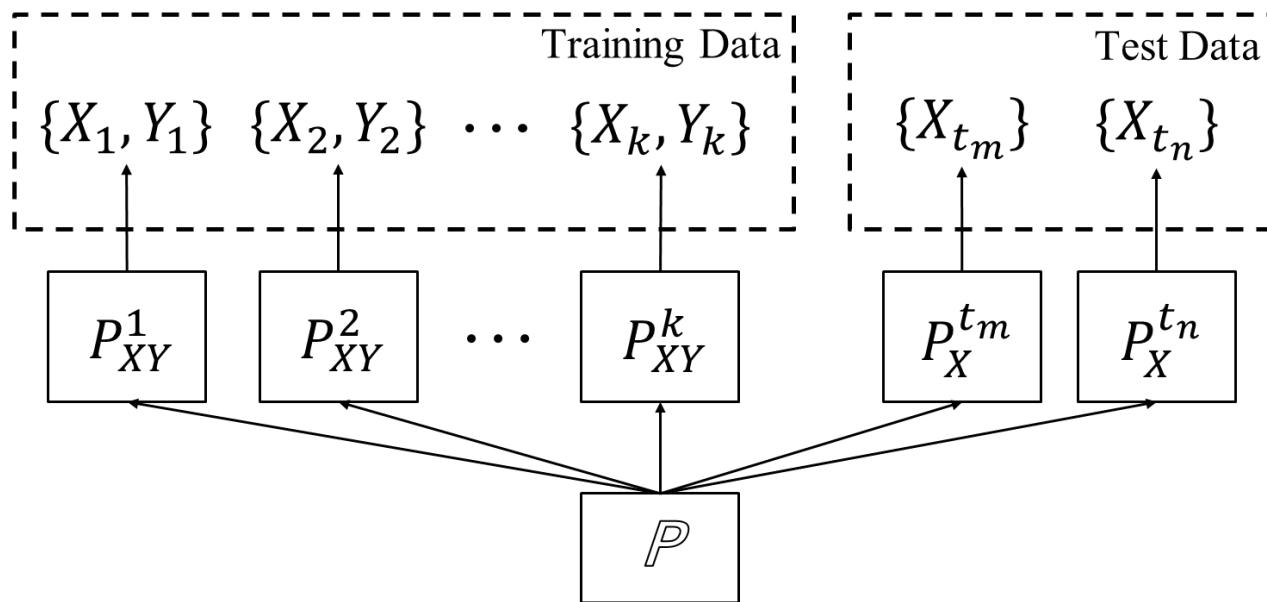
- Subject-independent (Domain Generalization)
 - No data from the new user
 - A robust feature-extraction method



Wei C S, Lin Y P, Wang Y T, et al. A subject-transfer framework for obviating inter-and intra-subject variability in EEG-based drowsiness detection. NeuroImage, 2018, 174: 407-419.

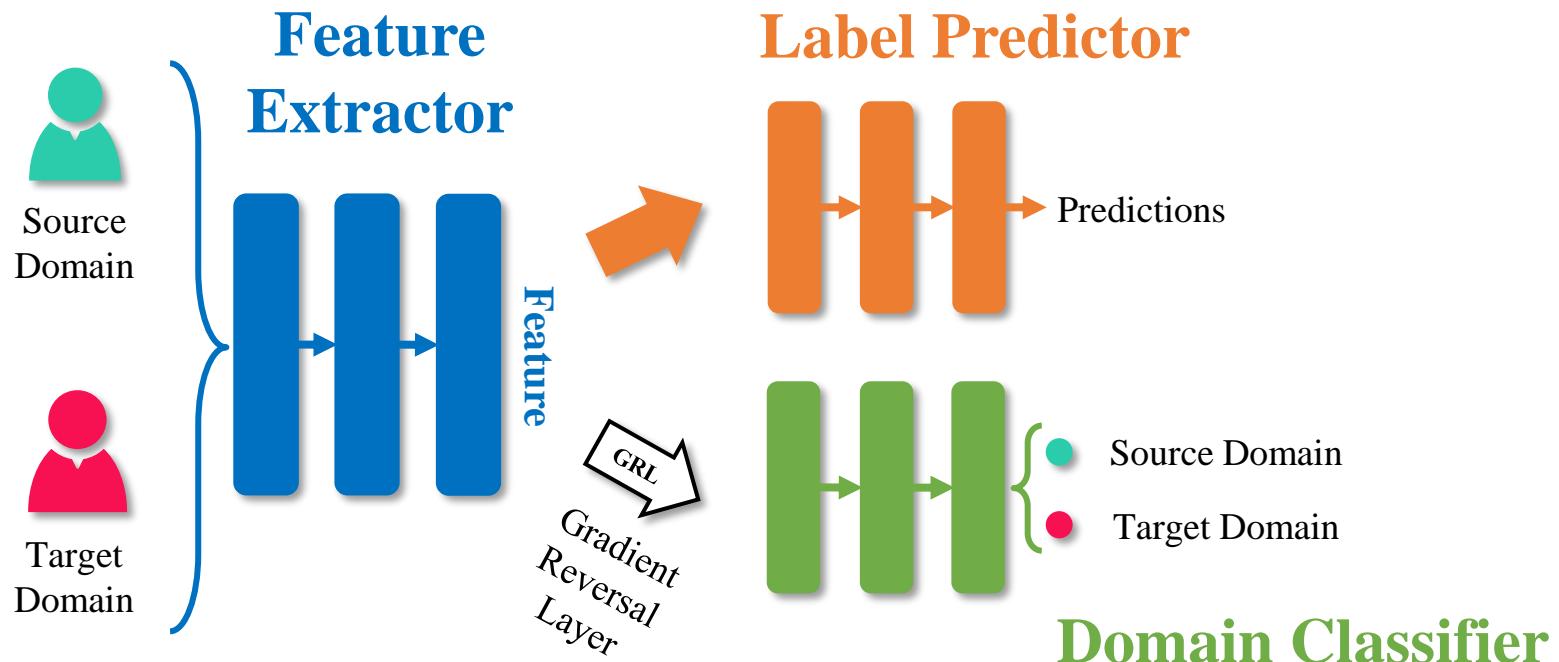
Domain Generalization

- Domain generalization considers how to take knowledge acquired from an arbitrary number of related domains, and apply it to previously unseen domains
- Estimate a functional relationship that handles changes in the marginal $P(X)$ or conditional $P(Y | X)$ well



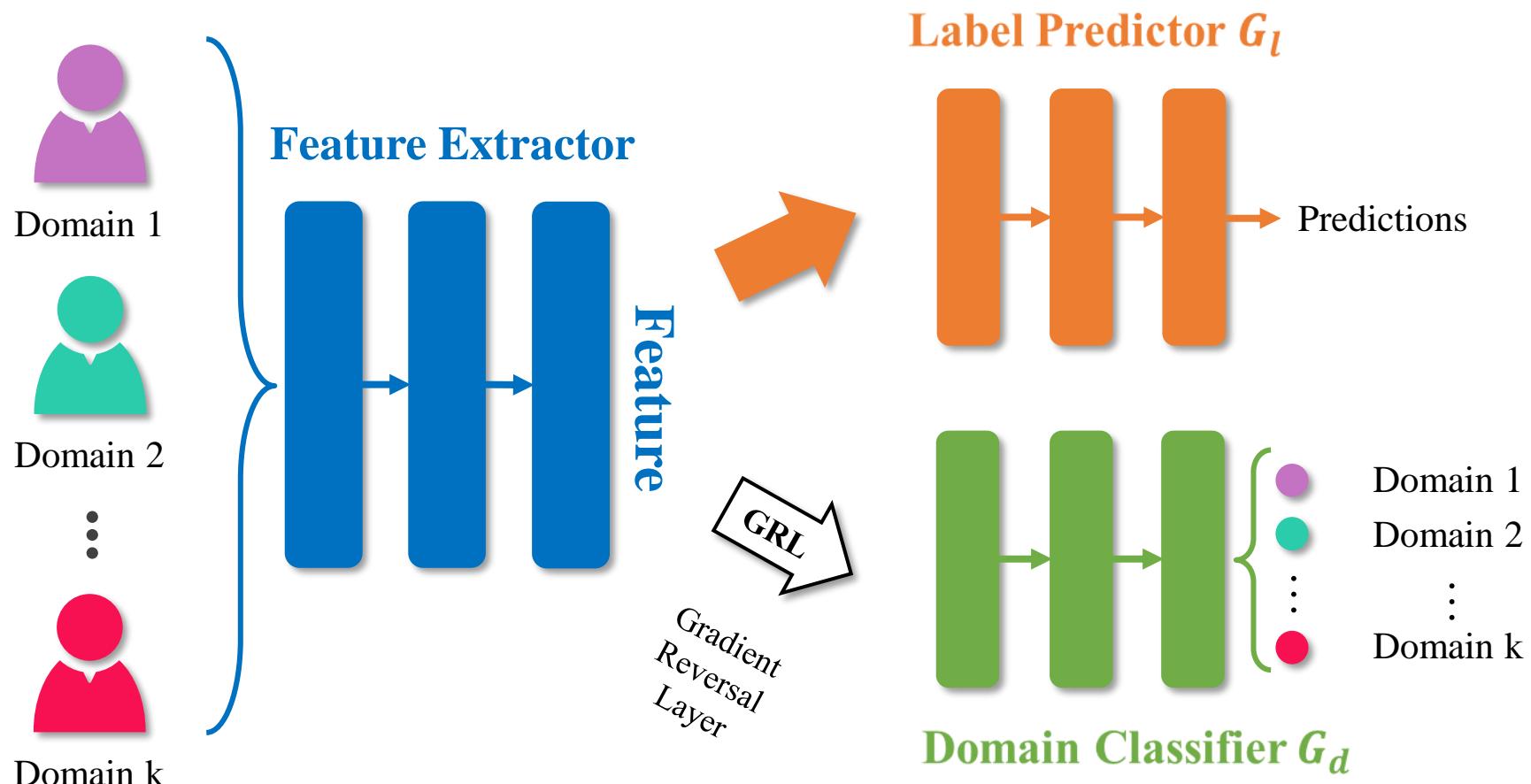
Domain Generalization with Deep Adversarial Network

- Our goal is to train a robust feature extractor (without the data from the new subject) which can :
 - Keep essential information for label prediction
 - Reduce the domain shift (features from both domains share similar distributions in the new feature space)



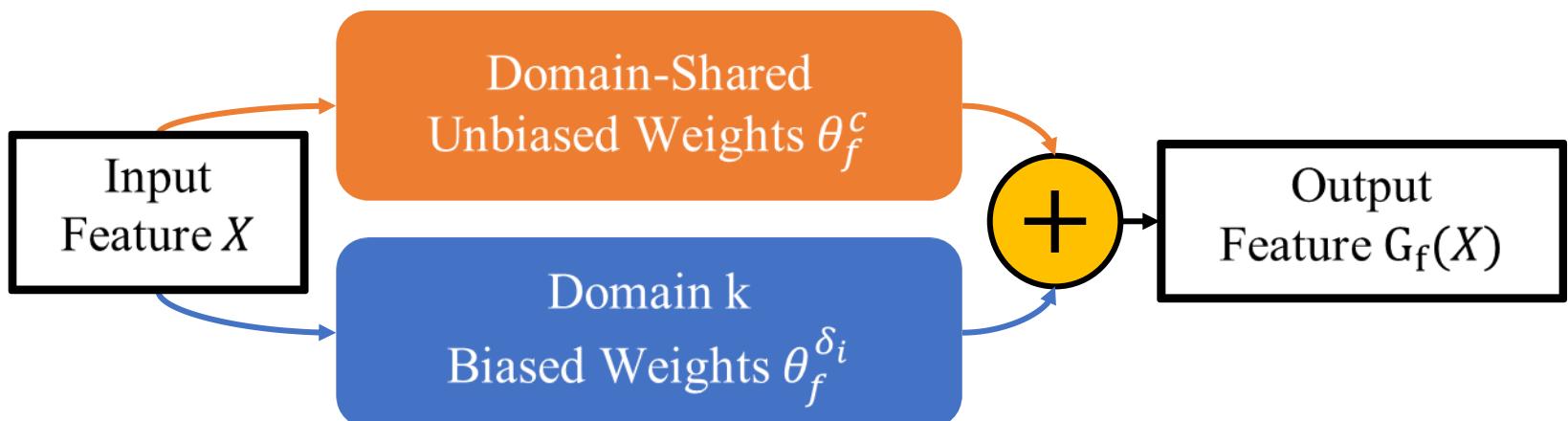
Subject-independent Approach: DG-DANN

- One way is to find the domain-invariant feature space with the information of multiple source domains

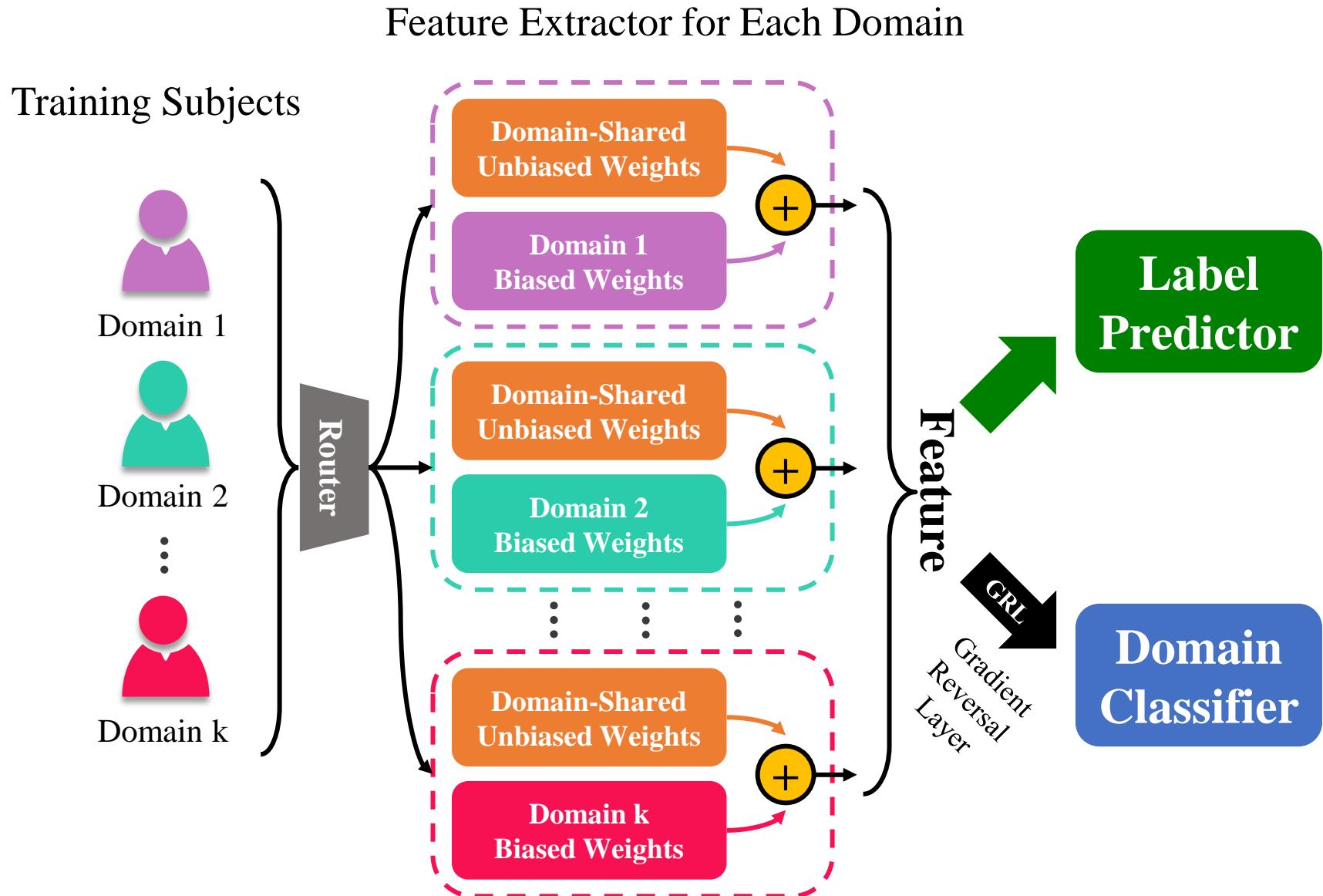


Modifying the Feature Extractor

- Another way is to regulate the model parameters
- We assume that the features consist of two parts :
 1. Common space information, which is learnt by the shared unbiased weights θ_f^c
 2. Domain-specific information, which is explicitly defined



Domain Residual Network: Training



Domain Residual Network: Training

Training Subjects



Domain 1



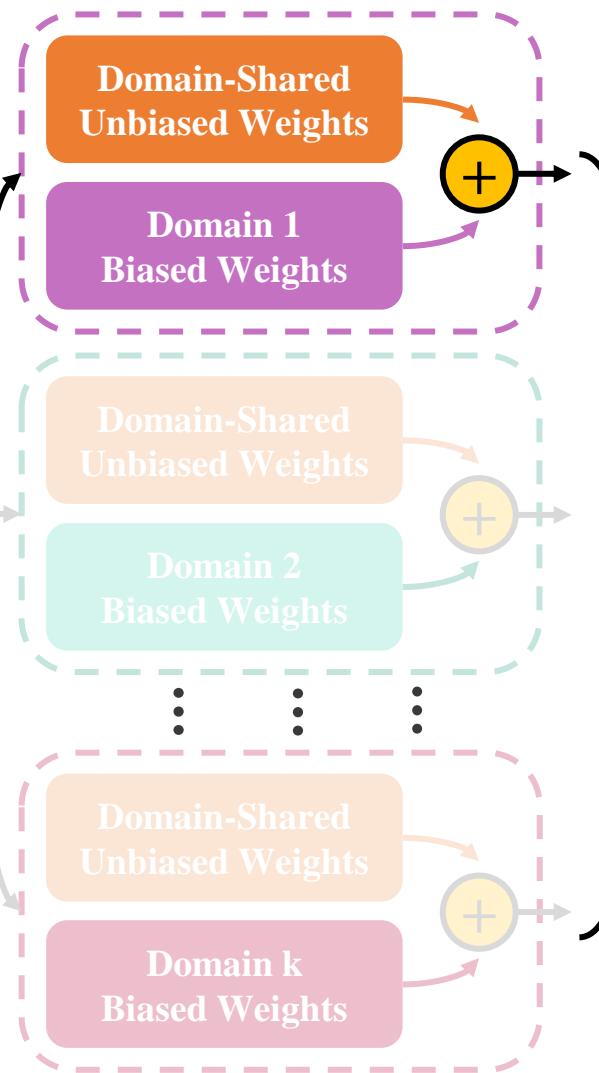
Domain 2

⋮



Domain k

Feature Extractor for Each Domain



Feature

GRL
Gradient
Reversal
Layer

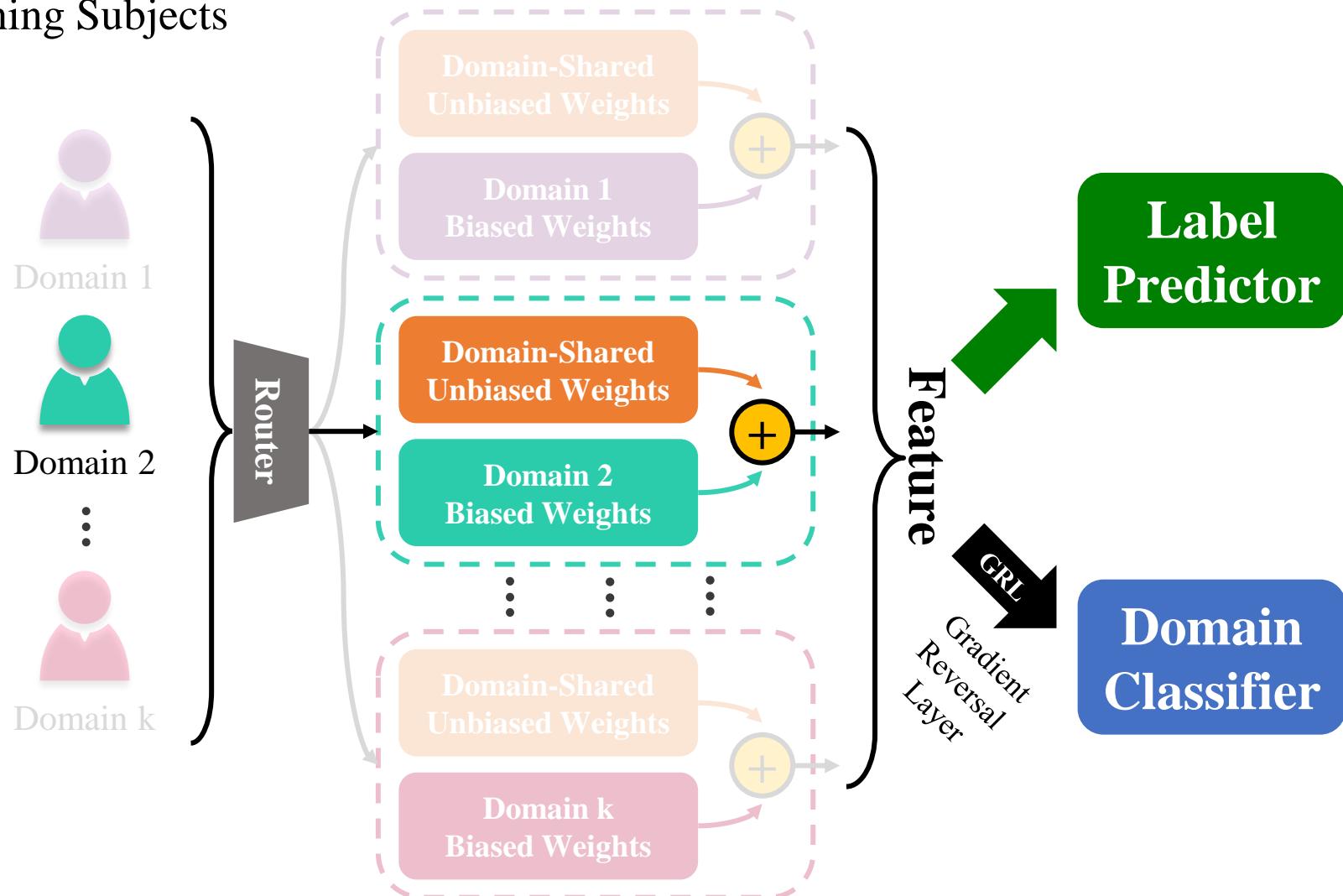
Label Predictor

Domain Classifier

Domain Residual Network: Training

Training Subjects

Feature Extractor for Each Domain

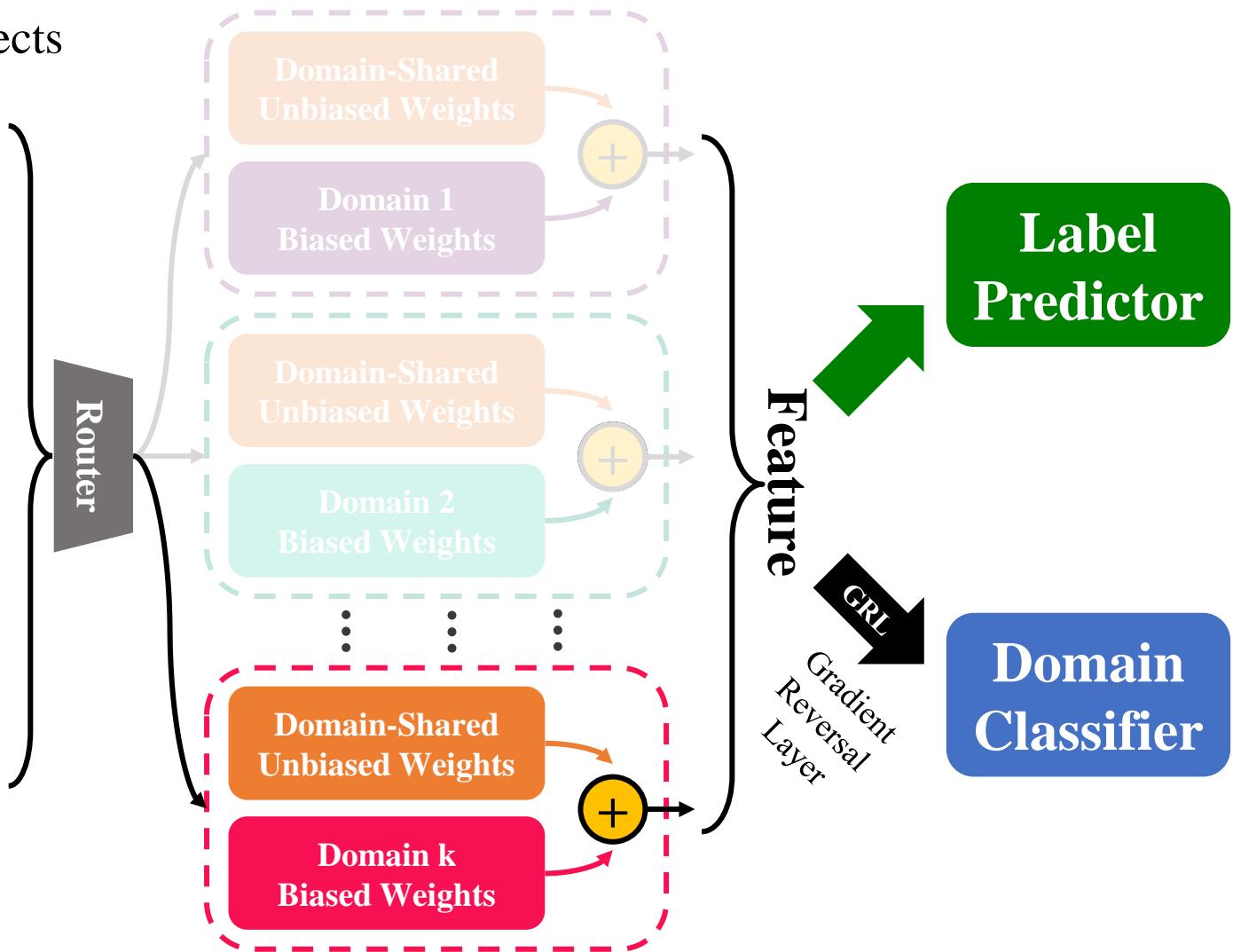


Domain Residual Network: Training

Training Subjects

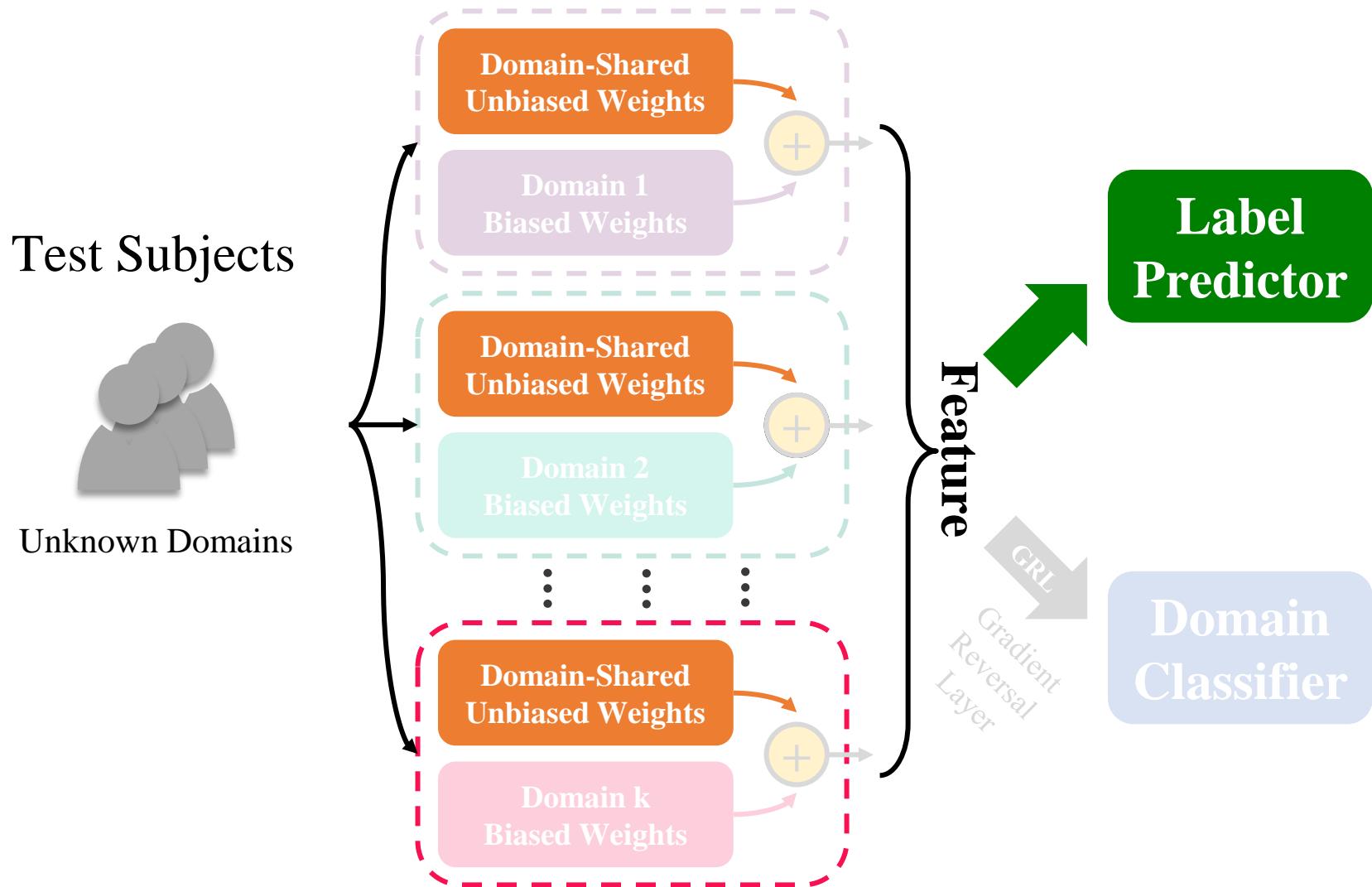


Feature Extractor for Each Domain



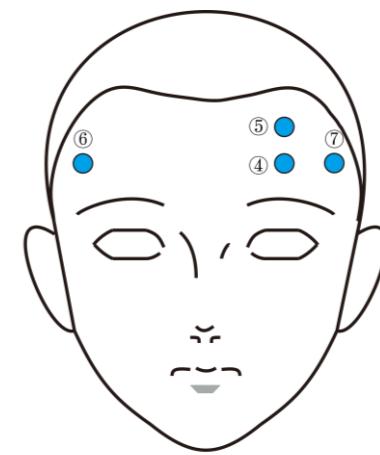
Domain Residual Network: Testing

Feature Extractor for Each Domain



Experiment: Vigilance Estimation

- The SEED-VIG dataset is a publicly available dataset. It contains 23-subject EEG and EOG signals recorded during simulated driving procedures.
- The vigilance levels are labeled in form of PERCLOS indices
- 885*136 features for each subject



(b) Forehead EOG

Wei-Long Zheng and Bao-Liang Lu, A multimodal approach to estimating vigilance using EEG and forehead EOG. Journal of Neural Engineering, 14(2): 026017, 2017.

Result : Vigilance Estimation

□ Leave-one-subject-out cross validation:

- 22 subjects for training
- 1 subject for testing

		Baseline	Domain Adaptation				Domain Generalization		
		SVR	TCA	GFK	DANN	ADDA	DICA	DG-DANN	DResNet
PCC	Avg	0.7606	0.7786	0.7907	0.8402	0.8442	0.7733	0.8320	0.8440
	Std	0.2314	0.2152	0.1260	0.1535	0.1336	0.1382	0.1000	0.0935
RMSE	Avg	0.1689	0.1596	0.1910	0.1427	0.1405	0.2007	0.1470	0.1420
	Std	0.0673	0.0544	0.0636	0.0588	0.0514	0.0674	0.0444	0.0402

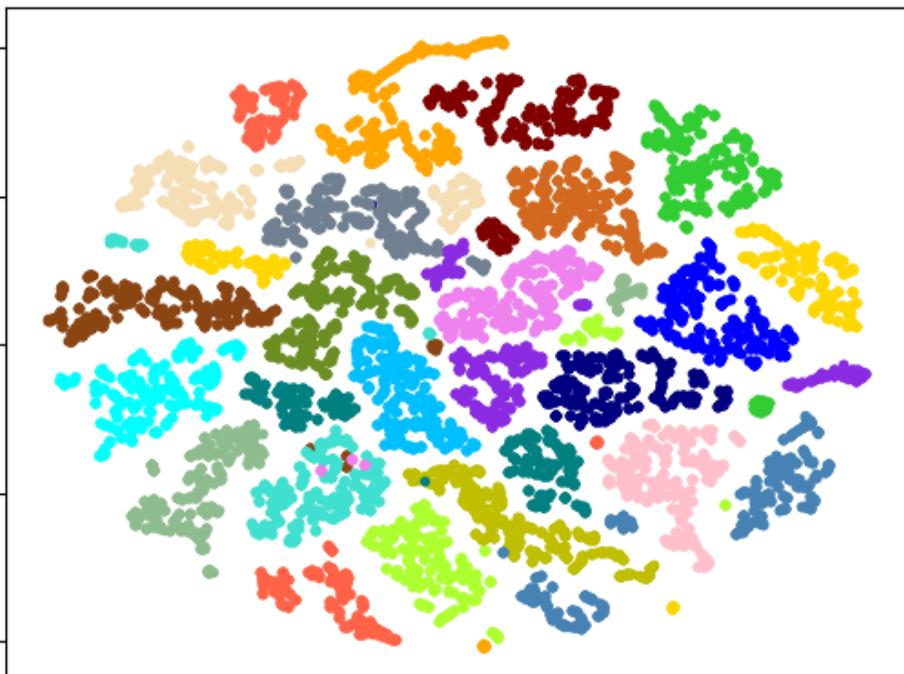
Result : Vigilance Estimation

- Leave-multiple-random-subjects-out evaluation:
 - 2/3 of the subjects (15) for training
 - 1/3 of the subjects (8) for testing

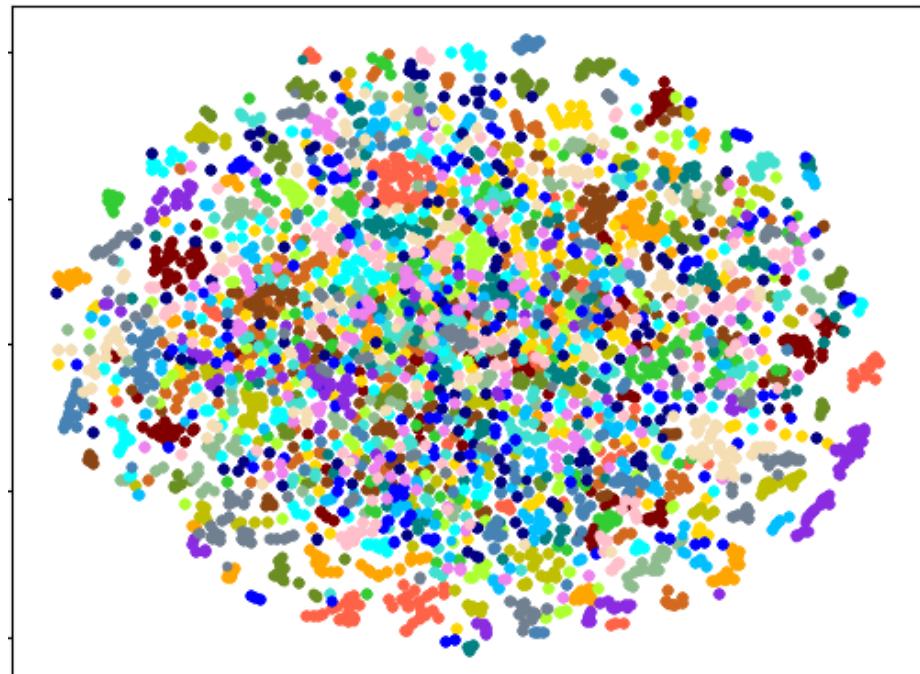
		Baseline	Domain Generalization		
		SVR	DICA	DG-DANN	DResNet
PCC	Avg	0.7499	0.7719	0.8294	0.8386
	Std	0.1980	0.1841	0.1541	0.1532
RMSE	Avg	0.2068	0.1735	0.1604	0.1569
	Std	0.0587	0.0468	0.0782	0.0735

Visualized Feature for SEED-VIG

Raw Feature

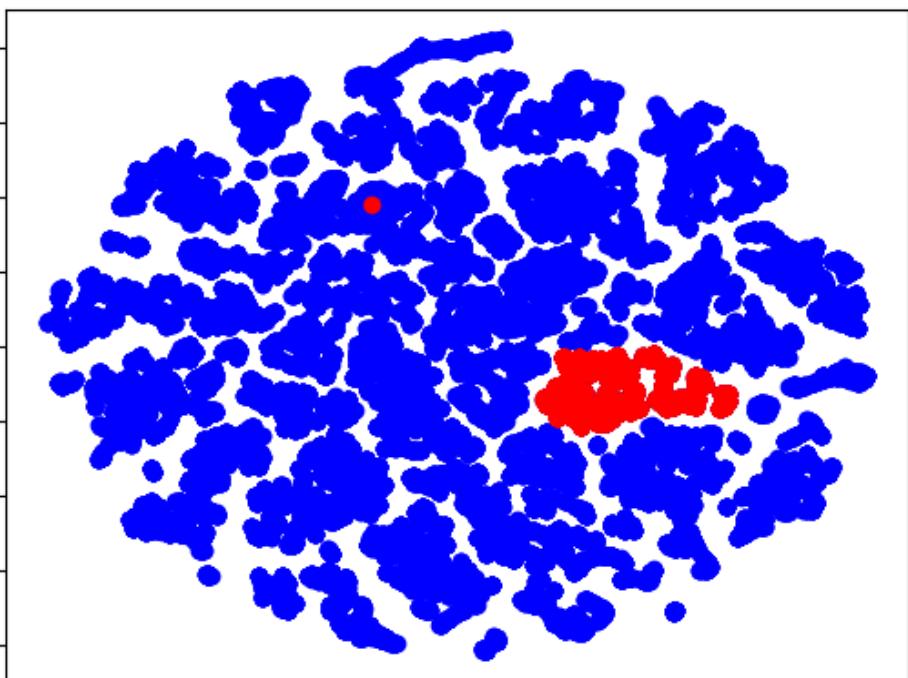


DResNet Feature

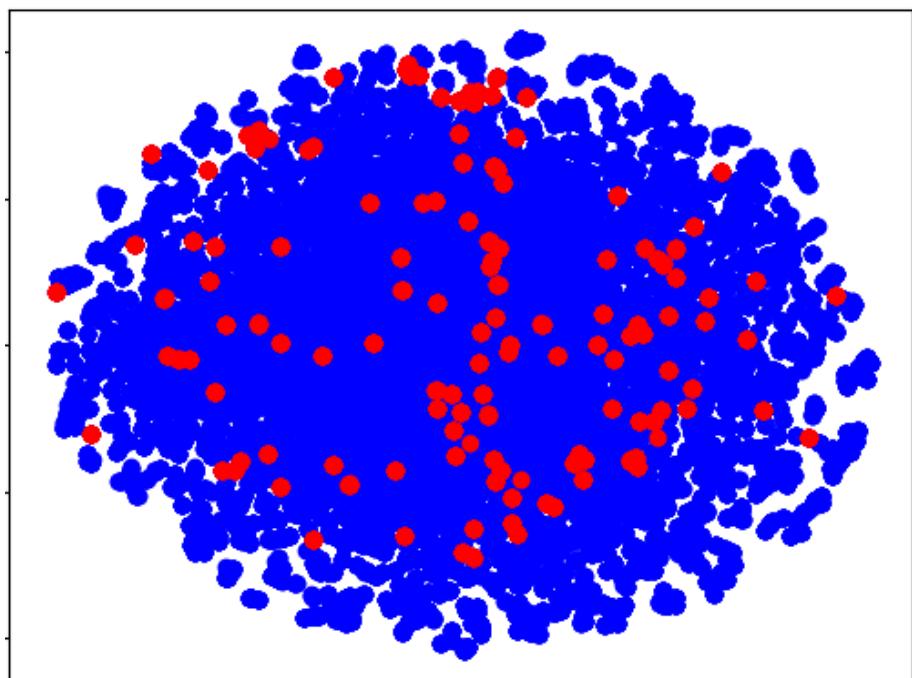


Visualized Feature for SEED-VIG

Raw Feature

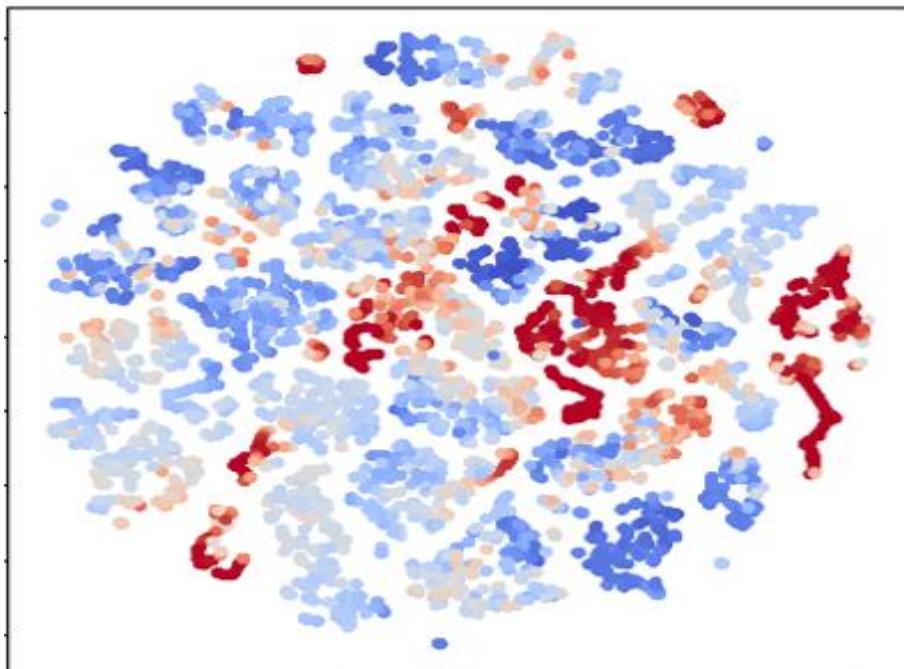


DResNet Feature

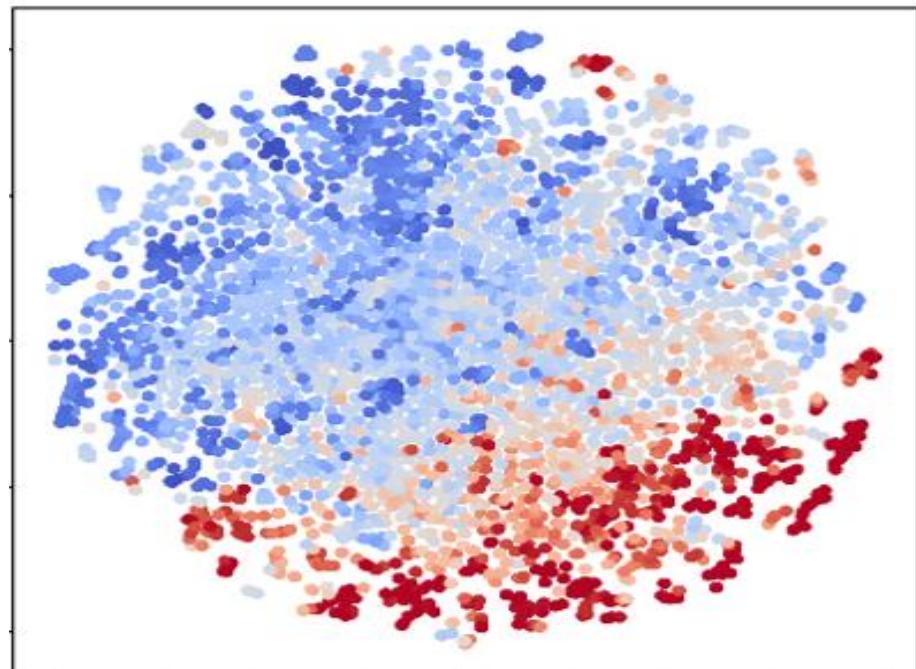


Visualized Feature for SEED-VIG

Raw Feature

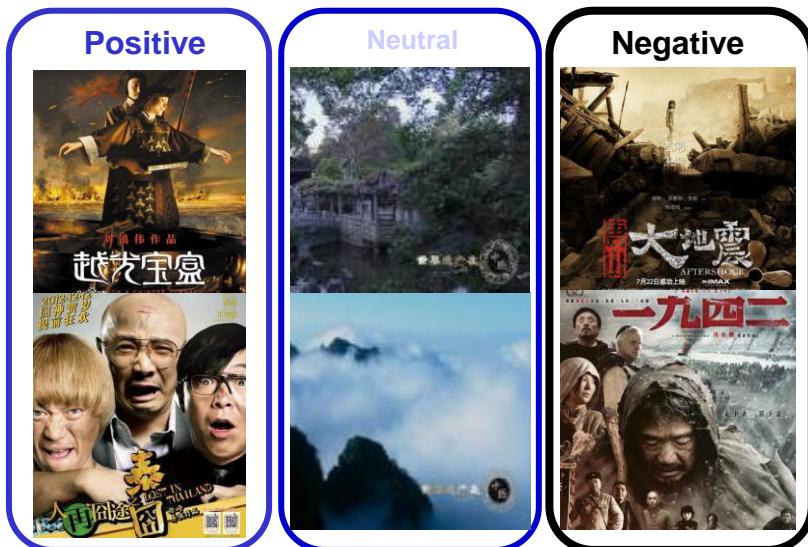


DResNet Feature



Experiment: Emotion Recognition

- 15 Chinese movie clips were used in each experiment, 5 for each emotion (positive, negative, neutral)
- 15 subjects (7 females), each participates in the experiment for three times
- 3394*310 features for each subject



Zheng, W.L., Lu, B.L.: Investigating Critical Frequency Bands and Channels for EEG-based Emotion Recognition with Deep Neural Networks. IEEE Transactions on Autonomous Mental Development 7(3), 162{175 (2015)

Result: Emotion Recognition

□ Leave-one-subject-out cross validation:

- 14 subjects for training
- 1 subject for testing

Base line	Domain Adaptation					Domain Generalization				
	SVM	TCA	TPT	DANN	DAN	WGA NDA	DICA	SCA	DG- DANN	DResNet
Avg	0.5818	0.6400	0.7517	0.7919	0.8381	0.8707	0.6941	0.6633	0.8430	0.8530
Std	0.1385	0.1466	0.1283	0.1314	0.0856	0.0714	0.0779	0.1060	0.0832	0.0797

Luo Y, Zhang S Y, Zheng W L, et al. WGAN Domain Adaptation for EEG-Based Emotion Recognition[C]//International Conference on Neural Information Processing. Springer, Cham, 2018: 275-286.

Li H, Jin Y M, Zheng W L, et al. Cross-Subject Emotion Recognition Using Deep Adaptation Networks[C]//International Conference on Neural Information Processing. Springer, Cham, 2018: 403-413.

Result: Emotion Recognition

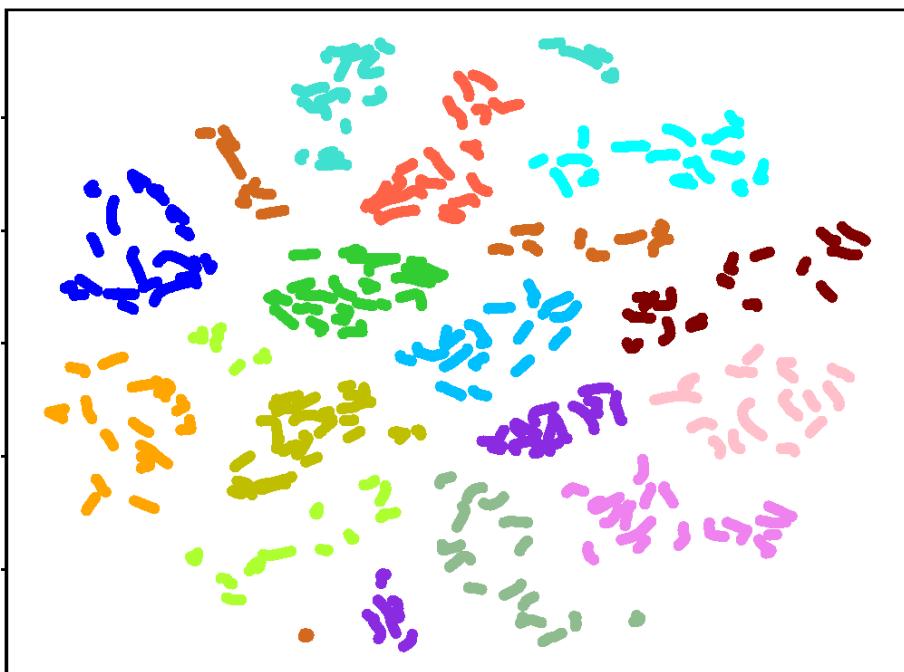
□ Leave-multiple-subject-out evaluation:

- Two-thirds of the subjects (10) for training
- One-third of the subjects (5) for testing

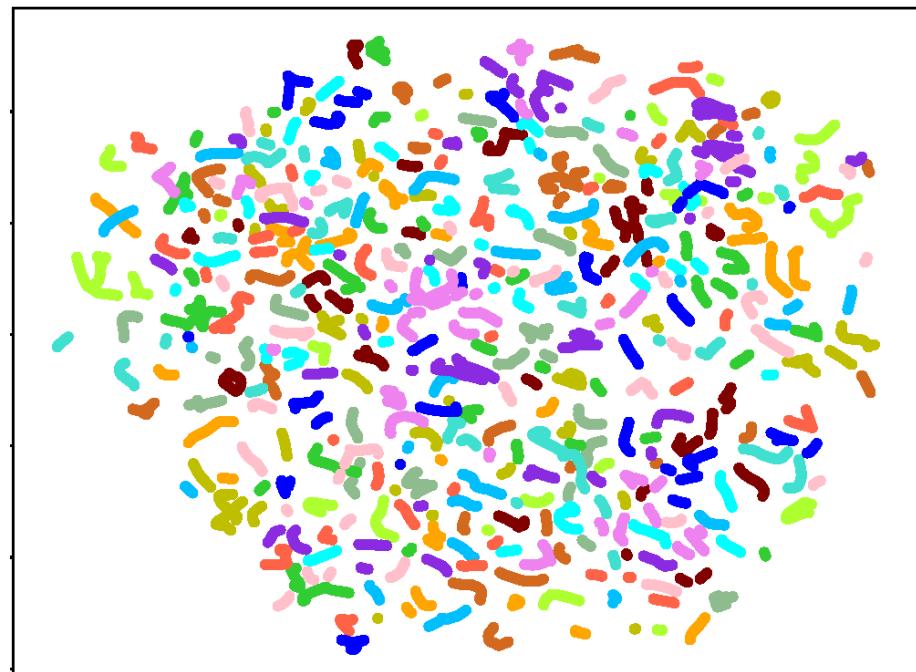
Baseline		Domain Generalization			
	SVM	DICA	SCA	DG-DANN	DResNet
Avg	0.5413	0.6435	0.6083	0.8146	0.8170
Std	0.1348	0.0896	0.0505	0.0788	0.0737

Visualized Feature for SEED

Raw Feature

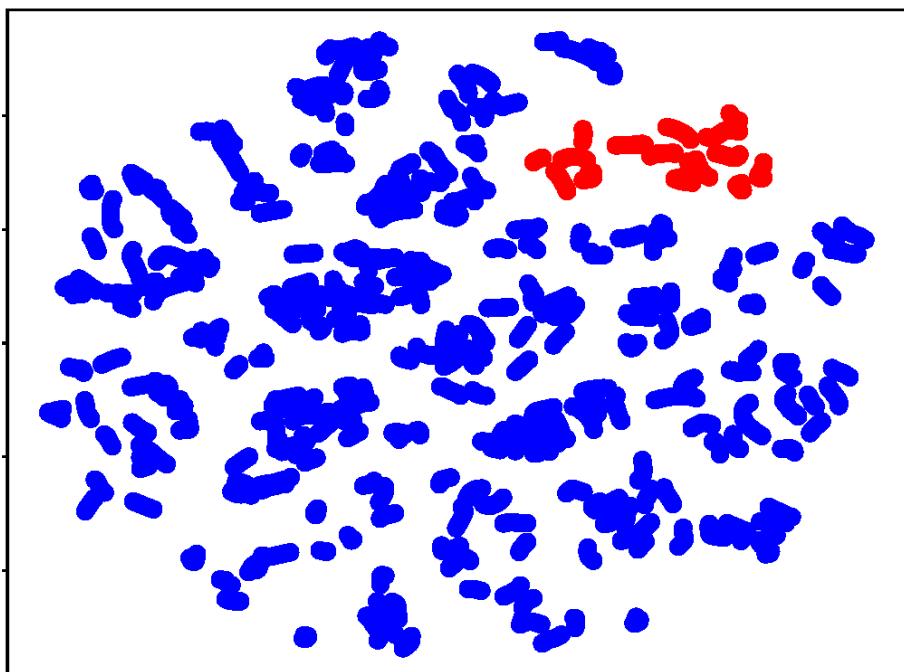


DResNet Feature

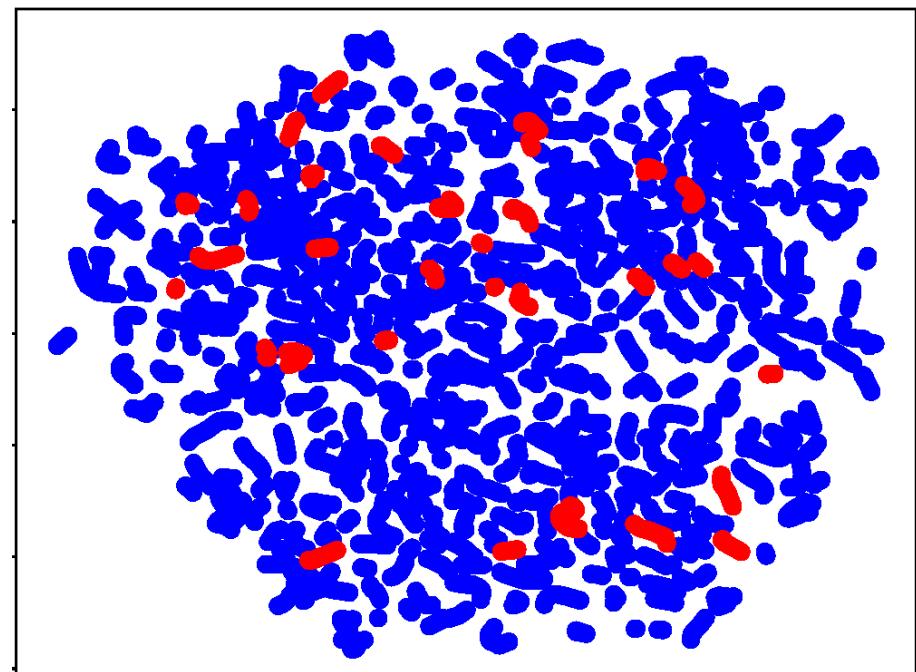


Visualized Feature for SEED

Raw Feature

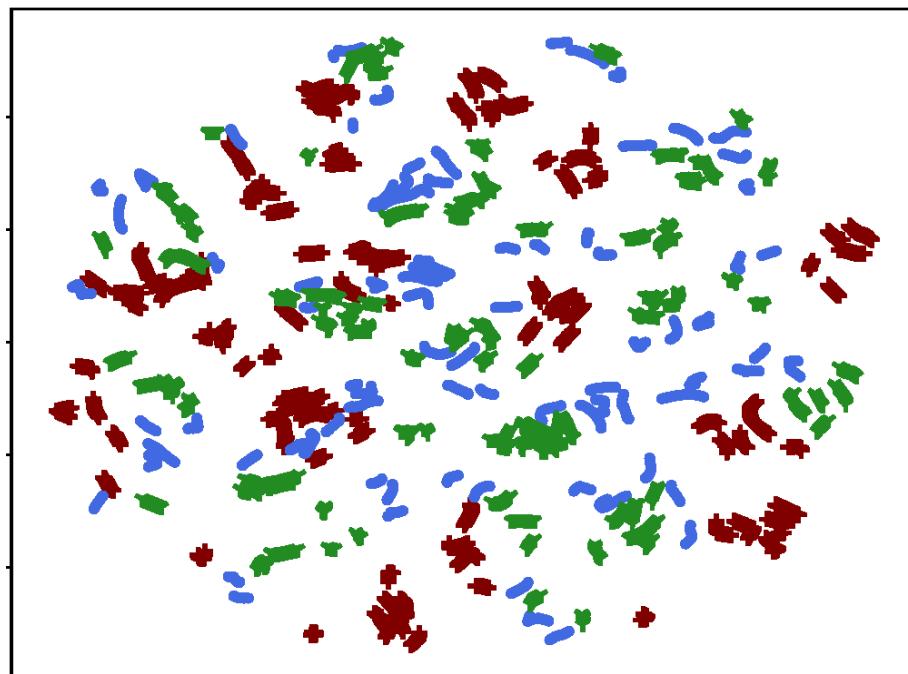


DResNet Feature

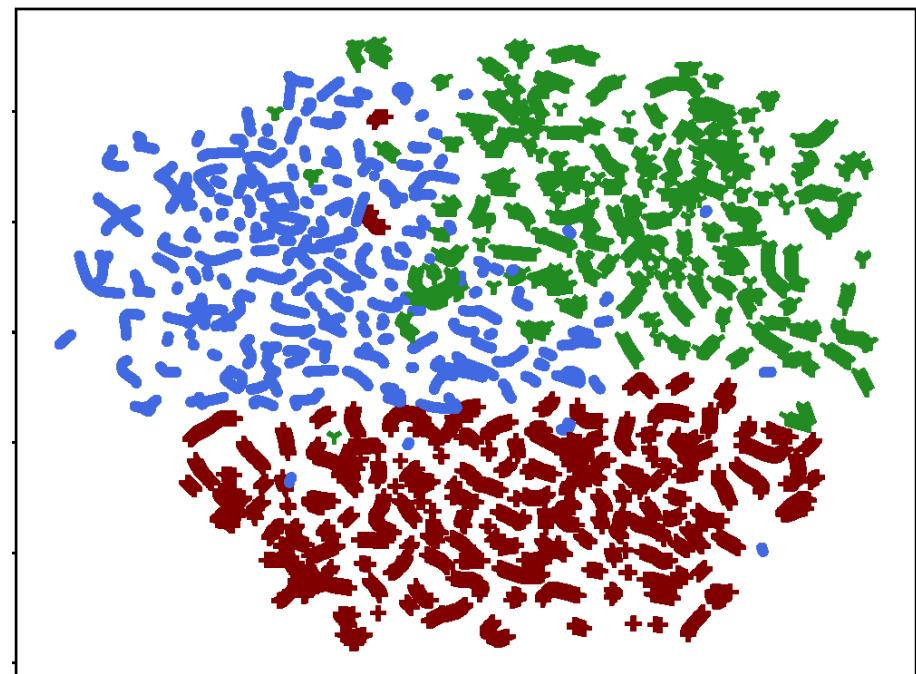


Visualized Feature for SEED

Raw Feature



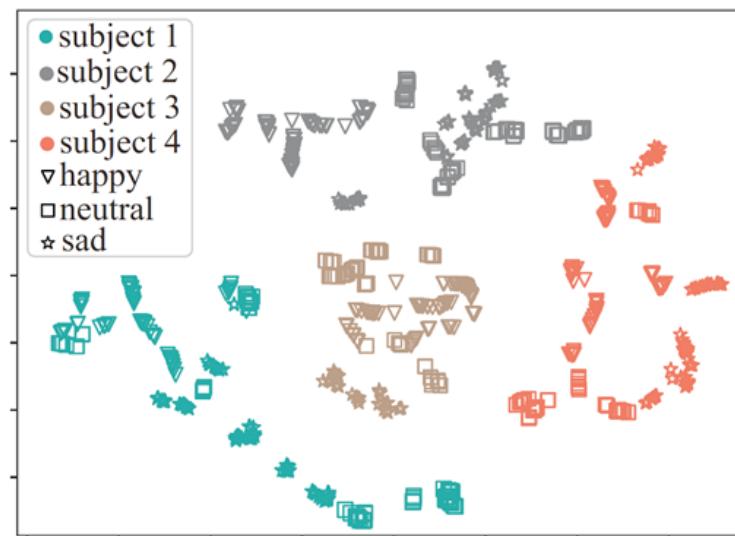
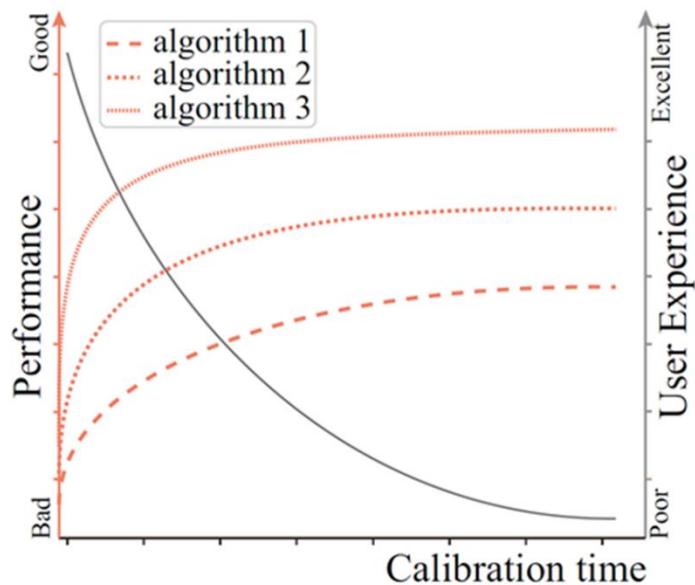
DResNet Feature



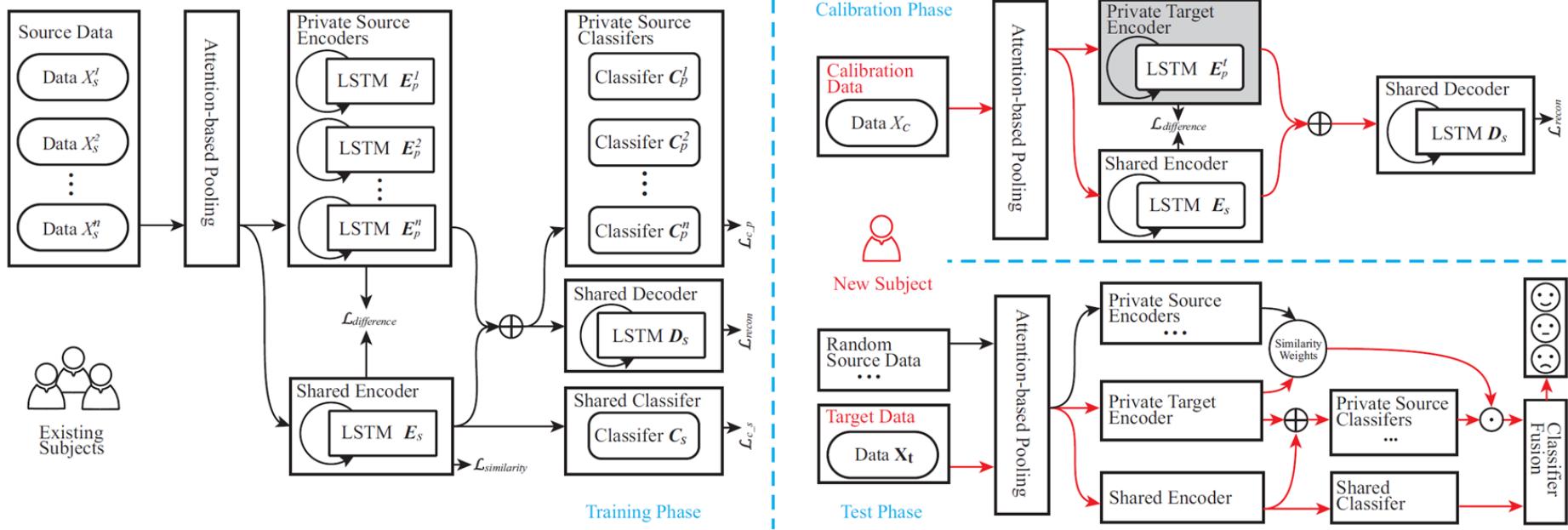
Conclusion

- We focused on reducing the influence of EEG subject variability on BCI systems for unknown subjects by introducing the idea of domain generalization
- Following different approaches to domain generalization, two novel deep adversarial DG models have been proposed
- Evaluations under two settings on public datasets of different topics have indicated that our proposed DG methods are effective for reducing subject variability on depersonalized cross-subject vigilance estimation and emotion recognition problem

跨被试脑电情绪识别模型性能与系统校准时间的平衡



即插即用域适应 (Plug-and-Play Domain Adaptation)

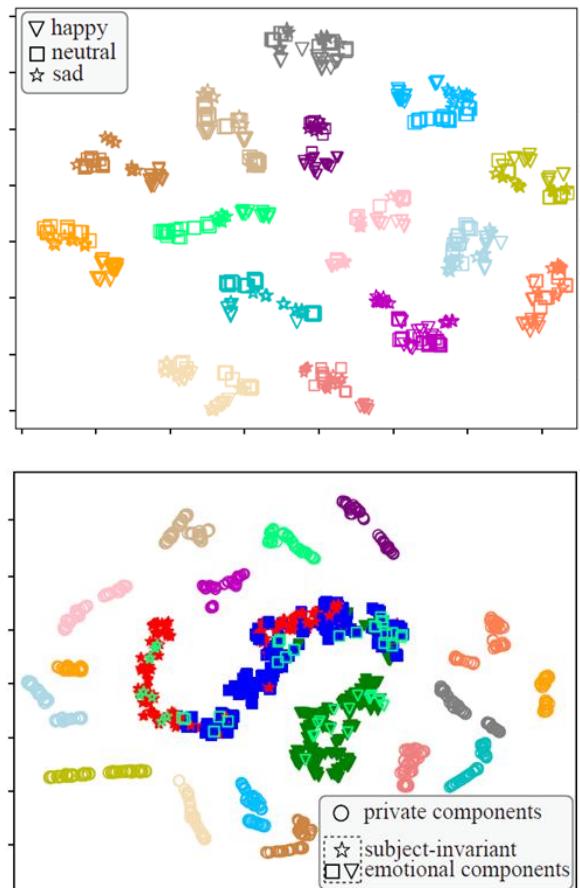


- 1) 模型分为三个部分：训练、校准与测试。
- 2) 注意力机制用于寻找与情绪识别相关的关键脑区与关键频段。

即插即用域适应PPDA的实验结果

Methods	#ATD	Avg.	Std.
Baseline SVM (Zheng and Lu 2016)	None	0.567	0.163
DICA (Ma et al. 2019)	None	0.694	0.078
DResNet (Ma et al. 2019)		0.853	0.080
TCA (Zheng and Lu 2016)		0.640	0.146
TPT (Zheng and Lu 2016)		0.752	0.128
DANN (Li et al. 2018)	All	0.792	0.131
DAN (Li et al. 2018)		0.838	0.086
WGANDA (Luo et al. 2018)		0.871	0.071
PPDA_NC	None	0.854	0.071
PPDA	Few	0.867	0.071

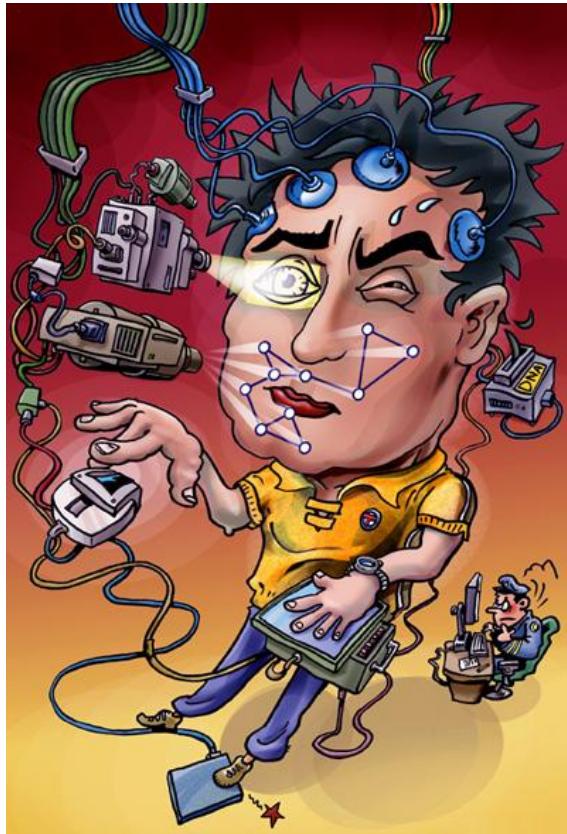
注：在本实验中ALL表示使用53分钟的目标数据，Few表示使用45秒的目标。



Useful Sites on Transfer Learning

- Resources:
- Open source program: <http://www.cse.ust.hk/TL/>
- Qiang Yang: <http://www.cs.ust.hk/~qyang/>
- Sinno Jialin Pan:
<http://www.ntu.edu.sg/home/sinnopan/>
- Wenyuan Dai:
<http://www.4paradigm.com/homepage.html>

- Survey:
- A survey on Transfer Learning.
- Transfer learning for activity recognition: A survey.
- Transitive Transfer Learning.
- Fuzzy Transfer Learning: Methodology and application.



谢谢！下周见！