

# Domain Adversarial Neural Networks (DANN)

Hana Ajakan, Pascal Germain, Hugo Larochelle, Francois  
Laviolette, Mario Marchhand

Département d'informatique et de génie logiciel, Université Laval, Québec, Canada  
Département d'informatique, Université de Sherbrooke, Québec, Canada

[rayjang111@gmail.com](mailto:rayjang111@gmail.com)  
HyunSukJang



# CONTENTS

---

01

Introduction

02

DANN Model  
Architecture

03

Train&test  
method

04

Evaluation



# 01 Introduction

.

# 01. Introduction

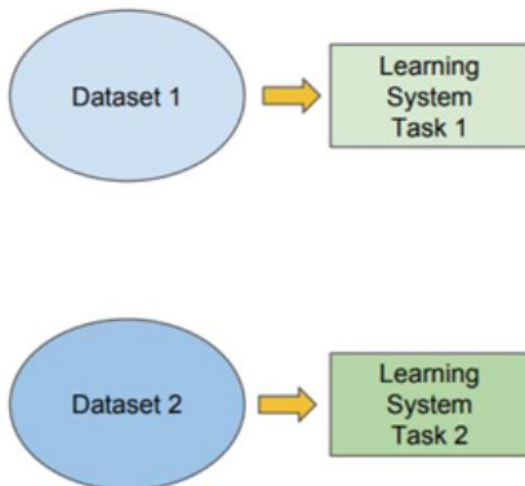
## What is Adaptation learning???

### Traditional ML

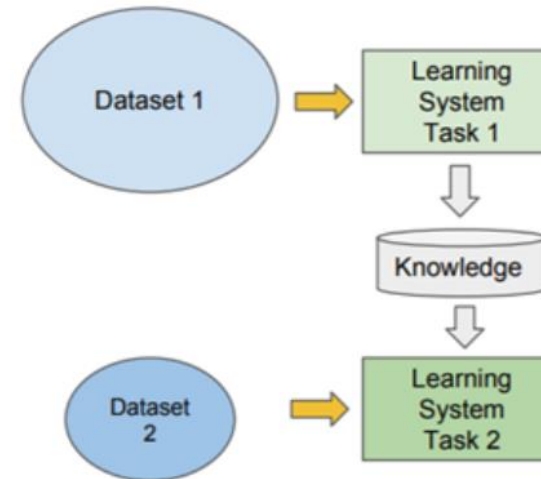
vs

### Transfer Learning

- Isolated, single task learning:
  - Knowledge is not retained or accumulated. Learning is performed w.o. considering past learned knowledge in other tasks

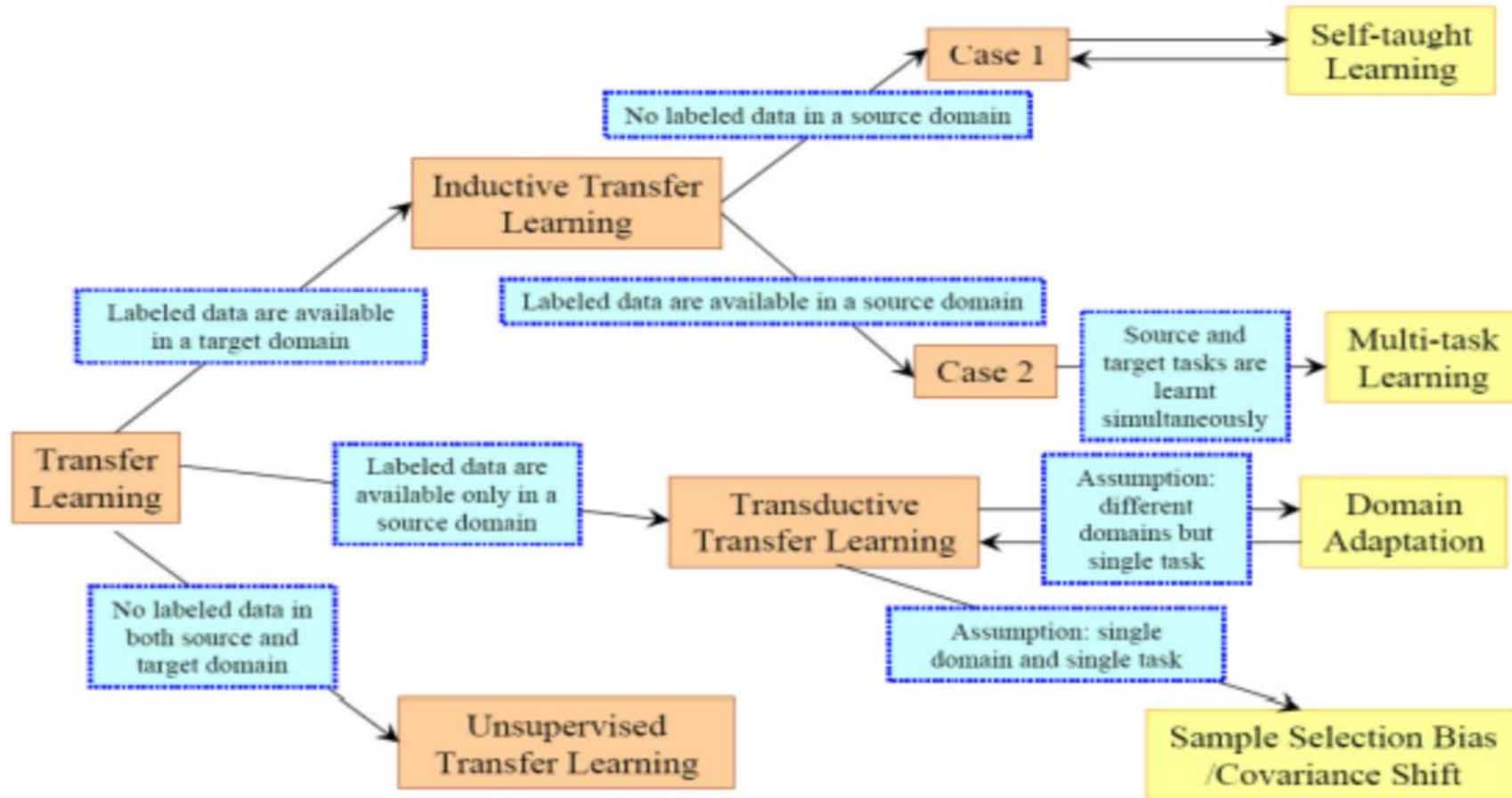


- Learning of a new tasks relies on the previous learned tasks:
  - Learning process can be faster, more accurate and/or need less training data



Traditional Learning vs Transfer Learning

# 01. Introduction



Transfer Learning Strategies

# 01. Introduction

*Source Domain*



*Target Domain*



Adaptation

## 02 DANN Model architecture



## 02. DANN Model architecture

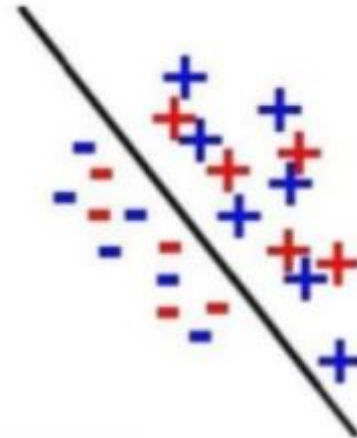
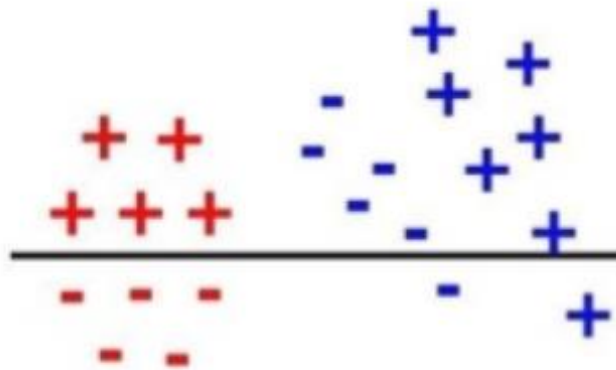
- **Labeled Source** Sample

$S = \{(\mathbf{x}_i, y_i)\}_{i=1}^{m_s}$  **Source** sample drawn i.i.d. from  $P_S$

- **Unlabeled Target** Sample

$T = \{\mathbf{x}_j\}_{j=1}^{m_t}$  **Target** Sample drawn i.i.d. from  $D_T$   
optionnal: a few labeled target examples

If  $h$  is learned from **source** domain, how does it perform on **target** domain?





## 02. DANN Model architecture

### Reweighting/Instance-based methods

Correct a sample bias by reweighting source labeled data: source instances close to target instances are more important.



### Feature-based methods/Find new representation spaces

Find a common space where source and target are close (projection, new features, etc)

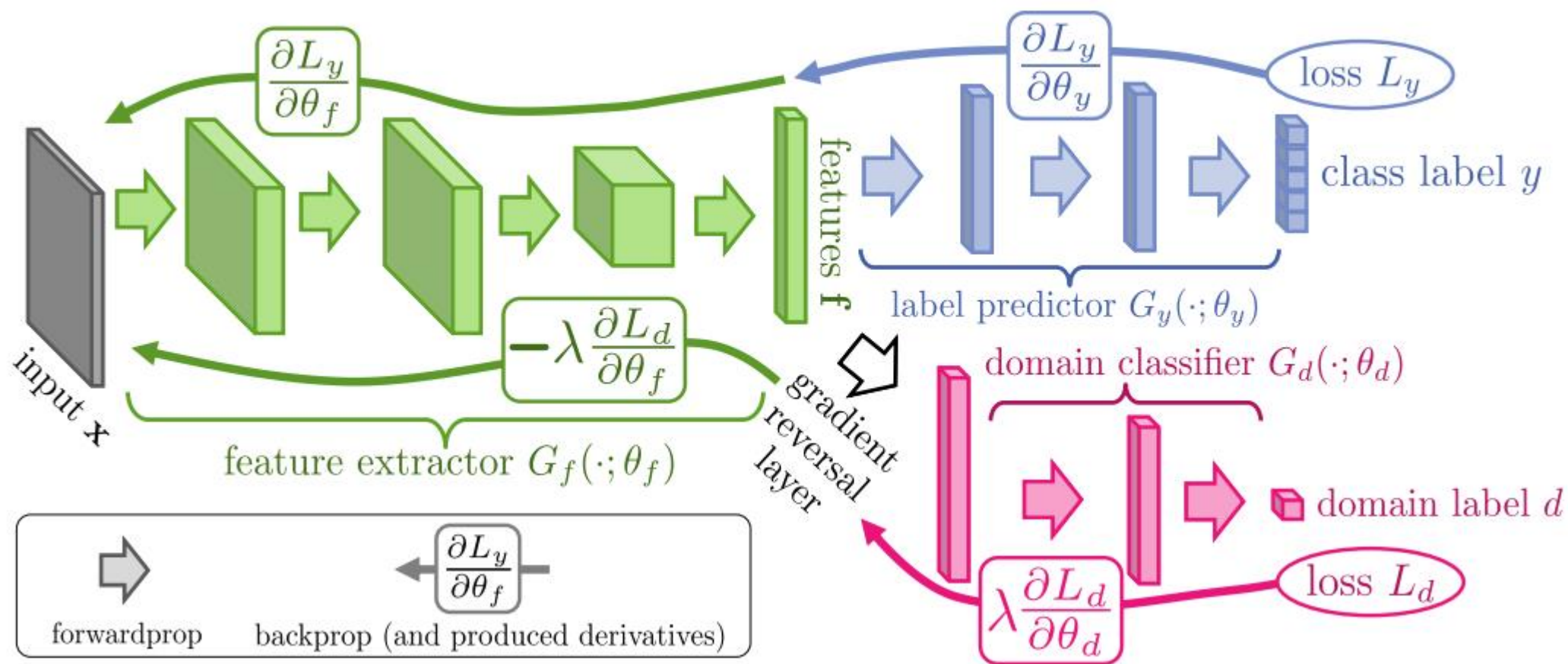


### Ajustement/Iterative methods

Modify the model by incorporating pseudo-labeled information



## 02. DANN Model architecture

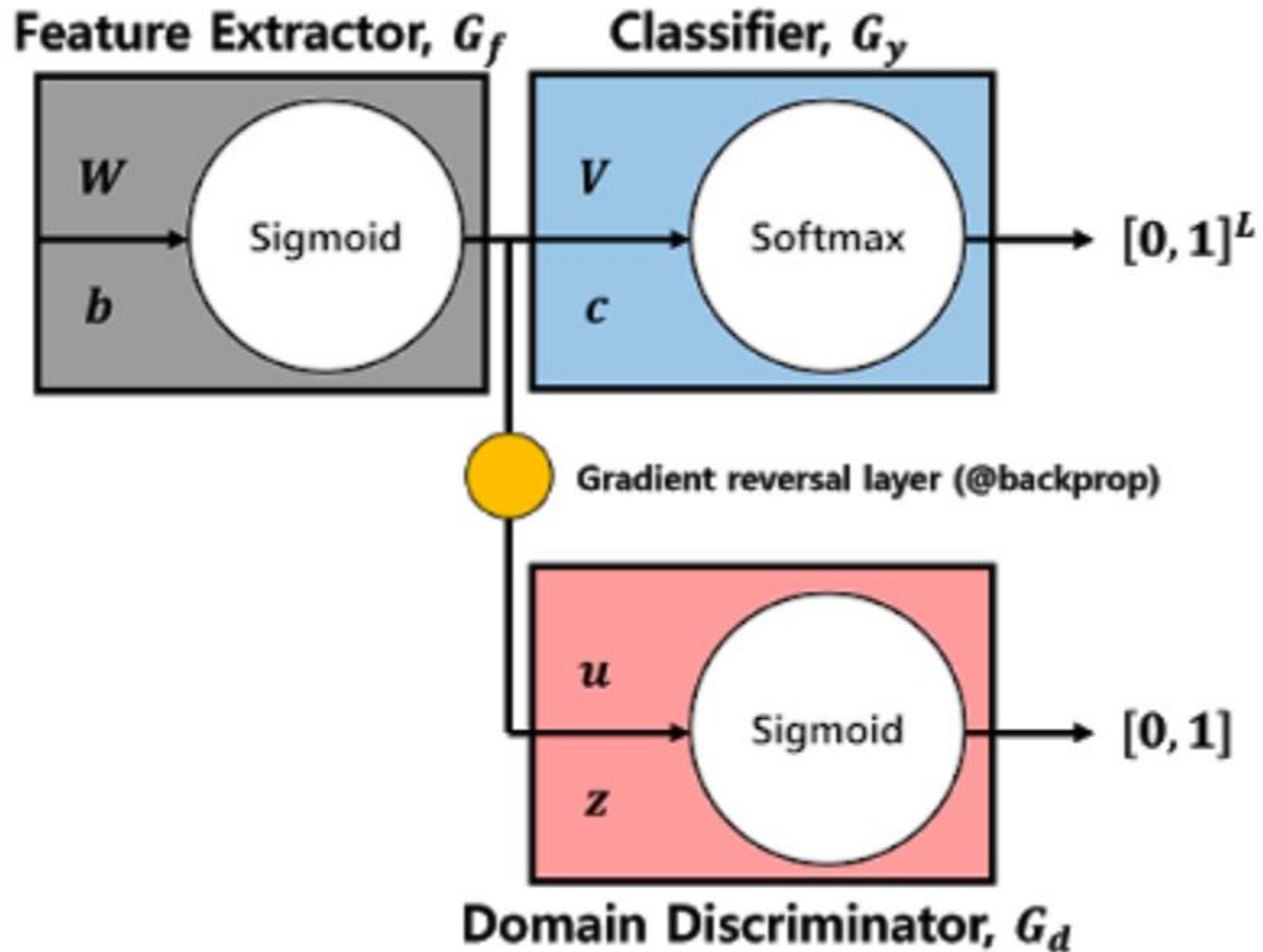


## 02. DANN Model architecture

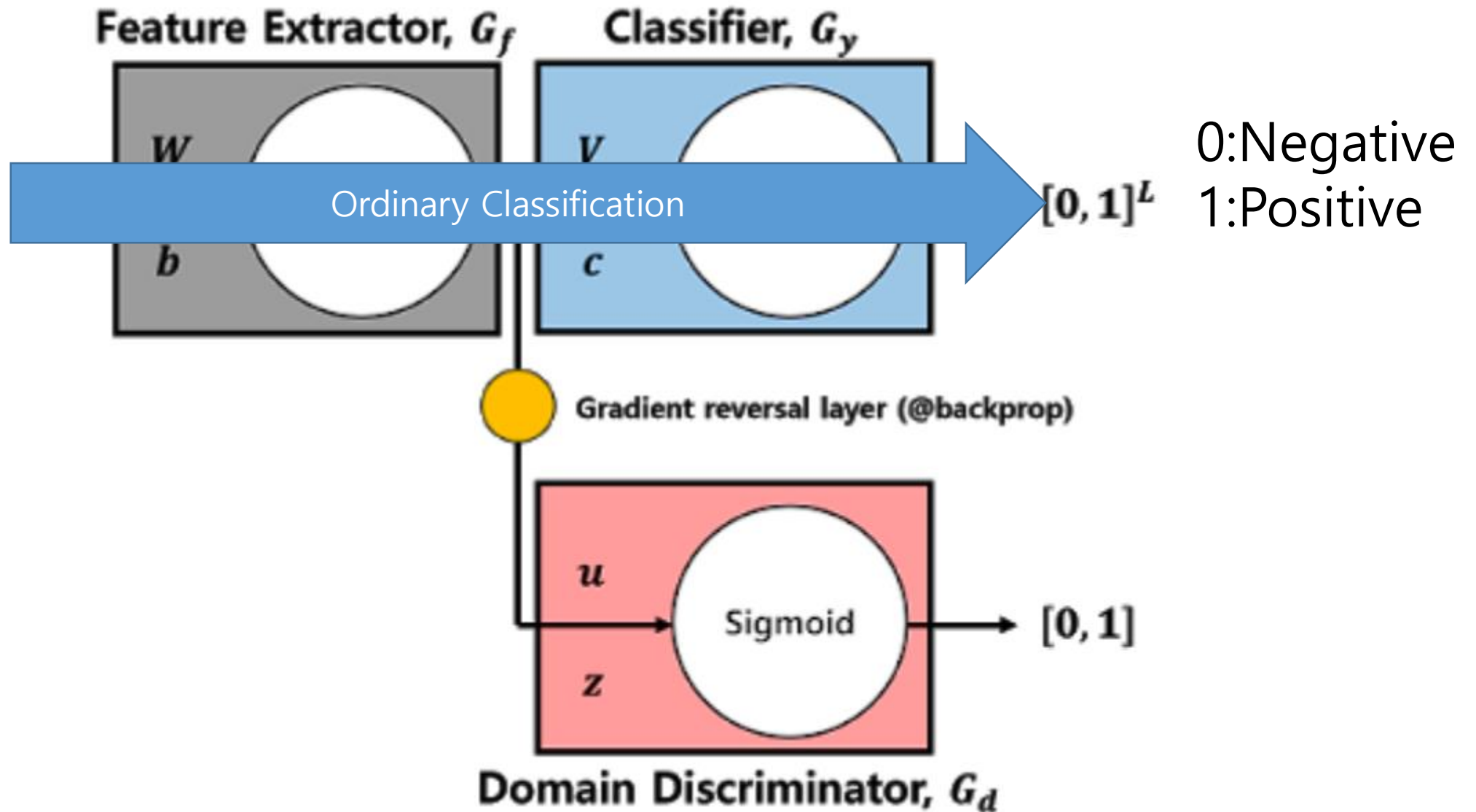
	Electronics	Video games
✓	(1) <u>Compact</u> ; easy to operate; very good picture quality; looks <u>sharp</u> !	(2) A very <u>good</u> game! It is action packed and <u>full</u> of excitement. I am very much <u>hooked</u> on this game.
✓	(3) I purchased this unit from Circuit City and I was very <u>excited</u> about the quality of the picture. It is really <u>nice</u> and <u>sharp</u> .	(4) Very <u>realistic</u> shooting action and good plots. We played this and were <u>hooked</u> .
✗	(5) It is also quite <u>blurry</u> in very dark settings. I will <u>never_buy</u> HP again.	(6) It is so boring. I am extremely <u>unhappy</u> and will probably <u>never_buy</u> UbiSoft again.

- Source specific: *compact, sharp, blurry*.
- Target specific: *hooked, realistic, boring*.
- Domain independent: *good, excited, nice, never\_buy, unhappy*.

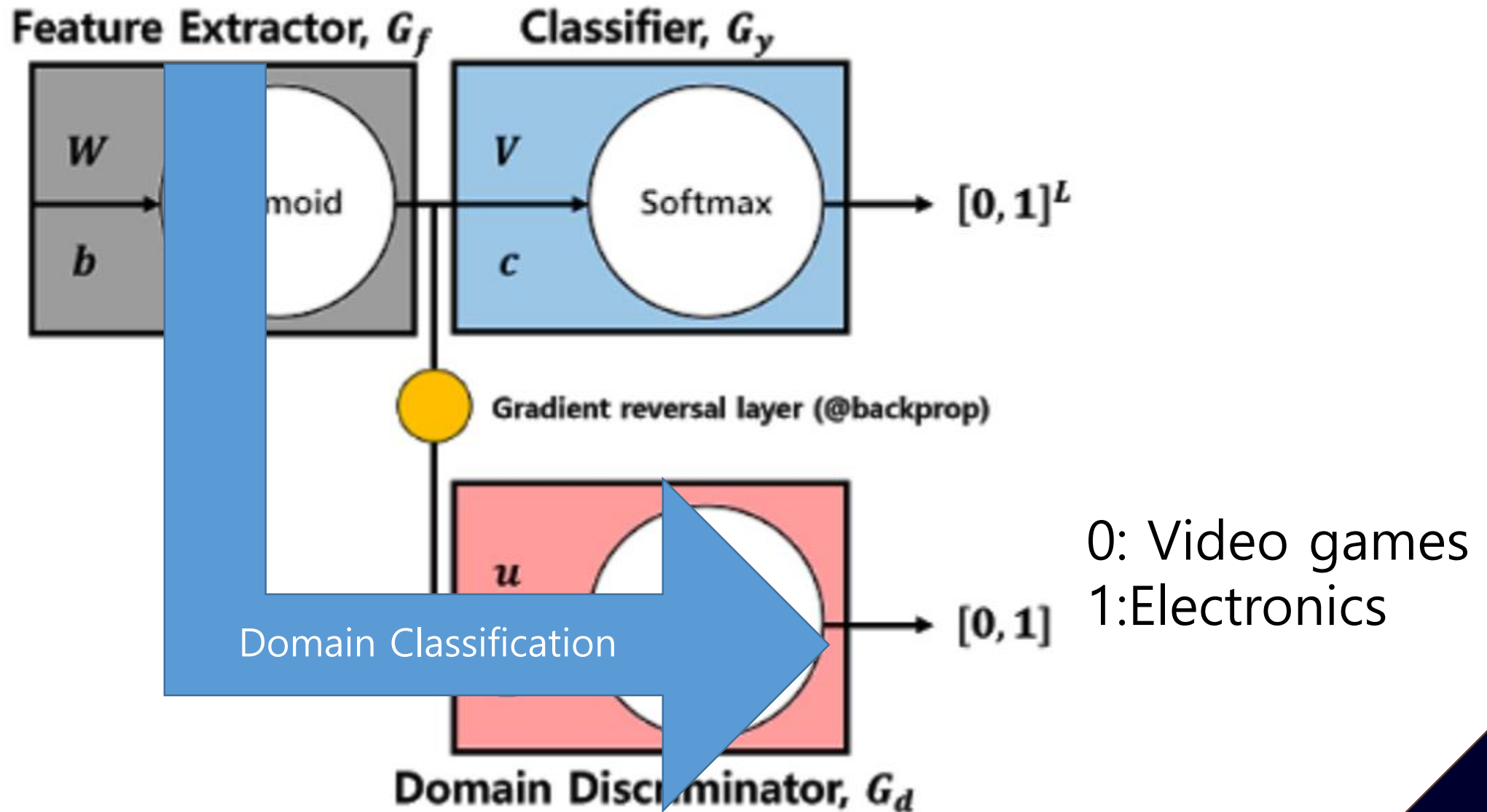
## 02. DANN Model architecture



## 02. DANN Model architecture



## 02. DANN Model architecture





## 02. DANN Model architecture

By minimizing classification error of source data & divergence between source domain and target domain

=> Features will be extracted to classify both source domain and target domain

# 03training&test method





### 03. training&test method

#### Classical Test Error:

$$\epsilon_{\text{test}} \leq \hat{\epsilon}_{\text{train}} + \sqrt{\frac{\text{complexity}}{n}}$$

Measured on the  
same distribution!

#### Adaptation Target Error:

$$\epsilon_{\text{test}} \leq ??$$

Measured on a  
**new** distribution!

### 03. training&test method

#### Theorem 1.

(Ben-David et al., 2006) Let  $\mathcal{H}$  is a hypothesis class of VC dimension  $d$ . With probability  $1 - \delta$  over the choice of samples  $S \sim (\mathbb{D}_S)^n$  and  $T \sim (\mathbb{D}_T^X)^n$ , for every  $\eta \in \mathcal{H}$ :

$$R_{\mathbb{D}_T}(\eta) \leq R_S(\eta) + \sqrt{\frac{4}{m} \left( d \log \frac{2em}{d} + \log \frac{4}{\delta} \right)} \quad (3)$$

$$+ \hat{d}_{\mathcal{H}}(S, T) + 4\sqrt{\frac{1}{m} \left( d \log \frac{2m}{d} + \log 4\delta \right)} + \beta, \quad (4)$$

with  $\beta \geq \inf_{\eta^* \in \mathcal{H}} [R_{\mathbb{D}_S}(\eta^*) + R_{\mathbb{D}_T}(\eta^*)]$ , and

$$R_S(\eta) = \frac{1}{m} \sum_{i=1}^m I[\eta(x_i) \neq y_i]$$

is the empirical source risk.

### 03. training&test method

#### Definition 1.

(Ben-David et al., 2006, 2010; Kifer et al., 2004) Given two domain distributions  $\mathbb{D}_S^X$  and  $\mathbb{D}_T^X$  over  $X$ , and a hypothesis class  $\mathcal{H}$ , the  $\mathcal{H}$ -divergence between  $\mathbb{D}_S^X$  and  $\mathbb{D}_T^X$  is

$$\hat{h}_{\mathcal{H}}(\mathbb{D}_S^X, \mathbb{D}_T^X) = 2 \sup_{\eta \in \mathcal{H}} \left| \Pr_{x \sim \mathbb{D}_S^X} [\eta(x_i) = 1] - \Pr_{x \sim \mathbb{D}_T^X} [\eta(x_i) = 1] \right|.$$

### 03. training&test method

$$\begin{aligned}\hat{h}_{\mathcal{H}}(\mathbb{D}_S^X, \mathbb{D}_T^X) &= 2 \sup_{\eta \in \mathcal{H}} \left| \Pr_{x \sim \mathbb{D}_S^X} [\eta(x_i) = 1] - \Pr_{x \sim \mathbb{D}_T^X} [\eta(x_i) = 1] \right| \\ &= 2 \sup_{\eta \in \mathcal{H}} \left| \Pr_{x \sim \mathbb{D}_S^X} [\eta(x_i) = 1] + \Pr_{x \sim \mathbb{D}_T^X} [\eta(x_i) = 0] - 1 \right|\end{aligned}$$



$$\hat{d}_{\mathcal{H}}(S, T) = 2 \left( 1 - \min_{\eta \in \mathcal{H}} \left[ \frac{1}{n} \sum_{i=1}^n I[\eta(x_i) = 1] + \frac{1}{n'} \sum_{i=n+1}^N I[\eta(x_i) = 0] \right] \right),$$

$$\hat{h}_{\mathcal{H}}(S, T) = 2(1 - 2\epsilon)$$



### 03. training&test method

$$E(W, V, b, c, u, z) \\ = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_y^i(W, b, V, c) - \lambda \left( \frac{1}{n} \sum_{i=1}^n \mathcal{L}_d^i(W, b, u, z) + \frac{1}{n'} \sum_{i=n+1}^N \mathcal{L}_d^i(W, b, u, z) \right),$$

where we are seeking the parameters  $\hat{W}, \hat{V}, \hat{b}, \hat{c}, \hat{u}, \hat{z}$  that deliver a saddle point given by

$$(\hat{W}, \hat{V}, \hat{b}, \hat{c}) = \arg \min_{W, V, b, c} E(W, V, b, c, \hat{u}, \hat{z}), \quad (1)$$

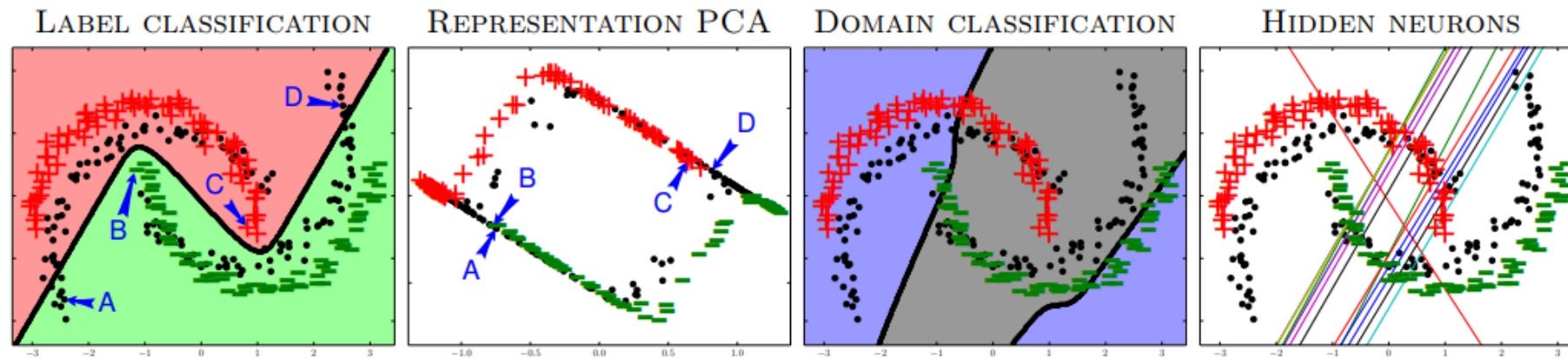
$$(\hat{u}, \hat{z}) = \arg \max_{\hat{u}, \hat{z}} E(\hat{W}, \hat{V}, \hat{b}, \hat{c}, u, z). \quad (2)$$

# 04 evaluation

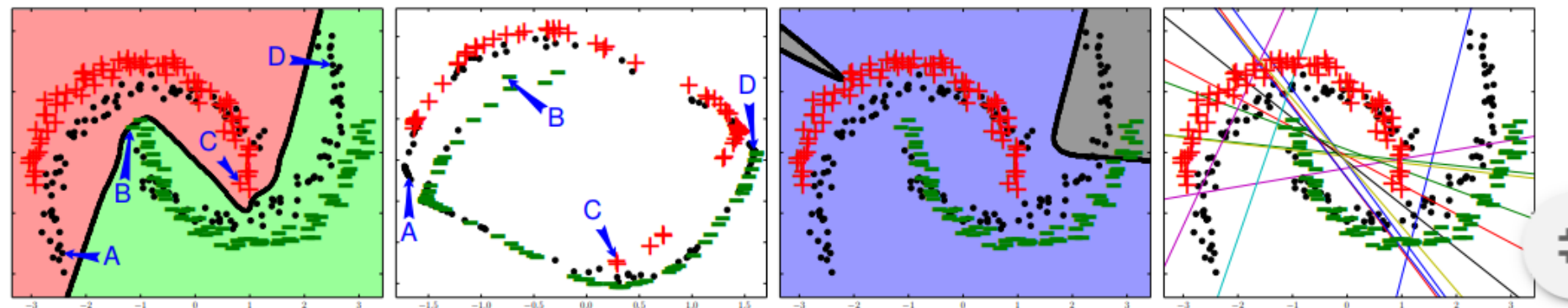


## 04. evaluation

Toy problem



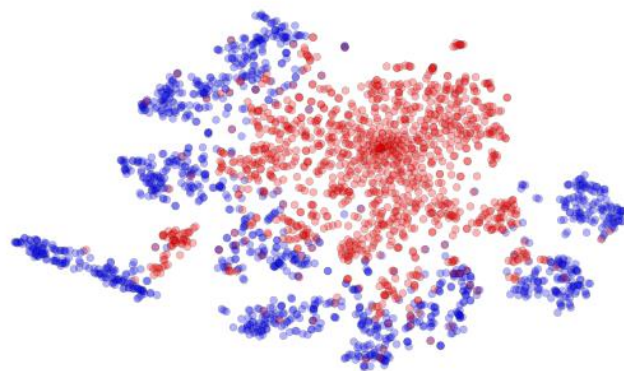
(a) Standard NN. For the “domain classification”, we use a *non adversarial* domain regressor on the hidden neurons learned by the Standard NN. (This is equivalent to run Algorithm 1, without Lines 22 and 31)



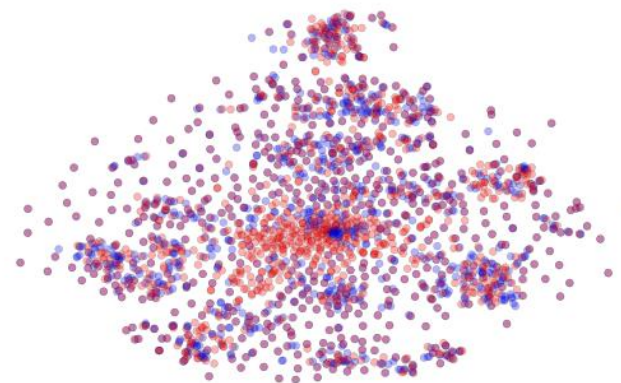
(b) DANN (Algorithm 1)

## 04. evaluation

MNIST  $\rightarrow$  MNIST-M: top feature extractor layer

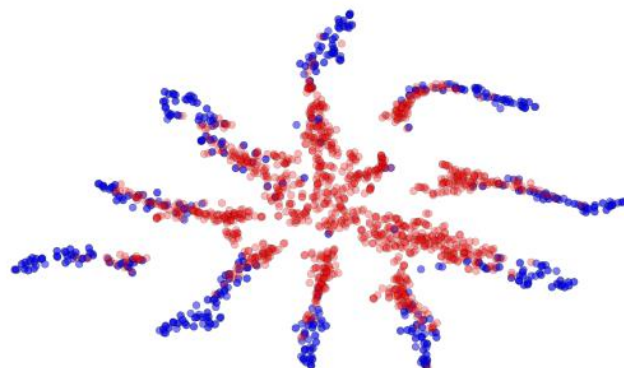


(a) Non-adapted

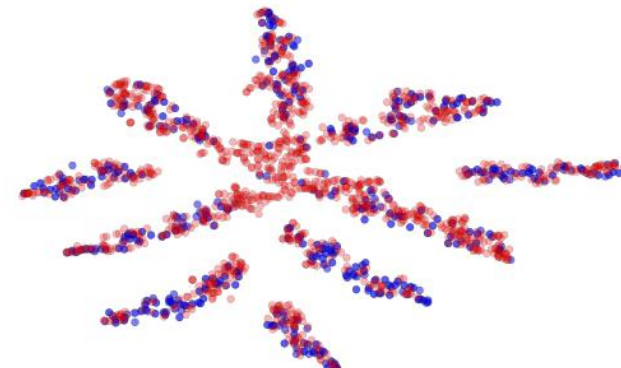


(b) Adapted

SYN NUMBERS  $\rightarrow$  SVHN: last hidden layer of the label predictor



(a) Non-adapted



(b) Adapted



## 04. evaluation

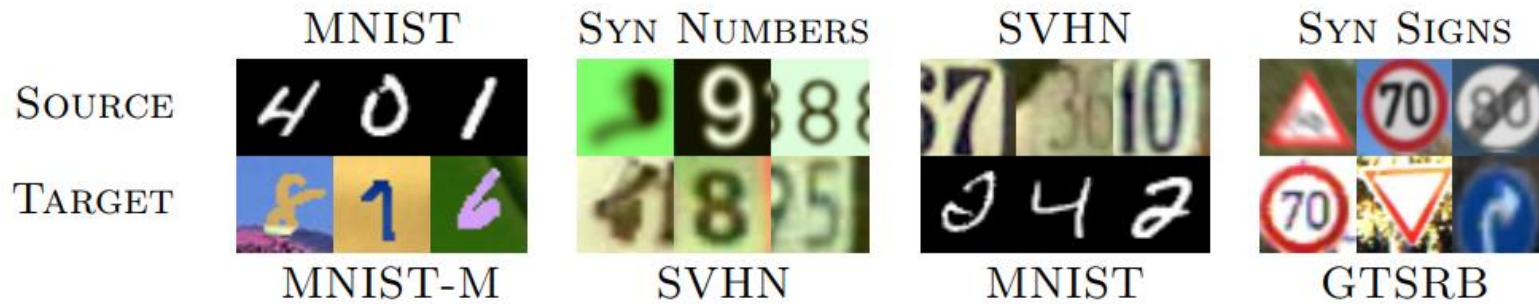


Figure 6: Examples of domain pairs used in the experiments. See Section 5.2.4 for details.

METHOD	SOURCE	MNIST	SYN NUMBERS	SVHN	SYN SIGNS
	TARGET	MNIST-M	SVHN	MNIST	GTSRB
SOURCE ONLY		.5225	.8674	.5490	.7900
SA (Fernando et al., 2013)		.5690 (4.1%)	.8644 (−5.5%)	.5932 (9.9%)	.8165 (12.7%)
DANN		<b>.7666</b> (52.9%)	<b>.9109</b> (79.7%)	<b>.7385</b> (42.6%)	<b>.8865</b> (46.4%)
TRAIN ON TARGET		.9596	.9220	.9942	.9980

## 04. evaluation

SOURCE	TARGET	Original data			mSDA representation		
		DANN	NN	SVM	DANN	NN	SVM
BOOKS	DVD	.784	.790	<b>.799</b>	.829	.824	<b>.830</b>
BOOKS	ELECTRONICS	.733	.747	<b>.748</b>	<b>.804</b>	.770	.766
BOOKS	KITCHEN	<b>.779</b>	.778	.769	<b>.843</b>	.842	.821
DVD	BOOKS	.723	.720	<b>.743</b>	.825	.823	<b>.826</b>
DVD	ELECTRONICS	<b>.754</b>	.732	.748	<b>.809</b>	.768	.739
DVD	KITCHEN	<b>.783</b>	.778	.746	.849	<b>.853</b>	.842
ELECTRONICS	BOOKS	<b>.713</b>	.709	.705	<b>.774</b>	.770	.762
ELECTRONICS	DVD	<b>.738</b>	.733	.726	<b>.781</b>	.759	.770
ELECTRONICS	KITCHEN	<b>.854</b>	<b>.854</b>	.847	.881	<b>.863</b>	.847
KITCHEN	BOOKS	<b>.709</b>	.708	.707	.718	.721	<b>.769</b>
KITCHEN	DVD	<b>.740</b>	.739	.736	<b>.789</b>	<b>.789</b>	.788
KITCHEN	ELECTRONICS	<b>.843</b>	.841	.842	.856	.850	<b>.861</b>

(a) Classification accuracy on the Amazon reviews data set



# Thank you

