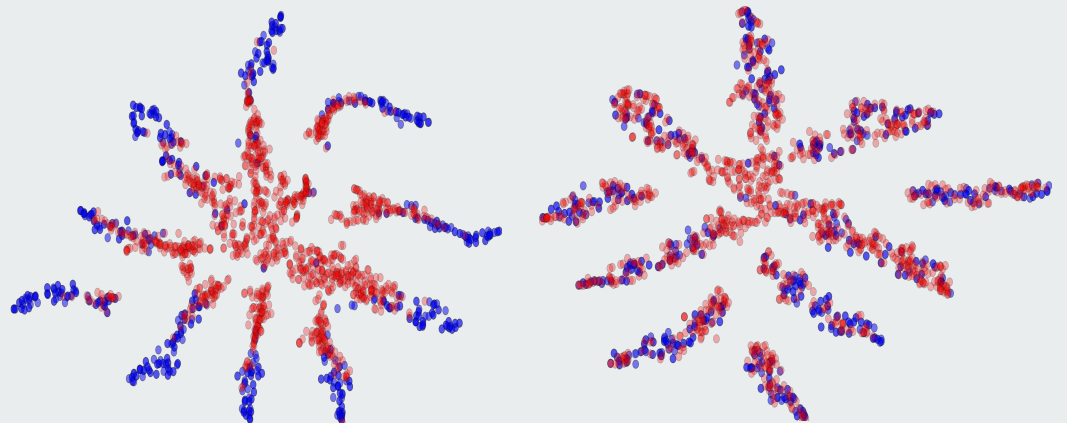




Unsupervised Domain Adaptation by Backpropagation

Chih-Hui Ho, Xingyu Gu, Yuan Qi



Outline



- Introduction
- Related works
- Proposed solution
- Experiments
- Conclusions

Problems



Deep network: requires massive **labeled** training data.

Labeled data:

- **Available** sometimes:
 - Image recognition
 - Speech recognition
 - Recommendation
- **Difficult to collect** sometimes:
 - Robotics
 - Disaster
 - Medical diagnosis
 - Bioinformatics

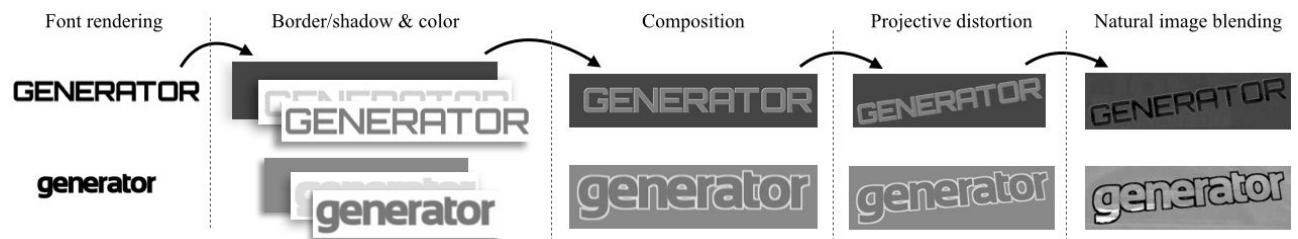
Problems

Test time failure: distribution of actual data is different from training data.

Example: Model is

- Trained on synthetic data (abundant and fully labeled), but
- Tested on real data.

MJSynth
(synthetic)



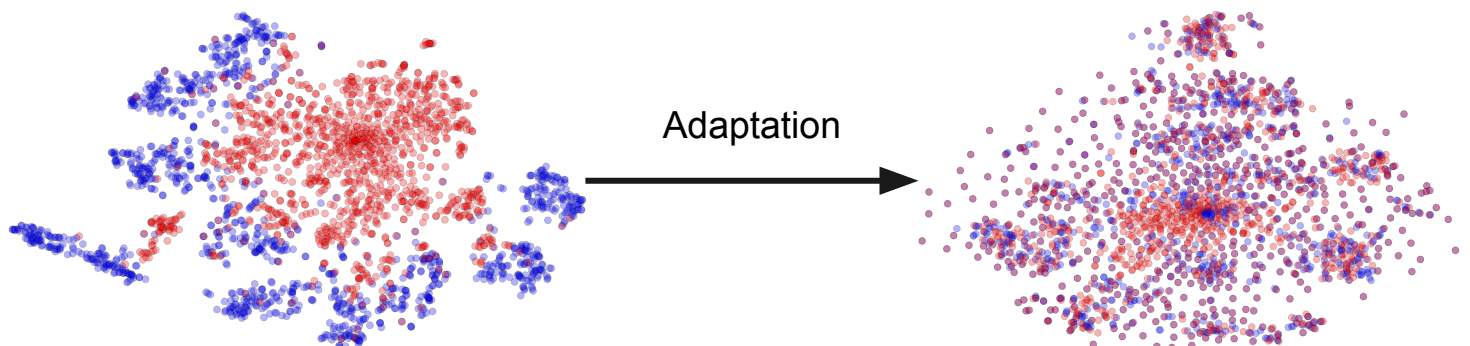
IIIT5K
(real)



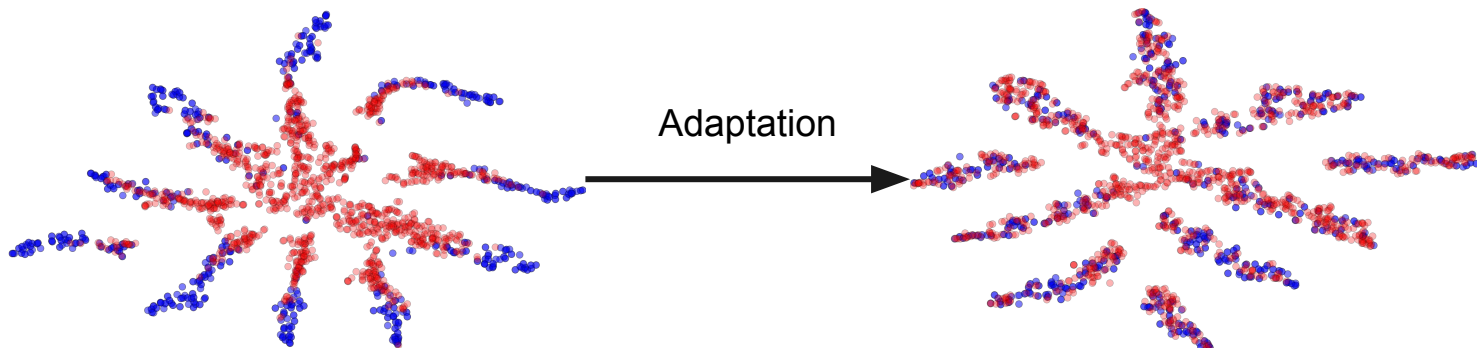
Results

- Source datapoint
- Target datapoint

- MNIST → MNIST-M (extracted features)



- SYN NUMBERS → SVHN (label classifier's last hidden layer)



Objective



Given:

- Lots of **labeled** data in the **source** domain (e.g. synthetic images)
- Lots of **unlabeled** data in the **target** domain (e.g. real images)

Domain Adaptation (DA):

In the presence of a *shift* between source and target domain,
Train a network on **source** domain that performs well on **target** domain.

Objective

Example: Office dataset

- **Source:**
Amazon photos of office objects
(on white background)
- **Target:**
Consumer photos of office objects
(taken by DSLR camera / webcam)



Amazon



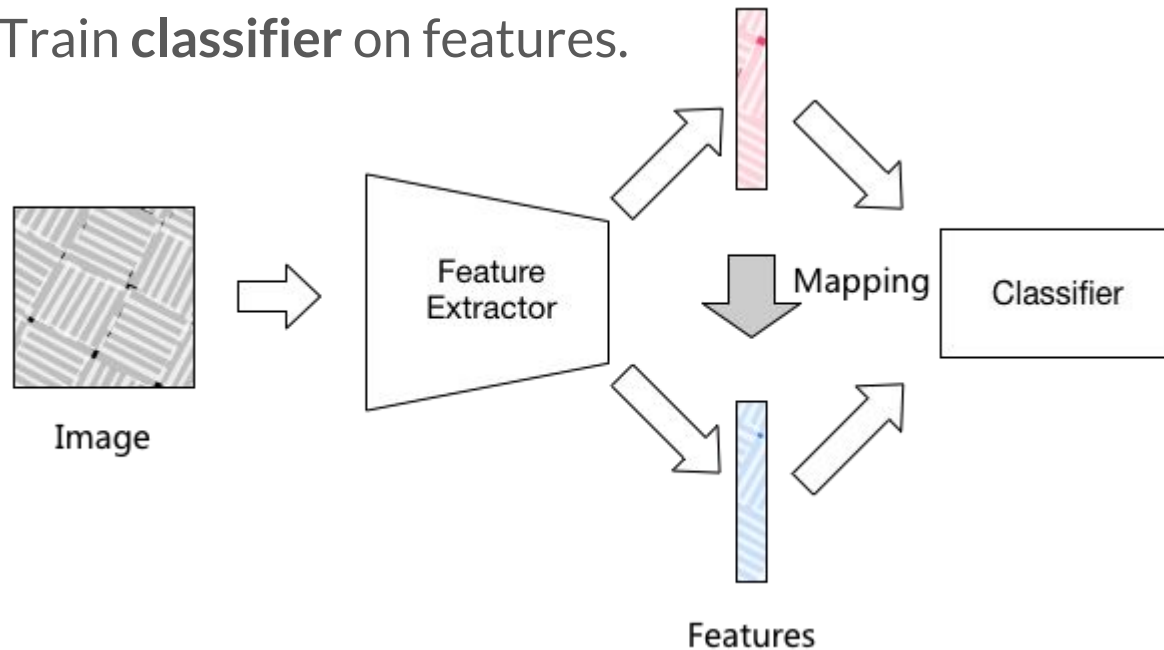
DSLR

Webcam

Previous Approaches - DLID

Deep Learning by Interpolating between Domains

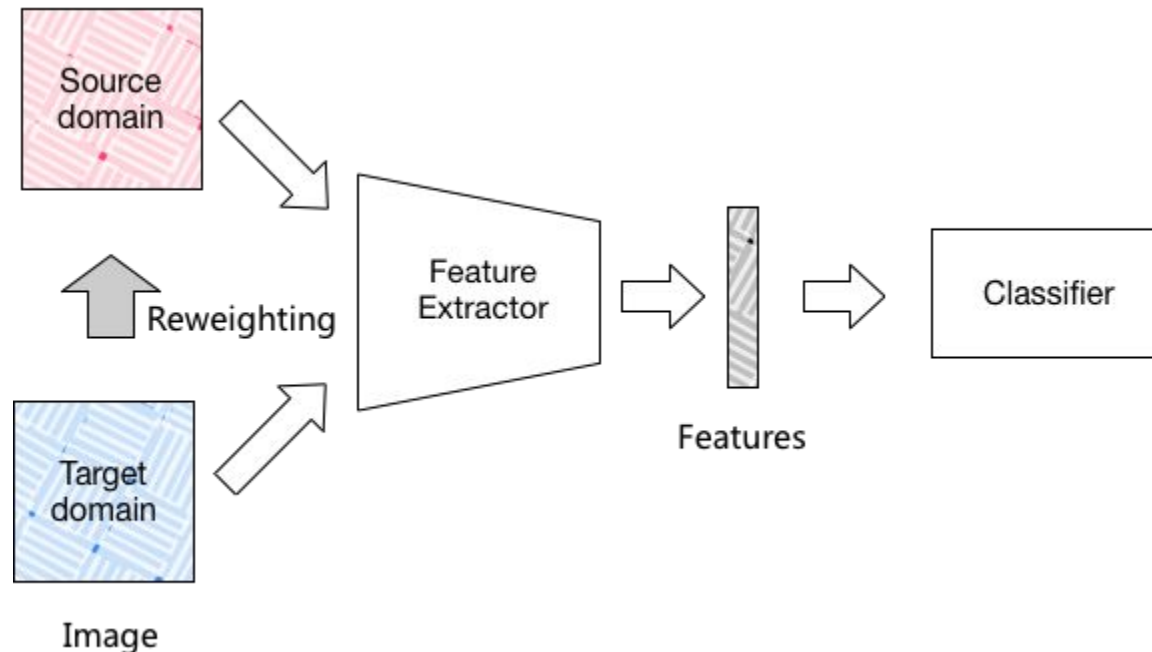
- Feature transformation mapping **source** into **target**.
 - Train feature extractor layer-wise.
 - Gradually replacing **source** samples with **target** samples.
 - Train classifier on features.



Previous Approaches - MMD

Maximum Mean Discrepancy (measures domain-distance)

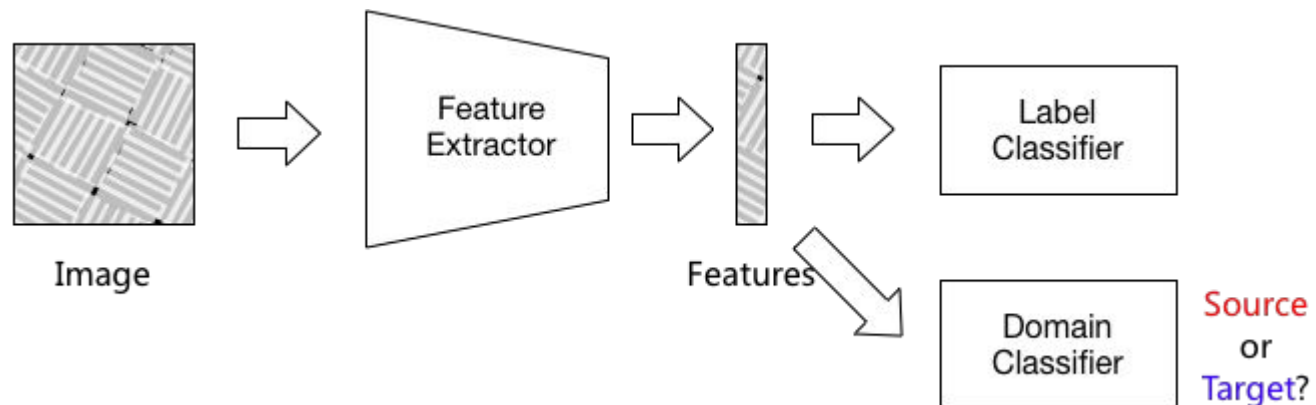
- Reweighting **target** domain images.
 - Distance between **source** and **target** distributions.
 - Explicit distance measurement (e.g. kernel Hilbert space).



Proposed Solution - Deep Domain Adaptation (DDA)

Standard CNN + **domain classifier**.

- An **implicit** way to measure similarity between **source** and **target**.
 - If domain classifier performs **good**: **dissimilar** features.
 - If domain classifier performs **bad**: **similar** features.
- Objective: feature is **best** for label classifier, and **worst** for domain classifier.



Improvement



	Previous approaches	Proposed solution
Measurement of similarity between domains	Explicit (distance in Hilbert space)	Implicit (performance of domain classifier)
Training steps	Separate feature extractor and label classifier	Jointly trained by backpropagation
Architecture	Complicated	Simple (standard CNN + domain classifier)

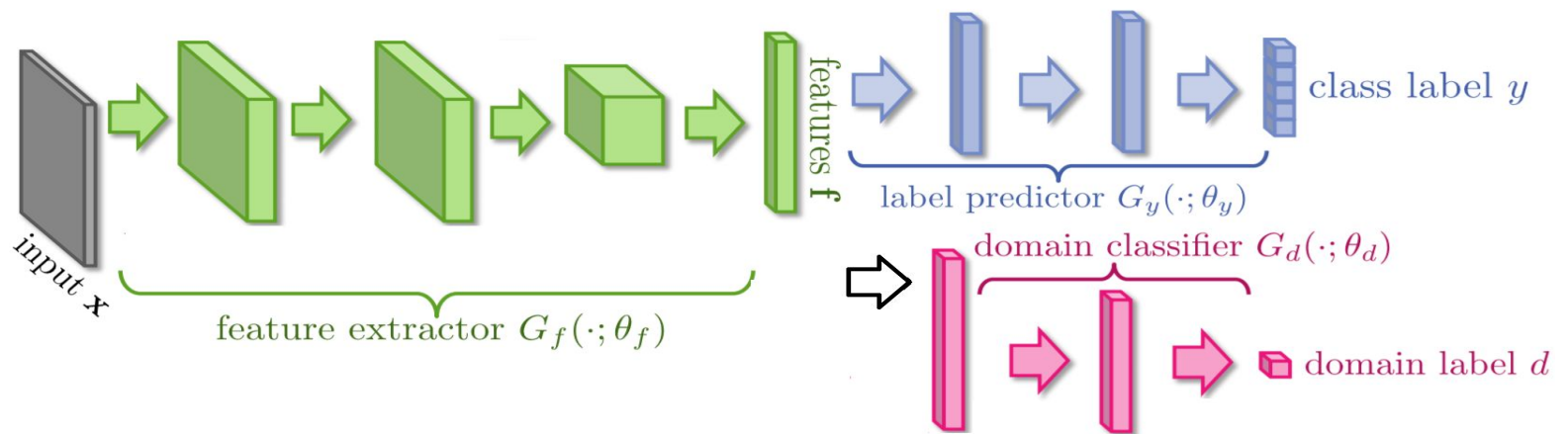
Proposed Solution



- Notation
 - x_i : training samples (from both source and target domain)
 - y_i : class label (only source domain has labels)
 - d_i : 0 (source domain) or 1 (target domain)
- x_i in source domain has $d_i = 0$ and y_i
- x_i in target domain has $d_i = 1$ (target domain has no label)

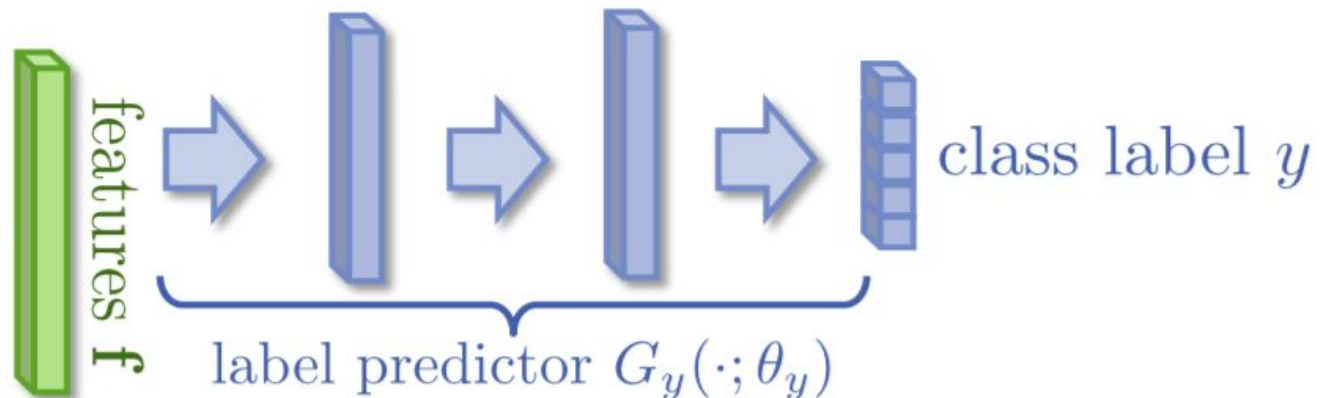
Proposed Solution

- G_f : feature extractor
- G_y : label predictor
- G_d : domain classifier



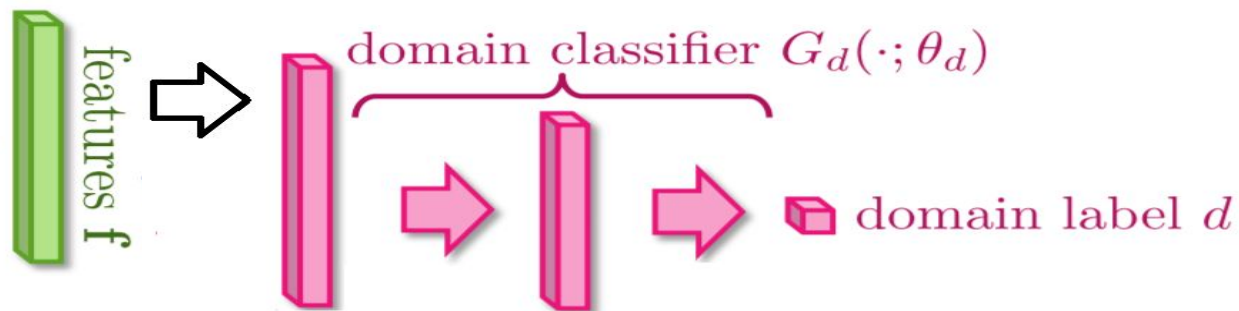
Proposed Solution – Label predictor

- θ_y denotes the parameters in label classifier
- Given a feature, label classifier tries to predict the label of the feature
- Multi-class classification task
- The loss is categorical cross entropy
- θ_y minimize the loss the label classifier
- θ_y is updated as $\theta_y \leftarrow \theta_y - \mu \frac{\partial L_y^i}{\partial \theta_y}$
- μ is the learning rate



Proposed Solution

- θ_d denotes the parameters in domain classifier
- Given a feature, domain classifier tries to predict whether the feature comes from source or target domain
- Binary classification task
- The loss is binary cross entropy
- θ_d minimizes the loss the domain classifier
- θ_d is updated as $\theta_d \leftarrow \theta_d - \mu \frac{\partial L_d^i}{\partial \theta_d}$
- μ is the learning rate



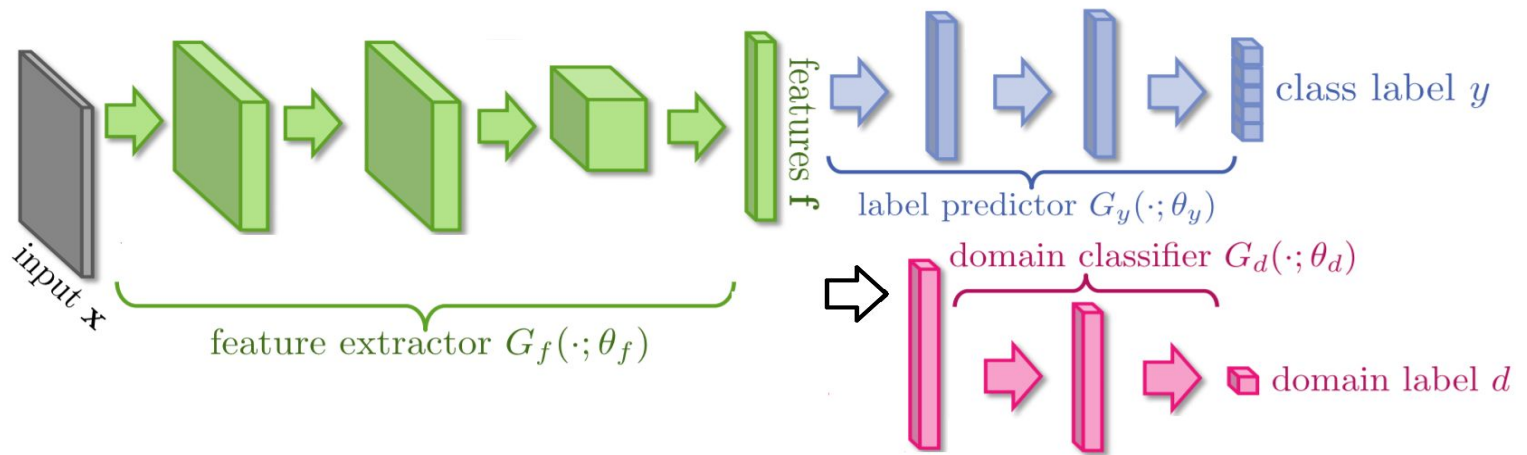
Proposed Solution

- θ_f denotes the parameters in feature extractor
- Given an image, feature extractor generates a deep feature for the image
- θ_f minimizes the loss the label classifier
- θ_f maximize the loss the domain classifier
- θ_d is updated as $\theta_f \leftarrow \theta_f - \mu \left(\frac{\partial L_y^i}{\partial \theta_f} - \lambda \frac{\partial L_d^i}{\partial \theta_f} \right)$
- μ is the learning rate
- The parameter λ controls the trade-off between the two objectives



Proposed Solution

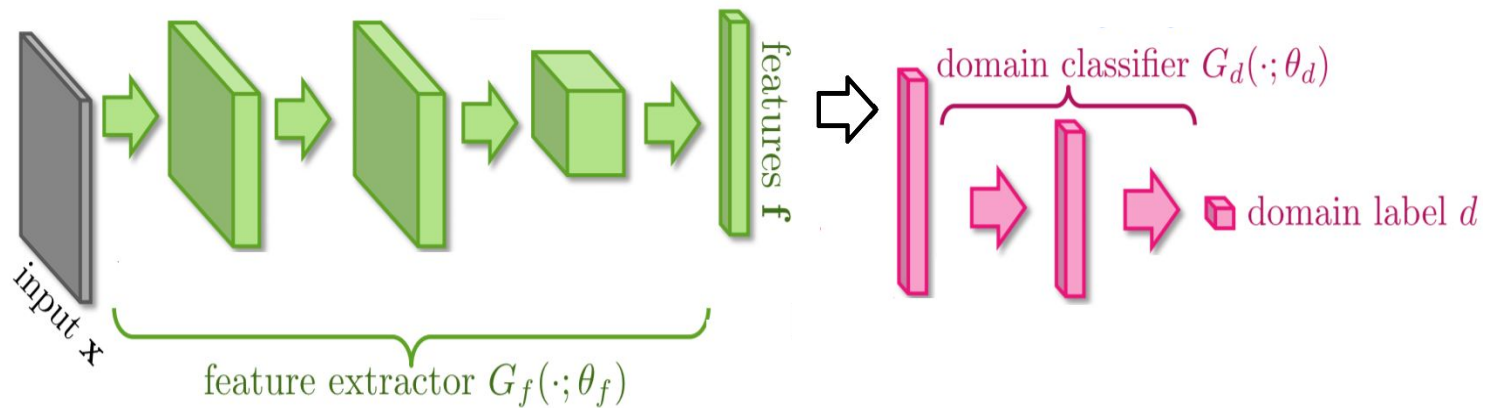
Consider an image from source domain



$$E(\theta_f, \theta_y, \theta_d) = \sum_{\substack{i=1..N \\ d_i=0}} L_y(\hat{y}_i, y_i) - \lambda \sum_{i=1..N} L_d(\hat{d}_i, d_i) =$$

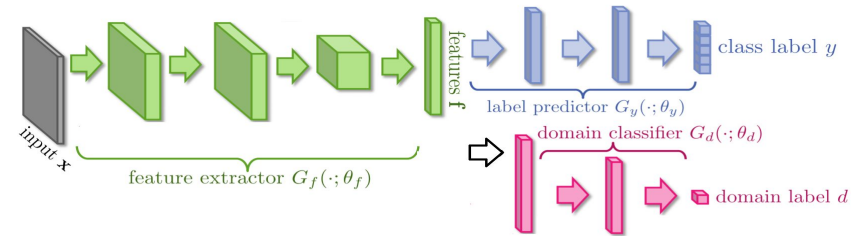
Proposed Solution

Consider an image from target domain



$$E(\theta_f, \theta_y, \theta_d) = \lambda \sum_{i=1..N} L_d (G_d(G_f(\mathbf{x}_i; \theta_f); \theta_d), \underline{y_i}^{d_i})$$

Proposed Solution



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

$$\lambda \sum_{i=1..N} L_d(G_d(G_f(\mathbf{x}_i; \theta_f); \theta_d), y_i) =$$

$$= \sum_{\substack{i=1..N \\ d_i=0}} L_y^i(\theta_f, \theta_y) - \lambda \sum_{i=1..N} L_d^i(\theta_f, \theta_d)$$

$$(\hat{\theta}_f, \hat{\theta}_y) = \arg \min_{\theta_f, \theta_y} E(\theta_f, \theta_y, \hat{\theta}_d)$$

$$\hat{\theta}_d = \arg \max_{\theta_d} E(\hat{\theta}_f, \hat{\theta}_y, \theta_d).$$

Proposed Solution

- At saddle point
 - θ_d minimizes domain classification loss
 - θ_y minimizes label prediction loss
 - θ_f minimizes label prediction loss and maximize domain classification loss

$$E(\theta_f, \theta_y, \theta_d) = \sum_{\substack{i=1..N \\ d_i=0}} L_y^i(\theta_f, \theta_y) - \lambda \sum_{i=1..N} L_d^i(\theta_f, \theta_d)$$

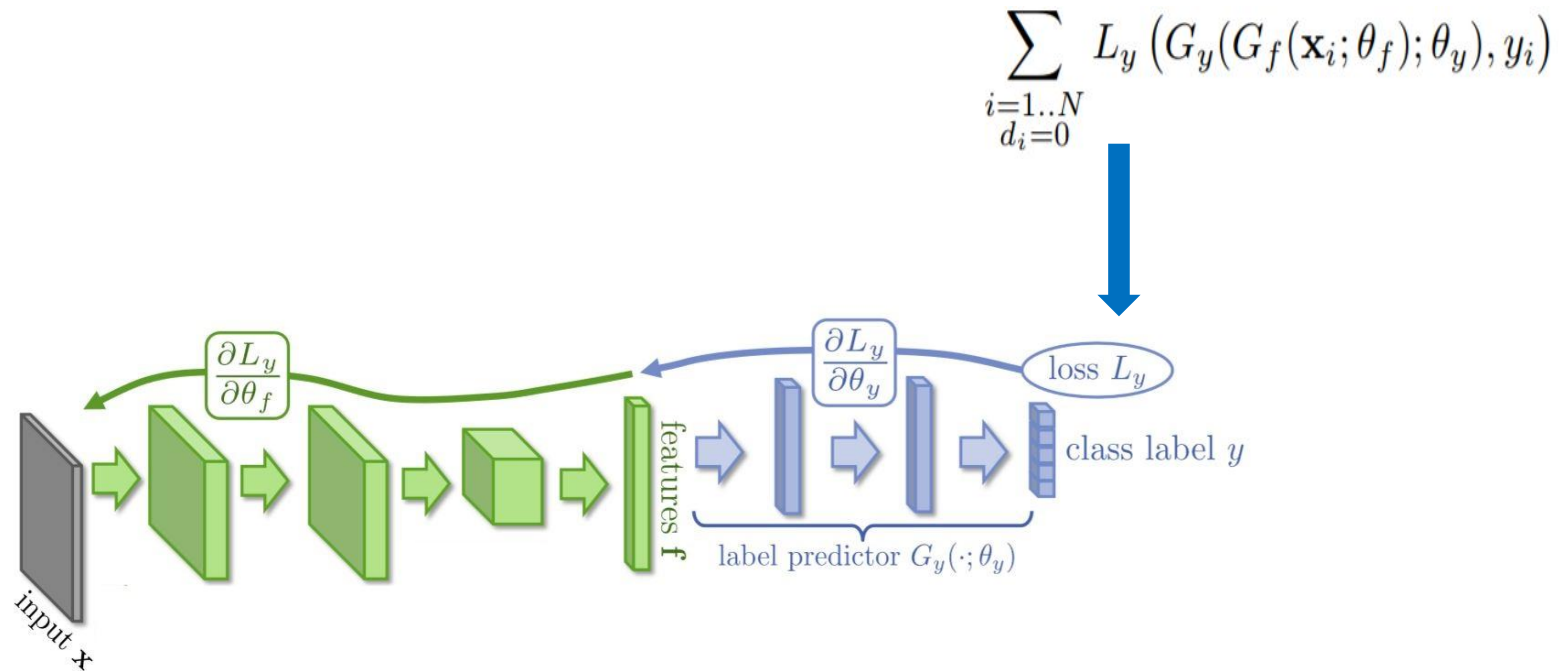
$$\theta_f \longleftarrow \theta_f - \mu \left(\frac{\partial L_y^i}{\partial \theta_f} - \lambda \frac{\partial L_d^i}{\partial \theta_f} \right)$$

$$\theta_y \longleftarrow \theta_y - \mu \frac{\partial L_y^i}{\partial \theta_y}$$

$$\theta_d \longleftarrow \theta_d - \mu \frac{\partial L_d^i}{\partial \theta_d}$$

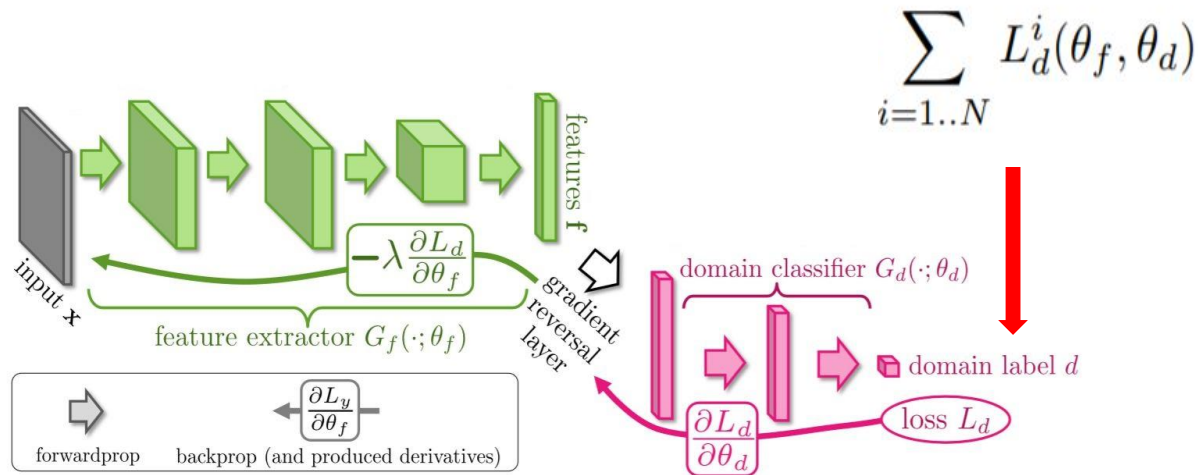
Proposed Solution

- How to backpropagate the label classifier loss?
- Consider only the upper architecture
- This is typical backpropagation



Proposed Solution

- How to backpropagate the domain classifier loss?
- Consider only the upper architecture
- Define gradient reversal layer (GRL)



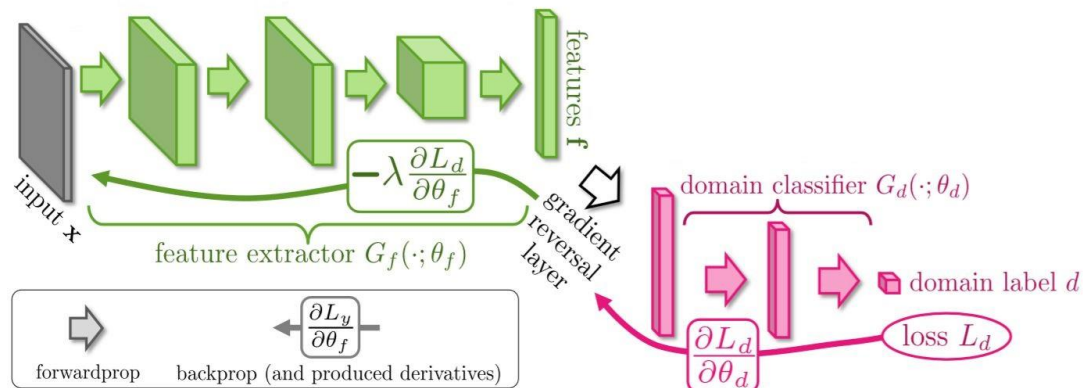
$$\theta_f \leftarrow \theta_f + \lambda \frac{\partial L_d^i}{\partial \theta_f} \quad \theta_d \leftarrow \theta_d - \mu \frac{\partial L_d^i}{\partial \theta_d}$$

Proposed Solution

- Forward : GRL is an identity transformation
- Backward: GRL takes gradient from subsequent level, multiply by λ and pass it to previous layer
- Treat GRL as a pseudo function $R_\lambda(x)$

$$\text{Forward} \quad R_\lambda(\mathbf{x}) = \mathbf{x}$$

$$\text{Backward} \quad \frac{dR_\lambda}{d\mathbf{x}} = -\lambda \mathbf{I}$$



Proposed Solution



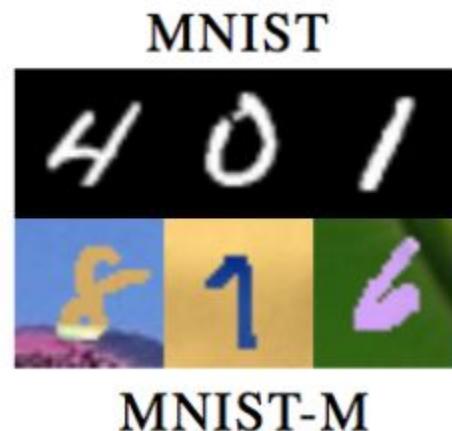
- After training, the label predictor can be used to predict labels for samples from either source or target domain
- Experiment results

Source & Target Datasets



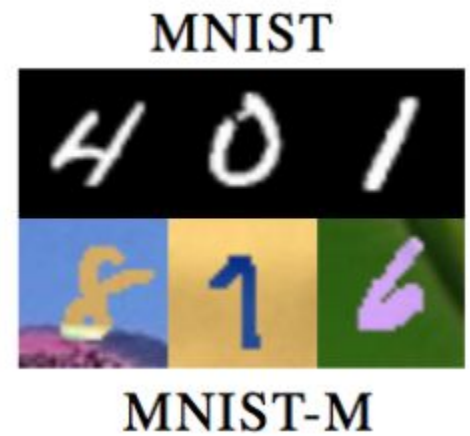
	MNIST	SYN NUMBERS	SVHN	SYN SIGNS
SOURCE				
TARGET				
	MNIST-M	SVHN	MNIST	GTSRB

MNIST → MNIST-M

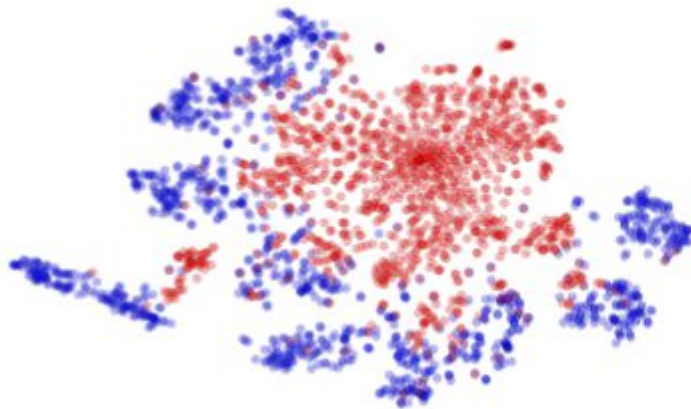


METHOD	SOURCE	MNIST
	TARGET	MNIST-M
SOURCE ONLY		.5749
SA (FERNANDO ET AL., 2013)		.6078 (7.9%)
PROPOSED APPROACH		.8149 (57.9%)
TRAIN ON TARGET		.9891

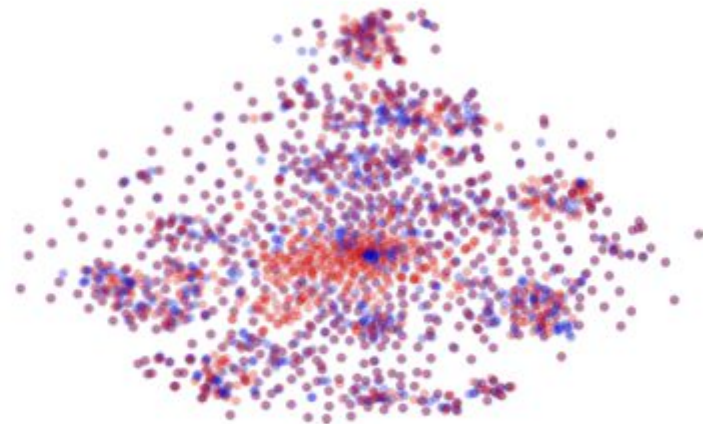
MNIST \rightarrow MNIST-M



MNIST \rightarrow MNIST-M: top feature extractor layer



(a) Non-adapted



(b) Adapted

SYN NUMBERS



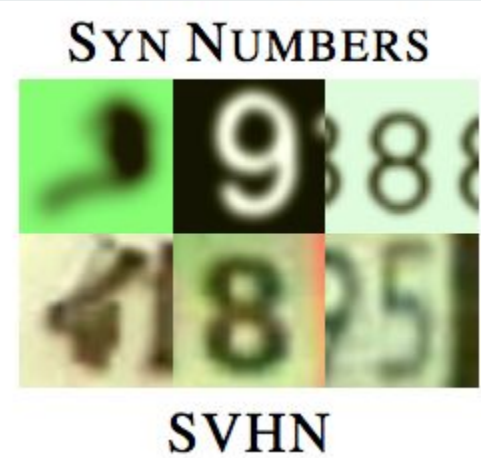
SVHN

Synthetic numbers → SVHN

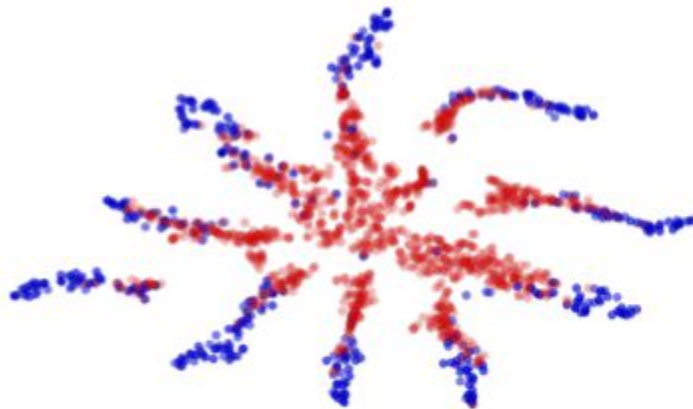


METHOD	SOURCE	MNIST	SYN NUMBERS
	TARGET	MNIST-M	SVHN
SOURCE ONLY		.5749	.8665
SA (FERNANDO ET AL., 2013)		.6078 (7.9%)	.8672 (1.3%)
PROPOSED APPROACH		.8149 (57.9%)	.9048 (66.1%)
TRAIN ON TARGET		.9891	.9244

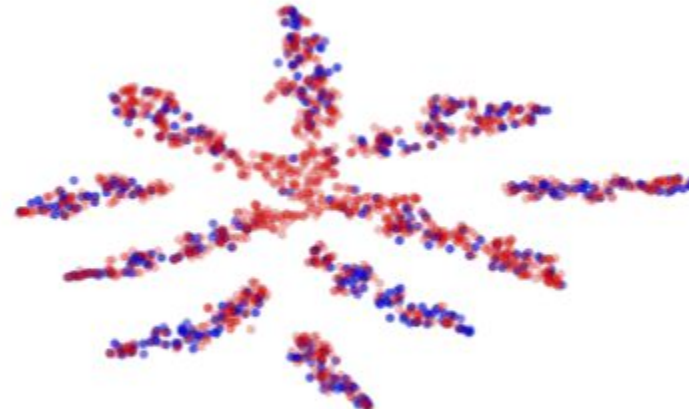
Synthetic numbers \rightarrow SVHN



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

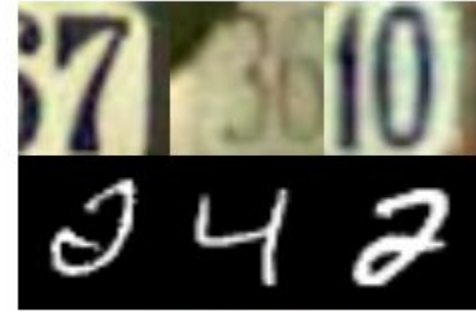


(a) Non-adapted



(b) Adapted

SVHN



MNIST

MNIST ↔ SVHN

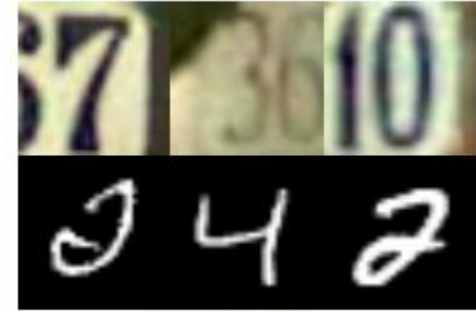


The two directions (MNIST → SVHN and SVHN → MNIST) are not equally difficult.

SVHN is more diverse, a model trained on SVHN is expected to be more generic and to perform reasonably on the MNIST dataset.

Unsupervised adaptation from MNIST to SVHN gives a failure example for this approach.

SVHN



MNIST

SVHN → MNIST



METHOD	SOURCE	MNIST	SYN NUMBERS	SVHN
	TARGET	MNIST-M	SVHN	MNIST
SOURCE ONLY		.5749	.8665	.5919
SA (FERNANDO ET AL., 2013)		.6078 (7.9%)	.8672 (1.3%)	.6157 (5.9%)
PROPOSED APPROACH		.8149 (57.9%)	.9048 (66.1%)	.7107 (29.3%)
TRAIN ON TARGET		.9891	.9244	.9951

Synthetic Signs → GTSRB

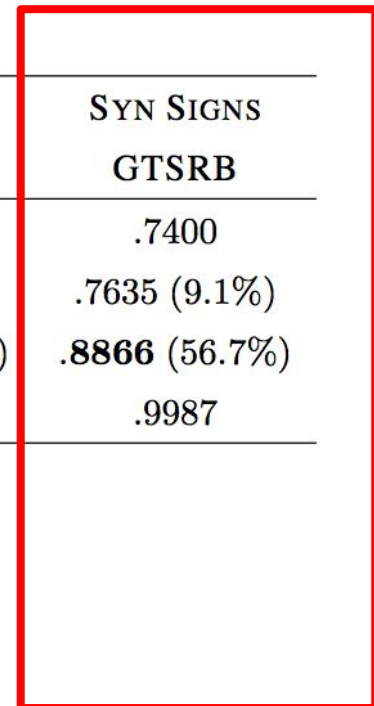


SYN SIGNS



GTSRB

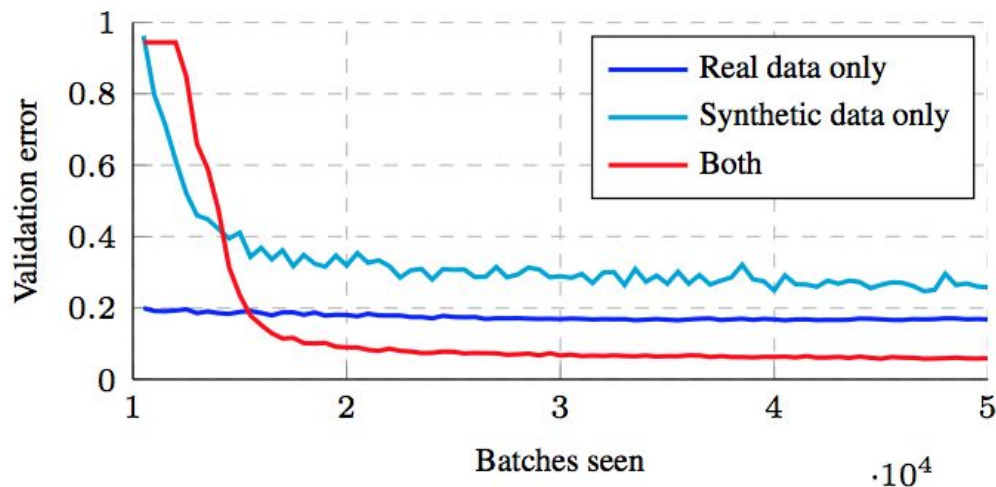
METHOD	SOURCE	MNIST	SYN NUMBERS	SVHN	SYN SIGNS
	TARGET	MNIST-M	SVHN	MNIST	GTSRB
SOURCE ONLY		.5749	.8665	.5919	.7400
SA (FERNANDO ET AL., 2013)		.6078 (7.9%)	.8672 (1.3%)	.6157 (5.9%)	.7635 (9.1%)
PROPOSED APPROACH		.8149 (57.9%)	.9048 (66.1%)	.7107 (29.3%)	.8866 (56.7%)
TRAIN ON TARGET		.9891	.9244	.9951	.9987



Synthetic Signs → GTSRB



This paper also evaluates the proposed algorithm for semi-supervised domain adaptation, i.e. when one is additionally provided with a small amount of labeled target data.



Office dataset



METHOD	SOURCE	AMAZON	DSLR	WEBCAM
	TARGET	WEBCAM	WEBCAM	DSLR
GFK(PLS, PCA) (GONG ET AL., 2012)		.464 ± .005	.613 ± .004	.663 ± .004
SA (FERNANDO ET AL., 2013)		.450	.648	.699
DA-NBNN (TOMMASI & CAPUTO, 2013)		.528 ± .037	.766 ± .017	.762 ± .025
DLID (S. CHOPRA & GOPALAN, 2013)		.519	.782	.899
DECAF ₆ SOURCE ONLY (DONAHUE ET AL., 2014)		.522 ± .017	.915 ± .015	–
DANN (GHIFARY ET AL., 2014)		.536 ± .002	.712 ± .000	.835 ± .000
DDC (TZENG ET AL., 2014)		.594 ± .008	.925 ± .003	.917 ± .008
PROPOSED APPROACH		.673 ± .017	.940 ± .008	.937 ± .010

Conclusions



- Proposed a new approach to unsupervised domain adaptation of deep feed-forward architectures;
- Unlike previous approaches, this approach is accomplished through standard backpropagation training;
- The approach is scalable, and can be implemented using any deep learning package.