



**KOREA**  
UNIVERSITY

# Autoencoder based Domain Adaptation for Speaker Recognition under Insufficient Channel Information

**Suwon Shon**, Seongkyu Mun<sup>\*</sup>, Wooil Kim<sup>\*\*</sup>, Hanseok Ko<sup>\*</sup>

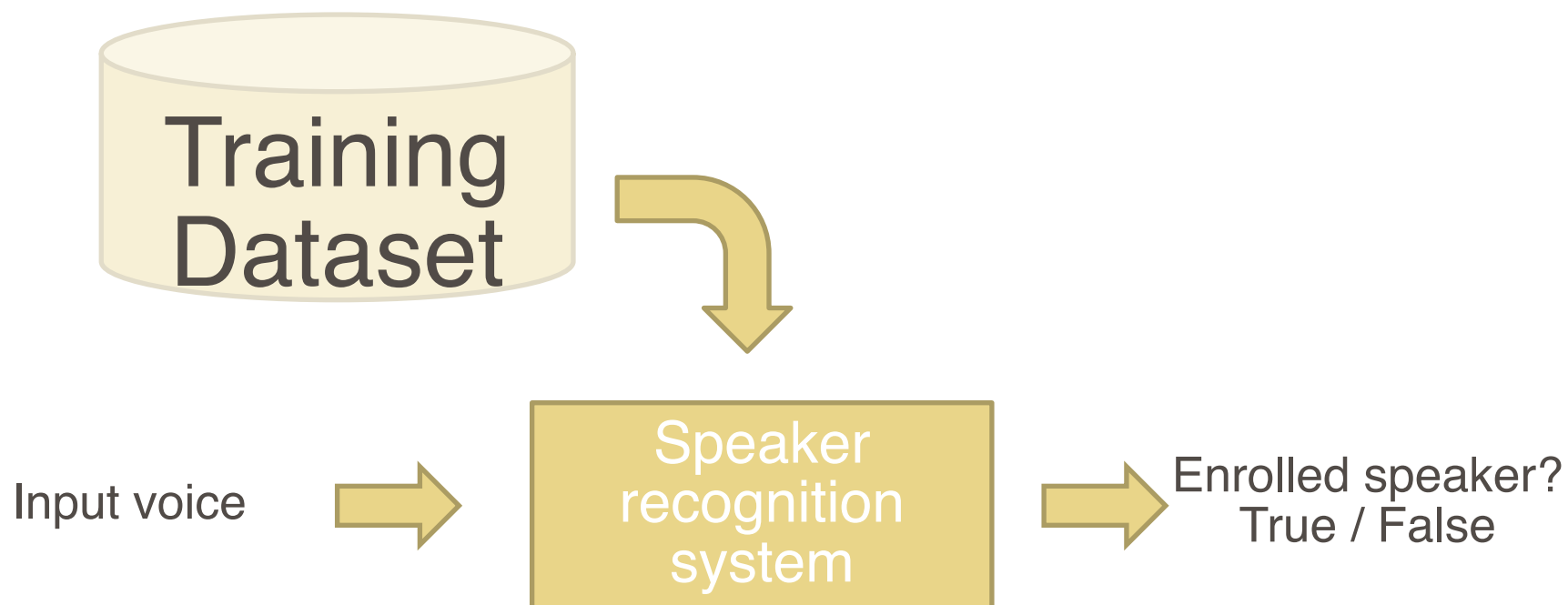
MIT **C**omputer **S**cience and **A**rtificial **I**ntelligence **L**aboratory,  
Cambridge, MA, USA

Korea University, Seoul, South Korea<sup>\*</sup>

Incheon National University, South Korea<sup>\*\*</sup>

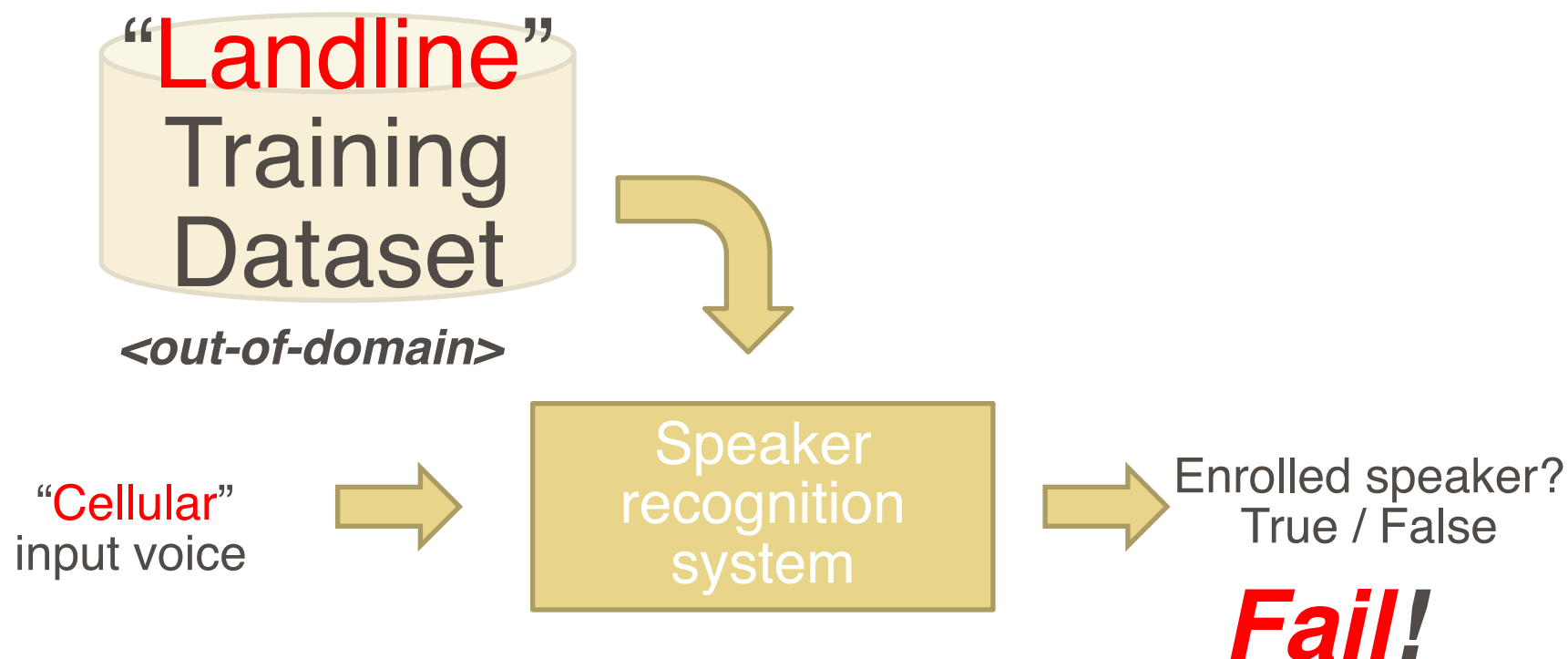
# Introduction

- **Speaker recognition task**



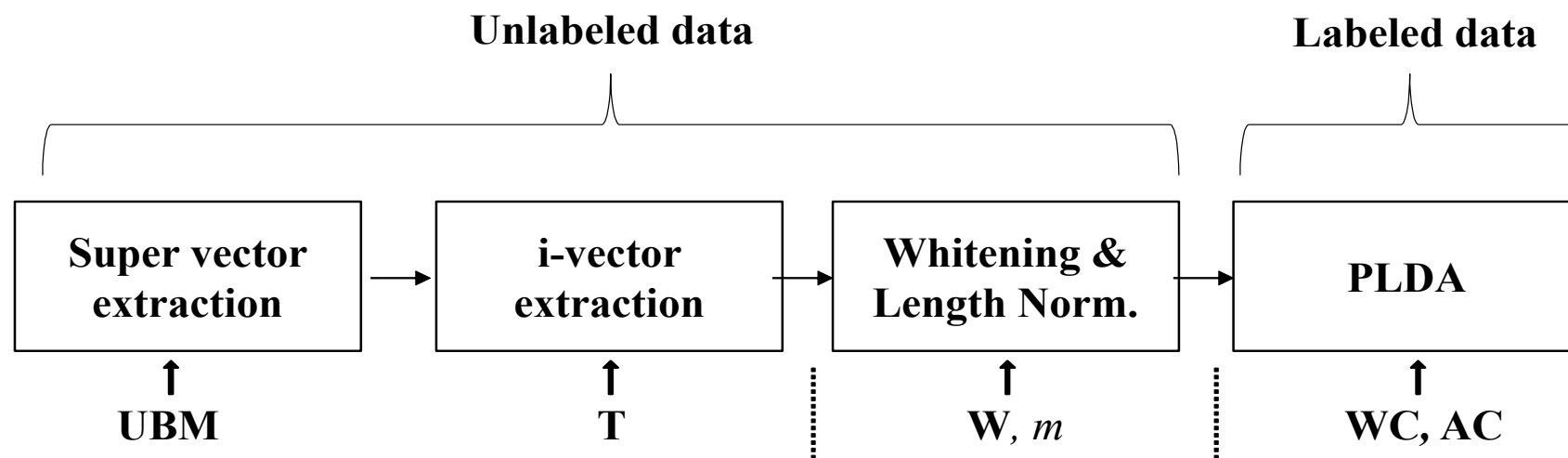
# Introduction

- Channel domain mismatched condition



# Introduction

- **Domain adaptation challenge 2013 @ JHU workshop**
  - **SRE10** (evaluation) collected in 2010 (mostly cellular)
    - \* 7,169 target and 408,950 non-target trials
  - **SWB** collected from 1992-2000 (mostly landline), mismatched
  - **SRE** collected from 2004-2008 (mostly cellular), matched
    - \* Suppose we don't have labels on **SRE**



# Introduction

- **Domain adaptation challenge 2013 @ JHU workshop**
  - **SRE10** (evaluation) collected in 2010 (mostly cellular)
    - \* 7,169 target and 408,950 non-target trials
  - **SWB** collected from 1992-2000 (mostly landline), mismatched
  - **SRE** collected from 2004-2008 (mostly cellular), matched
    - \* Suppose we don't have labels on **SRE**

System #	Unlabeled data		Labeled data		EER	
	UBM, T	W, <i>m</i>	WC,AC			
0*	<b>SRE</b>	<b>SRE</b>	<b>SRE</b>		2.43	→ Domain matched benchmark
1	<b>SWB</b>	<b>SRE</b>	<b>SRE</b>		2.33	
2	<b>SWB</b>	<b>SRE</b>	<b>SWB</b>		5.70	
3*	<b>SWB</b>	<b>SWB</b>	<b>SWB</b>		6.92	→ Domain mismatched

Table 2: *SRE10* Test using *DAC13* *i*-vector set.

# Motivation

- **Insufficient Channel Information**

	SWB	SRE	SRE-1phn
#spkrs	3114	3790	3787
#calls	33039	36470	25640
Avg. #calls/spkrs	10.6	9.6	6.77
Avg. #phone_num/spkr	3.8	2.8	1

*<Statistics in DAC 13 i-vector Dataset>*

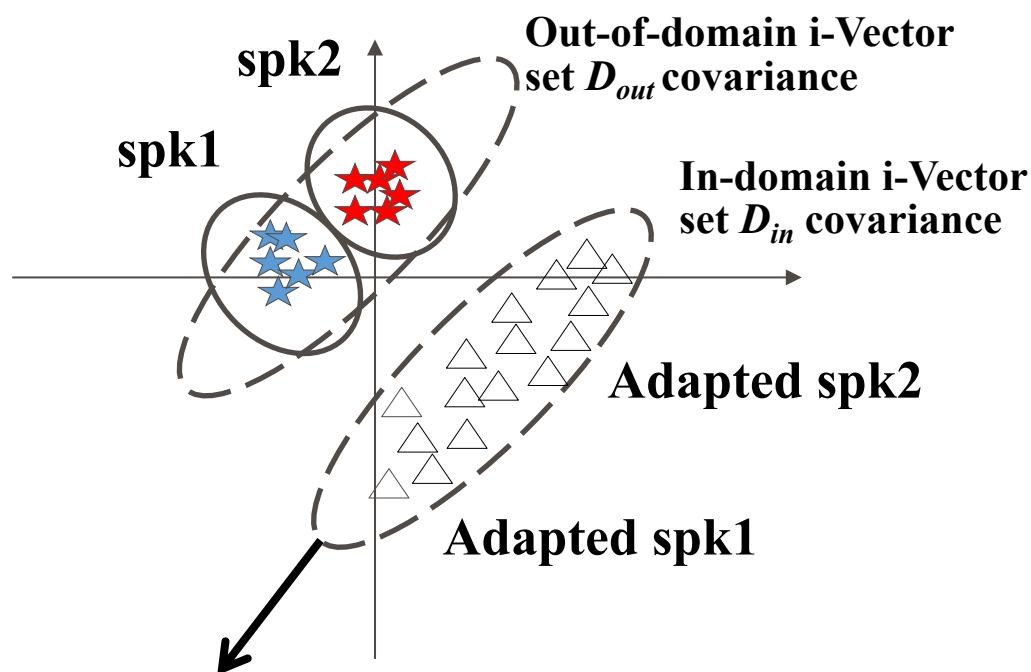
System #	UBM, T	W <sub>m</sub>	WC,AC	EER	
1	<b>SWB</b>	<b>SRE</b>	<b>SRE</b>	2.33	better
2	<b>SWB</b>	<b>SRE</b>	<b>SWB</b>	5.70	
3	<b>SWB</b>	<b>SRE-1phn</b>	<b>SRE-1phn</b>	9.34	worse
4	<b>SWB</b>	<b>SRE-1phn</b>	<b>SWB</b>	5.66	

*<SRE10 Test using DAC13 i-vector set>*

**Performance degraded by Insufficient channel information although it is matched domain dataset**

# Proposed Approach

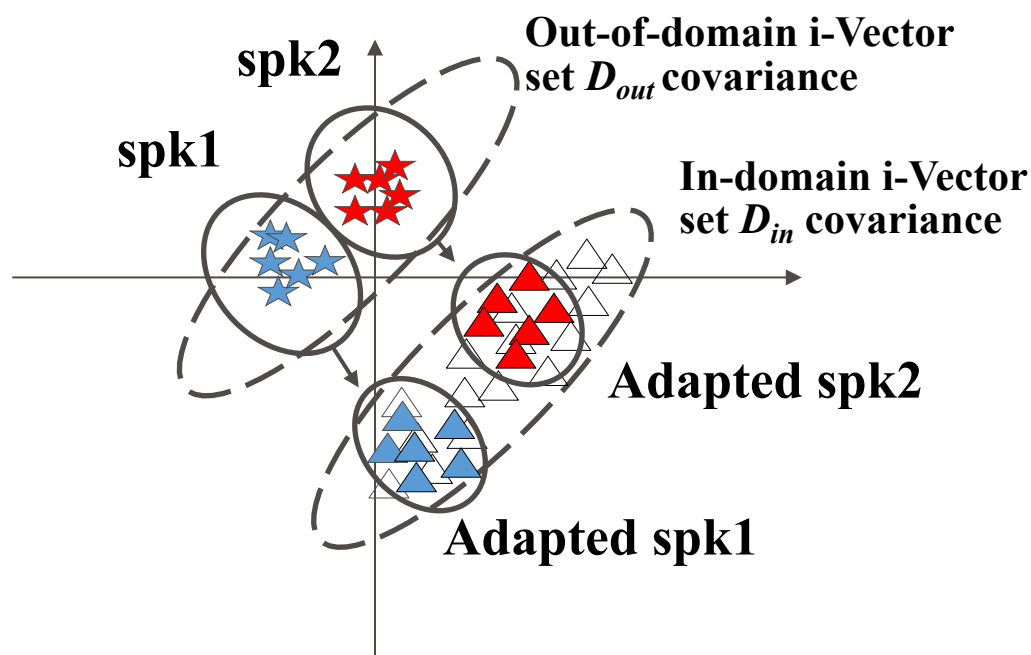
- **Auto-encoder based Domain Adaptation (AEDA)**



***Useless because of insufficient channel information***

# Proposed Approach

- **Auto-encoder based Domain Adaptation (AEDA)**

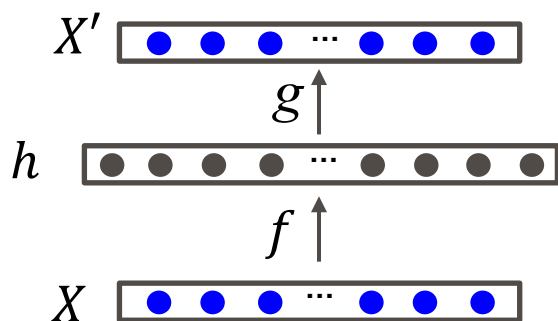


***Transferring labeled out-of-domain dataset to in-domain***



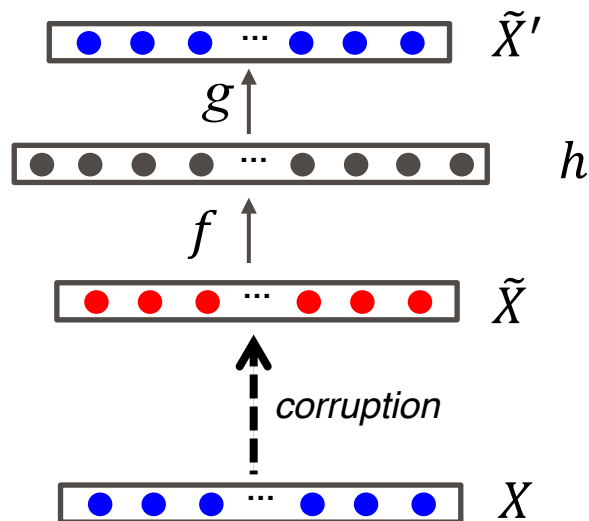
# Proposed Approach

- Autoencoder and Denoising Autoencoder



$$\mathcal{L}(X, X') = \|X - X'\|^2$$

<Autoencoder>

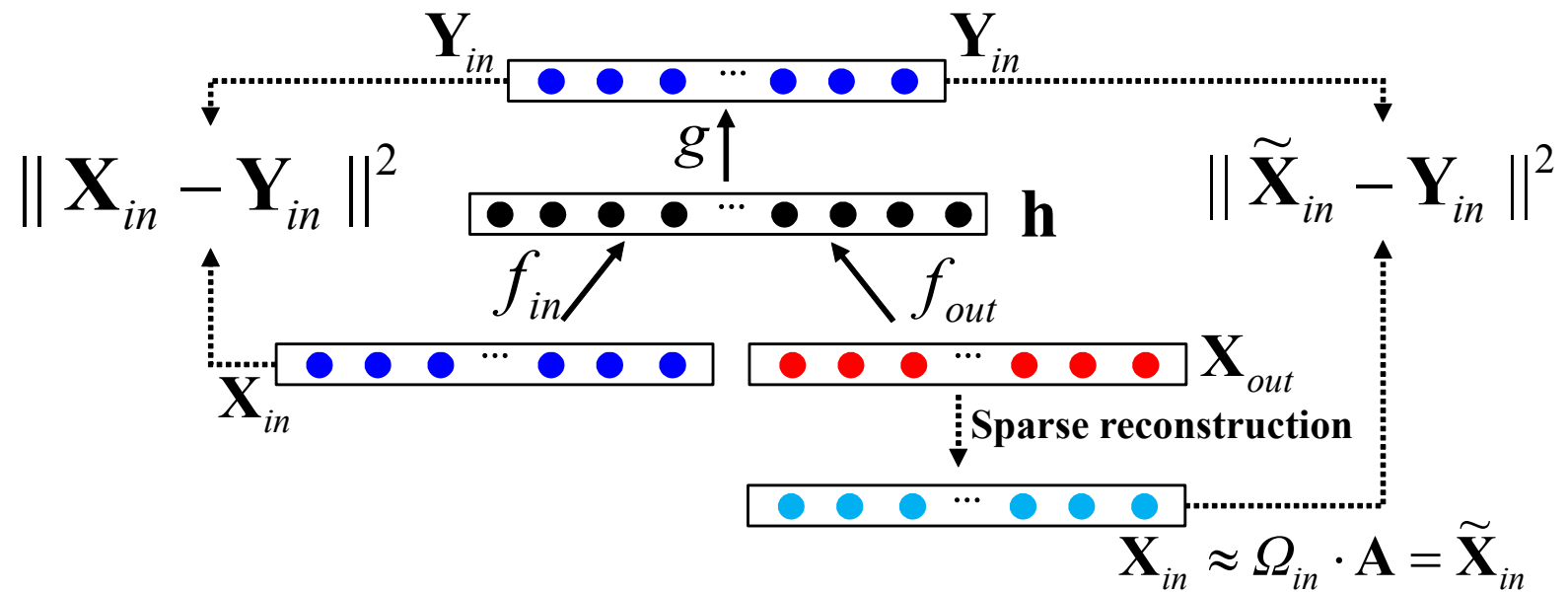


$$\mathcal{L}(X, X') = \|X - \tilde{X}'\|^2$$

<Denoising Autoencoder>

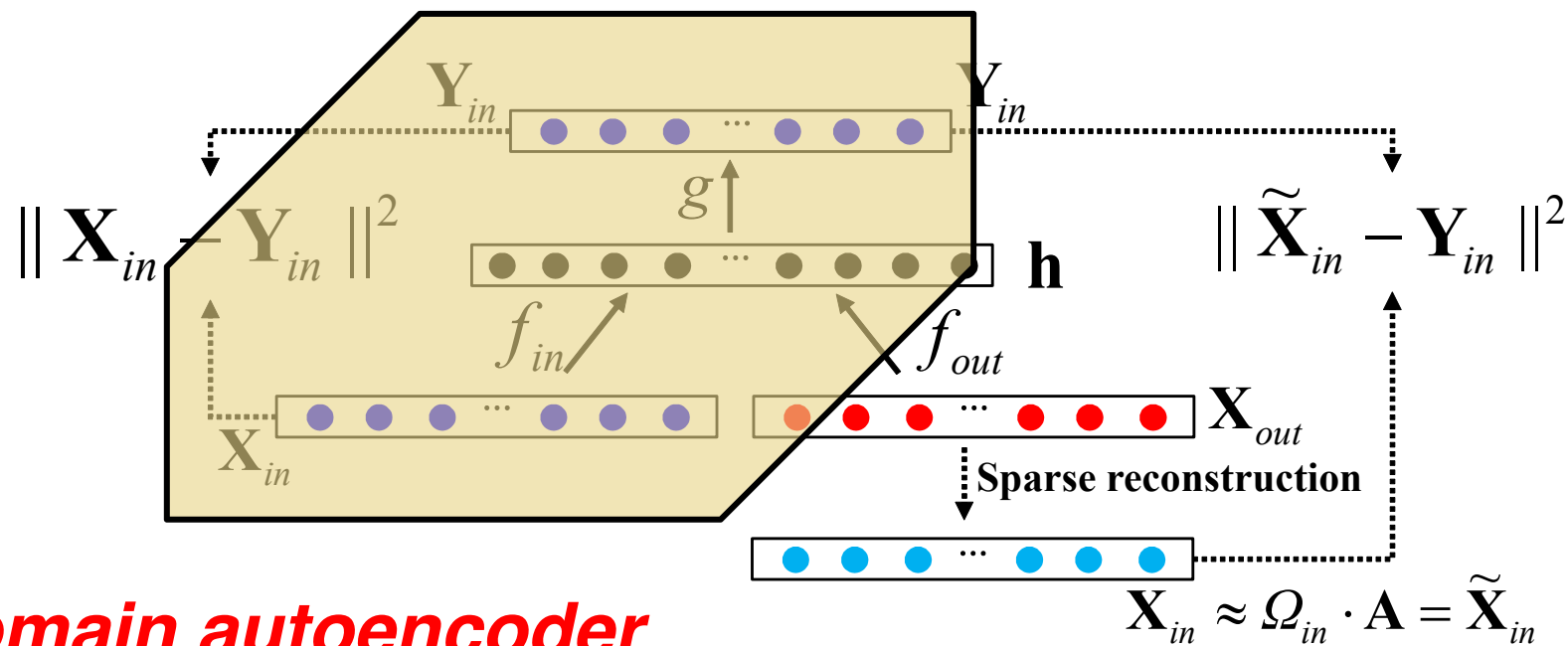
# Proposed Approach

- Auto-encoder based Domain Adaptation (AEDA)



# Proposed Approach

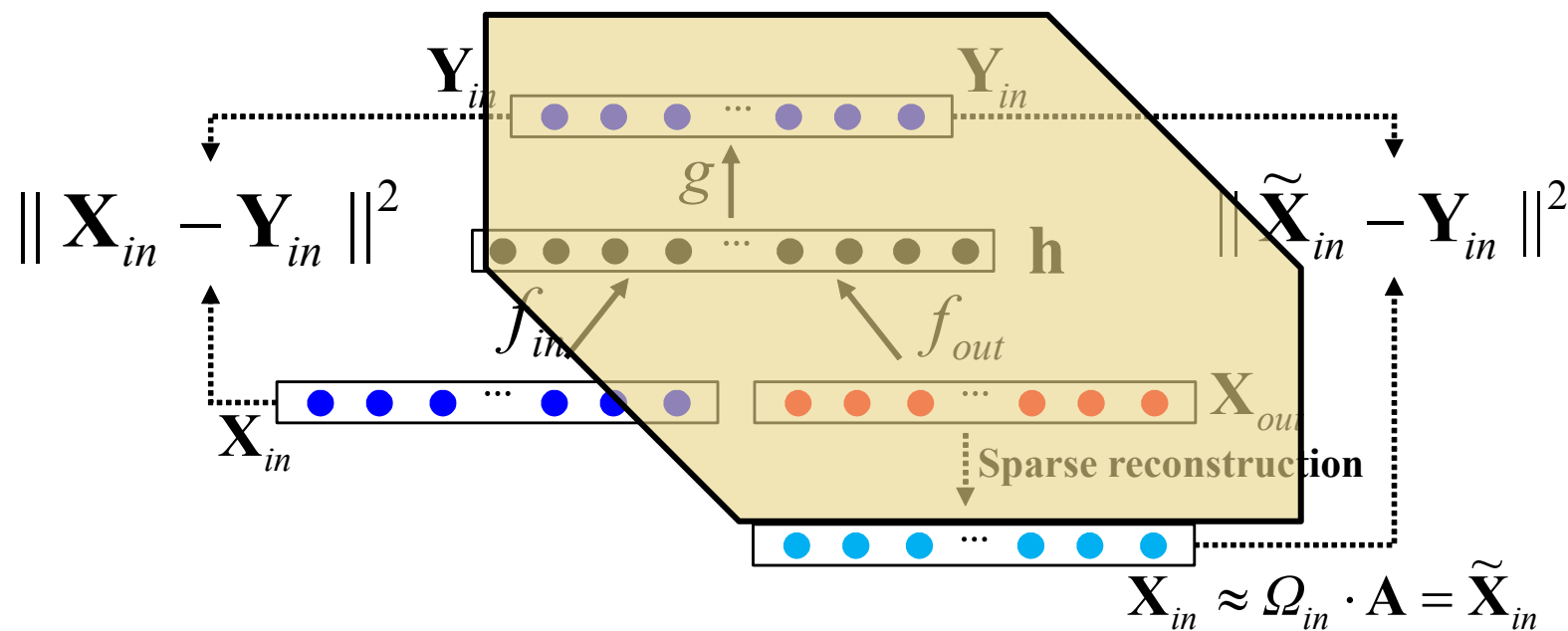
- Auto-encoder based Domain Adaptation (AEDA)



***In-domain autoencoder  
(using unlabeled in-domain dataset)***

# Proposed Approach

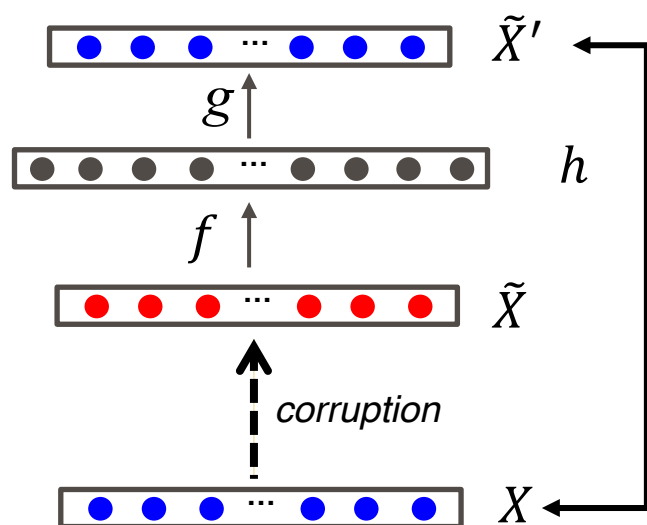
- Auto-encoder based Domain Adaptation (AEDA)



***domain transferring  
autoencoder  
(using labeled out-of-domain dataset)***

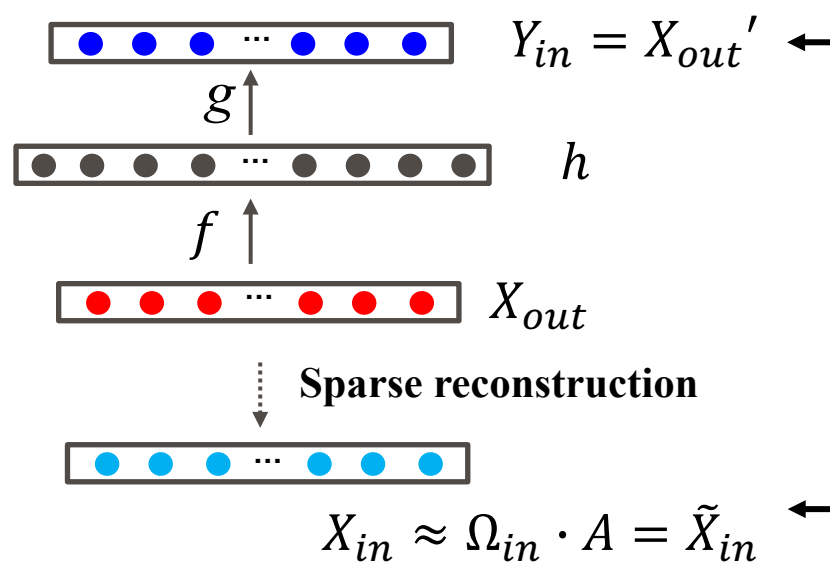
# Proposed Approach

- Sparse reconstruction



$$\mathcal{L}(X, X') = \|X - \tilde{X}'\|^2$$

<Denoising Autoencoder>



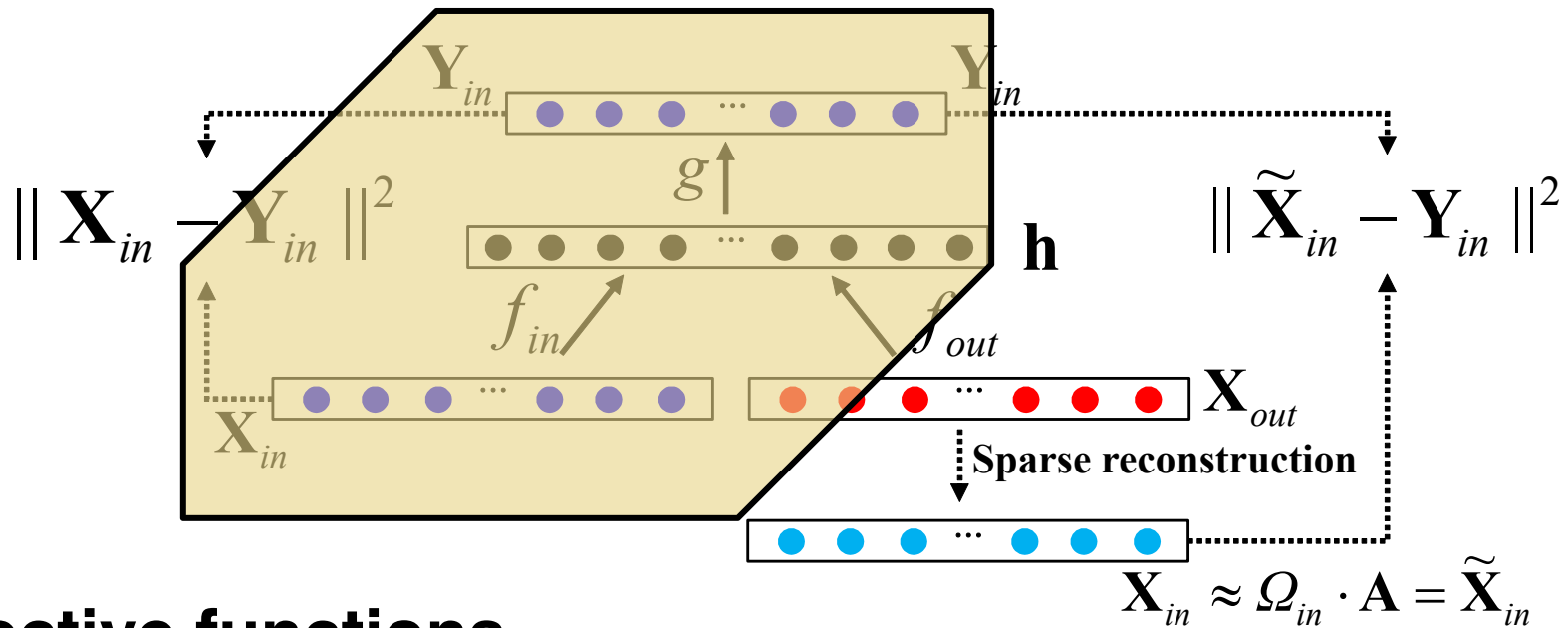
$$\text{Objective function : } \min_{\alpha_j} \|\Omega_{in} \alpha_j - \mathbf{y}_j^{in}\|^2 + \gamma |\alpha_j|^2$$

$$\begin{aligned} \mathcal{L}(X_{in}, Y_{in}) &= \|X_{in} - Y_{in}\|^2 \\ &= \|\tilde{X}_{in} - Y_{in}\|^2 \end{aligned}$$

<Out-of-domain transferring autoencoder>

# Proposed Approach

- Structure of Autoencoder which sharing hidden layer  $h$



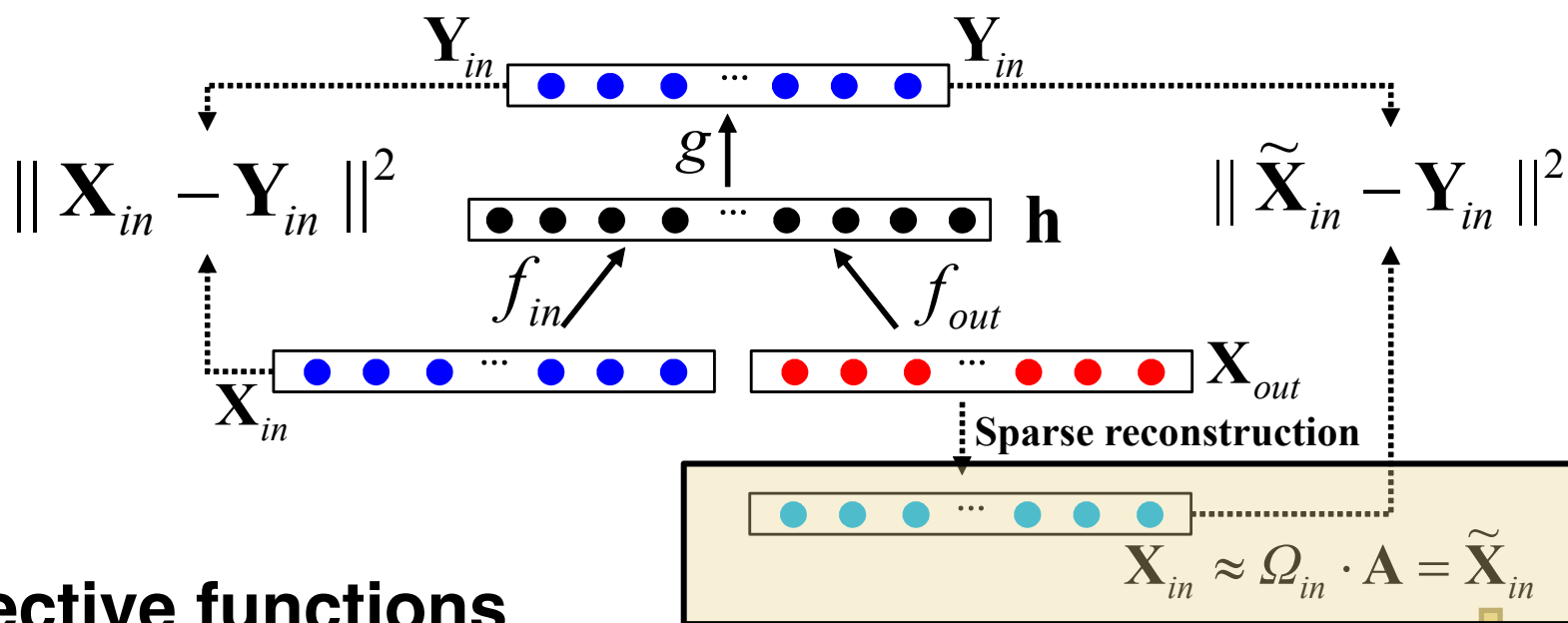
- Objective functions

– AE part 
$$\min_{f_{in}, g} \|\mathbf{X}_{in} - \mathbf{Y}_{in}\|^2 = \min_{f_{in}, g} \|\mathbf{X}_{in} - g(f_{in}(\mathbf{X}_{in}))\|^2$$

$\mathbf{X}_{in} \approx \Omega_{in} \cdot \mathbf{A} = \tilde{\mathbf{X}}_{in}$   
 ↓  
**Least angle regression**

# Proposed Approach

- Structure of Autoencoder which sharing hidden layer  $h$



- Objective functions

– AE part  $\min_{f_{in}, g} \|X_{in} - Y_{in}\|^2 = \min_{f_{in}, g} \|X_{in} - g(f_{in}(X_{in}))\|^2$

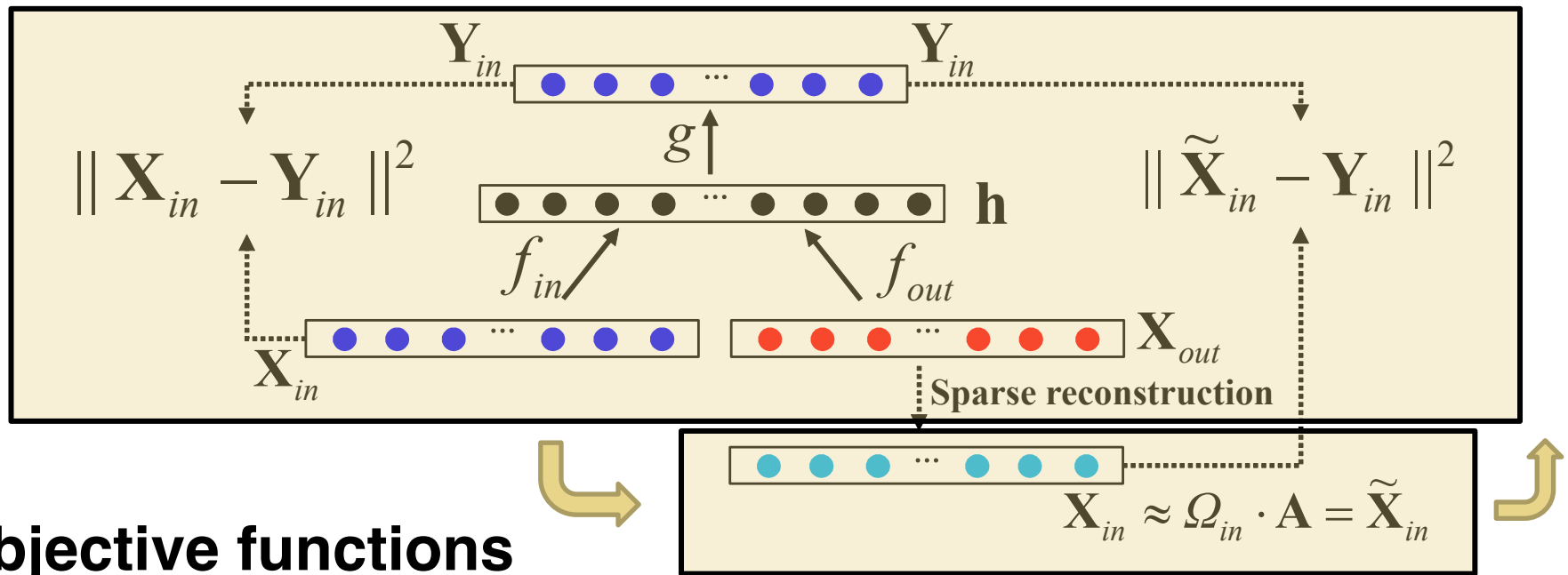
Least angle regression

- Least angle regression for sparse reconstruction

$$\min_{\alpha_j} \|\Omega_{in} \alpha_j - \mathbf{y}_j^{in}\|^2 + \gamma |\alpha_j|^2$$

# Proposed Approach

- Structure of Autoencoder which sharing hidden layer  $h$



- Objective functions

– AE part  $\min_{f_{in}, g} \|X_{in} - Y_{in}\|^2 = \min_{f_{in}, g} \|X_{in} - g(f_{in}(X_{in}))\|^2$

– DAE part  $\min_{f_{out}, g} \|X_{in} - Y_{in}\|^2 = \min_{f_{out}, g} \|X_{in} - g(f_{out}(X_{out}))\|^2$

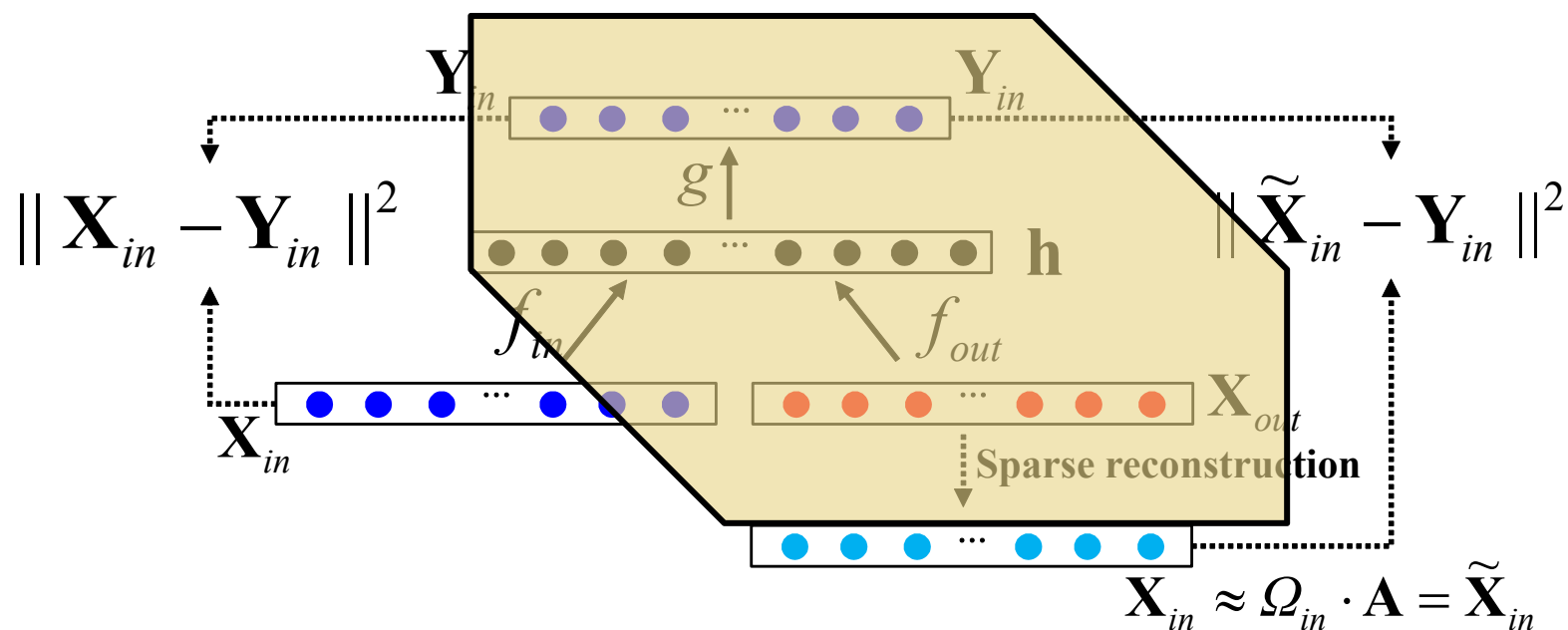
– AEDA  $\min_{f_{in}, f_{out}, g} \|X_{in} - g(f_{in}(X_{in}))\|^2 + \|\tilde{X}_{in} - g(f_{out}(X_{out}))\|^2$

Training network  $\min_{\alpha_j} \|\Omega_{in} \alpha_j - y_j^{in}\|^2 + \gamma |\alpha_j|^2$  Sparse reconstruction



# Proposed Approach

- **Structure of Autoencoder which sharing hidden layer h**



- **AEDA**
  - 600 dim i-vector with 1000 hidden node with learning rate 0.005
- **Sparse reconstruction**
  - Least Angle Regression(LARS)
  - Sparsity 0.01
  - Random 1500 spk i-vector for in-domain dictionary  $\Omega_{in}$
- **Performance**
  - Using PLDA with 400 eigenvoice after 400 dim LDA transform
  - EER, DCF10, DCF08

# Experimental result

- **Auto-encoder based Domain Adaptation (AEDA)**

#	Adaptation & Compensation	WC,AC	EER	DCF10	DCF08
3	-	SRE-1phn	9.34	0.721	0.520
4	-	SWB	5.66	0.633	0.426
5	Interpolated [13]	SWB + SRE-1phn	6.55	0.652	0.454
6	IDV [15]	IDV-SWB	6.15	0.676	0.476
7	DICN [16]	DICN-SWB	4.99	0.623	0.416
8	DAE [23]	DAE-SWB	4.81	0.610	0.398
<b>9</b>	<b>AEDA</b>	<b>AEDA-SWB</b>	<b>4.50</b>	<b>0.589</b>	<b>0.362</b>

*<SRE10 evaluation result with DAC 13 Dataset when Unlabeled In-Domain Dataset is Available >*

# Conclusion

- **Only small subset of unlabeled in-domain dataset is used for domain adaptation**
- **Insufficient channel information dataset is effectively used for transferring knowledge of in-domain**
- **Domain transferring autoencoder part of AEDA can be trained using sparse reconstruction without actual pair of in-domain and out-of-domain**

# Q & A

- **Thanks!**
  
- **Domain related paper :**  
**Suwon Shon, Seongkyu Mun and Hanseok Ko,**  
**“Recursive whitening transformation for speaker**  
**recognition on Language Mismatched Condition”**  
**@ 4.9 Evaluation of Speaker and language identification**  
**systems session, Wednesday 10:00~12:00**