# OVER-VIEW OF TRANSFER LEARNING

HAPPY BUZAABA

ABI-RIKEN AIP, University of Tsukuba

# Foundations of Transfer Learning

[1] As humans, we find it easy to transfer knowledge we have learned from one domain or task to another. When we encounter a new task, we don't have to start from scratch. Instead, we use our previous experience to learn and adapt to that new task faster and more accurately.

[2] **Transfer Learning**: refers to a situation where what has been learned from one setting (eg., distribution $P1$) is exploted in another setting (say, distribution $P2$).

**Assumption**: Many of the factors that explain the variations in $P1$ are relevant to the variations that need to be captured to learn $P2$.

[1]. Pan & Yang,., A Survey on Transfer Learning, IEEE 2010
[2]. Book: DEEP LEARNING Ian Godfellow, Yoshua Bengio, and Aaron Courville page 526

# Fundamental Questions in Transfer Learning

1. What information is useful and transferable from source domain to the target domain?

2. What is the best way of transfering the information?

3. How to avoid transfering information that is detrimental to the desired outcome?

Note: To answer these questions, we consider similarities between the feature spaces, models & tasks of the source and target domains.

Pan & Yang: A survey on Transfer Learning. IEEE 2010

# Notation in Transfer Learning

- Domain $D$: $D = \{X, P(X)\}$,

    $X$ is the feature space,

    $P(X)$ is marginal probability distribution, where $X = \{x_1, \dots, x_n\} \in X$.

> If two domains are different, either $X_s \neq X_T$, or $P(X_s) \neq P(X_T)$

- Task $\mathcal{T}$: Given a specific $D$, a task $\mathcal{T} = \{Y, f(\cdot)\}$

    $Y$ is label space, & $f(\cdot)$ objective predictive function.

    $f(\cdot)$ can be learned from training data $\{(x_i, y_i) | i \in \{1, 2, \dots, N\}\}$, where $x_i \in X$ & $y_i \in Y$

> From probabilistic view, $f(x_i)$ can be written as $P(y_i | x_i)$, and the task as $\mathcal{T} = \{Y, P(Y|X)\}$

In general, if two tasks are different, they may have different;

label spaces $Y_s \neq Y_T$, or $P(Y_s | X_s) \neq P(Y_T | X_T)$

Pan & Yang: A survey on Transfer Learning. IEEE 2010
Weiss et al., A survey of transfer Learning. Journal of Big Data 2016

# Definition of Transfer Learning

Given $D_s$ and $T_s$, $D_T$ and $T_T$,

Transfer learning aims to improve the learning of the target predictive function $f_T(\cdot) \sim P(Y_T|X_T)$ in $D_T$ using the knowledge from $D_s$ & $T_s$, where $D_s \neq D_T$, or $T_s \neq T_T$.

1. A domain is a pair $D = \{\mathcal{X}, P(X)\}$ thus the condition;

$$D_s \neq D_T \text{ implies that either } \mathcal{X}_s \neq \mathcal{X}_T \text{ or } P(X_s) \neq P(X_T)$$

2. A task is a pair $\mathcal{T} = \{\mathcal{Y}, P(Y|X)\}$ thus the condition;

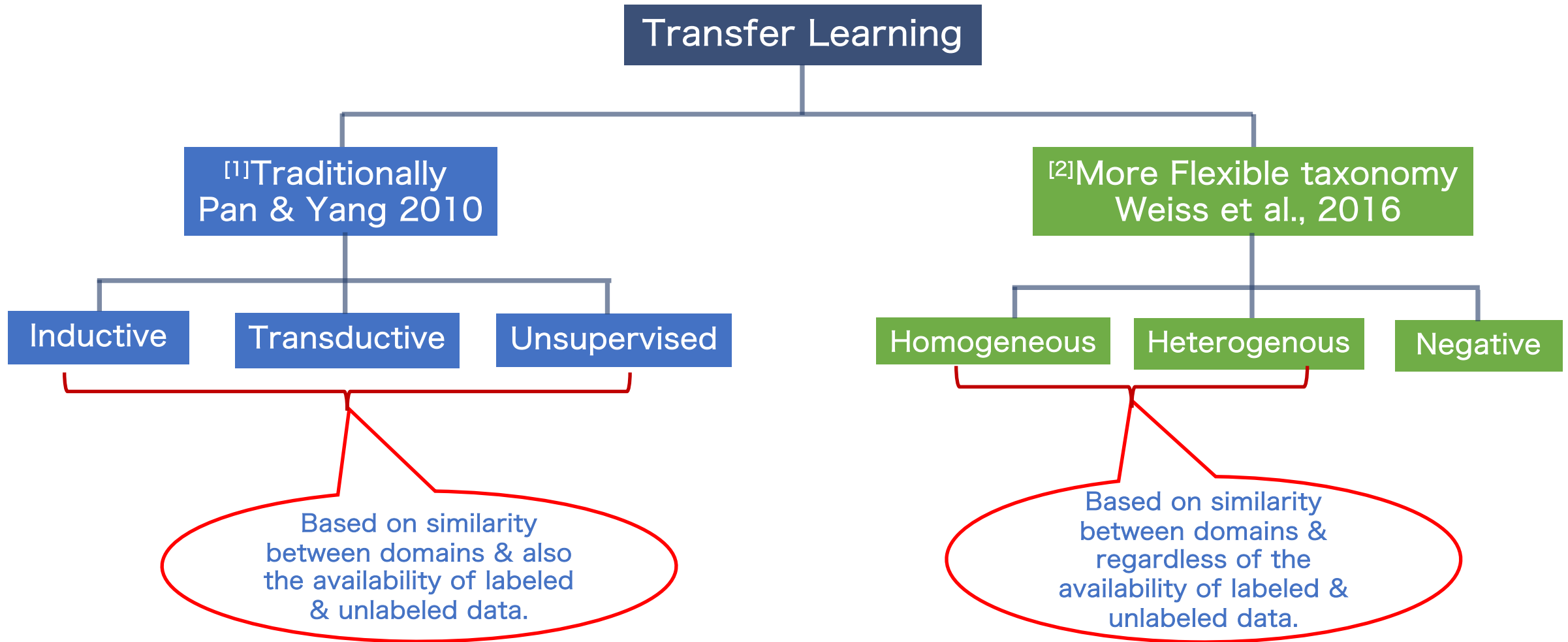$$T_s \neq T_T \text{ implies that either } y_s \neq y_T \text{ or } P(Y_s|X_s) \neq P(Y_T|X_T)$$

# Definition of Transfer Learning

Table 1: All possible combinations for domain and task pair

| Scenario | Example in Sentiment classification |
|---|---|
| $X_s \neq X_t$ | The source domain could be English and the target could be Arabic. |
| $P(X_s) \neq P(X_t)$ | The review could be written in the topic of hotels in the first domain while on restaurants on the target domain. |
| $y_s \neq y_t$ | As an example, the reviews in the source task might be binary while in the target task is categorical. |
| $P(y_s \mid x_s) \neq P(y_t \mid x_t)$ | For example, given a specific review in the domain task might have a label negative while in the target task it has a label neutral. |

Table Source; Zaid et al., A survey on Transfer Learning in Natural Language Processing
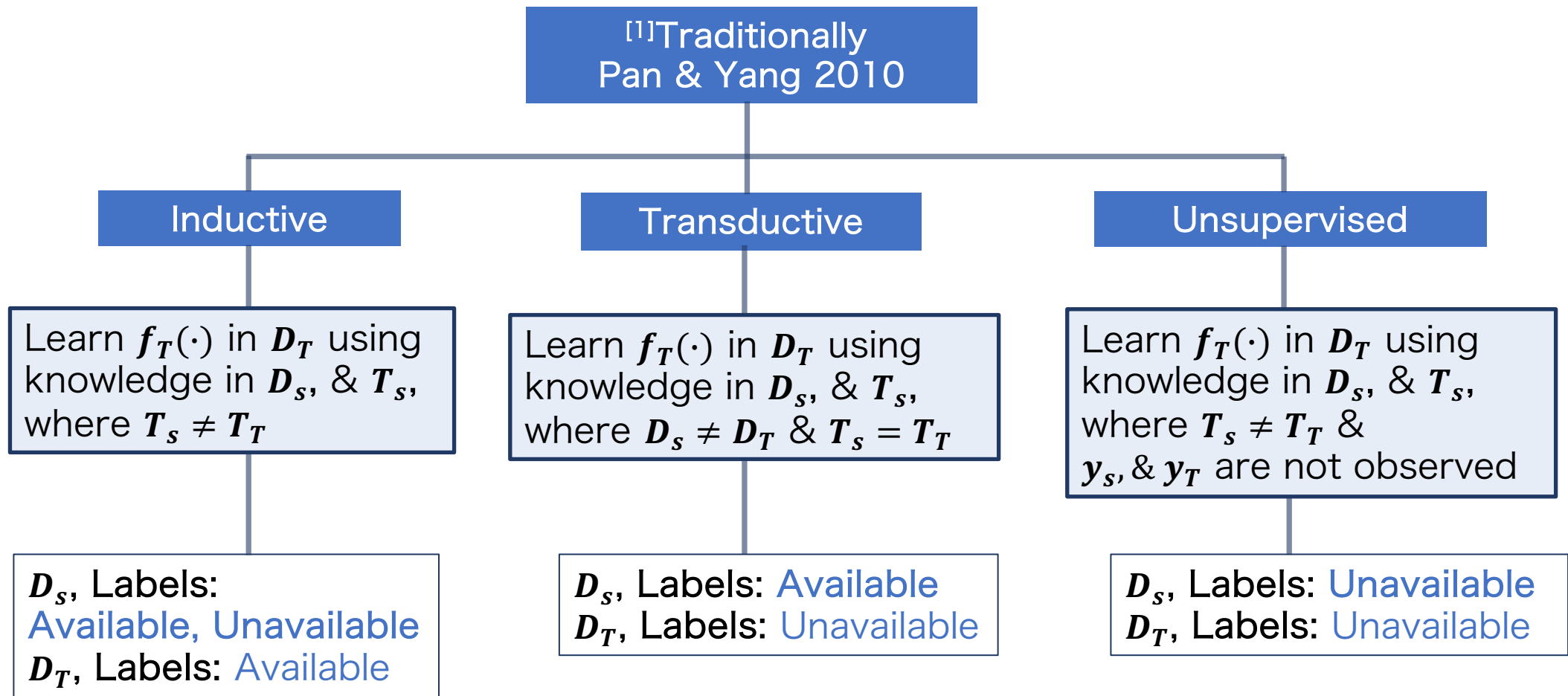
# Categorization of Transfer Learning Problems



Pan & Yang: A survey on Transfer Learning. IEEE 2010
Weiss et al., A survey of transfer Learning. Journal of Big Data 2016

# Categorization of Transfer Learning

[1]Traditionally
Pan & Yang 2010

**Inductive**

Learn $f_T(\cdot)$ in $\boldsymbol{D_T}$ using knowledge in $\boldsymbol{D_s}$, & $\boldsymbol{T_s}$, where $\boldsymbol{T_s} \neq \boldsymbol{T_T}$

$\boldsymbol{D_s}$, Labels: Available, Unavailable
$\boldsymbol{D_T}$, Labels: Available

**Transductive**

Learn $f_T(\cdot)$ in $\boldsymbol{D_T}$ using knowledge in $\boldsymbol{D_s}$, & $\boldsymbol{T_s}$, where $\boldsymbol{D_s} \neq \boldsymbol{D_T}$ & $\boldsymbol{T_s} = \boldsymbol{T_T}$

$\boldsymbol{D_s}$, Labels: Available
$\boldsymbol{D_T}$, Labels: Unavailable

**Unsupervised**

Learn $f_T(\cdot)$ in $\boldsymbol{D_T}$ using knowledge in $\boldsymbol{D_s}$, & $\boldsymbol{T_s}$, where $\boldsymbol{T_s} \neq \boldsymbol{T_T}$ & $\boldsymbol{y_s}, \& \boldsymbol{y_T}$ are not observed

$\boldsymbol{D_s}$, Labels: Unavailable
$\boldsymbol{D_T}$, Labels: Unavailable

Dai et al.,  Boosting for Transfer Learning ICML' 07
Arnold et al., A Comprehensive study of Methods for Transductive Transfer Learning IEEE' 07
Dai et al., Self-taught Clustering ICML' 08

# Categorization of Transfer Learning Problems

## 1. Homogenous Transfer Learning

$(X_s = X_T$ and $Y_s = Y_T)$
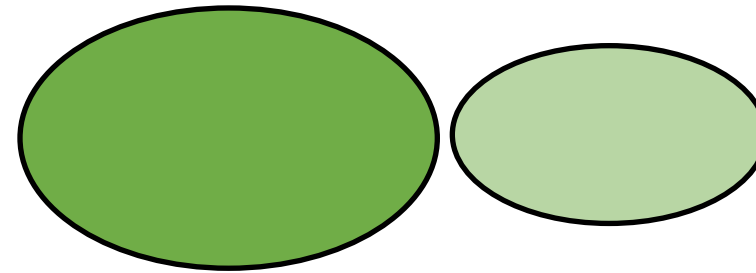
$$X_s \approx X_T$$



Task: Bridge the gap between the source and target data distributions i.e.,

$$P(X_s) \neq P(X_T) \text{ and/or } P(Y_s|X_s) \neq P(Y_T|X_T)$$

## 2. Heterogenous Transfer Learning

$(X_s \neq X_T$ and/or $Y_s \neq Y_T)$

$$X_s \quad \neq \quad X_T$$



Task: Bridge the gap between feature spaces and reduce the problem to homegenous

## 3. Negative Transfer:
If the source domain is not very similar to the target domain, the information learned from the source can have a detrimental effect on a target learner.

Weiss et al., A survey of transfer Learning. Journal of Big Data 2016

# Categorization of Transfer Learninng Solutions

Table 2. Homogenous Transfer learning Approaches

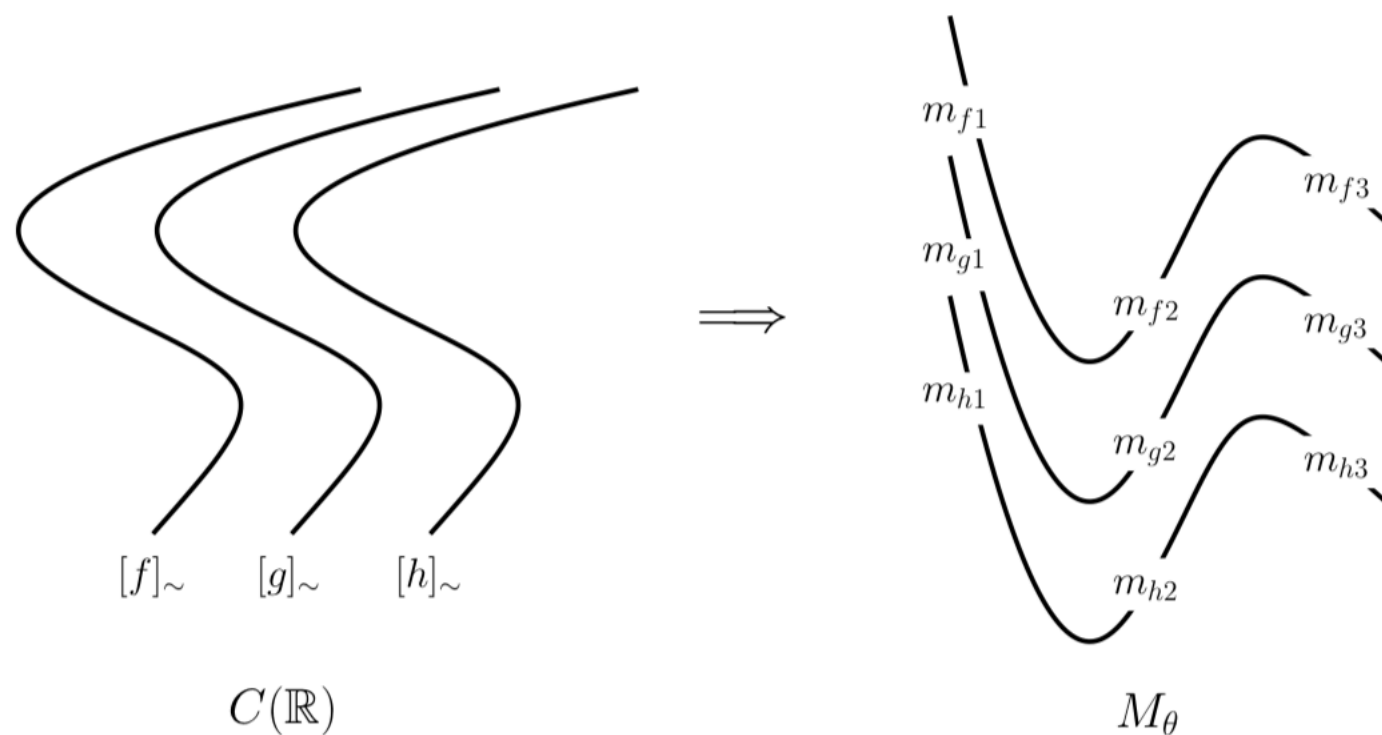| Transfer Learning Approaches | Description |
|---|---|
| Instance-transfer | Try to re-weight samples in the source domain for use in the target domain. *Sugiyama et al., 2008, Yao et al., 2010, Asgarian et al., 2018* |
| Feature-based-transfer | Aim to reduce gap between marginal and conditional distribution between source and target domains. *Long et al., 2014, Oquab et al., 2014, Pan et al., 2011.*<br>Two transformation groups:<br>1. **Asymmetric:** Transforms one of the domain into the other *[Hoffman et al., 2014]*.<br>2. **Symmetric:** Transforms both domains to a common latent space *[Ganin et al., 2014]* |
| Parameter-Transfer | Discover shared parameters or priors between the source domain and target domain models, which can benefit for transfer learning. *Duan et al., 2012, Yao et al., 2010,* |
| Relational-knowledge-transfer | Transfer knowledge through learning a common relationship between source and target domain. *Li et al., 2012, Yang et al., 2018* |
| Hybrid-based | Transfer through both instance and shared parameters *[Xia et al., 2013]* |

# DISCUSSION ON TRANSFER LEARNING

Fariz Ikhwantri

ABI-RIKEN AIP, Tokyo Institute of Technology

# Current Work on Transfer Learning

**Formalization of Relation Between Task**

- Relatedness
  - From several equivariance of Continuous function C(R)
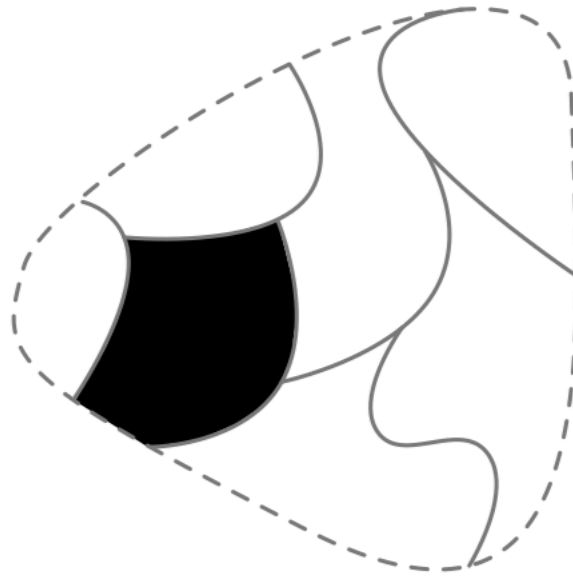


Petangoda, Janith C. et al. "A Foliated View of Transfer Learning." 2020.

# Current Work on Transfer Learning

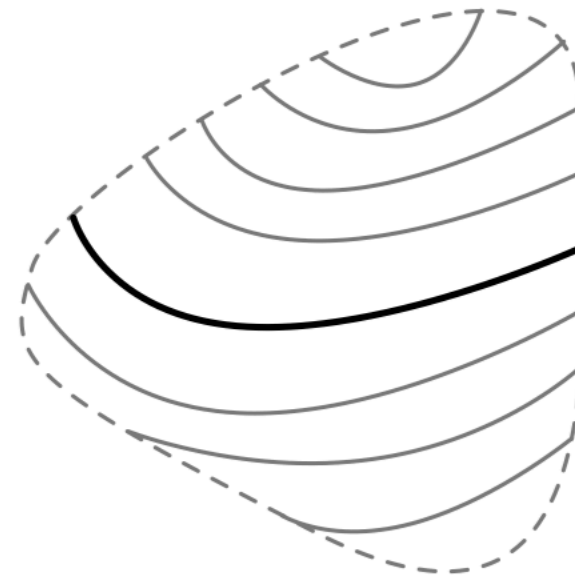## Formalization of Relation Between Task
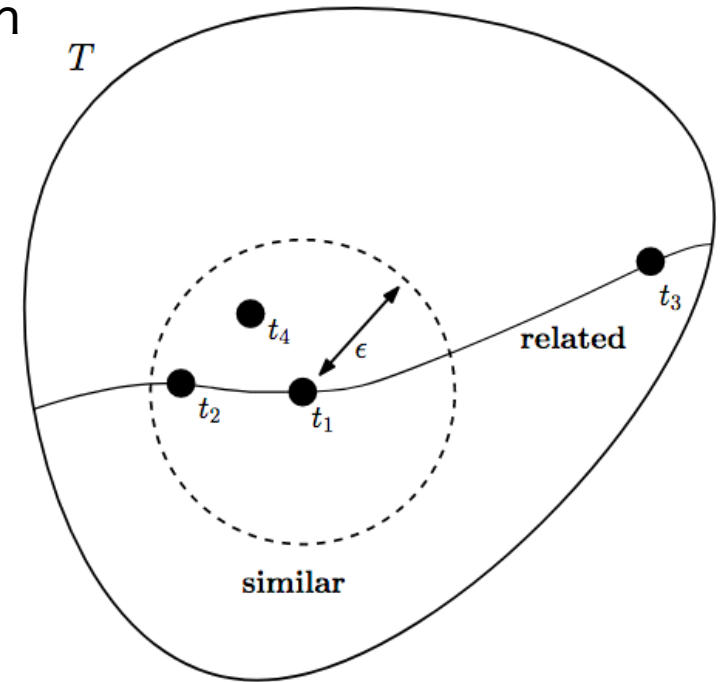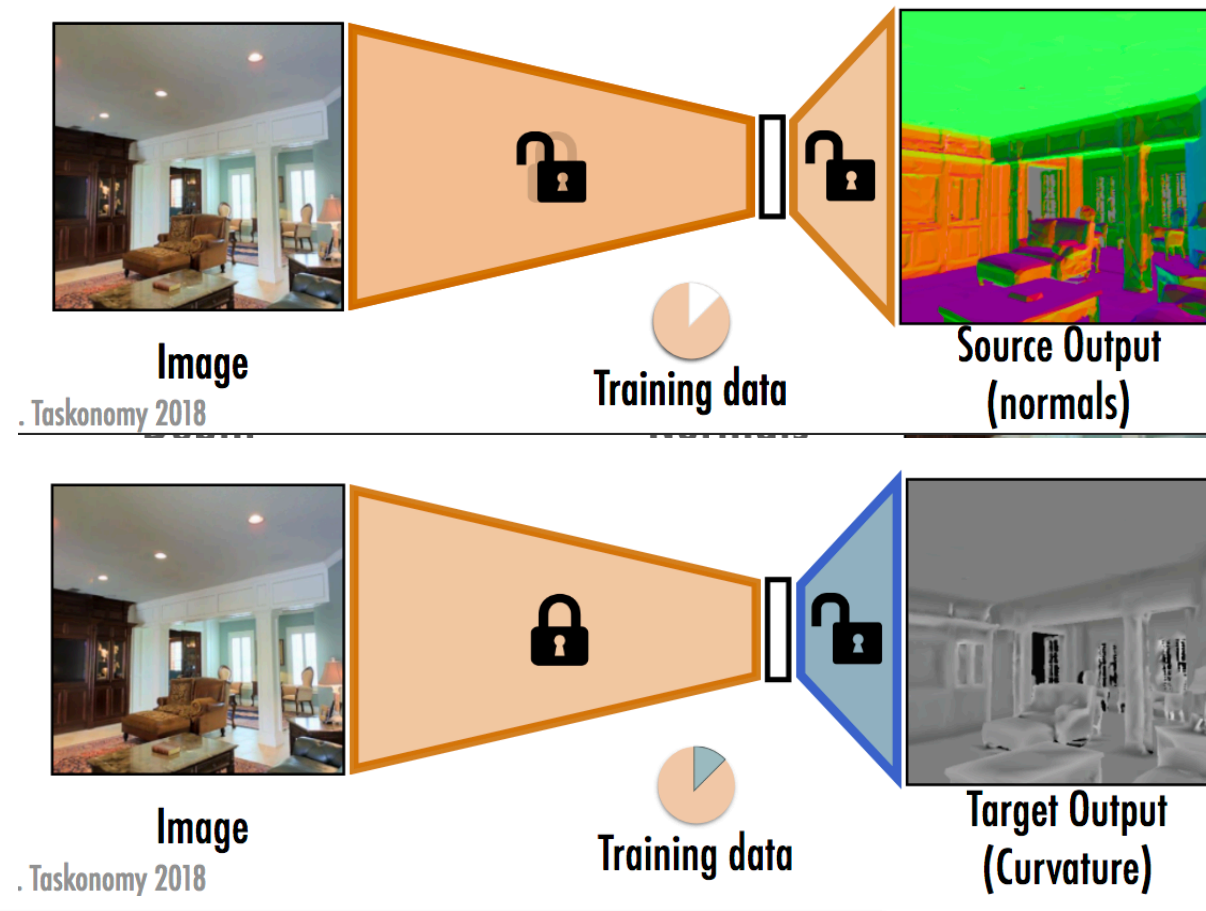
- Relatedness vs Similarity in a Topological



Tessellation    vs    Parallel Spaces

Petangoda, Janith C. et al.  "A Foliated View of Transfer Learning."  2020.

**Formalization of Relation Between Task**

- Related
  - Defined in transformation of function f to g
- Similarity
  - Defined in geometric distance between two task f and g
  - f and g are similar iff $\rho$ (f, g) < ε, where ε is to be chosen



Petangoda, Janith C. et al. "A Foliated View of Transfer Learning." 2020.

**Relation Between Task**

- Task taxonomy example in Computer Vision (Zamir et al., 2018)
    1. Specific Task Training
    2. Transfer Learning Model
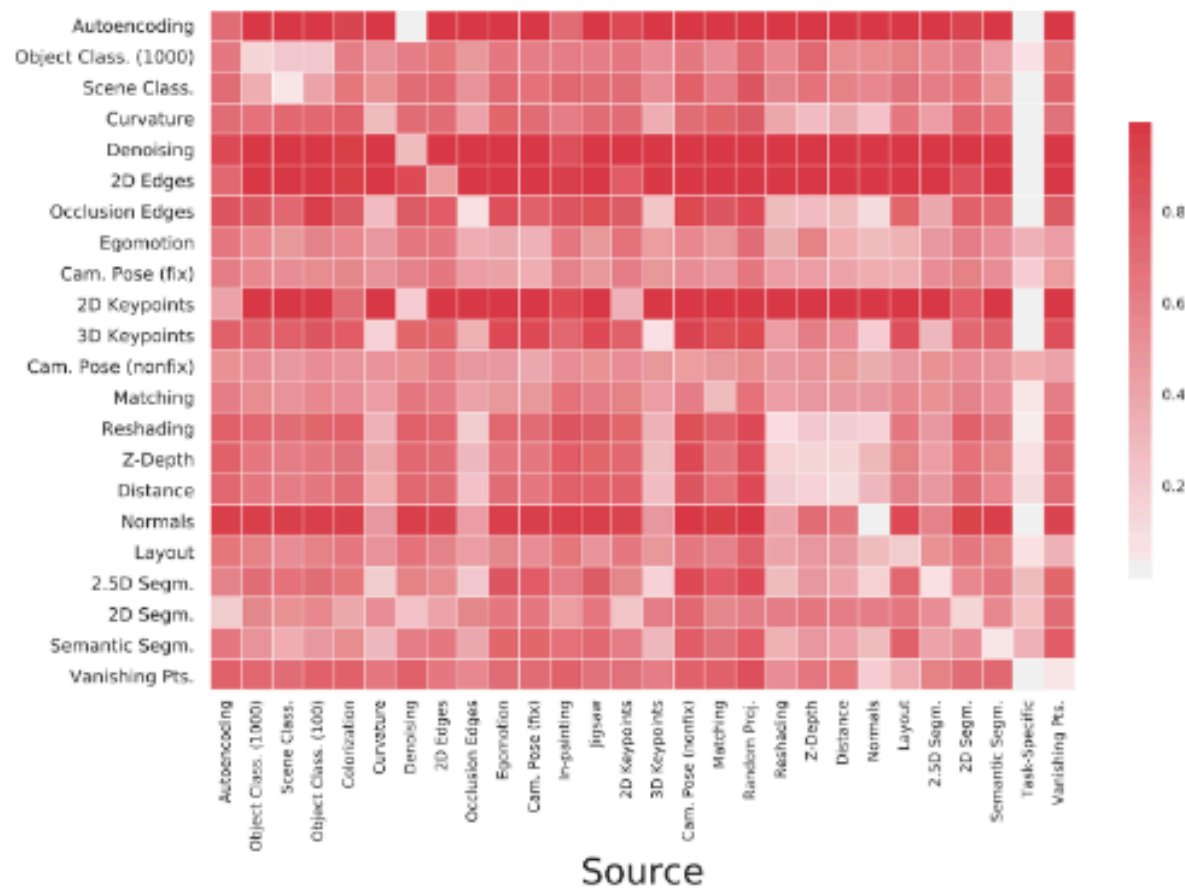        a. One to one task
        b. Many to one task



Image
. Taskonomy 2018

Training data

Source Output (normals)

Image
. Taskonomy 2018

Training data

Target Output (Curvature)

**Zamir et al. "Taskonomy: Disentangling Task Transfer Learning." 2018 CVPR.**

**Relation Between Task**

- Task taxonomy example in Computer Vision (Zamir et al., 2018)
    1. Specific Task Training
    2. Transfer Learning Model
        a. One to one task
        b. Many to one task
    3. Aggregate normalized (ordinal) raw-loss/evaluation between transfer model in Pairwise Matrix
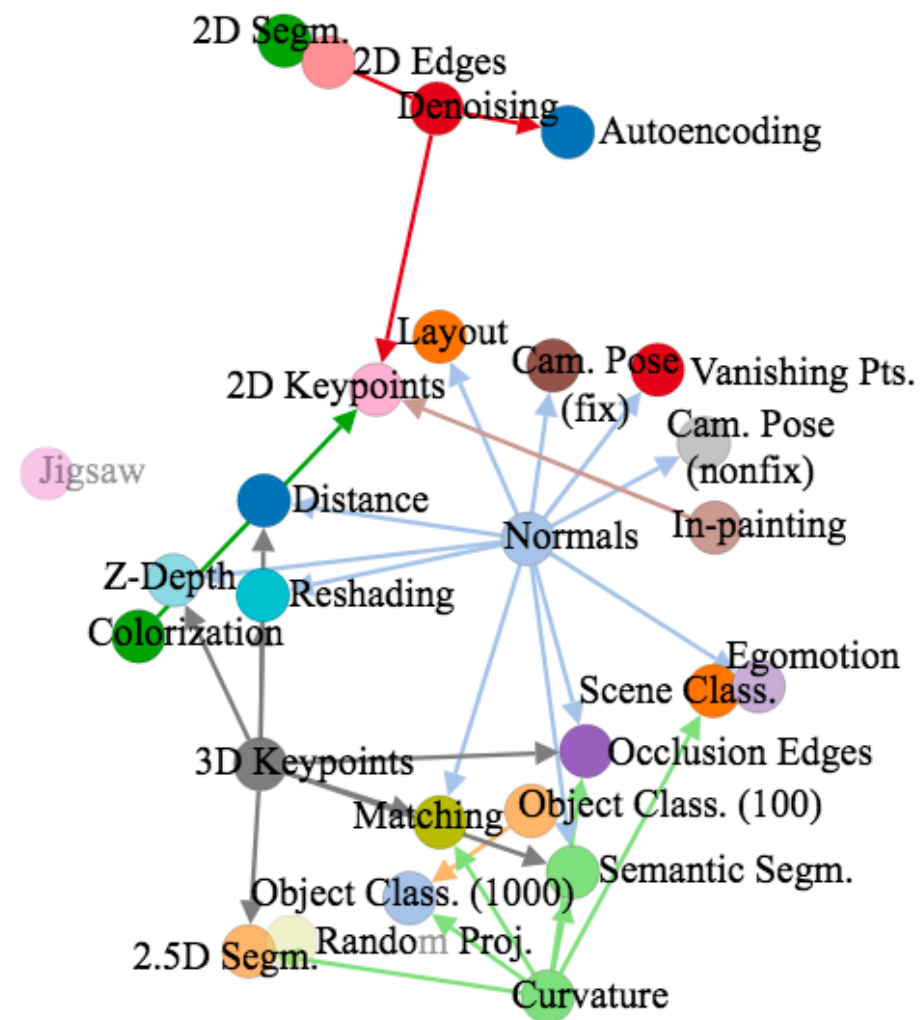


Zamir et al. "Taskonomy: Disentangling Task Transfer Learning." 2018 CVPR.

## Relation Between Task

- Task taxonomy example in Computer Vision (Zamir et al., 2018)

  1. Specific Task Training

  2. Transfer Learning Model

     a. One to one task

     b. Many to one task

  3. Aggregate normalized (ordinal) raw-loss/evaluation between transfer model in Pairwise Matrix

  4. Compute Global Taxonomy



**Zamir et al. "Taskonomy: Disentangling Task Transfer Learning." 2018 CVPR.**

# Current Work on Transfer Learning

**Analyzing Model Weight in Transfer Learning**

- Inductive Bias
  - What information that source task contains ??
- Transferred Knowledge: weight Similarity between Random and Pre-train (**Neyshabur et al. 2020**)
  - Centered Kernel Alignment (**Kornblith et al., 2019**)

$$\text{CKA}(XX^{\text{T}}, YY^{\text{T}}) = \frac{||Y^{\text{T}}X||_{\text{F}}^2}{||X^{\text{T}}X||_{\text{F}}||Y^{\text{T}}Y||_{\text{F}}}$$

$$= \frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \lambda_X^i \lambda_Y^j \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2}{\sqrt{\sum_{i=1}^{p_1} (\lambda_X^i)^2} \sqrt{\sum_{j=1}^{p_2} (\lambda_Y^j)^2}}.$$

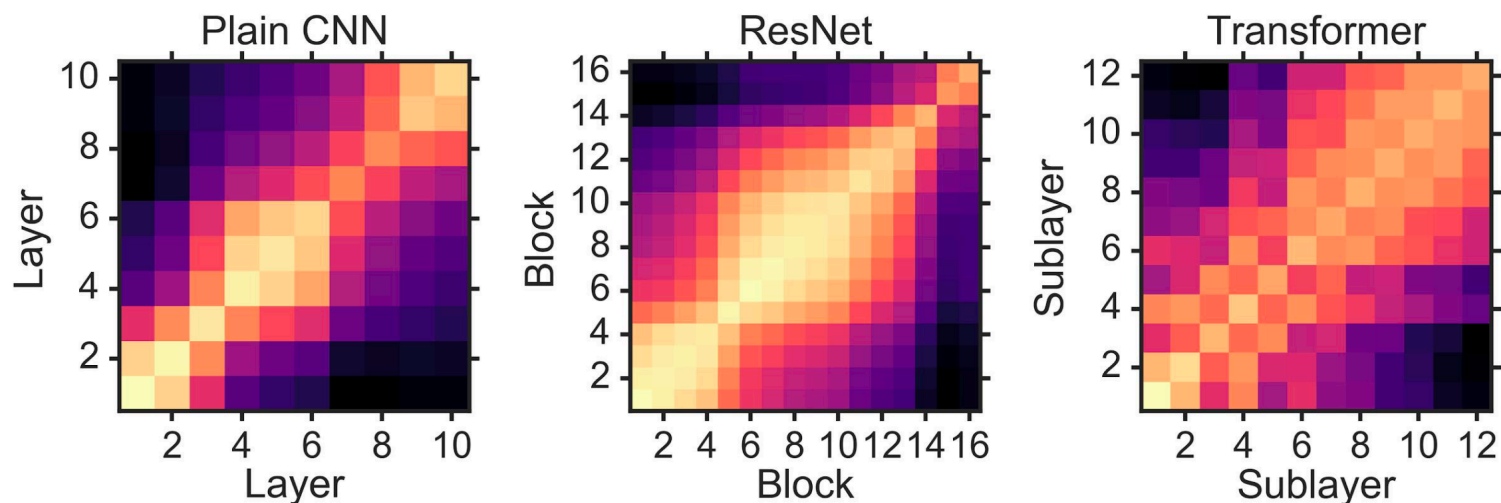Neyshabur et al. "What is being transferred in transfer learning?" (2020).
Kornblith et al. "Similarity of Neural Network Representations Revisited." ICML (2019).
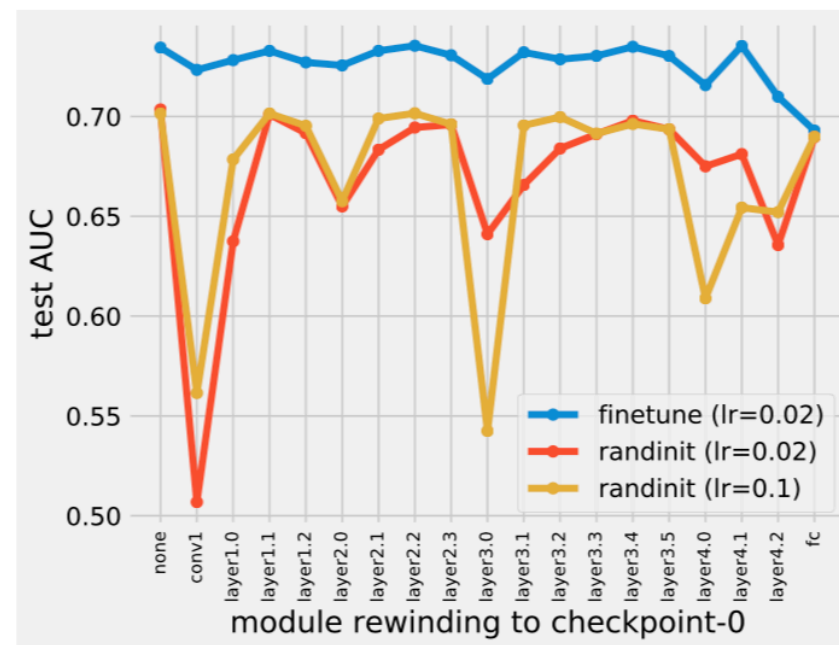
**Analyzing Model Weight in Transfer Learning**

- Inductive Bias
  - What information that source task contains ??
- Transferred Knowledge: weight Similarity between Random and Pre-train (**Neyshabur et al. 2020**)
  - Centered Kernel Alignment (**Kornblith et al., 2019**)

**A Sanity Check for Similarity**



Neyshabur et al. "What is being transferred in transfer learning?" (2020).
Kornblith et al. "Similarity of Neural Network Representations Revisited." ICML (2019).

**Analyzing Model Weight in Transfer Learning**

- Inductive Bias
  - What information source task contains ??
- Transferred Knowledge : weight Similarity between Random and Pre-train (**Neyshabur et al. 2020**)



Table 1: Feature similarity for different layers of ResNet-50, target domain CHEXPERT

| models/layer | conv1 | layer 1 | layer 2 | layer 3 | layer 4 |
|---|---|---|---|---|---|
| P-T & P | 0.6225 | 0.4592 | 0.2896 | 0.1877 | 0.0453 |
| P-T & P-T | 0.6710 | 0.8230 | 0.6052 | 0.4089 | 0.1628 |
| P-T & RI-T | 0.0036 | 0.0011 | 0.0022 | 0.0003 | 0.0808 |
| RI-T & RI-T | 0.0016 | 0.0088 | 0.0004 | 0.0004 | 0.0424 |

RI (random initialization), P (pre-trained model),
RI-T (model trained on target domain from random initialization),
P-T (model trained/fine-tuned on target domain starting from pre-trained weights).

**Neyshabur et al. "What is being transferred in transfer learning?" (2020).**

# Transfer Learning vs Continual Learning

| Difference | Transfer learning | Continual Learning |
|---|---|---|
| **Task Boundaries** | Source → Target | No Boundaries |
| **Learning Goal** | Target Task | Past (Source) and Future (Target) Task |
| **Access to Past Data/Task** | Directly (Instance, Parameter re-use/sharing) | Indirectly (Memory in weight/function space) |
| **??? (Comment welcome)** | ??? | ??? |
| **Similarity** | Feature / Parameter Re-use ??? (Comment welcome) | |

# Example: Transfer Learning vs Continual Learning

## Representation Across Task

- Cross-Entropy loss over large number of Classes (e.g. |Vocabulary|)
- Generative Model



McCann, B. et al. "The Natural Language Decathlon: Multitask Learning as Question Answering." 2019
Sun, Fan-Keng et al. "LAMOL: LAnguage MOdeling for Lifelong Language Learning." ICLR 2020

# Example: Transfer Learning vs Continual Learning

## Representation Across Task

| Methods | SST SRL WOZ | SST WOZ SRL | SRL SST WOZ | SRL WOZ SST | WOZ SST SRL | WOZ SRL SST | Average | Std |
|---|---|---|---|---|---|---|---|---|
| Fine-tuned | 50.2 | 24.7 | 62.9 | 31.3 | 32.8 | 33.9 | 39.3 | 12 |
| EWC | 50.6 | 48.4 | 64.7 | 35.5 | 43.9 | 39.0 | 47.0 | 8.7 |
| MAS | 36.5 | 45.3 | 56.6 | 31.0 | 49.7 | 30.8 | 41.6 | 8.9 |
| GEM | 50.4 | 29.8 | 63.3 | 32.6 | 44.1 | 36.3 | 42.8 | 11 |
| LAMOL$_{GEN}^{0}$ | 46.5 | 36.6 | 56.6 | 38.6 | 44.9 | 45.2 | 44.8 | 6.0 |
| LAMOL$_{GEN}^{0.05}$ | 79.6 | 78.9 | 73.1 | 73.7 | 68.6 | 75.7 | 74.9 | 3.4 |
| LAMOL$_{GEN}^{0.2}$ | 80.0 | 80.7 | 79.6 | 78.7 | 78.4 | 80.5 | **79.7** | 0.8 |
| LAMOL$_{TASK}^{0}$ | 41.0 | 33.5 | 50.1 | 41.9 | 49.3 | 41.5 | 42.9 | 5.2 |
| LAMOL$_{TASK}^{0.05}$ | 77.3 | 76.9 | 78.1 | 74.7 | 73.4 | 75.8 | 76.0 | 1.5 |
| LAMOL$_{TASK}^{0.2}$ | 79.4 | 79.9 | 80.1 | 78.7 | 79.8 | 79.0 | 79.5 | **0.5** |
| LAMOL$_{REAL}^{0.05}$ | 81.0 | 78.9 | 80.1 | 80.9 | 77.7 | 78.0 | 79.4 | 1.2 |
| LAMOL$_{REAL}^{0.02}$ | 81.8 | 80.6 | 81.6 | 81.2 | 80.4 | 80.5 | 81.0 | 0.5 |
| Multitasked | | | | 81.5 | | | | |

Sun, Fan-Keng et al. "LAMOL: LAnguage MOdeling for Lifelong Language Learning." ICLR 2020

# Possible Direction on Transfer for Continual Learning

**Analyzing weight space in transfer learning**

- On different architecture

- Continual Learning

  - Could useful for identifying general parameter across task

- On heterogenous transfer learning (output space is different)

  - e.g. classification vs structured prediction