

# Learning Twofold Heterogeneous Multi-Task by Sharing Similar Convolution Kernel Pairs

Quan Feng<sup>a,b</sup>, Songcan Chen<sup>a,b,\*</sup>

<sup>a</sup>*College of Computer Science & Technology, Nanjing University of Aeronautics & Astronautics, Nanjing, Jiangsu, 211106, China*

<sup>b</sup>*MIIT Key Laboratory of Pattern Analysis and Machine Intelligence, Nanjing University of Aeronautics & Astronautics, Nanjing, Jiangsu, 211106, China*

---

## Abstract

Heterogeneous multi-task learning (HMTL) is an important topic in multi-task learning (MTL). Most existing HMTL methods usually solve either scenario where all tasks reside in the same input (feature) space yet unnecessarily the consistent output (label) space or scenario where their input (feature) spaces are heterogeneous while the output (label) space is consistent. However, to the best of our knowledge, there is limited study on twofold heterogeneous MTL (THMTL) scenario where the input and the output spaces are both inconsistent or heterogeneous. In order to handle this complicated scenario, in this paper, we design a simple and effective multi-task adaptive learning (MTAL) network to learn multiple tasks in such THMTL setting. Specifically, we explore and utilize the inherent relationship between tasks for knowledge sharing from similar convolution kernels in individual layers of the MTAL network. Then in order to realize the sharing, we weightedly aggregate any pair of convolutional kernels with their similarity greater than some threshold  $\rho$ , consequently, our model effectively performs cross-task learning while suppresses the intra-redundancy of the entire network. Finally, we conduct end-to-end training. Our experimental results demonstrate the effectiveness of our method in comparison with the state-of-the-art counterparts.

---

\*Corresponding author

*Email address:* `s.chen@nuaa.edu.cn` (Songcan Chen)

## 1. Introduction

MTL aims to improve the generalization performance of individual tasks by learning multiple related tasks simultaneously. It has been successfully applied in computer vision [1], natural language processing [2], speech recognition [3], medical images processing [4], autonomous vehicles [5] and so on. Existing MTL methods assume that the inputs (feature) spaces and outputs (label) spaces of individual tasks are homogeneous (e.g., they have the same feature space and label space) [6], [7], [8]. This assumption is only applicable in limited real world scenarios.

In practice, general HMTL methods also likewise assume that 1) the inputs (feature) spaces are homogeneous yet the outputs (label) spaces are inconsistent, for example, [9] predicts the heterogeneous attributes of face images; [10] evaluates poses of persons; 2) the input (feature) spaces are heterogeneous while the output (label) space is consistent, for example, [11] uses the features of heterogeneous domains for image classification; [12] learns heterogeneous domain metrics for text classification. However, there are more general scenarios in reality, where the input and output spaces are heterogeneous (i.e., twofold heterogeneity), thus posing a big challenge for HMTL. To the best of our knowledge, there has been almost no study yet to focus on the scenario.

In order to deal with such a complicated scenario, firstly, let us clarify the two unique issues involved that will affect the performance of the HMTL methods. The first is how to model the relationships among heterogeneous tasks as done for homogeneous tasks. Due to the existence of the heterogeneity among those tasks, there do more likely exist strong similarity, weak similarity, which in turn leads to positive, negative correlations [13]. Therefore, modeling such unknown relationships is important yet challenging. A recent work [14] has shown that we can capture such relationship among tasks by embedding a shared unit into the multi-layer perceptron to reflect the specific information of tasks. While by

sharing the hyper-parameters, [15] designs a single pipeline HMTL network for heterogeneous pose evaluation and action recognition of humans.

Secondly, how to obtain useful shared information during the learning process. Since various information in heterogeneous tasks may have an implicit interrelationship with a single task or all tasks, thus undoubtedly affecting the performance of HMTL [16]. Therefore, it is very important to obtain effective common information for individual tasks by the twofold heterogeneous multi-task (THMT) learning. Recent research work has designed a feature matching network to capture these shared information in heterogeneous tasks for HMTL [17].

Even so, most previous HTML works relied on manually designing a shared layer for tasks including the latter layers until a branching layer. For example, [18] uses a BlitzNet shared layer for object detection and semantic segmentation. [19] designs a dual-attribute-aware hierarchical network (DARN) for cross-domain image retrieval at fully connected layers. These constraints bring about comprehensive cross-task knowledge sharing. To overcome these, we designed a MTAL network. Different from existing methods, the designed network automatically selects similar convolution kernel pairs across-tasks to obtain common knowledge. Specifically, as shown in Fig (1): in the learning process, we automatically capture the inherent relationship between tasks by exploiting the similarity between convolution kernels. Next, we use a soft threshold to select the suitable convolution kernel pairs for aggregation to form a set of new kernel bank to learn individual tasks. Finally, we train the entire network in an end-to-end manner. Our experimental results demonstrate the effectiveness of our method in comparison with the state-of-the-art counter-parts.

In summary, our contribution can be summarized as four folds below:

- 1) We propose a set of novel MTAL networks, which provides a new way of solving the information sharing problem in THMTL.
- 2) We design a soft threshold to select the suitable kernel pairs to prevent possible negative transfer caused by existing MTAL network learning.
- 3) We design a novel sharing strategy, which utilizes the similarity between

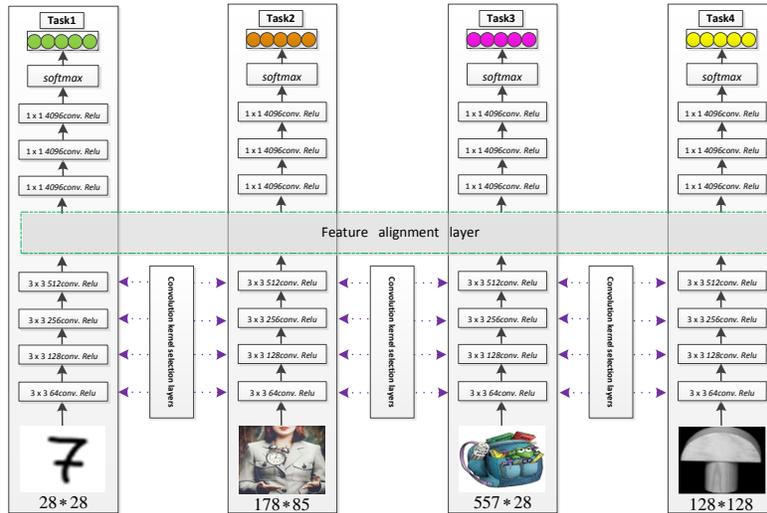


Figure 1: Our proposed MTAL networks. The network is composed of four identical individual task networks, and the input/output task (feature/label) spaces of each task network are heterogeneous, respectively. The black solid line mullion represents the convolution kernel selection layer. Green, orange, purple, and yellow solid circles indicate different output tasks.

convolution kernels in the same individual layers to formulate the intrinsic relationship between tasks and aggregate them to suppresses the intra-redundancy of the entire network.

4) We conduct experiments on eight public datasets and compare the state-of-the-art methods to validate the effectiveness of our method.

The rest of this article is organized as follows: in Section 2, we review the related methods of homogeneous and heterogeneous multi-task learnings respectively; in Section 3, we introduced the specific implementation of our method; in Section 4, we conduct experiments on the eight public benchmark datasets and compare the state-of-the-art methods to show the effectiveness of our method. Finally, Section 5 concludes this paper with future research directions. The code is available at <http://parnec.nuaa.edu.cn/3021/list.htm>.

## 2. Related Work

At present, deep learning has achieved remarkable results in computer vision and MTL. Our work is also based on deep learning framework, thus in this section, we mainly review deep learning-based homogeneous MTL and heterogeneous MTL in recent years, respectively.

### 2.1. Homogeneous MTL

Such methods generally assume that the input (feature) space and the output (label) space are homogeneous between tasks (i.e., the same feature space and label space). In such a scenario, existing methods share the same features between tasks, and they are applied in object detection, image classification, medical image segmentation, and so on. To date, there have had approaches proposed, we divide them into three categories according to the homogeneous MTL sharing architectures: hard sharing, soft sharing, and hybrid sharing.

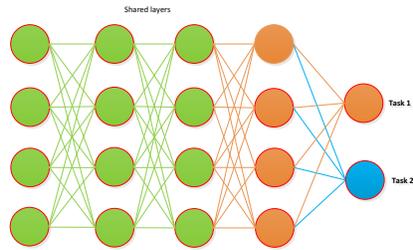


Figure 2: Homogeneous MTL based on hard sharing.

1) Hard sharing MTL. As shown in Fig.(2), in the architecture, all tasks share their knowledge in the same hidden space. For example, [20] uses a mutual representation of hard parameters shared for personalized stress recognition. [21] predicts the mortality of diverse rare diseases by initializing shared parameters in hidden space. The MTL methods of this kind can not only assist in effective learning between tasks but also minimize the risk of over fitting during the training process [22]. However, it is difficult to handle loosely related tasks [23].

2) Soft sharing MTL. As shown in Fig.(3), in its architecture, all task models and parameters are independent, and the distance between the model parame-

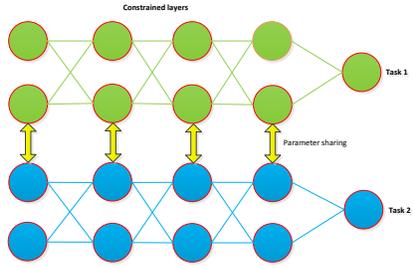


Figure 3: Homogeneous MTL based on soft sharing.

ters is regularized to obtain similar parameters for joint learning. For example, [24] uses a neural network parser to share parameters between the source domain language and the target domain language for natural language processing. [25] utilizes tensor trace norms to regularize parameters in multiple networks for image recognition. [26] learns multiple different tasks by jointly learning transferable features and multi-linear relationships between tasks and features in a fully connected layer. However, this method obviously relies on a predefined shared structure, and the model has poor generalization performance for the new tasks.

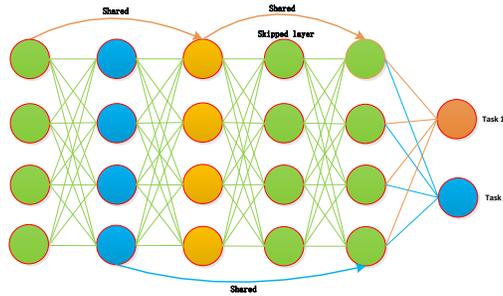


Figure 4: Homogeneous MTL based on hybrid sharing.

3) Hybrid sharing MTL. As shown in Fig.(4), in the hybrid shared architecture, specific task strategies are used to select which layer of multiple task network models can perform shared learning. For example, [27] uses task-specific strategies to learn shared patterns for image semantic and normal segmentation. The advantage of this method is that the number of parameters does not in-

crease as the number of tasks does. The disadvantage is that it cannot handle heterogeneous tasks.

## 2.2. Heterogeneous MTL

HMTL usually assumes that the input (feature) space is the same, but the output (label) space is unnecessarily consistent, or the input (feature) space is heterogeneous and the output (label) space is consistent. Such methods share heterogeneous features to make some predictions such as the heterogeneous attributes of human faces and human poses, and classify various images or texts. For these existing works, we can also divide them into two categories to the HMTL sharing architectures: hierarchical sharing and sparse sharing.

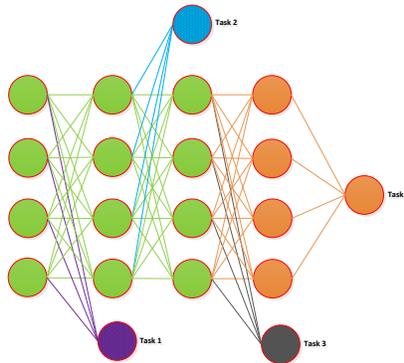


Figure 5: HMTL based on hierarchical sharing.

1) Hierarchical sharing HMTL. As shown in Fig.(5), in the architecture, various tasks perform hierarchical sharing learning in multiple task networks. For example, [28] is used for image classification by weighting and sharing similar features among different levels in the multi-task network. [29] learns a set of shared semantic representations from the bottom of the supervised model to multiple task hierarchies at the top of the model for natural language processing.

2) Sparse sharing HMTL. As shown in Fig.(6), in the sparse sharing architecture, an HMT network is composed of various task networks and is sparse. For example, [23] extracts sub-networks of different tasks from an over-parameterized base network and uses masks to sparsify the features of different individual task

networks to retain partly shareable features and delete irrelevant features for the Part-of-Speech , Named Entity Recognition and Chunking. The advantage of this method is to retain some useful information to help a specific task, and remove the useless information. Obviously, sparse sharing from the perspective of information sharing can be regarded as an example of hard sharing and hierarchical sharing.

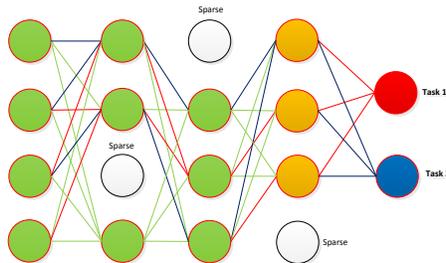


Figure 6: HMTL based on sparse sharing.

On the other hand, there are also typical works not using deep learning architecture in dealing with HMTL problems. For example, [30] uses a linear discriminant analysis of multi-task expansion algorithm (MTDA) for processing tasks with different data representations. The method learns a separate heterogeneous feature transformation for each task. It’s purpose is to alleviate the problem of insufficient label data during learning, and can jointly handle binary and multi-class problems for each task. [31] uses non-negative matrix factorization to learn sharing common semantic features in the feature space across heterogeneous tasks for image classification. [32] learns a shared graph Laplacian matrix in unified image features for visual clustering of different modalities.

### 3. Our Method

Our method is inspired by the fact that the activation maps generated by similar convolution kernels in each convolution layer are also similar [33] [34]. For this reason, our method aims to automatically select similar convolution kernel pairs across-task to obtain common knowledge for THMTL. As shown in

Fig.(7), our method consists of measuring the similarity of convolution kernels, selecting suitable kernel pairs, and aggregating them to form a set of new kernel bank. Specifically, Section 3.1 formally formulates the problem. Section 3.2 illustrates how to find similar convolutional kernels in the convolutional layer of the network and utilize a soft threshold to select suitable kernel pairs. How to aggregate these kernel pairs is discussed in Section 3.3. Section 3.4 derives the objective function of THMTL to deploy such a mechanism.

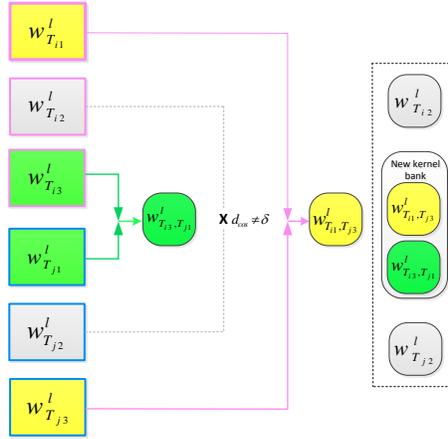


Figure 7: Convolution kernel pairs sharing mechanism. In the  $l$ -th layer of the MTAL network, the model measures the similarity of convolution kernels in task networks  $T_1$  and  $T_2$  while selects the suitable convolution kernel pairs through the threshold for aggregation to form a set of new kernel bank. The purple and blue rectangles represent the convolution kernels of task networks  $T_1$  and  $T_2$ , respectively. The yellow and cyan octagons represent the new kernels aggregated. The solid line box on the right represents the new kernel bank. The gray octagons represent independent kernels that are not shared to continually learning their corresponding tasks.

### 3.1. Problem formulation

Suppose we are given  $N$  tasks  $\{T_i\}_{i=1}^N$ , some of which are related or unrelated. The training dataset  $D_i = \{(x_h^i, y_h^i)\}_{h=1}^{n_i}$  for  $T_i$  contains  $n_i$  samples with  $x_h^i \in \mathbb{R}^{d_i}$  and its corresponding label  $y_h^i \in \{1, \dots, c_i\}$ , where  $d_i$  and  $c_i$  are the numbers of dimensions and classes in the dataset  $D_i$ , respectively. We do not

assume that the datasets from different tasks share the same feature space, so the dimensions of feature spaces can be different. This makes MTAL applicable under more general settings than most existing handling heterogeneous task methods.

### 3.2. Measure the similarity of the convolution kernels

The goal of HMTL is to improve the generalization performance of individual tasks by sharing their relevant information. However, such information is generally unavailable and implicit. To mine this information, there have been some works proposed, they can mainly be summarized as HMTLs based on shallow and deep networks according to the architecture. Here we only focus on the deep HMTL under the same network structures since the deep networks has been more satisfactorily applied in the MTL fields, which can be subdivide into two types according to sharing way:

1) *feature-sharing-based*. This type of method is to learn a common representation from different task features in the same hypothetical space, thereby effectively helping to learn each task. Typically, [35] uses tensors to represent the feature interactions from different tasks in the same shared subspace for inductive transfer of related information, thereby providing better generalization performance for multi-task models. [36] uses the  $\ell_{\{1,2\}}$  norm to regularize the weight matrix to extract relevant features between tasks for learning multi-tasks with different feature dimensions. These methods will lower their generalization performance when the tasks are unrelated or the distributions of the data are different.

2) *parameter/weight-sharing-based*. This type of methods mainly learns multiple tasks jointly by sharing common parameters hidden in the weights of different task models. According to different implementation manners, we further divide it into the following three sub-types: 1) sharing the parameters between task models based on the same space assumption: for example, [37] shares weight parameters in different tasks in the same subspace for face detection, key point positioning, pose estimation, and gender prediction. 2) parameter

matrix factorization based: e.g., [38] uses the matrix tri-factorization to process collective matrices associated with different tasks and performs joint learning to predict two types of drug-disease associations. 3) equal prior shared assumptions based: e.g., [39] uses a kind of meta data (i.e., contextual attributes) as a priori information to capture the relationships between different tasks for multiple tasks clustering. The above methods mainly utilize model parameters to associate different tasks. However, it is a great challenge to design an appropriate implementation manner in different tasks to obtain shared parameters. To demonstrate the importance of these parameter sharing manners in HMTL, we further studied the current two excellent deep network frameworks based on parameter sharing that are typically capable of handling heterogeneous tasks. E.g., the cross stitch network [40] connects the low-level layers of different single-task networks by learning the parameter  $\alpha$  in the task features, which is defined as

$$\begin{bmatrix} \tilde{x}_A^{ij} \\ \tilde{x}_B^{ij} \end{bmatrix} = \begin{bmatrix} \alpha_{AA} & \alpha_{AB} \\ \alpha_{BA} & \alpha_{BB} \end{bmatrix} \begin{bmatrix} x_A^{ij} \\ x_B^{ij} \end{bmatrix} \quad (1)$$

where  $x_A^{ij}$  and  $x_B^{ij}$  represent the activation-maps learned in the two sub-networks, and are linearly combined to realize the information interaction in the neurons of the hidden layer, thereby outputting new hidden features  $\tilde{x}_A^{ij}$  and  $\tilde{x}_B^{ij}$ . The SubNetwork Routing (SNR) [41] decomposes the shared-bottom module of the MMOE [42] network into three subnetworks and shares the parameters  $z$  to learn relevant information among different tasks, which is defined as

$$\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} z_{11}W_{11} & z_{12}W_{12} & z_{13}W_{13} \\ z_{21}W_{21} & z_{22}W_{22} & z_{23}W_{23} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} \quad (2)$$

where  $u_1$ ,  $u_2$  and  $u_3$  represent the outputs by the hidden layer of each lower-level sub-networks,  $W$  is the transformation matrix,  $v_1$  and  $v_2$  represent the inputs of the higher-level sub-networks of the next layer. Instead of the above methods of sharing neuron units in the hidden layers of the network, our method is to

share the set of neurons in (similar) convolution kernel pairs to improve the generalization and efficiency of the MTAL network.

In order to implement the above sharing, we assume that in the  $l$ -th layer for an individual task, there is a group of  $m$  convolution kernels  $\mathbf{W}_{T_i}^l = \{w_{T_{i1}}^l, w_{T_{i2}}^l \dots, w_{T_{im}}^l\}$ , ( $i = 1, 2, \dots, N$ ), which form  $\mathbf{W}^l = \{\mathbf{W}_{T_1}^l, \mathbf{W}_{T_2}^l, \dots, \mathbf{W}_{T_N}^l\}$  in the MTAL network. Then, we measure the similarity between convolution kernel pairs by the following cosine similarity to capture the relationship between tasks:

$$d_{\cos}(\text{vec}(w_{T_i}^l), \text{vec}(w_{T_j}^l)) = \frac{\text{vec}(w_{T_i}^l)^\top \cdot \text{vec}(w_{T_j}^l)}{\|\text{vec}(w_{T_i}^l)\|_2 \cdot \|\text{vec}(w_{T_j}^l)\|_2} \quad (3)$$

where  $\text{vec}(\cdot)$  represents the vectorization operator.

Next, we just consider sharing the kernel pairs that satisfy formulation (4) while retaining the remaining independent kernels to prevent possible negative transfer caused by existing MTAL network learning.

$$d_{\cos}(\text{vec}(w_{T_i}^l), \text{vec}(w_{T_j}^l)) \geq \delta, \quad \delta \in [0.1, 0.9] \quad (4)$$

### 3.3. Aggregation of Convolution Kernels

To further perform the sharing, we use the weighted aggregation of pairwise  $w_{T_i}^l$  and  $w_{T_j}^l$  to model the shared representation between individual tasks as follows:

$$\left\{ w_{T_i, T_j}^l \mid w_{T_i, T_j}^l = \varphi_{i,j}^l w_{T_i}^l + \varphi_{j,i}^l w_{T_j}^l, \quad i \neq j, \quad \varphi_{j,i}^l + \varphi_{i,j}^l = 1 \right\} \quad (5)$$

where  $\varphi_{i,j}^l$  and  $\varphi_{j,i}^l$  are weight coefficients,  $w_{T_i, T_j}^l$  represents the aggregated kernels. In this way, we reduce the number of weights in the entire network while suppressing intra-redundancy (due to over-parameterization [43]). Our experiments show that the memory of the entire network can be saved by about 13.1% compared to the original structures, as empirically analyzed in 4.6.

For joint learning of multiple tasks in the MTAL network, we gather (5) to form a kernel bank and perform a simple average aggregation as follows:

$$\hat{w}_{T_i}^l = \frac{\sum_{w \in \mathbf{M}_{T_i}^l} w_{\mathbf{M}_{T_i}^l}^l}{|\mathbf{M}_{T_i}^l|} \quad (6)$$

where  $\hat{w}_{T_i}^l$  represents the averaged kernel,  $\mathbf{M}_{T_i}^l$  is the new kernel bank (i.e.,  $\mathbf{M}_{T_i}^l = \{w_{T_i, T_j}^l, \dots, w_{T_i, T_N}^l\}, (j = 2, \dots, N, j \neq i)$ ),  $w_{\mathbf{M}_{T_i}^l}^l$  is the kernel of the bank.

#### 3.4. Objective Function of individual & total task

In the MTAL learning, we use the inputs  $X_{T_i}$  ( $i = 1, 2, \dots, N$ ) and the labels  $Y_{T_i}$  ( $i = 1, 2, \dots, N$ ) to minimize following the individual task loss function  $\mathcal{L}_{T_i}$  as follow

$$\mathcal{L}_{T_i} = - \sum_{h=1}^{n_i} y_h (\log y'_h) + \lambda \|\mathbf{W}_{T_i}\|_F^2 \quad (7)$$

where the first term is the cross-entropy loss of individual tasks, the second term is the  $\ell_2$  norm regularization.  $\lambda$  is an adjustment hyper-parameter. Then, we define the total objective function of the entire network as

$$\mathcal{L}_{T_{total}} = \mathcal{L}_{T_1} + \mathcal{L}_{T_2} + \dots + \mathcal{L}_{T_N} \quad (8)$$

The whole process of the proposed method to solve THMT is summarized in Algorithm (1).

## 4. Experiments

In this section, we use VGG [44] network as a base-model but also other networks such as ResNets [45]. We conduct two sets of experiments on eight public datasets to verify the performance, that is, one set of experiments uses the prior relationships between tasks (e.g., obtained through cross-validation experiments) while the other set does not.

### 4.1. Datasets

We use the following datasets for experiments, and divide 70% for training, and the remaining 30% for testing, which is detailed in Table (1).

The *Office-Caltech* dataset<sup>1</sup> contains the Office-Caltech10 dataset and the Office-Caltech31 dataset, each of which has a total of 2533 samples and is composed of a subset of image datasets from three different databases: Caltech,

<sup>1</sup><https://people.eecs.berkeley.edu/~jhoffman/domainadapt/>

---

**Algorithm 1:** MTAL

---

**Notations:**  $T_1, T_2, \dots, T_N$  denote input heterogeneous tasks;

$\mathbf{W}^l = \{W_{T_1}^l, W_{T_2}^l, \dots, W_{T_N}^l\}$  denotes the convolution kernels of the  $l$ -th layer;  $\eta$  denotes learning rate;  $\delta$  is the threshold;  $\theta$  are the weights of the entire network;  $\mathbf{M}_{T_i}^l$  denotes the new kernel bank.

**Input:**  $T_1, T_2, \dots, T_N, \eta, \delta$

**Output:**  $\theta$

```
1 random initialization  $\mathbf{W}^l$ 
2 repeat
3   for all  $\{(x_h, y_h) \in D_{T_j}\}_{j=1}^N$  do
4     for each convolution layer  $l$  do
5       for each  $T_i$  do
6          $\mathbf{M}_{T_i}^l = \{\}$ 
7         for each  $T_j (j \neq i)$  do
8           The convolution kernel pairs similarity is calculated
9             by (3)
10          If  $d_{\cos}(\text{vec}(w_{T_i}^l), \text{vec}(w_{T_j}^l)) \geq \delta$  : put  $w_{T_i, T_j}^l$  into  $\mathbf{M}_{T_i}^l$ 
11          Convolution kernel sharing by using (5)
12          end for
13          Update weight  $\hat{w}_{T_i}^l$  by (6)
14          end for
15        end for
16        Calculate the output of the current sample sequence  $\{y'_h\}_{h=1}^N$ 
17        Update weight  $\theta$  by using back propagation algorithm
18 until convergence;
```

---

*Amazon*, and Webcam. We select a set of *Amazon* from the Office-Caltech31 dataset, and randomly select 10 categories, with a minimum image size of  $200 * 150$  and a maximum image size of  $900 * 557$ .

The *Office Home* dataset<sup>2</sup> is composed of subsets of image datasets from different fields of *Art*, *Clipart*, *Product*, *Real-World*. Each subset has 65 different categories and 15,500 images. We randomly select 10 classes in the *Art* subset of the *Office Home* with the image size of  $117 * 85$  and  $4384 * 2686$ .

The *Coil-20* dataset<sup>3</sup> is a 20-object grayscale image dataset, consisting of a set of 720 unprocessed images of 10 objects and another set of 1,440 normalized image datasets of 20 objects. The image acquisition comes from placing the object on the electric turntable against a black background, rotate the turntable by 360 degrees to capture the pose of the object with a fixed camera or take an image of the object by rotating the turntable by 5 degrees. We select the first set of unprocessed image datasets for the experiment.

The *Chars74K* dataset<sup>4</sup> is composed of two datasets of English characters and Kannada characters, where the English characters include 3 datasets of *A* (*A-Z*), *a* (*a-z*), and 0-9 handwritten digits (*HD*) with a total of 62 categories, 3410 images, handwritten by 55 volunteers. We use the *A*, *a*, *HD* subsets in the English character datasets, and randomly select 10 categories from the *A*, *a* subset for the experiment.

The *Typographic* dataset<sup>5</sup> is a classic dataset used for machine learning, image classification, and image recognition. It is mainly composed of 0-9 typographic numbers, a total of 10,000 pictures, and the picture size is  $12 * 16$ .

#### 4.2. Comparison Methods

We compare our method with the following five representative methods.

---

<sup>2</sup><http://hemantdv.org/OfficeHome-Dataset/>

<sup>3</sup><https://www.cs.columbia.edu/CAVE/software/softlib/coil-20.php>

<sup>4</sup><http://www.ee.surrey.ac.uk/CVSSP/demos/chars74k/>

<sup>5</sup>[http://www.catalina.com.cn/info\\_249972.html](http://www.catalina.com.cn/info_249972.html)

Table 1: Parameters and settings for each datasets.

Datasets	Number of categories	Size of image	Dataset partition (%)
<i>Amazon</i>	10	150 * 900/557 * 28	70/30
<i>Art</i>	10	117 * 85/4384 * 2686	70/30
<i>Coil-20</i>	10	128 * 128	70/30
<i>HD</i>	10	28 * 28	70/30
<i>Typographic</i>	10	12 * 16	70/30
<i>a</i>	10	28 * 28	70/30
<i>A</i>	10	28 * 28	70/30

Single task<sup>6</sup> baseline: This method uses a single VGG network to learn the predictive model for each independent task.

Multi-task<sup>7</sup> baseline: This method uses multiple identical VGG networks to jointly learn a multi-task prediction model.

MTDA<sup>8</sup> [30]: This method uses linear discriminant analysis to handle multiple tasks represented by different data.

Cross-Stitch network<sup>9</sup> [40]: This method introduces the cross stitch unit in the convolutional neural network of two different tasks for knowledge sharing, thereby improving the learning performance of the network.

NDDR-CNN<sup>10</sup> [46]: This method performs feature fusion at each layer on different tasks to obtain shared information, thereby improving the predicting accuracy of the model.

MTAL<sup>11</sup>: This is our proposed method, which mainly aims at automatically selecting similar convolution kernel pairs across tasks to obtain common information for THMTL.

<sup>6</sup><https://github.com/machrisaa/tensorflow-vgg>

<sup>7</sup>[https://github.com/luntai/VGG16\\_Keras\\_TensorFlow](https://github.com/luntai/VGG16_Keras_TensorFlow)

<sup>8</sup><https://yuzhanghk.github.io/>

<sup>9</sup><https://github.com/helloyide/Cross-stitch-Networks-for-Multi-task-Learning>

<sup>10</sup><https://github.com/ethanygao/NDDR-CNN>

<sup>11</sup><http://parnec.nuaa.edu.cn/3021/list.htm>

MTAL<sub>R</sub><sup>12</sup>: This is our proposed method, which uses the prior relationship between tasks in the experiment.

#### 4.3. Hyper-Parameter Tuning

In the contrasted deep neural network methods, we adjust the hidden units, learning rate, and the number of training steps in each layer according to the parameter settings of the corresponding references. In MTAL, we adjusted the hyper-parameters in the same way and chose the stochastic gradient descent method as the network optimizer. Specifically, we set the learning rate  $\eta$  to 0.01 and  $\lambda$  to 0.1. For  $\delta$ , we have verified through multiple experiments that its value is 0.4 when the tasks are related, but it is 0.55 when the tasks are unrelated. All the deep models are implemented by Tensorflow.

#### 4.4. Results of Model Performance

We respectively show the performance of various methods on the datasets *HD*, *a*, *Typographic*, *A*, *Coil-20*, *Art*, and *Amazon*. The detailed analysis is show as follows:

Firstly, when using the relationship between tasks, we find from Table(2) and Table(4), that the accuracy of MTAL<sub>R</sub> is better than most methods. However, when the relationship between tasks is not used, we find from Table(3) and Table(5) that the accuracy of MTAL is better than other methods.

Secondly, in Table(2) and Table(4) we find that most of the MTL methods are better than the single-task learning method, which proves the effectiveness of jointly learning multiple heterogeneous tasks by exploring the relationship between tasks. In particular, Table(2) shows that all methods are better than single-task learning methods, which indicates that the more similar the relationship between tasks, the better the generalization performance of all methods. In addition, we observe that the results of different multi-task methods are different, which is caused by the differences among tasks.

---

<sup>12</sup><http://parnec.nuaa.edu.cn/3021/list.htm>

Table 2: Performance comparison among various methods that use task relationships on *HD*, *a*, *typographic*, and *A* datasets. Among them, the bold numbers are the best classification results, and the underlined numbers are the second-best classification results.

Methods	Group 1		Group 2	
	<i>HD</i>	<i>a</i>	<i>Typographic</i>	<i>A</i>
Single-task	0.76 ± 0.025	0.80 ± 0.023	0.95 ± 0.011	0.84 ± 0.020
Multi-task	<b>0.83 ± 0.015</b>	0.85 ± 0.015	<u>0.97 ± 0.015</u>	0.86 ± 0.020
MTDA	0.78 ± 0.032	0.82 ± 0.014	0.95 ± 0.015	0.85 ± 0.040
Cross-Stitch	0.81 ± 0.017	0.86 ± 0.012	<u>0.97 ± 0.015</u>	<u>0.88 ± 0.015</u>
NDDR-CNN	0.80 ± 0.040	<u>0.88 ± 0.035</u>	<b>0.97 ± 0.014</b>	0.87 ± 0.033
MTAL <sub>R</sub>	<u>0.83 ± 0.035</u>	<b>0.92 ± 0.029</b>	<b>0.97 ± 0.014</b>	<b>0.98 ± 0.013</b>

Again, the results in Tables (2) to (5) show that most MTL methods using task relationships are better than those without task relationships. It suggests that most of the current MTL methods rely on the relationship between tasks. However, from these results, we find that MTAL is equivalent to the best method of the first set of experiments, and it improves the accuracy on *HD* and *Amazon* in Tables (3) and (5). This further reflects the excellent performance of the convolution kernel pairs sharing mechanism.

Table 3: Performance comparison among various methods in *HD*, *A*, *typographic*, and *a* datasets without the task relationships. Among them, the bold numbers are the best classification results, and the underlined numbers are the second-best classification results.

Models	<i>HD</i>	<i>A</i>	<i>Typographic</i>	<i>a</i>
Single-task	0.76 ± 0.025	0.84 ± 0.020	0.95 ± 0.011	0.80 ± 0.023
Multi-task	0.82 ± 0.030	0.85 ± 0.090	0.92 ± 0.012	0.84 ± 0.090
MTDA	0.75 ± 0.047	0.78 ± 0.025	0.90 ± 0.024	0.82 ± 0.012
Cross-Stitch	<u>0.82 ± 0.019</u>	0.86 ± 0.020	0.95 ± 0.015	<u>0.87 ± 0.090</u>
NDDR-CNN	0.80 ± 0.026	<u>0.87 ± 0.017</u>	<u>0.96 ± 0.005</u>	0.85 ± 0.023
MTAL	<b>0.86 ± 0.038</b>	<b>0.98 ± 0.014</b>	<b>0.97 ± 0.016</b>	<b>0.92 ± 0.025</b>

Finally, as shown in Figure (8), we find that the overall performance of the MTAL and MTAL<sub>R</sub> methods is better than other methods. The above experimental results are consistent with our theoretical analysis.

Table 4: Performance comparison among various methods that use task relationships on *HD*, *Coil-20*, *Art*, and *Amazon* datasets. Among them, the bold numbers are the best classification results, and the underlined numbers are the second-best classification results.

Methods	Group 1		Group 2	
	<i>HD</i>	<i>Coil-20</i>	<i>Art</i>	<i>Amazon</i>
Single-task	0.76 ± 0.025	<b>1</b>	0.59 ± 0.041	0.76 ± 0.032
Multi-task	0.81 ± 0.041	<b>1</b>	0.57 ± 0.059	0.77 ± 0.042
MTDA	0.80 ± 0.007	<u>0.97 ± 0.018</u>	0.59 ± 0.013	0.72 ± 0.043
Cross-Stitch	<u>0.81 ± 0.040</u>	<b>1</b>	0.60 ± 0.055	0.77 ± 0.045
NDDR-CNN	0.79 ± 0.041	<b>1</b>	<u>0.60 ± 0.051</u>	<u>0.80 ± 0.044</u>
MTAL <sub>R</sub>	<b>0.86 ± 0.034</b>	<b>1</b>	<b>0.67 ± 0.050</b>	<b>0.80 ± 0.043</b>

Table 5: Performance comparison among various methods in *HD*, *Art*, *Coil-20*, and *Amazon* datasets without task relationships. Among them, the bold numbers are the best classification results, and the underlined numbers are the second-best classification results.

Models	<i>HD</i>	<i>Art</i>	<i>Coil-20</i>	<i>Amazon</i>
Single-task	0.76 ± 0.025	0.59 ± 0.041	<b>1</b>	0.76 ± 0.032
Multi-task	0.79 ± 0.047	<u>0.61 ± 0.060</u>	<b>1</b>	0.77 ± 0.047
MTDA	0.77 ± 0.051	0.58 ± 0.025	0.95 ± 0.042	0.70 ± 0.044
Cross-Stitch	0.78 ± 0.037	0.59 ± 0.058	<b>1</b>	0.78 ± 0.047
NDDR-CNN	<u>0.78 ± 0.038</u>	0.58 ± 0.063	<b>1</b>	<u>0.78 ± 0.044</u>
MTAL	<b>0.82 ± 0.035</b>	<b>0.65 ± 0.056</b>	<b>1</b>	<b>0.82 ± 0.040</b>

#### 4.5. Threshold selection analysis

To prevent the negative transfer caused by existing MTAL network learning, we conduct threshold experiments on related and unrelated scenarios between tasks. We set the threshold in the range of 0.1-0.9, train 10 epochs for each value, and finally compute the mean and standard deviation of classification accuracy. As shown in Fig.(9), we can obtain: 1) When the tasks are related and  $\delta \geq 0.4$ , the generalization performance of the MTAL network is the best. 2) When the tasks are unrelated and  $\delta \geq 0.55$ , the generalization performance of the MTAL network is the best.

#### 4.6. Model parameters compression and sharing strategy visualization

We have obtained the sharing rate (i.e., network redundancy compression) of the convolution kernels in the MTAL network and each convolution layer by experiments. As shown in Table(6), the kernel pairs shared ratios of each

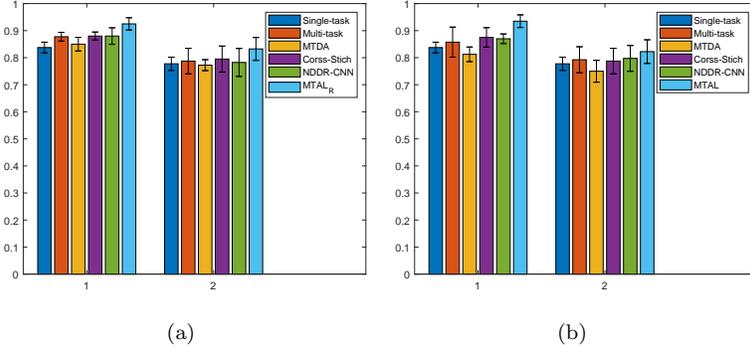


Figure 8: Performance comparison of various methods on mean and mean square error. In Fig.8 (a), 1 and 2, the performance comparison of various methods using task relationships when the tasks are related and unrelated, respectively. In Fig.8 (b), 1 and 2, the performance comparison of various methods that do not use task relationships when tasks are related and unrelated, respectively.

convolutional layer in the MTAL network are 12.5%, 12.5%, 16.7%, 12.5%, 8.3%, 8.3%, 12.5%, 16.7%, 12.5%, 16.7%, 12.5%, 16.7%, 16.7% respectively. The solution space of the entire network can be compressed to 13.1%.

To verify the feasibility of the convolution kernel pairs sharing mechanism proposed in this paper, we respectively visualized the activation maps in the first and second convolution layer of the MTAL network. Fig.(10) shows the activation maps generated in the first layer of the convolution layer when the task is partially related. We find that more shapes and features information can be shared in the active maps area. Fig.(11) shows the activation maps generate in the second convolution layer when the tasks are unrelated. We further find that more textures and contours information can be shared in the activation map area. The above visualization results show that it is feasible to perform cross-task learning by sharing similar convolution kernel pairs.

#### 4.7. Model convergence analysis

In this section, we show the loss functions of the two sets of tasks in Fig.(12) (a) and (b) to analyze the convergence of the model. In Figure 12(a), we find that the loss function of task 2 tends to converge after about 40 iterations,

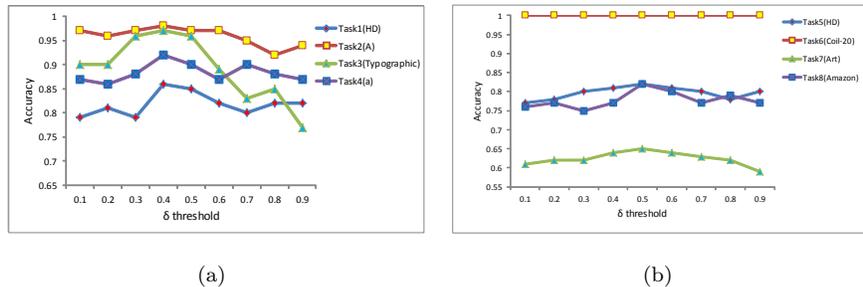


Figure 9: The selection result of  $\delta$  in various tasks. Fig.9 (a) and (b) show the corresponding learning value  $\delta$  when the tasks are related and unrelated.

while tasks 4, 3, and 1 tend to converge after about 60, 70, and 80 iterations, respectively. In Figure 12 (b), we observe that the loss function of task 6 tends to converge after about 40 iterations, while tasks 5, 8, and 7 converge after about 80, 120, and 150 iterations, respectively. The above shows that our model converge relatively fast.

## 5. Conclusion

In this paper, we provide a deep learning framework MTAL for processing THMT. Compared with the previous MTL method, the MTAL network explores and uses the inherent relationship between tasks to share knowledge of similar convolution kernel pairs in each of their layers to learn THMT. The network not only effectively performs cross-task learning but also suppresses the intra-redundancy of the entire network. Meanwhile, MTAL can handle related heterogeneous tasks well and achieve great performance when they are unrelated. At the same time, the designed sharing strategy in MTAL can be flexibly embedded in other deep multi-task learning frameworks. To evaluate the proposed MTAL, we conduct experiments on eight public datasets and compare with the state-of-the-art HMTL methods. Experimental results show superiority of our MTAL. In summary, our work can enrich HMTL research from three aspects: 1) an adaptive THMT learning mechanism that can avoid negative transfer caused by jointly learning multiple tasks due to incorrect pre-defined

Table 6: Kernels sharing ratio

MTAL network	Kernels sharing ratio (%)
Conv1_1 layer	12.5
Conv1_2 layer	12.5
Conv2_1 layer	16.7
Conv2_2 layer	12.5
Conv3_1 layer	8.3
Conv3_2 layer	8.3
Conv3_3 layer	12.5
Conv4_1 layer	16.7
Conv4_2 layer	12.5
Conv4_3 layer	16.7
Conv5_1 layer	12.5
Conv5_2 layer	16.7
Conv5_3 layer	16.7
Total network	13.1

task relationships; 2) a new method for solving THMT with no relation among tasks; 3) a new THMT sharing strategy for learning multiple heterogeneous tasks. However, in this work, we do not solve the problem of unbalanced and interpretable shared learning among multiple heterogeneous tasks and will devote ourselves to solving these problems next.

### Acknowledgments

This work is supported in part by Key Program of NSFC under Grant No. 61732006 and the NSFC under Grant No. 62076124.

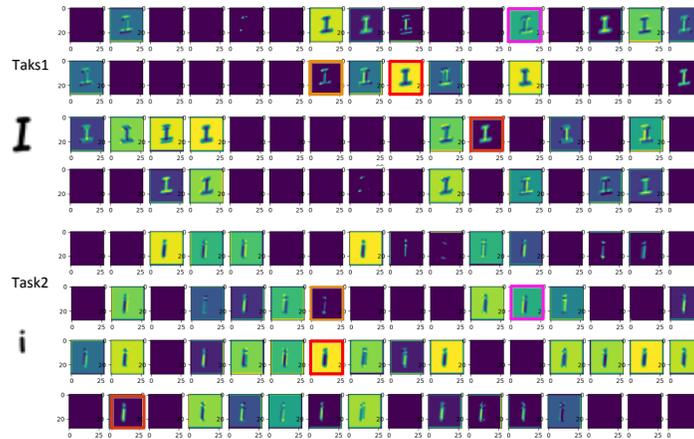


Figure 10: Visual activation maps when the tasks are relevant. The activation maps are extracted from the first convolutional layer in the MTAL network. The red, purple, brown, and orange solid line boxes denote the activation maps generated from different convolution kernels in task 1 and task 2, respectively.

## References

## References

- [1] A. Kendall, Y. Gal, R. Cipolla, Multi-task learning using uncertainty to weigh losses for scene geometry and semantics, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 7482–7491.
- [2] D. Dong, H. Wu, W. He, D. Yu, H. Wang, Multi-task learning for multiple language translation, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), 2015, pp. 1723–1732.
- [3] R. Giri, M. L. Seltzer, J. Droppo, D. Yu, Improving speech recognition in reverberation using a room-aware deep neural network and multi-task learning, in: 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2015, pp. 5014–5018.

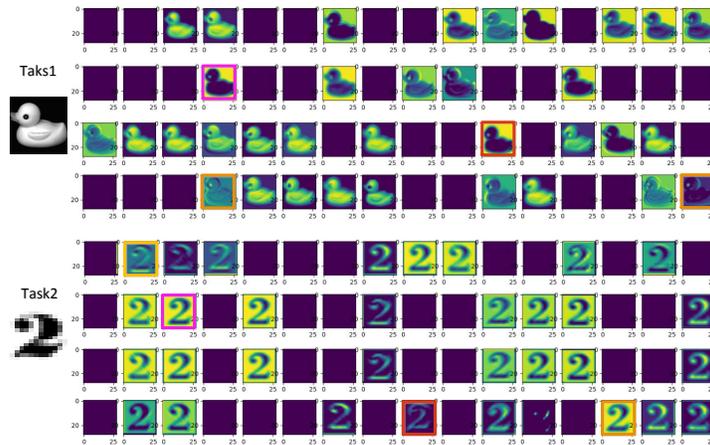


Figure 11: Visual activation maps when the tasks are unrelated. The activation maps are extracted from the second convolutional layer in the MTAL network. The red, purple, brown, and orange solid line boxes denote the activation maps generated from different convolution kernels in task 1 and task 2, respectively.

- [4] D. Zhang, D. Shen, Multi-modal multi-task learning for joint prediction of multiple regression and classification variables in alzheimer’s disease, *Neuroimage* 59 (2) (2012) 895–907.
- [5] Y. Chen, D. Zhao, L. Lv, Q. Zhang, Multi-task learning for dangerous object detection in autonomous driving, *Information Sciences* 432 (2018) 559–571.
- [6] M. Long, Z. Cao, J. Wang, S. Y. Philip, Learning multiple tasks with multi-linear relationship networks, in: *Advances in neural information processing systems*, 2017, pp. 1594–1603.
- [7] X. Zhang, F. Zhou, Y. Lin, S. Zhang, Embedding label structures for fine-grained feature representation, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1114–1123.
- [8] T.-Y. Lin, A. RoyChowdhury, S. Maji, Bilinear cnn models for fine-grained visual recognition, in: *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1449–1457.

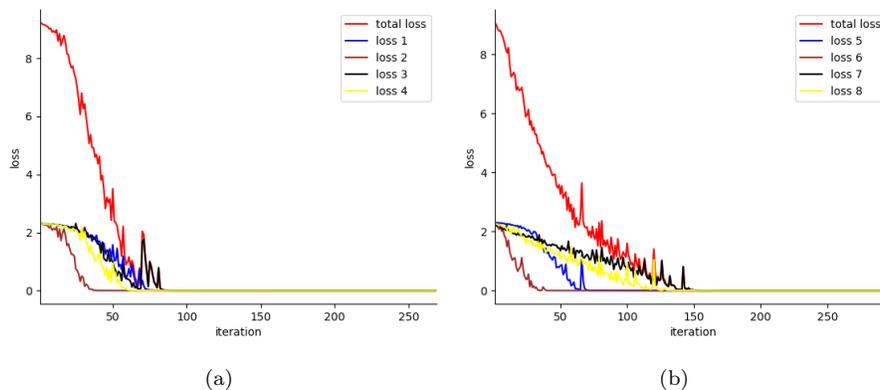


Figure 12: Model convergence experiment. Fig. 12 (a) and (b) show the corresponding convergence curves of the model when the heterogeneous tasks are related and unrelated.

- [9] H. Han, A. K. Jain, F. Wang, S. Shan, X. Chen, Heterogeneous face attribute estimation: A deep multi-task learning approach, *IEEE transactions on pattern analysis and machine intelligence* 40 (11) (2017) 2597–2609.
- [10] S. Li, Z.-Q. Liu, A. B. Chan, Heterogeneous multi-task learning for human pose estimation with deep convolutional neural network, in: *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2014, pp. 482–489.
- [11] Y. Luo, Y. Wen, D. Tao, J. Gui, C. Xu, Large margin multi-modal multi-task feature extraction for image classification, *IEEE Transactions on Image Processing* 25 (1) (2015) 414–427.
- [12] Y. Luo, Y. Wen, D. Tao, Heterogeneous multitask metric learning across multiple domains, *IEEE Transactions on Neural Networks and Learning Systems* 29 (9) (2018) 4051–4064.
- [13] J. Zhao, B. Du, L. Sun, F. Zhuang, W. Lv, H. Xiong, Multiple relational attention network for multi-task learning (2019) 1123–1131.

- [14] J. Schreiber, B. Sick, Emerging relation network and task embedding for multi-task regression problems, ArXiv abs/2020.14034.
- [15] D. C. Luvizon, D. Picard, H. Tabia, Multi-task deep learning for real-time 3d human pose estimation and action recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence (2020) 1–1.
- [16] S. Wu, H. R. Zhang, C. Ré, Understanding and improving information transfer in multi-task learning (2020). arXiv:2005.00944.
- [17] Q. Liu, X. Li, Z. He, N. Fan, D. Yuan, W. Liu, Y. Liang, Multi-task driven feature models for thermal infrared tracking, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, 2020, pp. 11604–11611.
- [18] N. Dvornik, K. Shmelkov, J. Mairal, C. Schmid, Blitznet: A real-time deep network for scene understanding (2017) 4174–4182.
- [19] J. Huang, R. S. Feris, Q. Chen, S. Yan, Cross-domain image retrieval with a dual attribute-aware ranking network, CoRR abs/1505.07922.
- [20] A. Saeed, S. Trajanovski, Personalized driver stress detection with multi-task neural networks using physiological signals, arXiv preprint arXiv:1711.06116.
- [21] L. chen Liu, Z. Liu, H. Wu, Z. Wang, J. Shen, Y. Song, M. Zhang, Multi-task learning via adaptation to similar tasks for mortality prediction of diverse rare diseases, ArXiv abs/2004.05318.
- [22] J. Baxter, A bayesian/information theoretic model of learning to learn via multiple task sampling, Machine Learning 28 (1) (1997) 7–39.
- [23] T. Sun, Y. Shao, X. Li, P. Liu, H. Yan, X. Qiu, X. Huang, Learning sparse sharing architectures for multiple tasks, in: AAAI, 2020.
- [24] L. Duong, T. Cohn, S. Bird, P. Cook, Low resource dependency parsing: Cross-lingual parameter sharing in a neural network parser, in: Proceed-

- ings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), 2015, pp. 845–850.
- [25] Y. Yang, T. M. Hospedales, Trace norm regularised deep multi-task learning, arXiv preprint arXiv:1606.04038.
- [26] M. Long, Z. CAO, J. Wang, P. S. Yu, Learning multiple tasks with multilinear relationship networks, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), *Advances in Neural Information Processing Systems 30*, Curran Associates, Inc., 2017, pp. 1594–1603.
- [27] X. Sun, R. Panda, R. S. Feris, Adashare: Learning what to share for efficient deep multi-task learning, ArXiv abs/2019.12423.
- [28] X. Cai, F. Nie, W. Cai, H. Huang, Heterogeneous image features integration via multi-modal semi-supervised learning model, in: *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 1737–1744.
- [29] V. Sanh, T. Wolf, S. Ruder, A hierarchical multi-task approach for learning embeddings from semantic tasks, in: *AAAI*, 2019.
- [30] Y. Zhang, D.-Y. Yeung, Multi-task learning in heterogeneous feature spaces, in: *Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011.
- [31] F. Zhuang, X. Li, X. Jin, D. Zhang, L. Qiu, Q. He, Semantic feature learning for heterogeneous multitask classification via non-negative matrix factorization, *IEEE transactions on cybernetics* 48 (8) (2017) 2284–2293.
- [32] X. Cai, F. Nie, H. Huang, F. Kamangar, Heterogeneous image feature integration via multi-modal spectral clustering, in: *CVPR 2011, IEEE*, 2011, pp. 1977–1984.
- [33] Y. He, P. Liu, L. Zhu, Y. Yang, Meta filter pruning to accelerate deep convolutional neural networks, arXiv preprint arXiv:1904.03961.

- [34] O. Slizovskaia, E. Gómez, G. Haro, A case study of deep-learned activations via hand-crafted audio features, arXiv preprint arXiv:1907.01813.
- [35] K. Lin, J. Xu, I. M. Baytas, S. Ji, J. Zhou, Multi-task feature interaction learning, in: Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016, pp. 1735–1744.
- [36] J. Zhang, J. Miao, K. Zhao, Y. Tian, Multi-task feature selection with sparse regularization to extract common and task-specific features, *Neurocomputing* 340 (2019) 76–89.
- [37] R. Ranjan, V. M. Patel, R. Chellappa, Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41 (1) (2019) 121–135. doi:10.1109/TPAMI.2017.2781233.
- [38] F. Huang, Y. Qiu, Q. Li, S. Liu, F. Ni, Predicting drug-disease associations via multi-task learning based on collective matrix factorization, *Frontiers in Bioengineering and Biotechnology* 8 (2020) 218.
- [39] Z. Zheng, Y. Wang, Q. Dai, H. Zheng, D. Wang, Metadata-driven task relation discovery for multi-task learning., in: *IJCAI*, 2019, pp. 4426–4432.
- [40] I. Misra, A. Shrivastava, A. Gupta, M. Hebert, Cross-stitch networks for multi-task learning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 3994–4003.
- [41] J. Ma, Z. Zhao, J. Chen, A. Li, L. Hong, E. H. Chi, Snr: Sub-network routing for flexible parameter sharing in multi-task learning, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33, 2019, pp. 216–223.
- [42] J. Ma, Z. Zhao, X. Yi, J. Chen, L. Hong, E. H. Chi, Modeling task relationships in multi-task learning with multi-gate mixture-of-experts, in: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2018, pp. 1930–1939.

- [43] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, N. De Freitas, Predicting parameters in deep learning, in: *Advances in neural information processing systems*, 2013, pp. 2148–2156.
- [44] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.
- [45] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [46] Y. Gao, J. Ma, M. Zhao, W. Liu, A. L. Yuille, Nddr-cnn: Layerwise feature fusing in multi-task cnns by neural discriminative dimensionality reduction, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3205–3214.