

# Rapid Performance Gain through Active Model Reuse

Feng Shi  
Yu-Feng Li

IJCAI 2019

# Introduction

- Learning models work well on the specific tasks, but have to be discarded once the target task changes.
- **Model reuse** tries to reduce the learning resources for a new target task with the exploitation of pre-trained models.
- Previous model reuse studies usually assume that the labeled data for the target task are **passively** collected.
- **Active Model Reuse** (AcMR) construct queries through pre-trained models to facilitate the active learner when labeled examples are insufficient.
- Different from Active Learning.

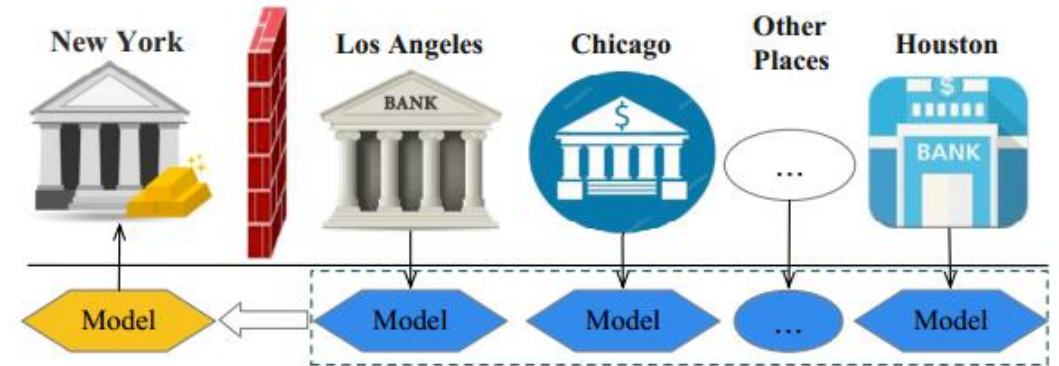


Figure 1: A practical example of model reuse. It is often challenging to share bank data due to the privacy issue. Instead we can reuse pre-trained models of banks to help improve the performance.

# Notation

Table 1: Summary of Notation

Notation	Meaning
$N$	number of training data
$y_t \in \{+1, -1\}$	target label of instance $\mathbf{x}_t$
$\mathcal{L}$	labeled data in the target task
$\mathcal{U}$	unlabeled data in the target task
$\Rightarrow \{f_1, f_2, \dots, f_k\}$	$k$ pre-trained models
$\Rightarrow f_j^{(t)}$	the prediction of $f_j$ on $\mathbf{x}_t$
$\Rightarrow \boldsymbol{\eta} = [\eta_1, \eta_2, \dots, \eta_k]$	the weight vector of $k$ models
$l_j^{(t)} = (1 - y^{(t)} f_j^{(t)})_+$	the hinge loss of $j$ -th model for $\mathbf{x}_t$
$L_j = \sum_{\mathbf{x}_t \in \mathcal{L}} l_j^{(t)}$	the empirical loss of $j$ -th model
$\Rightarrow f_{\mathcal{L}}$	the active learner built on $\mathcal{L}$
$\Rightarrow f_{\mathcal{L}}^{(t)}$	the prediction of $f_{\mathcal{L}}$ on $\mathbf{x}_t$

# Deficiencies of Baseline Approaches

$$\hat{y}^{(t)} = \arg \max_{c \in \{-1, +1\}} \sum_{j=1}^k \eta_j \cdot \mathbb{I}(f_j^{(t)} = c)$$

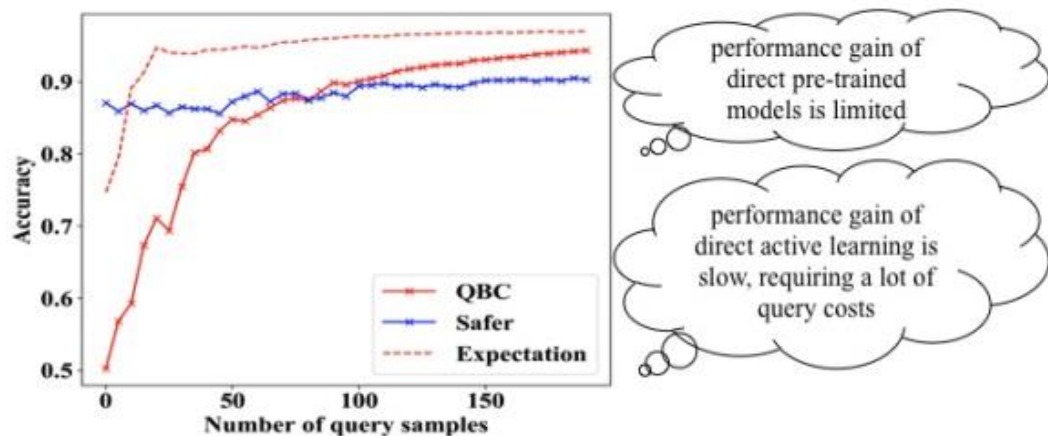


Figure 2: A true example of learning curves of two methods on a classification task. We expect that with the help of pre-trained models, the performance of active learner can get rapid improvement.

This prediction may be risky since we can not get accurate prior weights especially when there are very few labeled examples in target tasks.

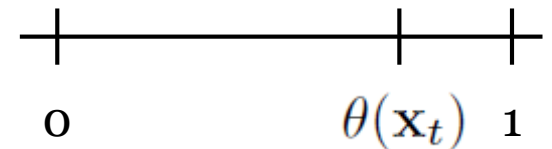
## Strategies:

- Rather than fully trusting pre-trained models, we consider trusting the pre-trained models partially and leverage pre-trained models to **filter out not very necessary queries** generated by direct active learning to facilitate an effective active learner.
- Moreover, once the label capacity is enhanced, we continually **update the confident weights of pre-trained models** to improve their confidence for target task.

# Actively Reuse Pre-trained Models

- Predicted probability of pre-trained model  $f_j$  for instance  $\mathbf{x}_t$   $P(\hat{y}^{(t)}|\mathbf{x}_t) = \sum_{j=1}^k \eta_j P_j(\hat{y}^{(t)}|\mathbf{x}_t) \cdot \mathbb{I}(f_j^{(t)} = \hat{y}^{(t)})$
- Consistency of pre-trained model and active learner  $\alpha(\mathbf{x}_t) = (1 - \mathbb{I}(\hat{y}^{(t)} \neq f_{\mathcal{L}}^{(t)})) P(\hat{y}^{(t)}|\mathbf{x}_t)$
- Query necessity  $\theta(\mathbf{x}_t) = (1 + \alpha(\mathbf{x}_t))^{-1}$
- Randomly generate a real number  $R$  within 0 to 1, and then the **decision function** is defined as

$$\mathbb{F}(\mathbf{x}_t) = \begin{cases} 0, & \text{if } R > \theta(\mathbf{x}_t) \\ 1, & \text{otherwise} \end{cases}$$



Label the instance  $\mathbf{x}_t$  by the pre-trained models with probability  $1 - \theta(\mathbf{x}_t)$

# Sampling Error

错误地相信了 pre-trained models

Let  $\varepsilon_p$  and  $\varepsilon_a$  denote the expected error of the pre-trained models and the active learner  $f_{\mathcal{L}}$  respectively, and let  $\delta = \varepsilon_p + \varepsilon_a$ .

**Theorem 1.** *In the algorithm ACMR, we assume that  $\varepsilon_a \leq \varepsilon_p$ , then the sampling error  $\varepsilon$  for ACMR satisfies:*

$$\varepsilon \leq \frac{\varepsilon_p^2}{1 + (1 - \varepsilon_p)} \quad (7)$$

**Proof.** *According to the analysis of the decision function  $\mathbb{F}(\mathbf{x})$  described above, ACMR makes wrong decision only when both the pre-trained models and the active learner  $f_{\mathcal{L}}$  agree on the wrong label. In this case, ACMR has probability  $1 - \theta(\mathbf{x})$  to trust the classification result given by the pre-trained models, where  $\theta(\mathbf{x})$  is defined in Eq.(5). Thus, the sampling error of ACMR can be written as  $\varepsilon = \varepsilon_p \varepsilon_a (1 - \theta(\mathbf{x})) \leq \varepsilon_p^2 (1 - \theta(\mathbf{x}))$ . Moreover, in this situation,  $\theta(\mathbf{x}) = \frac{1}{1 + (1 - \varepsilon_p)}$ . Thus,*

$$\varepsilon \leq \frac{\varepsilon_p^2 \times (1 - \varepsilon_p)}{1 + (1 - \varepsilon_p)} \leq \frac{\varepsilon_p^2}{1 + (1 - \varepsilon_p)} \quad (8)$$

# Querying Rate

**Theorem 2.** *In the algorithm ACMR, for an unlabeled instance, the probability that ACMR queries the label from the experts (with cost) satisfies:*

$$P(Q) \leq \delta + \frac{1 - \delta}{1 + (1 - \varepsilon_p)} \quad (9)$$

**Proof.** *According to the analysis of the decision function  $\mathbb{F}(\mathbf{x})$ , ACMR will query the experts to label the instance when the pre-trained models and the active learner hold different predictions on the classification result. And when the two classifiers agree on the result, it still has probability  $\theta(\mathbf{x})$  to query the experts. Thus:*

$$\begin{aligned} P(Q) &= \varepsilon_a(1 - \varepsilon_p) + [\varepsilon_p\varepsilon_a + (1 - \varepsilon_p)(1 - \varepsilon_a)]\theta(\mathbf{x}) \\ &\quad + (1 - \varepsilon_a)\varepsilon_p \quad (10) \\ &= \theta(\mathbf{x}) + (\varepsilon_p + \varepsilon_a - 2\varepsilon_p\varepsilon_a)(1 - \theta(\mathbf{x})) \\ &\leq \delta + (1 - \delta)\theta(\mathbf{x}) \\ &\leq \delta + \frac{1 - \delta}{1 + (1 - \varepsilon_p)} \end{aligned}$$

- The more accurate the pre-trained models and active learner become, the less necessarily will we query the experts to label the instance.
- Reduce  $\varepsilon_p$  to improve performance in the whole learning process.

# Update Weights for Pre-Trained Models

At time  $m$ , the target task receives a training instance  $\mathbf{x}_{t(m)}$ , the pre-trained models make a prediction and suffer a loss after  $\mathbf{y}_{t(m)}$  is revealed.

$$\boldsymbol{\eta}^{(m+1)} = \arg \min_{\boldsymbol{\eta} \in \Theta} \sum_{j \in [k]} \eta_j l_j^{(t(m))} + \lambda \mathcal{D}_{KL}(\boldsymbol{\eta} \| \boldsymbol{\eta}^{(m)})$$

$$\eta_j^{(m+1)} = \frac{\eta_j^{(m)} \exp(-l_j^{(m)} / \lambda)}{\sum_{j'=1}^k \eta_{j'}^{(m)} \exp(-l_{j'}^{(m)} / \lambda)}, \quad j \in [k]$$

**Proposition 1.** (*Weight Concentration*). *During the weight update procedure in the whole learning process, the weights will concentrate on those pre-trained models who suffer a small cumulative loss on the target task.*

**Proof.** *Through the update rule, we know that the weight associated with the  $j$ -th previous model is equal to  $\eta_j = \frac{\exp(-L_j / \lambda)}{\sum_{j'=1}^k \exp(-L_{j'} / \lambda)}$ ,  $j \in [k]$ . The smaller the loss  $L_j$  of the pre-trained model, the higher the weight  $\eta_j$ .*

# AcMR

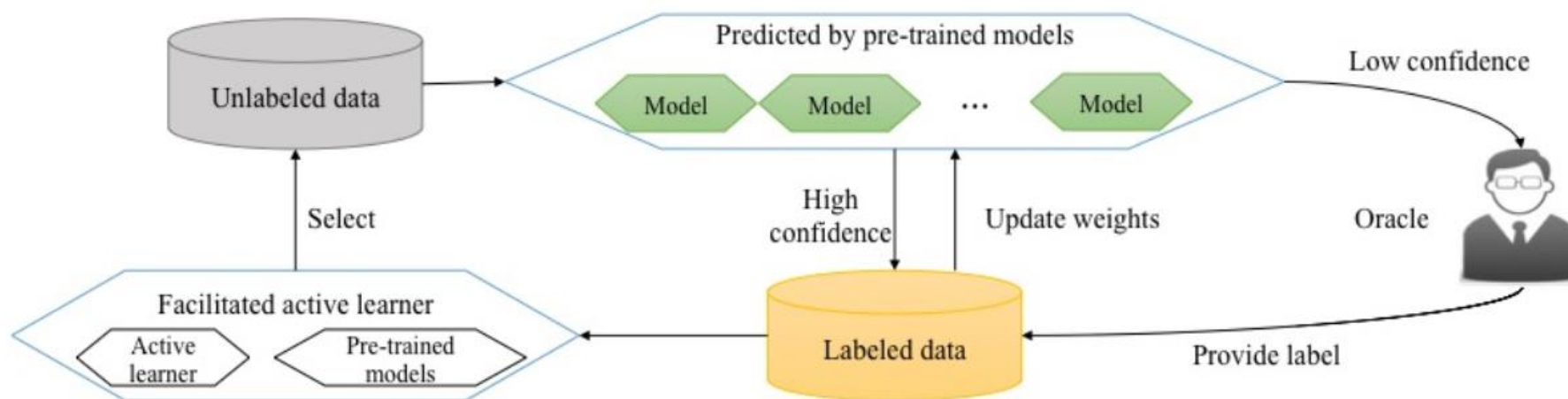


Figure 3: Model architecture of AcMR. In contrast to traditional active learning, we propose to facilitate the active learner with pre-trained models, and we also use pre-trained models to filter out not very necessary queries. During the learning process, the weights of  $k$  models are continually updated such that they can predict unlabeled samples more accurately.

---

**Algorithm 1** The learning algorithm for ACMR

---

**Input:** labeled dataset  $\mathcal{L}$ , unlabeled dataset  $\mathcal{U}$ , sampling size  $N$ ,  $k$  pre-trained models  $\{f_1, f_2, \dots, f_k\}$  and  $\lambda > 0$

**Output:** the model  $f_{\mathcal{L}}$  for the target task

- 1: Initialize weight vector  $\boldsymbol{\eta} = [1/k, \dots, 1/k]$  and  $f_{\mathcal{L}}$ .
  - 2: **for**  $m = 1, 2, \dots, N$  **do**
  - 3:   Select an instance  $\mathbf{x}_{t(m)}$  from  $\mathcal{U}$  by traditional active learning
  - 4:   Predict  $\mathbf{x}_{t(m)}$  by pre-trained models, and calculate the decision function  $\mathbb{F}(\mathbf{x}_{t(m)})$  via Eq.(6)
  - 5:   **if**  $\mathbb{F}(\mathbf{x}_{t(m)}) = 0$  **then**
  - 6:      $y_{t(m)} \leftarrow \hat{y}^{(t(m))}$
  - 7:   **else**
  - 8:      $y_{t(m)} \leftarrow$  query from oracle
  - 9:     **if**  $y_{t(m)} \neq \hat{y}^{(t(m))}$  **then**
  - 10:       update the weights via Eq.(12)
  - 11:     **else**
  - 12:        $\boldsymbol{\eta}^{(m+1)} \leftarrow \boldsymbol{\eta}^{(m)}$
  - 13:     **end if**
  - 14:   **end if**
  - 15:    $\mathcal{L} = \mathcal{L} \cup (\mathbf{x}_{t(m)}, y_{t(m)}), \mathcal{U} = \mathcal{U} / (\mathbf{x}_{t(m)})$ ;
  - 16:   Train the learner  $f_{\mathcal{L}}$  with  $\mathcal{L}$
  - 17:   Facilitate active learner  $f_{\mathcal{L}}$  via Eq.(1)
  - 18: **end for**
  - 19: **return**  $f_{\mathcal{L}}$
-

# Experiments

Baselines:

- 1) An actively transfer learning method **AcTraK** with QBC and random sample selection criteria
- 2) **Safer**: safely exploits pre-trained models
- 3) **Logistic Regression** with QBC and random
- 4) **AcMR** with QBC and random

Compare these methods in different number of samples labeled by experts (30, 60, 90).

Three tasks:

- 1) Text Classification
- 2) Sentiment Analysis Task
- 3) Spam Detection Task

# Text Classification

Dataset: 20 Newsgroups

6 text classification tasks

- comp vs. talk
- comp vs. sci
- comp vs. rec
- rec vs. talk
- rec vs. sci
- sci vs. talk

Table 3: Classification accuracy of compared methods with different queries on 6 classification tasks in 20 Newsgroups dataset. The boldfaces denote the best and the second best methods in terms of the accuracy, and Queries mean that the number of samples labeled by experts.

Task	Queries	Safer	QBC	AcTrak-QBC	AcMR-QBC	Random	AcTrak-Ran	AcMR-Ran
Task1	30	.876 ± .033	.792 ± .118	.817 ± .107	<b>.923 ± .016</b>	.806 ± .126	.839 ± .079	<b>.926 ± .014</b>
	60	.880 ± .029	.916 ± .029	.909 ± .052	<b>.932 ± .013</b>	.911 ± .030	.915 ± .034	<b>.933 ± .008</b>
	90	.891 ± .040	.930 ± .022	.929 ± .026	<b>.934 ± .011</b>	.928 ± .025	.934 ± .019	<b>.936 ± .011</b>
Task2	30	.737 ± .027	.739 ± .124	.711 ± .143	<b>.838 ± .039</b>	.689 ± .110	.756 ± .102	<b>.839 ± .049</b>
	60	.738 ± .036	.821 ± .076	.809 ± .091	<b>.870 ± .044</b>	.833 ± .059	.842 ± .071	<b>.873 ± .038</b>
	90	.749 ± .038	.856 ± .065	.851 ± .050	<b>.885 ± .032</b>	.871 ± .032	.861 ± .056	<b>.892 ± .027</b>
Task3	30	.944 ± .010	.788 ± .155	.803 ± .128	<b>.954 ± .012</b>	.748 ± .146	.839 ± .095	<b>.952 ± .012</b>
	60	.945 ± .011	.903 ± .080	.920 ± .049	<b>.956 ± .011</b>	.868 ± .115	.921 ± .050	<b>.956 ± .011</b>
	90	.945 ± .010	.937 ± .040	.948 ± .015	<b>.957 ± .012</b>	.928 ± .043	.951 ± .020	<b>.957 ± .009</b>
Task4	30	.689 ± .059	.705 ± .119	.664 ± .133	<b>.905 ± .027</b>	.786 ± .142	.757 ± .131	<b>.910 ± .027</b>
	60	.680 ± .049	.853 ± .085	.797 ± .126	<b>.928 ± .012</b>	.875 ± .063	.875 ± .083	<b>.927 ± .018</b>
	90	.703 ± .040	.907 ± .044	.859 ± .086	<b>.938 ± .011</b>	.914 ± .032	.922 ± .023	<b>.933 ± .017</b>
Task5	30	.850 ± .061	.765 ± .133	.753 ± .142	<b>.915 ± .059</b>	.765 ± .134	.758 ± .177	<b>.941 ± .015</b>
	60	.869 ± .053	.888 ± .081	.880 ± .065	<b>.941 ± .029</b>	.864 ± .085	.833 ± .128	<b>.937 ± .027</b>
	90	.871 ± .057	.925 ± .030	.915 ± .055	<b>.949 ± .015</b>	.909 ± .042	.886 ± .087	<b>.950 ± .011</b>
Task6	30	.770 ± .042	.645 ± .105	.646 ± .093	<b>.819 ± .029</b>	.643 ± .104	.676 ± .113	<b>.821 ± .040</b>
	60	.763 ± .042	.679 ± .089	.718 ± .094	<b>.838 ± .032</b>	.743 ± .103	.779 ± .090	<b>.850 ± .038</b>
	90	.788 ± .034	.747 ± .106	.781 ± .081	<b>.860 ± .031</b>	.841 ± .083	.823 ± .075	<b>.867 ± .031</b>

# Sentiment Analysis

Product reviews from Amazon dataset containing reviews from 4 domains: Book, DVD, Electronics and Kitchen.

Consider each domain as a binary classification task: reviews with rating  $> 3$  are labeled positive(+), those with rating  $< 3$  are labeled negative(-).

Table 4: Classification accuracy of compared methods with different queries on 4 classification tasks in Sentiment dataset. The boldfaces denote the best and the second best methods in terms of the accuracy, and Queries mean that the number of samples labeled by experts.

Task	Queries	Safer	QBC	AcTrak-QBC	AcMR-QBC	Random	AcTrak-Ran	AcMR-Ran
Task1	30	<b>.612 ± .022</b>	.529 ± .033	.578 ± .033	<b>.607 ± .028</b>	.577 ± .036	.582 ± .039	.598 ± .032
	60	.614 ± .018	.579 ± .048	.609 ± .035	<b>.631 ± .023</b>	.607 ± .025	.615 ± .037	<b>.631 ± .026</b>
	90	.614 ± .018	.615 ± .047	.627 ± .038	<b>.645 ± .023</b>	.625 ± .024	.637 ± .034	<b>.649 ± .020</b>
Task2	30	<b>.648 ± .020</b>	.562 ± .048	.582 ± .041	.610 ± .033	.594 ± .031	.594 ± .036	<b>.611 ± .046</b>
	60	<b>.648 ± .018</b>	.591 ± .049	.626 ± .038	.631 ± .030	.623 ± .033	.629 ± .034	<b>.643 ± .032</b>
	90	.648 ± .020	.621 ± .047	.649 ± .033	<b>.660 ± .031</b>	.650 ± .026	.657 ± .028	<b>.668 ± .029</b>
Task3	30	<b>.651 ± .044</b>	.578 ± .057	.629 ± .033	<b>.645 ± .036</b>	.612 ± .036	.629 ± .033	.642 ± .037
	60	.654 ± .041	.640 ± .046	.665 ± .024	<b>.676 ± .028</b>	.643 ± .032	.664 ± .032	<b>.676 ± .026</b>
	90	.665 ± .033	.672 ± .035	.685 ± .023	<b>.694 ± .025</b>	.671 ± .027	.679 ± .029	<b>.690 ± .025</b>
Task4	30	.621 ± .049	.628 ± .032	.634 ± .033	<b>.644 ± .026</b>	.624 ± .036	.628 ± .045	<b>.641 ± .035</b>
	60	.631 ± .051	.657 ± .035	.666 ± .024	<b>.680 ± .018</b>	.652 ± .032	.663 ± .035	<b>.684 ± .029</b>
	90	.643 ± .054	.684 ± .026	.687 ± .026	<b>.702 ± .022</b>	.678 ± .026	.681 ± .033	<b>.705 ± .024</b>

# Spam Detection

Dataset: ECML PAKDD Discovery challenge

Table 5: Classification accuracy of compared methods with different queries on 6 classification tasks in Spam dataset. The boldfaces denote the best and the second best methods in terms of the accuracy, and Queries mean that the number of samples labeled by experts.

Task	Queries	Safer	QBC	AcTrak-QBC	AcMR-QBC	Random	AcTrak-Ran	AcMR-Ran
Task1	30	.925 ± .028	.757 ± .152	.734 ± .144	<b>.956 ± .016</b>	.673 ± .148	.681 ± .149	<b>.951 ± .017</b>
	60	.935 ± .028	.778 ± .142	.750 ± .116	<b>.954 ± .021</b>	.731 ± .139	.755 ± .113	<b>.957 ± .018</b>
	90	.939 ± .021	.788 ± .117	.761 ± .092	<b>.954 ± .024</b>	.805 ± .143	.806 ± .088	<b>.957 ± .023</b>
Task2	30	.906 ± .037	.750 ± .163	.770 ± .161	<b>.952 ± .032</b>	.707 ± .170	.751 ± .159	<b>.949 ± .024</b>
	60	.908 ± .037	.824 ± .118	.791 ± .151	<b>.965 ± .017</b>	.787 ± .136	.849 ± .110	<b>.962 ± .021</b>
	90	.901 ± .038	.872 ± .097	.838 ± .144	<b>.968 ± .014</b>	.843 ± .099	.910 ± .078	<b>.960 ± .019</b>
Task3	30	.897 ± .051	.846 ± .116	.875 ± .098	<b>.970 ± .024</b>	.860 ± .064	.879 ± .047	<b>.965 ± .031</b>
	60	.914 ± .052	.895 ± .052	.916 ± .037	<b>.984 ± .014</b>	.895 ± .055	.895 ± .044	<b>.981 ± .015</b>
	90	.913 ± .042	.928 ± .039	.908 ± .042	<b>.986 ± .013</b>	.902 ± .046	.888 ± .039	<b>.983 ± .014</b>
Task4	30	<b>.962 ± .023</b>	.928 ± .072	.924 ± .081	.959 ± .023	.944 ± .030	.942 ± .022	<b>.961 ± .021</b>
	60	.964 ± .023	.964 ± .023	.968 ± .019	<b>.970 ± .019</b>	.959 ± .020	.967 ± .015	<b>.967 ± .016</b>
	90	.964 ± .023	.964 ± .023	.968 ± .019	<b>.970 ± .019</b>	.963 ± .021	.969 ± .012	<b>.973 ± .013</b>
Task5	30	.653 ± .137	.794 ± .076	.789 ± .097	<b>.896 ± .045</b>	.762 ± .093	.667 ± .066	<b>.882 ± .055</b>
	60	.664 ± .132	.802 ± .070	.811 ± .076	<b>.911 ± .039</b>	.781 ± .082	.702 ± .081	<b>.902 ± .044</b>
	90	.703 ± .118	.794 ± .058	.822 ± .068	<b>.924 ± .023</b>	.796 ± .069	.719 ± .056	<b>.906 ± .031</b>
Task6	30	<b>.782 ± .061</b>	.669 ± .097	.652 ± .091	<b>.778 ± .093</b>	.639 ± .112	.647 ± .107	.763 ± .095
	60	.793 ± .059	.768 ± .102	.711 ± .089	<b>.828 ± .058</b>	.700 ± .106	.682 ± .103	<b>.807 ± .048</b>
	90	.794 ± .062	.793 ± .090	.793 ± .076	<b>.854 ± .048</b>	.758 ± .099	.710 ± .087	<b>.820 ± .044</b>

# Conclusion

- Reusable model design becomes a desire for the rapid expansion of machine learning applications. However, previous model reuse studies assume that the target task receives labeled data passively. This leads to a slow performance improvement to the target task.
- In this paper, we study a kind of new model reuse problem, where the goal is that the model performance for the target task can be quickly improved.
- We propose the ACMR method which constructs queries through pre-trained models when labeled examples are insufficient for the target task, and leverages pre-trained models to filter out not very necessary queries so that considerable queries could be saved compared with direct active learning.