# ICS-Assist: Intelligent Customer Inquiry Resolution Recommendation in Online Customer Service for Large E-Commerce Businesses

Min Fu[1,2], Jiwei Guan[2], Xi Zheng[2], Jie Zhou[1], Jianchao Lu[2], Tianyi Zhang[3], Shoujie Zhuo[1], Lijun Zhan[1], and Jian Yang[2]

[1] Alibaba Group, Hangzhou, China
[2] Macquarie University, Sydney, Australia
[3] Harvard University, Cambridge, Massachusetts, USA
{hanhao.fm,zj236040,souljoy.zsj,zhanlijun.zlj}@alibaba-inc.com
{james.zheng,jian.yang}@mq.edu.au
{jiwei.guan,jianchao.lu}@hdr.mq.edu.au
{tianyi}@seas.harvard.edu

**Abstract.** Efficient and appropriate online customer service is essential to large e-commerce businesses. Existing solution recommendation methods for online customer service are unable to determine the best solutions at runtime, leading to poor satisfaction of end customers. This paper proposes a novel intelligent framework, called ICS-Assist, to recommend suitable customer service solutions for service staff at runtime. Specifically, we develop a generalizable two-stage machine learning model to identify customer service scenarios and determine customer service solutions based on a scenario-solution mapping table. A novel knowledge distillation network called "Panel-Student is proposed to derive a small yet efficient distilled learning model. We implement ICS-Assist and evaluate it using an over 6-month field study with Alibaba Group. In our experiment, over 12,000 customer service staff use ICS-Assist to serve for over 230,000 cases per day on average. The experimental results show that ICS-Assist significantly outperforms the traditional manual method, and improves the solution acceptance rate, the solution coverage rate, the average service time, the customer satisfaction rate, and the business domain catering rate by up to 16%, 25%, 6%, 14% and 17% respectively, compared to the state-of-the-art methods.

**Keywords:** Intelligent customer service · Natural language processing · Deep learning · Distilled learning.

## 1 Introduction

Large e-commerce businesses such as Alibaba and Amazon provide hundreds of thousands of customer services to end customers via conversations every day, and these customer service conversations contain several topics, such as refunds, delivering inquiries, and instructions for using lucky money [23]. When end customers make inquiries through online customer service, they usually demand

their requirements and intentions be addressed as fast as possible [14]. These requirements and intentions are usually versatile. As such, customer service solutions should be provided at runtime and should be able to correctly and timely address customers' requirements and intentions. For instance, when a customer calls in to complain about the poor quality of her newly bought shoes, we must recognize her intention of "returning the shoes" and provide her with the solution of how to return the shoes and apply for the refund [14].

Customer service solutions can be determined either manually or automatically. Determining customer service solutions manually is flexible and human-centric, and the representatives need to have enough expert knowledge to handle all types of customer problems [18]. Several existing automated mechanisms have required expert knowledge learned from rich transaction history data to target most customer requirements. However, these approaches are inaccurate, inefficient and unsatisfactory, and most critically they are unable to generalize for diverse business domains [20]. As such, end customers' satisfaction will be significantly affected, and business quality and profits will also be further influenced.

In this paper, we propose a novel machine learning-based approach, called ICS-Assist, to facilitate customer service staff to identify ideal customer service solutions at runtime. ICS-Assist uses a two-stage learning model, coarse-grained learning and fine-grained learning, to identify the proper service scenario of each query made by the end customer. Moreover, ICS-Assist uses multi-aspect features (i.e. multi-round conversations, customer profiles, staff profiles, and order details) as the inputs to train a deep learning model for fine-grained service scenario recognition. Then ICS-Assist further determines the final solutions based on the "scenario-solution" mapping table constructed by business operators. The main differences between our approach and existing methods are: 1) Our approach can achieve accurate customer service scenario recognition at runtime (i.e., while customer service staff are servicing end-customers); 2) We use a novel "Panel-Student" learning scheme to derive a much smaller yet efficient learning model which can recognize service scenario at a finer granularity, a significant improvement over the traditional "Teacher-Student" model [11]; 3) Our approach uses multi-aspect features instead of the commonly used language feature to train the "Panel-Student" learning scheme and recognize service scenarios.

We implement ICS-Assist and evaluate it using a real-world field study with Alibaba Group. The experiments are conducted for over 6 months. On average, over 12,000 customer service staff handle over 230,000 cases per day. We compare the performance of ICS-Assist with existing semantic and relevance matching methods, including HCAN [20], ESIM-seq [2], DAM [30], and DIIN [8]. The experimental results are two-fold: 1) Our method increases the solution acceptance rate by up to 16%, increases solution coverage rate by up to 25%, reduces average service time by up to 6%, increases customer satisfaction rate by up to 14%, and increases business domain catering rate by up to 17%, compared to the state-of-the-art methods; 2) Our method increases the solution acceptance rate by 24%, increases solution coverage rate by 34%, reduces average service

time by 8%, increases customer satisfaction rate by 19%, and increases business domain catering rate by 22%, compared to the traditional manual method.

The research contributions of this paper are 1) We propose a novel intelligent framework to recognize customer service scenarios and further determine appropriate customer service solutions at runtime. In this way, we extend the idea of the "Teacher-Student" model to propose a generalizable "Panel-Student" distilled learning method that determines suitable customer service scenarios and solutions for multiple e-commerce business domains. 2) We show a real-world field study to demonstrate the efficacy and validity of our proposed approach.

The remainder of this paper is as follows: Section 2 introduces the background; Section 3 illustrates our proposed approach; Section 4 describes the experimental evaluation; Section 5 discusses threats to validity; Section 6 provides related work; Section 7 provides the conclusion and future work.

## 2  Background

### 2.1  Intelligent Customer Service in E-Commerce

E-commerce customer service plays a significant role in business profit-making and customer satisfaction [23]. In contrast to traditional customers' service involving huge human efforts, organizations use intelligent customer service to promote effortless customers experiences and improve productivity. Specifically, the state-of-the-art intelligent customer service is not just multi-channel but omnichannel, which allows the organizations to facilitate effective interactions between them and their customers by unifying the experience across self-assisted and field-service channels [16]. In large e-commerce corporations, such as Alibaba, JD.Com and Amazon, intelligence customer service has been successfully used to save their customer service costs by over 20%. With these successful stories, many small to medium-sized e-commerce companies are starting to develop their intelligent customer service systems [14].

### 2.2  Business Requirements for Customer Service

As a critical component of the business chain, customer service has been regularized by standardized business requirements, which are formulated by several popular e-commerce corporations based on over 20 years' business exploration [5]. These requirements are 1) Customer service solutions should be correctly determined; 2) The customer service system should cover as many customer service solutions as possible; 3) The time spent on customer service dialogues should be minimized; 4) The satisfaction rate of end customers should be maximized; 5) The Customer service system should be able to cater for as many business domains as possible. Hence, the e-commerce industry uses the following business metrics to evaluate the quality of customer service: 1) Solution Acceptance Rate (SAR), which refers to the percentage of solutions that are accepted by end customers; 2) Solution Coverage Rate (SCR), which refers to the proportion of

the solutions that can be recalled from the overall solutions; 3) Average Service Time (AST), which refers to the average time spent on customer service conversations; 4) Customer Satisfaction Rate (CSR), which refers to the percentage of the customers who are satisfied with the customer service; 5) Business Domain Catering Rate (BCR), which refers to how many business domains can be catered for by the customer service system.

## 3    Our Proposed Method

Our approach is based on the following design decisions: 1) Service solutions should be mapped from recognized service scenarios based on the well-established "scenario to solution" mapping rules defined by the e-commerce business; 2) Customer service scenarios must be determined in a runtime manner; 3) The model can utilize a multi-stage paradigm in order to recognize the optimal customer service scenarios to determine the optimal service solutions. The overview of our proposed approach, named ICS-Assist, is shown in Fig. 1. When a customer inquires, the system selects an available customer service staff to start the conversation. After the customer makes each query, ICS-Assist recognizes the relevant service scenarios based on the two-stage machine learning (coarse-grained learning and fine-grained learning) scenario recognition model proposed by us. If scenarios are not found, the customer service staff responds to the customer on her own expert experience; otherwise, ICS-Assist determines the solutions based on the scenario-solution mapping table maintained by the business itself, and the customer service staff confirms and provides the solutions to the customer. If the problem of the end customer is solved, the customer service ends; otherwise, ICS-Assist waits for the end customer to make another query, and the aforementioned procedure repeats until the problem is solved.

### 3.1    Data Preparation & Preprocessing

The data processing pipeline for the service scenario recognition in ICS-Assist is shown in Fig. 2. The input data is generated from the historical customer service log, which contains customer utterances and staff operations (e.g. clicking, hovering and querying) in a service session. The service scenarios clicked or searched by the staff are paired with the customer utterances to form the positive samples in the dataset. We also manually check these automatically generated pairs.

However, the training set is imbalanced as some regular service scenarios have millions of cases, such as inquiries about a delivery, refunding, while others may only contain a few thousand. Thus the corpus of customer utterance can be too sparse to learn a well-generalized model. To address this, we apply the data augmentation method of up-sampling on the scarce cases to expand their original size to 100 times. Finally, we combine the augmented rare cases with the regular cases as the positive samples, and also randomly choose an equal number of irrelevant pairs of service scenarios and customer utterances as negative samples.
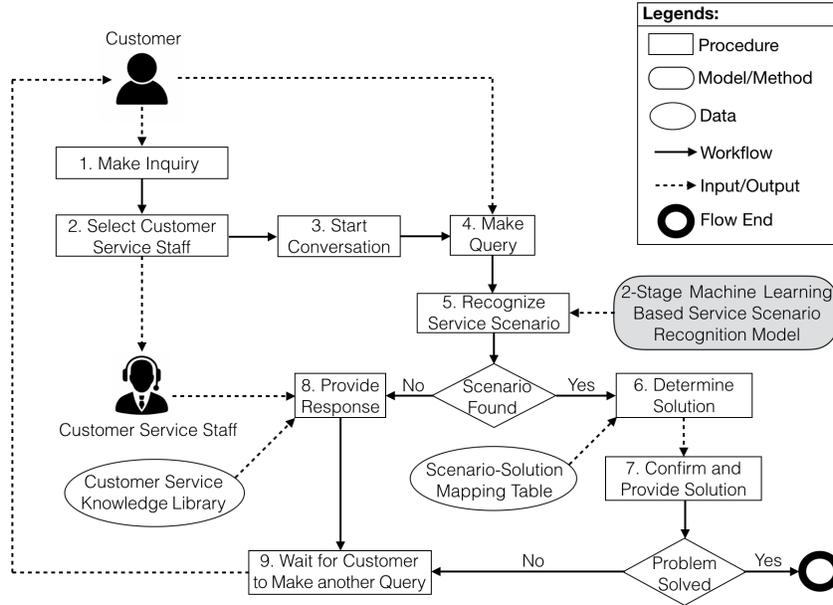
**Fig. 1.** Overview of ICS-Assist

The data structure for each training sample is a triplet consisting of the customer utterance $\mathcal{U}$, the description of standard service scenario $\mathcal{S}$ and label $y$, where $y \in \{0,1\}$, both $\mathcal{U}$ and $\mathcal{S}$ are text, $\mathcal{U} = (w_1^u, w_2^u, w_3^u, ...), \mathcal{S} = (w_1^s, w_2^s, w_3^s, ...)$, $w_i$ is the $i$-th word in the sequence.

The model learning follows a two-stage procedure that contains the coarse-grained ranking and the fine-grained ranking. At the coarse-grained ranking stage, we use a simple approach that narrows down the search range in the candidate set to filter out the irrelevant scenarios. At the fine-grained ranking stage, we propose a Panel-Student knowledge distillation approach to train a lighter model that is able to find out the most suitable service scenarios.

### 3.2    Coarse-Grained Learning Model

We describe the process for the coarse-grained model. First, we compute the representation $\mathbf{u}$, $\mathbf{s}$ for $\mathcal{U}$, $\mathcal{S}$:

$$\sum_i \mathrm{tf\_idf}(w_i) \times \mathrm{word2vec}(w_i) \tag{1}$$

where $u$, $s \in \mathbb{R}^{d\text{-}\{wv\}}$, d_{wv} is the dimensionality for Word2Vec [17]. The representation is exactly the weighted average of the word vector for the corresponding words in the text, where the weight we use here is tf-idf.

We get the top-$K$ suitable scenarios by comparing cosine similarity $cos\_sim(\mathbf{u}, \mathbf{s}_k)$, where $\mathbf{s}_k$ is the representation for a scenario in the candidate set. After that, the top-$K$ candidates would be fed into the fine-grained model.
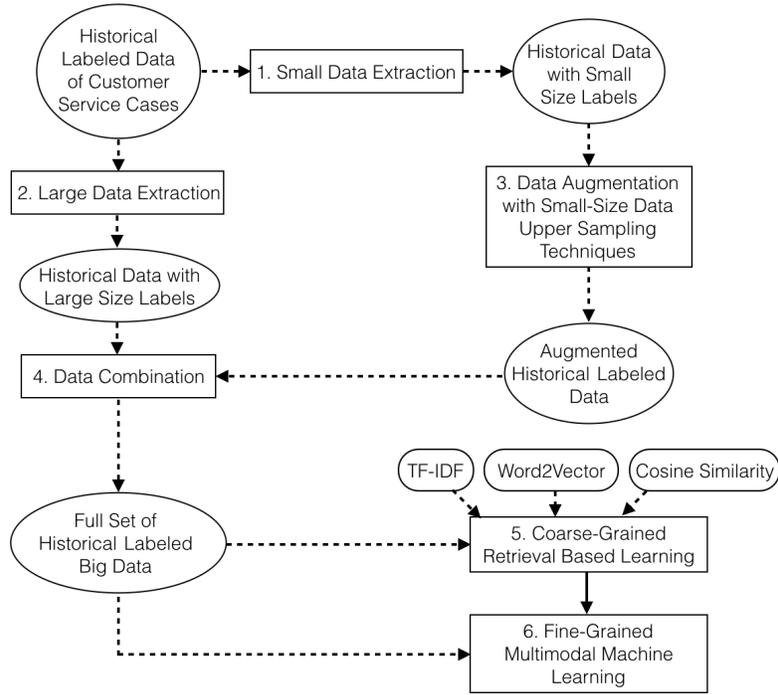
**Fig. 2.** Training Process for Service Scenario Recognition Model

### 3.3   Fine-Grained Learning Model

The fine-grained model learns complicated semantic relationships between customer utterances and service scenario descriptions and finds out the most suitable service scenarios. In our case, the ranking model requires very high precision. To achieve this goal, the simplest way is that we train a model as large as possible with strong fitting capacities, but by doing so the inference time would be slower, which is unfriendly for online recommendation at runtime.

Knowledge distillation [11, 29] is an effective way to distill the knowledge learned from the teacher model and builds an accurate lightweight student model. The teacher model is usually a large neural network or an ensemble of networks containing millions parameters. Hence, the state-of-the-art large-scale pre-trained language models, such as ELMO [19], BERT [4] and XLNet [28], can serve as the teacher network. These models are millstones in natural language processing filed and significantly improve the performance of many downstream tasks such as question answering, textual entailment, and text classification etc.

However, among the aforementioned pre-trained language models, using only one of them as a teacher network seems to be unable to completely train a generalizable student model that can achieve as good performance as the teacher network in diverse business domains. Besides, our empirical studies show that

ELMO has the best performance and is slightly better than BERT in the business domain of Alibaba Movie (which is an Alibaba business portal for watching online movies), while in the business domain of Tmall Global (is an Alibaba web portal for selling imported commodities), ELMO's performance is the worst among the three pre-trained language models. Thus, in order to cater for all types of business domains, we explore a Panel-Student knowledge distillation approach that combines all the three *teachers* to form a generalizable *panel*.

The full details of the fine-tuning "Panel-Student" learning scheme are illustrated in Fig. 3, which can be divided into four layers:

1. the input layer that maps the raw text data into the word embeddings;
2. the representation learning layer that encodes the word embeddings into a comprehensive contextualized representation;
3. the interaction learning layer that further processes the representation and extracts both semantic-oriented and relevance-oriented matching signals between the input utterance and service scenario;
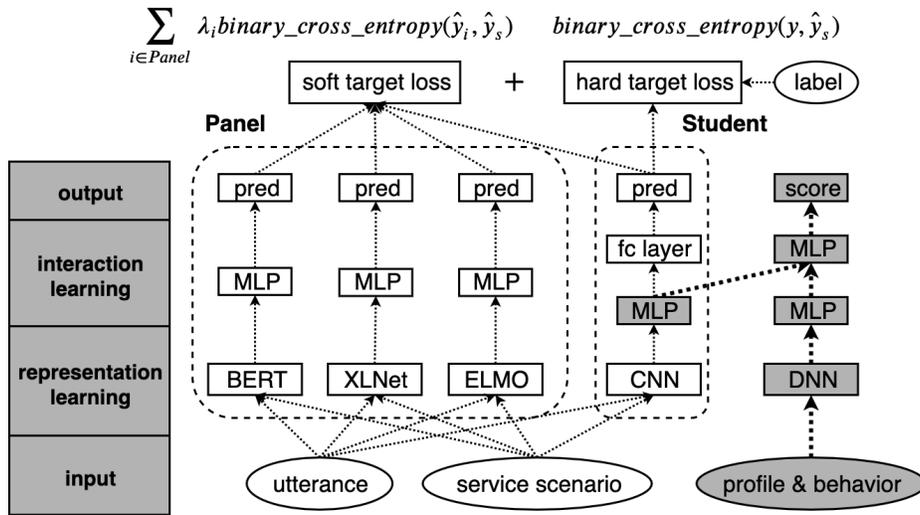4. the output layer that generates the final matching scores.



**Fig. 3.** Fine-Grained "Panel-Student" Learning Model

**Panel-Student Framework** We use three high-capacity pre-trained language models in our panel, including ELMO [19], BERT-LARGE [4], and XLNet [28]. As the core task of our ICS-Assist is to match the customer utterance with a suitable service scenario by leveraging semantic-oriented and relevance-oriented matching signals in the text pairs, we follow the fine-tuning setting of the text-entailment task described in each corresponding teacher model (text-entailment can be viewed as a special case of text match [20]). After the fine-tuning stage, we train the student model under the panel's supervision with the soft target loss

(shown in Figure 3) and the hard target loss with ground truth labels. We choose TextCNN [24] as the student model due to its lightweight and fast inference.

The input layer maps the words within $\mathcal{U}$ and $\mathcal{S}$ into the corresponding embeddings $U = [\mathbf{w}_1^u; \mathbf{w}_2^u; \mathbf{w}_2^u; ... \mathbf{w}_N^u]$ and $S = [\mathbf{w}_1^s; \mathbf{w}_2^s; \mathbf{w}_2^s; ... \mathbf{w}_N^s]$, where $\mathbf{w}_i \in \mathbb{R}^d$ is the corresponding embedding given a word $w_i$, $U, S \in \mathbb{R}^{N \times d}$. We pad the variable length sequence to fixed-length $N$.

At the representation learning stage, we first apply 1-d convolution over $U$ with $k$ different kernel size:

$$\bar{U}^k = \sigma(W_f^k * U + b^k), \tag{2}$$

where $W_f^k$ is the parameter for $k$-th convolution kernel, $b^k$ is the bias, $\sigma$ is the activation function, the output channel number for convolution is $d_o$, $\bar{U}^k \in \mathbb{R}^{N \times d_o}$. Then we take the maximum and average values over the sequence length dimension and concatenate them to form an overall representative semantic signal. Taking the maximum value can effectively extract the features for the occurrence of some keywords and taking the average can be more robust to the noise in the corpus.

$$\bar{\mathbf{u}}_{max}^k = \max(\bar{U}^k), \tag{3}$$

$$\bar{\mathbf{u}}_{mean}^k = \mathrm{mean}(\bar{U}^k), \tag{4}$$

$$\tilde{\mathbf{u}} = [\bar{\mathbf{u}}_{max}^1; ...; \bar{\mathbf{u}}_{max}^k; \bar{\mathbf{u}}_{mean}^1; ...; \bar{\mathbf{u}}_{mean}^k], \tag{5}$$

where $\bar{\mathbf{u}}_{max}^k, \bar{\mathbf{u}}_{mean}^k \in \mathbb{R}^{d_o}$, $\tilde{\mathbf{u}} \in \mathbb{R}^{2kd_o}$. The representation $\tilde{\mathbf{s}}$ for $\mathcal{S}$ is obtained in a similar way.

At the interaction learning stage, we enhance the interaction between $\mathcal{U}$ and $\mathcal{S}$ by applying more complicated arithmetic operations on the original signal obtained in the previous stage. The original signals $\tilde{\mathbf{u}}, \tilde{\mathbf{s}}$, the element-wise multiplication of the original signals, the element-wise square of the difference between the two signals are concatenated together and fed to an MLP to generate the final matching feature $\mathbf{m}$:

$$\mathbf{x} = [\tilde{\mathbf{u}}; \tilde{\mathbf{s}}; \tilde{\mathbf{u}} \otimes \tilde{\mathbf{s}}; (\tilde{\mathbf{u}} - \tilde{\mathbf{s}})^{\circ 2}], \tag{6}$$

$$\mathbf{m} = \mathrm{MLP}(\mathbf{x}), \tag{7}$$

where $\mathbf{x} \in \mathbb{R}^{8kd_o}$. After the student model is trained, the model can make the following prediction:

$$g^s = \sigma(W\mathbf{m} + b), \tag{8}$$

$$\hat{y}_s = \mathrm{sigmoid}(g^s), \tag{9}$$

where $\sigma$ is the activation function, and $\hat{y}_s$ is a real number between 0 and 1, which represents the probability of standard service scenario $\mathcal{S}$ matching the given utterance $\mathcal{U}$.

**Final Hybrid Model for Service Scenario Recognition**  As shown at the rightmost side in Figure 3, we also use multi-aspect features $\bar{\mathbf{m}}$ learned by a DNN model based on the customer profiles, the log of historical customer behavior and customer service staff. The intermediate multi-aspect feature $\bar{\mathbf{m}}$ will be combined with $\mathbf{m}$ (see equation 7) generated by the student TextCNN, and then fed to an MLP to make the final prediction as follows:

$$g^h = \text{MLP}([\mathbf{m}; \bar{\mathbf{m}}]), \tag{10}$$

$$\hat{y}_h = \text{sigmoid}(g^h). \tag{11}$$

The training of the final hybrid model which utilizes the output from our student model has three phases. The first two phases are used for training the student model alone, and the third phase is used for training the hybrid model:

1. Each teacher model in the panel is fine-tuned;
2. The student model (i.e., TextCNN) is trained within the "Panel-Student" scheme using the loss function below:

$$\sum_{i \in Panel} \lambda_i \text{binary\_cross\_entropy}(\hat{y}_i, \hat{y}_s) + \text{binary\_cross\_entropy}(y, \hat{y}_s), \tag{12}$$

   The lost function consists of two types of loss: 1) the soft-target loss between the student model's predictions $\hat{y}_s$ and each teacher model's predictions $\hat{y}_i$; 2) the hard target loss between the predictions of the student model $\hat{y}_s$ and the ground truth labels $y$.
3. All layers within TextCNN up to the MLP are extracted and combined with the output from the MLP layer within the DNN model at the rightmost in Figure 3 to construct the final hybrid model for scenario prediction. The hybrid model is trained using this loss function:

$$\text{binary\_cross\_entropy}(y, \hat{y}_h), \tag{13}$$

   where $\hat{y}_h$ is the prediction for the hybrid model.

### 3.4   Scenario Recognition & Solution Mapping

Customer service scenarios are determined by our two-stage scenario recognition model, which are represented as the "parameters" for determining solutions. It matches the customer utterance with the best service scenario by solving a pairwise text-match task. Taking the case of "complaining the bought shoes" as an example, the recognized best scenario is "returning the commodities (shoes)".

Given a determined customer service scenario determined by our two-stage scenario recognition framework, the customer service solution can be determined and selected based on the scenario-solution mapping table formulated by the e-commerce company itself according to its business strategies. The customer service solution can be a customized one-to-one mapping from the customer service scenario. Specifically, the solution can be a standard service manual, a road-map, or just a predefined answer. For instance, the above-determined scenario is mapped to the solution of how to apply for the refund of the shoes.

## 4    Experimental Evaluation

Our experiment was conducted in the real customer service of Alibaba Group. We implemented ICS-Assist as an enterprise-level service system. Over 12,000 customer service staff use ICS-Assist to serve for over 230,000 cases per day on average, and this procedure lasted for over 6 months. In our experimental environment, the queries made by end customers are dispensed to dedicated query processing servers by the query router. Each query processing server encapsulates and passes the query to the service scenario recognition model in our ICS-Assist to predict the service scenarios. The recognized scenarios are then mapped with the service solutions, which are sent to customers by customer service staff.

### 4.1    Experimental Procedure

The experimental procedure consists of three parts:

1. We apply our proposed "Panel-Student" model and the baseline models [20, 2, 30, 8] on the public dataset Quora [13], and compare the performance among them. For the baseline models, we reproduce them according to their best hyper-parameter settings. For our proposed model, we set the convolution kernel widths from 2 to 5, and the output channel numbers are all 64. The layer number of the MLP module in equation 7 is 3. The dropout rate is 0.2 and the L2 regularization coefficient is 0.05. The activation function for all the layers is ReLu. The optimizer is Adam [15] with a constant learning rate of 1e-4, decay rate $\beta_1$ of 0.9, and $\beta_2$ of 0.999.
2. We also conduct similar experiments on our historical dataset. The only difference is that, for our proposed model, we use an additional neural model to handle the handcrafted multi-aspect features to construct the final hybrid model, we adopt a similar architecture like a Wide-Deep model [3], and combine it with the text-based features from the TextCNN model. The training for the final model follows a two-stage paradigm: 1) We freeze the parameters within TextCNN model and train the DNN with a constant learning rate of 1e-3 until convergence; 2) We make TextCNN's parameters trainable and restart the training phase with an exponential decay learning rate (initial learning rate: 1e-4; decay rate: 0.95; decay step: 10000).
3. Since the purpose of the two steps above is to demonstrate the superiority of our proposed "Panel-Student" model, we replace this hybrid model in ICS-Assist with each of the state-of-the-art baseline models [20, 2, 30, 8] to create several variants of ICS-Assist and compare our proposed ICS-Assist (with the "Panel-Student" model) with these ICS-Assist variants as well as the manual method against the 5 business evaluation metrics (SAR, SCR, AST, CSR, and BCR) mentioned in Section 2.2.

### 4.2    Experimental Results

Table 1 shows the comparison between the performance (accuracy, precision, recall, f1-score, and latency) of our proposed Panel-Student model (PS model

henceforth) and the other existing models on the Quora dataset. Due to the less execution complexity of our proposed model compared to other models, the accuracy, precision, recall and f1-score of our PS model are slightly less than hcan-hybrid, hcan-only rm [20], dam [30], diin [8] and esim-seq [2], and slightly better than hcan-only sm [20], but our PS models latency is much lower than these models. We also implement the panel-student model with a single teacher (i.e. BERT, XLNet, and ELMO), and obtain three Teacher-Student models (TS models henceforth). The performance of these three TS models is also worse than our PS model. As such, our model is the best one among all the models.

**Table 1.** Model performance comparison on Quora dataset

| model | accuracy | precision | recall | f1 | latency(ms) |
|---|---|---|---|---|---|
| hcan - hybrid | 0.831 | 0.832 | 0.830 | 0.831 | 81 |
| hcan - only sm | 0.791 | 0.791 | 0.791 | 0.791 | 73 |
| hcan - only rm | 0.821 | 0.824 | 0.817 | 0.820 | 21 |
| dam | 0.855 | 0.856 | 0.854 | 0.855 | 109 |
| diin | 0.873 | 0.877 | 0.867 | 0.872 | 151 |
| esim - seq | 0.843 | 0.846 | 0.839 | 0.842 | 95 |
| TS - BERT | 0.791 | 0.783 | 0.792 | 0.787 | 15 |
| TS - XLNet | 0.807 | 0.795 | 0.809 | 0.802 | 15 |
| TS - ELMO | 0.781 | 0.769 | 0.759 | 0.764 | 13 |
| our PS model | 0.811 | 0.807 | 0.819 | 0.813 | 11 |

Table 2 shows the comparison between the performance of our proposed PS model and the existing state-of-the-art models using our dataset. Again, our PS model outperforms all the baseline models in terms of the overall model performance due to the less execution complexity of our model than others.

**Table 2.** Model performance comparison on our dataset

| model | accuracy | precision | recall | f1 | latency(ms) |
|---|---|---|---|---|---|
| hcan - hybrid | 0.878 | 0.876 | 0.879 | 0.877 | 198 |
| hcan - only sm | 0.845 | 0.847 | 0.841 | 0.844 | 186 |
| hcan - only rm | 0.850 | 0.848 | 0.852 | 0.850 | 53 |
| dam | 0.914 | 0.914 | 0.914 | 0.914 | 265 |
| diin | 0.894 | 0.899 | 0.887 | 0.893 | 387 |
| esim - seq | 0.930 | 0.932 | 0.928 | 0.930 | 241 |
| TS - BERT | 0.871 | 0.874 | 0.877 | 0.875 | 28 |
| TS - XLNet | 0.878 | 0.884 | 0.889 | 0.886 | 26 |
| TS - ELMO | 0.853 | 0.851 | 0.861 | 0.856 | 25 |
| our PS model | 0.894 | 0.892 | 0.895 | 0.893 | 25 |

Table 3 shows the improvement rate of the business metrics of our proposed ICS-Assist (ICS-Assist (PS)) over the manual method and the variants of ICS-

Assist with state-of-the-art models, including the teacher-student model using each single teacher model in our panel. Our ICS-Assist (PS) performs better than all the other variants of ICS-Assist by up to 16%, 25%, 6%, 14%, and 17%, in terms of SAR (Solution Acceptance Rate), SCR (Solution Coverage Rate), AST (Average Service Time), CSR (Customer Satisfaction Rate) and BCR (Business Domain Catering Rate). Although our PS model's performance (e.g. f1-score) is slightly worse than other models (e.g. dam and esim-seq), our approach still performs better in the business metrics. This is because our PS model has lower latency than other models, and it enables customer service staff to timely utilize the recommended solutions. ICS-Assist (PS) increases SAR by 24%, increases SCR by 34%, decreases AST by 8%, increases CSR by 19%, and increases BCR by 22%, compared to the manual method.

**Table 3.** Business performance improvement results

| ICS-Assist (PS) vs. | SAR | SCR | AST | CSR | BCR |
|---|---|---|---|---|---|
| manual | 24% | 34% | 8% | 19% | 22% |
| ICS-Assist (dam) | 13% | 19% | 4% | 5% | 7% |
| ICS-Assist (hcan) | 16% | 25% | 6% | 14% | 17% |
| ICS-Assist (diin) | 12% | 19% | 7% | 12% | 15% |
| ICS-Assist (esim) | 10% | 15% | 6% | 11% | 14% |
| ICS-Assist (TS-BERT) | 3% | 2% | 3% | 3% | 5% |
| ICS-Assist (TS-XLNet) | 1% | 2% | 2% | 2% | 7% |
| ICS-Assist (TS-ELMO) | 7% | 6% | 4% | 9% | 10% |

From the experimental results, we can conclude that our approach outperforms other automated state-of-the-art methods as well as the manual method in all the business metrics. The main reasons are as follows: 1) Our method assembles the three pre-trained language models (i.e. BERT, XLNet and ELMO) to distill a more generalizable model that creates better language representations for multiple business domains; 2) Our method takes multi-aspect features as the inputs for the scenario recognition model; 3) Our method employs a two-stage learning approach to maximize the validity of the recommended results, and it makes a reasonable prepossessing on the historical data to address its inevitable drawbacks related to quality, volume, and noise.

## 5   Threats to Validity

The threats to validity are as follows: 1) The historical customer service data provided by Alibaba Group largely focus on the east Asian and southeast Asian countries, and the countries from other continents are relatively few. 2) The model training parameters with the PS model can be further tuned. The current parameters may not yield an optimized deep learning model because they may cause a local minimum instead of a global minimum. 3) The three representation learning-oriented models (BERT, XLNet, and ELMO) constitute the Panel, but more pre-trained language models could have been investigated.

# 6   Related Work

## 6.1   Neural Text Matching Techniques

One line of work related to our system is Neural Text Matching. Text Matching is a core task in many NLP and information retrieval applications, the mainstream of which can be divided into Semantic Matching (SM) and Relevance Matching (RM). Although both SM and RM are modelling similarities between two pieces of texts, SM emphasizes the semantic understanding and reasoning while RM focuses more on keyword matching signals. Typical SM tasks includes question answering [1], paraphrase identification [27], and natural language inference [2, 8]. RM models, such as DRMM [9], Co-PACRR [12] and MP-HCNN [21] are frequently used in IR applications like search engines to rank documents by relevance given a user query. In our work, both semantic and relevance matching technologies are involved in our model to enable more comprehensive language understanding and identify suitable service scenarios.

## 6.2   Collaborative Filtering Techniques

Because our system aims to recommend suitable service scenarios and solutions to the customer service staff, in this way the research work on recommendation Systems is also related to our work. Most recommendation systems are based on collaborative filtering, which learns a representation of user and item based on the rating matrix, and then predict the rating assigned a user given an unseen item. Currently, many recommendation systems adopt neural networks [26, 10, 6] to learn a good dense representation and the interaction between the user and item and achieve the state-of-the-art performance. However in our scope, mechanically matching the user and service scenario would ignore the customer's intention and requirement thus impair our service quality.

## 6.3   Knowledge Distillation Techniques

Researchers from the University of Waterloo try to transfer deep language representation like BERT to a lightweight neural network such as single-layer BiLSTM [25]. But they do not employ multiple teachers knowledge to distil a simple student model. This experience motivates our multiple knowledge distilling. In addition, The model Fitnets [22] is proposed by A. Romero. This model has extended the model compression idea and introduces the intermediate-level hints techniques to simplify a deeper and thinner student network with fewer parameters and better generalization. The new loss function is imported in hidden layers feature maps, which helps to reduce parameters in our work. Recently, the IBM researchers have proposed to train the student model from an ensemble of multiple teachers [7]. They implemented different deep neural networks to train convolutional neural network acoustic models on a medium-sized speech corpus. The experimental results highlight that the proposed training techniques could increase a significant amount of knowledge to the student. Hence, our work also follows the idea of distilled learning by proposing a "Panel-Student" model.

## 7    Conclusion & Future Work

Identifying proper customer service solutions is critical to e-commerce businesses. Existing service solution determination methods are usually unsatisfactory to end customers. This is because they are of low efficiency and unable to achieve runtime solution determination. Hence, this paper proposes an innovative framework, called ICS-Assist, to determine customer service solutions at runtime. We designed a novel two-stage learning model to identify customer service scenarios, which are mapped to end solutions. We implemented ICS-Assist and evaluated it in a 6-month real-world field study at Alibaba Group. The experimental results show that ICS-Assist improves the five business evaluation metrics (solution acceptance rate, solution coverage rate, average service time, customer satisfaction rate and business domain catering rate) by up to 16%, 25%, 6%, 14%, and 17% respectively, compared to the state-of-the-art methods, and it outperforms the manual method by 24%, 34%, 8%, 19%, and 22% respectively, in terms of the five business evaluation metrics. Our future work includes: 1) Explore more representation learning models for determining the members in the panel; 2) Design robust light-weight pre-trained models for customer services; 3) Investigate different customer service application areas such as finance.

## References

1. Chen, D., Fisch, A., Weston, J., Bordes, A.: Reading wikipedia to answer open-domain questions. In: ACL. pp. 1870–1879 (2017)
2. Chen, Q., Zhu, X., Ling, Z.H., Wei, S., Jiang, H., Inkpen, D.: Enhanced lstm for natural language inference. In: ACL. pp. 1657–1668 (2017)
3. Cheng, H.T., Koc, L., Harmsen, J., Shaked, T., Chandra, T., Aradhye, H., Anderson, G., Corrado, G., Chai, W., Ispir, M., et al.: Wide & deep learning for recommender systems. In: Proceedings of the 1st workshop on deep learning for recommender systems. pp. 7–10 (2016)
4. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: ACL. pp. 4171–4186 (2019)
5. Eales-Reynolds, L.J., Clarke, C.: Impact of a novel training experience on the development of a customer service culture in a large hospital trust. International journal of health care quality assurance pp. 483–497 (2012)
6. Ebesu, T., Shen, B., Fang, Y.: Collaborative memory network for recommendation systems. In: SIGIR. pp. 515–524 (2018)
7. Fukuda, T., Suzuki, M., Kurata, G., Thomas, S., Cui, J., Ramabhadran, B.: Efficient knowledge distillation from an ensemble of teachers. In: Interspeech. pp. 3697–3701 (2017)
8. Gong, Y., Luo, H., Zhang, J.: Natural language inference over interaction space. In: ICLR (2018)
9. Guo, J., Fan, Y., Ai, Q., Croft, W.B.: A deep relevance matching model for ad-hoc retrieval. In: CIKM. pp. 55–64 (2016)
10. He, X., Liao, L., Zhang, H., Nie, L., Hu, X., Chua, T.S.: Neural collaborative filtering. In: WWW. pp. 173–182 (2017)

11. Hinton, G., Oriol, V., Jeff, D.: Distilling the knowledge in a neural network. In: NIPS Deep Learning and Representation Learning Workshop (2015)
12. Hui, K., Yates, A., Berberich, K., De Melo, G.: Co-pacrr: A context-aware neural IR model for ad-hoc retrieval. In: WSDM. pp. 279–287 (02 2018)
13. Iyer, S., Dandekar, N., Csernai, K.: First quora dataset release: Question pairs (2017), https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs
14. Kaufman, R.: Why your customer service training won't lead to happy customers (or inspired employees). The Journal for Quality and Participation p. 33 (2015)
15. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2015)
16. Kresch, M.: What is intelligent customer service (2016), https://cloudblogs.microsoft.com/dynamics365/bdm/2016/01/19/what-is-intelligent-customer-service Accessed May, 12, 2020
17. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: NIPS. pp. 3111–3119 (2013)
18. Mirchandani, K.: Learning racial hierarchies: Communication skills training in transnational customer service work. Journal of Workplace Learning pp. 338–350 (2012)
19. Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L.: Deep contextualized word representations. In: NAACL. pp. 2227–2237 (2018)
20. Rao, J., Liu, L., Tay, Y., Yang, W., Shi, P., Lin, J.: Bridging the gap between relevance matching and semantic matching for short text similarity modeling. In: EMNLP-IJCNLP. pp. 5373–5384 (2019)
21. Rao, J., Yang, W., Zhang, Y., Ture, F., Lin, J.: Multi-perspective relevance matching with hierarchical convnets for social media search. In: AAAI. pp. 232–240 (2019)
22. Romero, A., Ballas, N., Kahou, S.E., Chassang, A., Gatta, C., Bengio, Y.: Fitnets: Hints for thin deep nets. In: ICLR (2015)
23. Sari, P.K., Alamsyah, A., Wibowo, S.: Measuring e-commerce service quality from online customer review using sentiment analysis. In: Journal of Physics: Conference Series. p. 012053 (2018)
24. Sun, X., Ma, X., Ni, Z., Bian, L.: A new lstm network model combining textcnn. In: International Conference on Neural Information Processing. pp. 416–424. Springer (2018)
25. Tang, R., Lu, Y., Liu, L., Mou, L., Vechtomova, O., Lin, J.: Distilling task-specific knowledge from bert into simple neural networks. arXiv preprint arXiv:1903.12136 (March 2019)
26. Wang, H., Wang, N., Yeung, D.Y.: Collaborative deep learning for recommender systems. In: KDD. pp. 1235–1244 (2015)
27. Wang, Z., Hamza, W., Florian, R.: Bilateral multi-perspective matching for natural language sentences. In: IJCAI. pp. 4144–4150 (2017)
28. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R.R., Le, Q.V.: Xlnet: Generalized autoregressive pretraining for language understanding. In: NIPS. pp. 5754–5764 (2019)
29. You, S., Xu, C., Xu, C., Tao, D.: Learning from multiple teacher networks. In: KDD. pp. 1285–1294 (2017)
30. Zhou, X., Li, L., Dong, D., Liu, Y., Chen, Y., Zhao, W.X., Yu, D., Wu, H.: Multi-turn response selection for chatbots with deep attention matching network. In: ACL. pp. 1118–1127 (2018)