

Athens University of Economics and Business



School of Information Sciences and Technology
Department of Informatics
Athens, Greece

M.Sc. Thesis in Computer Science

OCaTS: an Online Cost-aware Teacher-Student Framework to Reduce the Calls to Large Language Models

Ilias Marios Stogiannidis

Supervisors: Ion Androutsopoulos

Prodromos Malakasiotis

Stavros Vassos

October, 2023

Ilias Marios Stogiannidis

OCaTS: an Online Cost-aware Teacher-Student Framework to Reduce the Calls to Large Language Models

October, 2023

Supervisors: Ion Androutsopoulos[‡], Prodromos Malakasiotis[‡], Stavros Vassos[†]

[‡] Athens University of Economics and Business

[†] helvia.ai

Athens University of Economics and Business

School of Information Sciences and Technology

Department of Informatics

Natural Language Processing Group, Information Processing Laboratory

Athens, Greece

Abstract

Recently, Deep Learning advancements have led to considerable progress in Natural Language Processing (NLP). One of these advancements is the emergence of Large Language Models (LLMs) which are highly capable of following instructions and can perform remarkably in downstream tasks in zero- and few-shot settings. Therefore, small and medium-sized businesses (SMEs) that cannot afford the cost of creating large task-specific training datasets, or the cost of training their own task-specific models, are increasingly adopting the use of such LLMs and particularly third-party services that allow them to prompt those LLMs without the need of having their own deployment. However, such services currently require a payment per call, which becomes a considerable operating expense (OpEx). Moreover, customer inputs (e.g., questions about company policies in Employee Experience chatbots) are often very similar over time, thus SMEs end up prompting LLMs with very similar instances. In this thesis, we propose a teacher-student framework that allows reducing the calls to LLMs (teacher) by caching previous LLM responses and using them to train a local inexpensive model (student) on the SME side. The framework includes criteria taken from active learning for deciding when to trust the local model or call the LLM, and a methodology to tune the criteria and measure the trade-off between performance and cost. We evaluate our framework with two of the most powerful LLMs at the moment, GPT-3.5 and GPT-4, and two inexpensive students, a k -NN algorithm and a Multi Layer Perceptron using two (very) common tasks for companies, Intent Recognition and Sentiment Analysis. Experimental results indicate that significant OpEx savings can be achieved with only slightly lower performance.

Περίληψη

Τα τελευταία χρόνια, η ανάπτυξη της βαθιάς μάθησης έχει προκαλέσει σημαντική πρόοδο στην Επεξεργασία Φυσικής Γλώσσας (ΕΦΓ). Μία από αυτές τις προόδους είναι η ανάδειξη των μεγάλων γλωσσικών μοντέλων (ΜΓΜ) τα οποία είναι ιδιαίτερα ικανά να ακολουθούν οδηγίες και μπορούν να αποδώσουν αξιοσημείωτα στη λύση προβλημάτων με τη χρήση λίγων ή καθόλου δεδομένων. Ως εκ τούτου, οι μικρές και μεσαίες επιχειρήσεις (ΜΜΕ) που δεν μπορούν να αντέξουν το κόστος δημιουργίας μεγάλων συνόλων δεδομένων εκπαίδευσης για συγκεκριμένες εργασίες ή το κόστος της εκπαίδευσης των δικών τους μοντέλων, στρέφονται όλο και περισσότερο στη χρήση τέτοιων ΜΓΜ και ιδιαίτερα σε υπηρεσίες τρίτων που τους επιτρέπουν να προτρέπουν (**prompt**) αυτά τα ΜΓΜ χωρίς να χρειάζεται να έχουν τα δικά τους σε περιβάλλοντα ανάπτυξης. Ωστόσο, οι εν λόγω υπηρεσίες απαιτούν επί του παρόντος πληρωμή ανά κλήση, η οποία γίνεται ένα σημαντικό λειτουργικό κόστος (ΛΚ). Επιπλέον, τα εισερχόμενα αιτήματα των πελατών (π.χ. ερωτήσεις σχετικά με τις πολιτικές της εταιρείας προς ένα ψηφιακό βοηθό εργαζομένων) είναι συχνά πολύ παρόμοια με την πάροδο του χρόνου, επομένως οι ΜΜΕ καταλήγουν να προτρέπουν τα ΜΓΜ με πολύ παρόμοιες περιπτώσεις. Στην παρούσα διπλωματική, προτείνουμε ένα πλαίσιο διδάσκοντα-μαθητή που επιτρέπει τη μείωση των κλήσεων σε ΜΓΜ (διδάσκων) με την αποθήκευση προηγούμενων απαντήσεων ΜΓΜ και τη χρήση τους για την εκπαίδευση ενός τοπικού χαμηλού κόστους μοντέλου (μαθητή) στην πλευρά της ΜΜΕ. Το πλαίσιο περιλαμβάνει κριτήρια που προέρχονται από την ενεργητική μάθηση για την απόφαση πότε να εμπιστευτούμε το τοπικό μοντέλο ή να καλέσουμε το ΜΓΜ, καθώς και μια μεθοδολογία για τη ρύθμιση των κριτηρίων και τη μέτρηση του συμβιβασμού μεταξύ απόδοσης και κόστους. Αξιολογούμε το πλαίσιο μας με δύο ισχυρά ΜΓΜ την συγκεκριμένη χρονική στιγμή, τα GPT-3.5 και GPT-4, και δύο πολύ φθηνά μοντέλα, όπως ένας αλγόριθμος k -NN και ένα Νευρωνικό Δίκτυο σε δύο πολύ συνηθισμένες εργασίες για εταιρείες, την Αναγνώριση Προθέσεων και την Ανάλυση Συναισθήματος. Τα πειραματικά αποτελέσματα δείχνουν ότι μπορεί να επιτευχθεί σημαντική μείωση λειτουργικών εξόδων με μόνοχσ ελαφρώς χαμηλότερες επιδόσεις.

Acknowledgements

I would like to sincerely thank my academic advisors Professor Ion Androutsopoulos and Prodromos Malakasiotis, as well as Stavros Vassos from helvia.ai, for their ongoing support and guidance throughout this thesis. I would like to express my appreciation to all the members of the Natural Language Processing Group in the Department of Informatics of *AUEB*, especially Dimitris Mamakas and Panos Kaliosis, and my long-time friend Odysseas Diamantopoulos for the many discussions we had about the issues encountered during this thesis. Lastly, my heartfelt appreciation goes to my parents, who have been a constant source of support and encouragement throughout both the highs and lows of this academic endeavor.

The work of this thesis was supported by Google's TPU Research Cloud (TRC) program.¹

¹<https://sites.research.google/trc/about/>

Contents

Abstract	iv
Acknowledgements	v
1 Introduction	1
1.1 Motivation and Problem Statement	1
1.2 Proposed Approach: OCATS	1
1.3 Contributions	2
1.4 Thesis Structure	3
2 Background and Related Work	5
2.1 Introduction	5
2.2 Background	5
2.2.1 Deep Learning	5
2.2.2 Transformers	6
2.2.3 Large Language Models (LLMs)	7
2.2.4 k -Nearest Neighbors	8
2.3 Related Work	10
2.3.1 Knowledge Distillation	10
2.3.2 Self-Labeling	12
2.3.3 Active Learning	13
2.3.4 Uncertainty Sampling	14
2.3.5 Diversity Sampling	15
2.3.6 Representativeness Sampling	16
3 System Design	17
3.1 Introduction	17
3.2 System Design	17
3.2.1 Cache	17
3.2.2 Teacher	18
3.2.3 Student	19
3.2.4 Criteria	19
3.2.5 Discounted evaluation metric	19
4 Experimental Setup	21

4.1	Introduction	21
4.2	Datasets	21
4.2.1	BANKING77 Dataset	21
4.2.2	Large Movie Review Dataset	22
4.2.3	Few-shot scenario	22
4.3	Experiments and System Evaluation	23
4.3.1	Online setting and evaluation measure	23
4.3.2	Criteria	23
4.3.3	Experiment on BANKING77 with a k -NN student and GPT-4 teacher	24
4.3.4	Experiment on BANKING77 with an MLP student and GPT-4 teacher	27
4.3.5	Experiment on LMR with the k -NN student and GPT-3.5 teacher	30
4.4	Related work on this topic	33
5	Conclusion and Future work	35
5.1	Conclusion	35
5.2	Future work	35
5.3	Limitations	36
5.4	Ethics statement	36
	Bibliography	37
	List of Acronyms	43
	List of Figures	45
	List of Tables	47
	List of Algorithms	49

Introduction

1.1 Motivation and Problem Statement

The rapid growth of Natural Language Processing (NLP) has led to the development of high-performing language models such as GPT-4 [Ope23], Llama-2 [Tou+23a], and others, which perform remarkably in most NLP tasks. These models have billions of parameters and do not necessitate the creation of large task-specific datasets for fine-tuning in order to perform downstream tasks. However, the cost of pre-training such models from scratch is significant. The use of these models in real-world applications is often hindered by the need for high operating costs for training and deployment. Prompting such models, which are tuned to follow instruction, performs remarkably well in zero- and few-shot settings. Therefore, most small and medium-sized enterprises (SMEs) turn to third-party services that allow them to prompt LLMs via the use of an API. For instance, an SME may offer customer care chatbots; for example, instruct LLMs such as GPT-4 to recognize user intents and drive the chatbot-customer interaction. However, the top LLMs now charge a fee per prompting call, which adds up to a major operating expense (OpEx) for SMEs. Furthermore, client inputs (e.g., dialog turns) are frequently quite similar over time, thus SMEs wind up calling LLMs to handle inputs that are very similar to inputs handled by the LLMs during earlier (already paid) calls. This thesis aims to explore a novel cost-aware Teacher-Student framework for Large Language Models and evaluate its potential to reduce the total cost required for inference while mitigating performance reduction by using previous responses of the LLM.

1.2 Proposed Approach: OCaTS

This thesis introduces OCaTS, an *Online Cost-aware Teacher-Student* framework that allows reducing calls to a commercial LLM treated as a teacher model by caching its previous responses and using it to train a local, low-cost student model. OCaTS includes criteria for deciding when to trust the student or call the teacher, as well as a methodology for tuning the criteria and measuring the trade-off between performance and cost. Unlike the usual teacher-student training for knowledge distillation [HVD15b], here the teacher does not train the student for all available

instances (in this case, all incoming customer input). Also, unlike teacher-student approaches to self-training [Li+21a; Mi+21b], the teacher is already reasonably effective (however costly). In this sense, this work is more similar to active learning [Set12; Mon21], but OCaTS trains the student using labels provided by an LLM teacher rather than humans, and there is no large pool of unlabeled instances (customer input) to select from, as instances (customer requests) arrive online.

1.3 Contributions

The main contributions of this thesis are:

1. We introduce a general teacher-student framework that helps SMEs reduce the prompting calls to commercial LLMs and the corresponding OpEx costs by caching the responses of the LLMs and training inexpensive local student models.
2. We present discounted versions of standard evaluation metrics that enable SMEs to quantify how much they prefer fewer LLM calls in comparison to increased user frustration (e.g., due to lower accuracy) and adjust the framework's criteria that decide when to trust the local student model or invoke the LLM teacher accordingly.
3. We instantiate the framework with GPT-3.5 or GPT-4 as teachers, and a k -NN or MLP classifier as students.
4. We conduct experiments on two common tasks for small and medium-sized enterprises (SMEs): intent recognition and sentiment analysis. Our results demonstrate that it is possible to achieve considerable cost savings while only slightly compromising performance. This is the initial step in exploring the advantages of the proposed framework with more datasets, models, and business scenarios.
5. We open-source the code of this thesis to facilitate further research on the subject.¹ Additionally, a shorter version of this thesis has been accepted to the Finding of The 2023 Conference on Empirical Methods in Natural Language Processing.

¹<https://github.com/stoyian/OCaTS>

1.4 Thesis Structure

Chapter 2.3 - Background and Related Work

In Chapter 2, we will provide an overview of similar work we draw inspiration from. This will involve a concise examination of active learning and distillation. We will discuss the advantages and limitations of each method and how they have been applied to NLP tasks. We will also highlight how we diverge from them.

Chapter 3 - System Design

In Chapter 3, we will present the theoretical architecture of the framework. We will also discuss the components of the framework and how these components interact with each other. Finally, we will discuss the criteria we use for these interactions and introduce a new cost-aware evaluation metric to measure the framework's performance.

Chapter 4 - Experimental Setup

In Chapter 4, we will introduce the first incarnation of our framework with a simple student and a powerful commercial LLM teacher. We will discuss the methodology for identifying the best criteria, and we will also evaluate our framework in the Intent recognition and Sentiment Analysis tasks incorporating both our proposed metric but also the conventional ones.

Chapter 5 - Conclusions and Future work

In the final chapter, we will summarize the key points of this thesis. Finally, we will also discuss possible future work that was not explored due to time restrictions.

Background and Related Work

2.1 Introduction

LLMs have become increasingly popular in the NLP world because of their impressive performance in adapting to various downstream tasks without requiring a lot of task-specific data. Therefore, more and more companies that cannot afford the cost of creating large task-specific training datasets are increasingly adopting the use of these models. However, their usage is often impeded by their large size, the amount of computing power, and the high monetary cost they require. To address these issues, researchers have investigated different techniques to reduce the resources needed while maintaining the same levels of performance, such as incorporating active learning [Set12] principles to acquire better demonstrators for in-context learning [Mar+23] or distillation [HVD15b] to reduce size. In this chapter, we will give an outline of related research in this area, starting with some basic information about the field, and then proceeding to examine active learning and distillation, two of the most commonly employed techniques for reducing the computing resources needed in neural networks.

2.2 Background

The use of deep learning techniques to learn intricate patterns and structures in language data has been a major area of research in NLP. This has enabled the generation of high-quality text and the execution of a wide range of NLP tasks. Nevertheless, these techniques require a high computational cost, which poses a major obstacle and has led to research in optimizing their efficiency.

2.2.1 Deep Learning

Deep learning [LBH15] has become a powerful technique for addressing intricate problems in a variety of fields, including computer vision, natural language processing, and speech recognition. At its core, deep learning involves training neural networks with multiple layers to recognize intricate patterns in data, which can be used for various NLP tasks such as text categorization and generation. The success

of deep learning is due to its capacity to automatically learn features from data, allowing it to manage complex and high-dimensional data with minimal human input. The increasing availability of large datasets and computing power has accelerated the progress of deep learning [Alo+19], resulting in breakthroughs in areas such as image recognition, speech synthesis, and natural language understanding. Deep learning is a very active area of research and there are ongoing efforts to develop new methods, architectures, and applications.

2.2.2 Transformers

Transformers are a type of deep learning model that has become increasingly popular in the field of natural language processing [Vas+17]. In contrast to recurrent neural networks (RNNs) [RHW86], which process input in a recursive manner, transformers, during training, are capable of processing entire sequences of data simultaneously, making them suitable for tasks such as language modeling, machine translation, and text generation. The success of transformers can be attributed to the use of self-attention mechanisms to weigh the importance of different input tokens, enabling them to capture long-range dependencies, which is critical for many natural language processing tasks. Several transformers employ the Encoder-Decoder architecture [Vas+17], which maps an input sequence of symbols to a sequence of continuous representations. The decoder then produces an output sequence of symbols one at a time based on this continuous representation. Figure 2.1 shows an analytical overview of the transformer architecture and Figure 2.2 shows the encoder-only and decoder-only architectures.

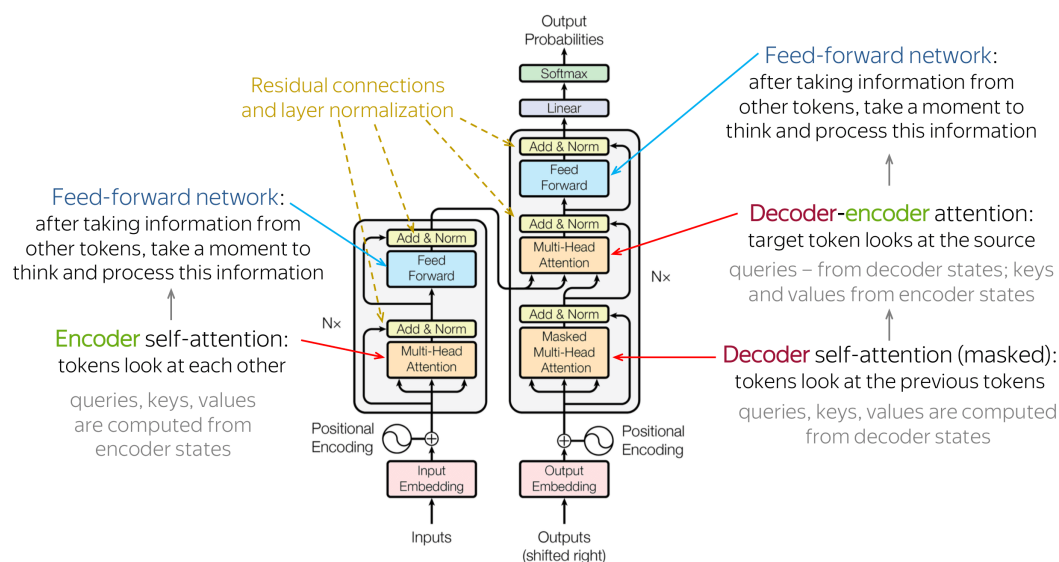


Fig. 2.1: The Transformer’s Encoder-Decoder architecture is illustrated on the left and right sides. The components of the encoder are displayed on the left, while the components of the decoder are shown on the right. | Source:[Vas+17]

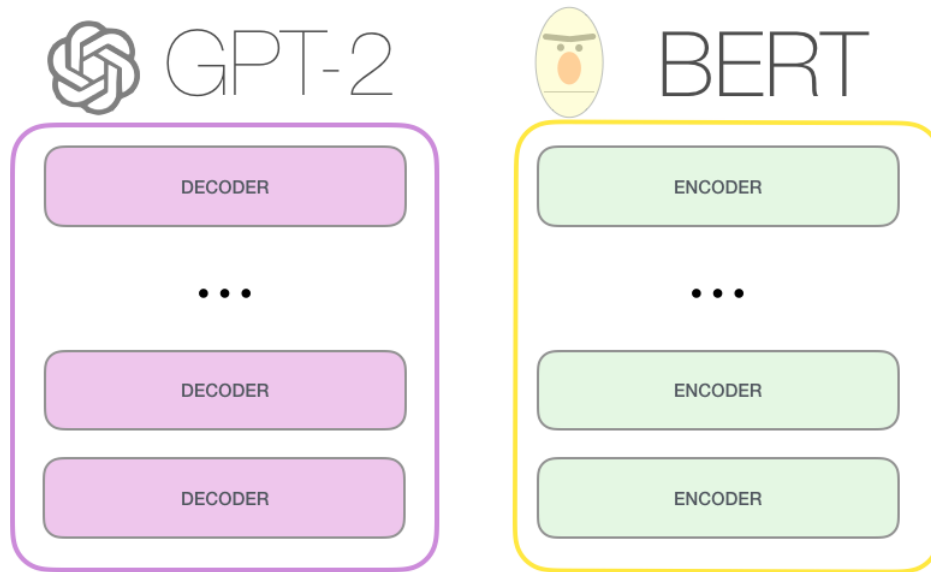


Fig. 2.2: Bert's encoder-only architecture and GPT-2's decoder-only architecture. | Source: <http://jalammr.github.io/illustrated-gpt2/>

Two of the most renowned transformer-based models are BERT [Dev+19], which is an encoder only transformer and the GPT (Generative Pretrained Transformer) series, which includes GPT-3, ChatGPT (or GPT 3.5 turbo), and GPT-4, which are decoder only models. GPT-3 and GPT-3.5 have a whopping 175 billion parameters, making them two of the most intricate and extensive models in existence [Bro+20]. The total number of parameters for GPT-4 is currently unknown. These models have achieved state-of-the-art performance in a variety of natural language processing tasks, such as text classification, named entity recognition, and sentiment analysis. However, like other large language models, transformers can be computationally expensive to train and deploy, and their decision-making processes can be hard to interpret.

2.2.3 Large Language Models (LLMs)

LLMs are powerful deep learning models pre-trained on vast amounts of data, capable of understanding and generating plausible language [Rad+18]. They are typically trained to predict the next word in a sentence, a process known as self-supervised learning. These models use a set of transformers, which consist of multiple decoder blocks with self-attention capabilities that extract meanings from a sequence of text and understand the relationships between words and phrases within it [Vas+17]. They can capture much of the syntax and semantics of human language because of their extensive training on large datasets and a large number of

parameters, which are the weights that the model learns during training. Currently, these models were also fine-tuned to follow instructions [Ouy+22], making them extremely capable of performing downstream tasks without the need for many task-specific data. GPT-4 is one of the most sophisticated models of this and is known for its impressive performance, sometimes even outperforming humans [Ope23]. The success of this instruction follow-up capability can be attributed to the use of reinforcement learning from human feedback (RLHF) [Ouy+22], where the model is trained with human-provided feedback to improve its training and its ability to follow instructions.

LLMs like GPT-4 have shown great potential in various NLP tasks when given an appropriate prompt. However, these models are not open-source to the public and require a cost per call to operate. This led researchers and companies to develop their own open-source models, like Llama [Tou+23a] and Llama-2 [Tou+23b] with the same capabilities. The Llama models have demonstrated comparable yet slightly inferior results in downstream tasks when given an appropriate prompt and few or no task-specific data. Despite this impressive performance, LLMs also have some limitations when the task requires extracting information from a prompt with large context [Liu+23a]. Another main challenge is the computational cost of training and deploying the model, which can be prohibitively high for some applications. Furthermore, these models are large and complex, making it difficult to interpret the underlying mechanisms of their decision-making process.

2.2.4 k -Nearest Neighbors

k -Nearest Neighbors (k -NN) [CH67] is a widely used machine learning algorithm that belongs to the family of instance-based or lazy learning algorithms. This technique is based on the concept that similar items tend to have similar characteristics. In k -NN, each item in the dataset is represented as a point in a high-dimensional space. To classify a new item, the algorithm finds the k nearest neighbors of that item in the feature space and assigns the class of the majority of those neighbors to the new item. The value of k is a parameter that needs to be specified and controls the number of neighbors considered. Figure 2.3 demonstrates how the k -NN algorithm works in a two-dimensional space. One advantage of k -NN is that it does not require assumptions about the underlying data distribution. However, k -NN can be computationally expensive at inference time, especially when dealing with large datasets. To optimize k -NN, data structures such as KD trees [AI08], locality-sensitive hashing [Ary+98], and parallelization [JDJ17] can be used to approximate nearest neighbors. In general, k -NN is a straightforward and effective

algorithm that can be used in a variety of machine learning tasks, particularly when the dataset is small or when interpretability is important.

Weighted k -Nearest Neighbors [Dud76] is a variation of the algorithm that assigns weights to the nearest neighbors based on their distance from the query point. This approach has been shown to be successful in improving the accuracy of classification problems. In a study exploring the use of weighted k -NN in ordinal classification [HS04], the authors found that it outperformed traditional k -NN in terms of accuracy. These findings suggest that weighted k -NN can be a powerful tool to improve the accuracy of classification problems and that selecting an appropriate distance metric is essential to achieve the best results.

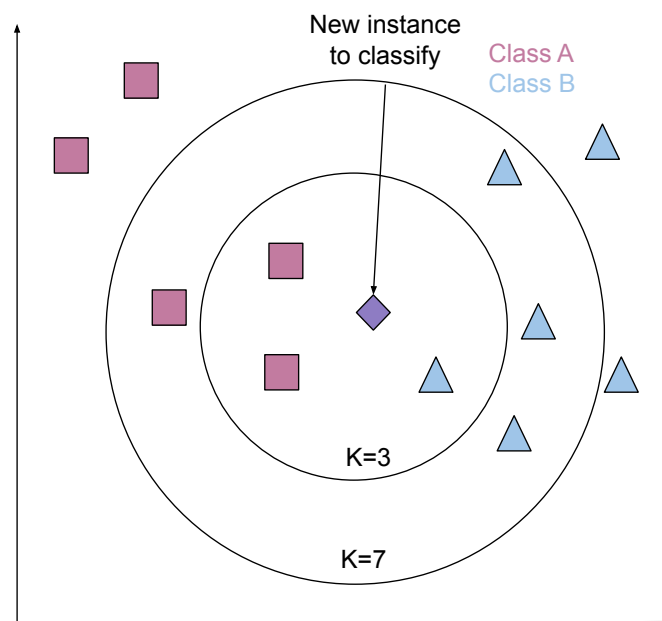


Fig. 2.3: The value of k affects how the k -NN algorithm works. In this example, when k is set to 3, the red point is classified as Class B, but when k is increased to 7, the same point is classified as Class A.

2.3 Related Work

As deep learning progresses, larger models with greater knowledge capacity are being developed. However, these models have millions or billions of parameters and require a great deal of computing power for inference, making their deployment difficult. Consequently, researchers have investigated various methods to reduce the computational resources needed for training and inference. Nevertheless, we want to train these models with quality data to maximize their performance. To do this, techniques such as knowledge distillation and active learning have emerged.

2.3.1 Knowledge Distillation

Knowledge distillation [HVD15b] is a technique that enables the capture and transfer of knowledge from a complex machine learning model or an ensemble of models into a single, smaller model. This process does not lead to a significant decrease in performance, and makes the model much easier to deploy. Essentially, it is a form of model compression that was first successfully introduced by Bucilua et al. [BCN06]. Knowledge distillation is performed on neural network models associated with complex architectures that include several layers and model parameters. Consequently, with the emergence of deep learning in the last 10 years and its success in various areas such as speech recognition, image recognition, and natural language processing, knowledge distillation techniques have become increasingly important for practical applications in the real world [Gou+21].

A knowledge distillation system consists of three main components: knowledge, the distillation algorithm, and the teacher-student architecture [Gou+21]. Figure 2.4 shows the knowledge distillation process, where a small *student* model learns to mimic a large *teacher* model and leverage the teacher's knowledge to obtain similar or higher precision.

Knowledge typically refers to the parameters learned in the neural network. At the same time, there is a rich diversity in the sources of knowledge in a large deep neural network. The different forms of knowledge are categorized into three different types: *response-based* knowledge, *feature-based* knowledge, and *relationship-based* knowledge. Response-based knowledge is concentrated on the last output layer of the teacher model. The idea is that the student model will learn to imitate the teacher model's predictions. This can be accomplished by utilizing a loss function, known as the distillation loss, which captures the discrepancy between the logits of the student model and the teacher model. As this loss is minimized through training, the student model will become more proficient at making the same predictions as

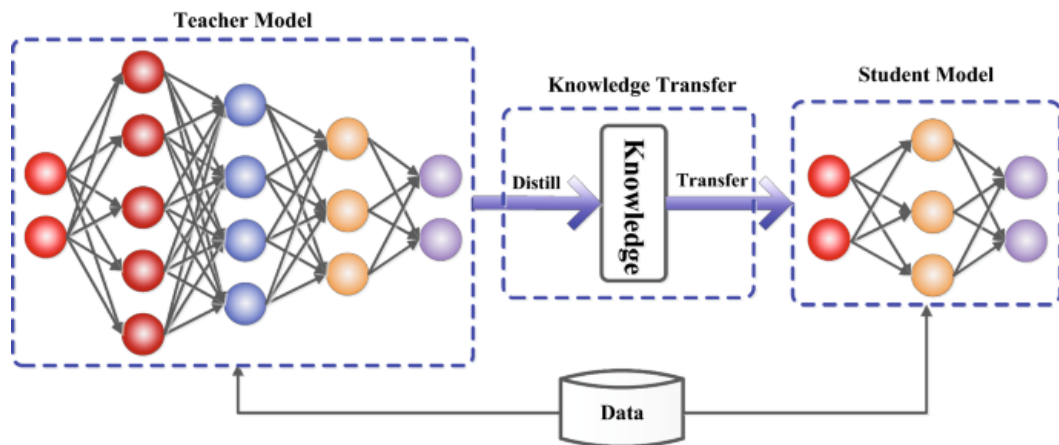


Fig. 2.4: The teacher-student framework for knowledge distillation

the teacher [Gou+21]. The trained teacher model captures the knowledge of the data in its intermediate layers, which is especially relevant for deep neural networks. Intermediate layers learn to differentiate particular features, and this information can be used to train a student model [Gou+21]. For feature-based distillation, the aim is to train the student model to acquire the same feature activations as the teacher model. The distillation loss function accomplishes this by reducing the gap between the feature activations of the teacher and student models. Both of the aforementioned methods use the outputs of specific layers in the teacher model. The focus of relation-based knowledge is to identify the links between different layers or data sets within a model. These connections can be established in a variety of ways, including between two layers, three layers, two samples, three samples, and more. In a text relevance task, the student model can be taught to recognize the connection between different sample pairs produced by the teacher model. If the teacher model establishes a certain degree of relevance between two text samples, the student model is trained to imitate that same degree of relevance [Yan+23].

The three main approaches to training the models of the student and teacher are offline, online, and self-distillation. The way in which the distillation training methods are classified is based on whether the teacher model is altered in tandem with the student model or not, as illustrated in Figure 2.5.

Offline distillation is the most common approach to knowledge transfer in deep learning, where a pre-trained teacher model is used to guide the student model. This technique requires that the teacher model be trained on a dataset prior to distilling its knowledge to the student model. With recent advances in deep learning, a wide variety of pre-trained neural network models are available to serve as the teacher depending on the use case. Offline distillation is a well-known method and is relatively easy to implement. The teacher model is usually a large-capacity deep neural network. However, in some cases, a pre-trained model may not be

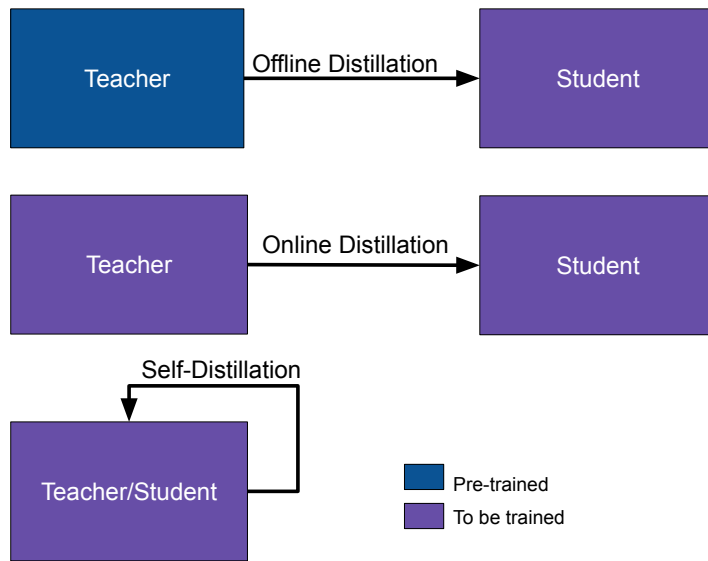


Fig. 2.5: Different knowledge distillation training schemes

available for offline distillation. To address this limitation, online distillation can be used, where both the teacher and student models are updated simultaneously in a single end-to-end training process. Online distillation can be implemented using parallel computing, making it a highly efficient method. Self-distillation is a type of online distillation in which the same model is used for both the teacher and the student. This process involves splitting the network into multiple parts and transferring knowledge from deeper sections to shallower ones, as opposed to traditional knowledge distillation, which requires the student model to approximate the softmax layer output of a pre-trained teacher model [Zha+19].

2.3.2 Self-Labeling

Self-labeling is a semi-supervised learning technique in which a model is initially trained on a limited number of labeled data points. This method involves an iterative process in which a Teacher model, trained on the labeled instances, is used to generate silver labels for the unlabeled dataset [Yar95; BM98]. A subset of the silver-labeled instances, along with the pseudo-labels, is then added to the labeled dataset. The Student, usually a model with the same architecture as the Teacher, is then trained on the augmented labeled dataset. This process can be repeated for a fixed number of iterations or until the evaluation metric starts to drop or it reaches convergence. This approach is referred to as Self-training when the Student model has a capacity similar to or greater than that of the Teacher, and knowledge distillation [HVD15a] when the Student model is smaller than the Teacher. Self-labeling

utilizes unlabeled data in a task-specific way during the pseudo-labeling process to train a more robust Student model and has been applied to a variety of NLP, computer vision, and speech recognition tasks. Examples of such tasks include sentiment analysis [Du+21], intent-classification and dialog-state tracking [Mi+21a], evidence-extraction [Niu+20], natural language inference [Li+21b], toxic span detection [SJ21], text classification, and rationale extraction [BSM21].

2.3.3 Active Learning

The goal of machine learning is to create systems that can improve their performance at some task as they acquire experience or data. In many natural learning tasks, this experience or data is gained interactively, by taking actions, making queries, or doing experiments. However, most machine learning research treats the learner as a passive recipient of the data to be processed. This approach disregards the fact that, in a lot of circumstances, the student's most potent resource is their capacity to take action, to acquire information, and to affect the environment they are attempting to comprehend. Active learning is the study of how to use this ability effectively. Formally, active learning represents a machine learning approach that aims to efficiently gather annotated data from a (typically extensive) set of unlabeled data [Set12], [CGJ96], [LC94]. As shown in Figure 2.6, active learning systems try to overcome the bottleneck of labeling by asking for queries in the form of unlabeled examples to be labeled by an expert (e.g., a human annotator). The goal of the active learner is to attain a high level of accuracy with the least amount of labeled data, thus reducing the cost of obtaining labeled information [Set12]. Algorithm 1 shows the typical procedure of an active learning paradigm.

In certain situations, active learners can pose questions, and a variety of strategies have been used to decide which cases are the most beneficial. The most popular methods to obtain active learning data are based on uncertainty and diversity, which are often referred to as the "two faces" of active learning [Das11].

Algorithm 1 Typical active learning procedure

- 1: **Input:** An unlabelled data pool \mathcal{U}
 - 2: **Output:** The final dataset \mathcal{D} and the model \mathcal{M}
 - 3: **Require:** Sampling algorithm \mathcal{A}
 - 4: **while** stop criteria not met **do**
 - 5: $\mathcal{I} \leftarrow$ Apply \mathcal{A} on \mathcal{M} and \mathcal{U}
 - 6: $\mathcal{I}' \leftarrow$ anotate(\mathcal{I})
 - 7: $\mathcal{U} \leftarrow \mathcal{U} - \mathcal{I}$; $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{I}'$
 - 8: $\mathcal{M} \leftarrow$ train(\mathcal{M}, \mathcal{D})
 - 9: **return** \mathcal{M}, \mathcal{D}
-

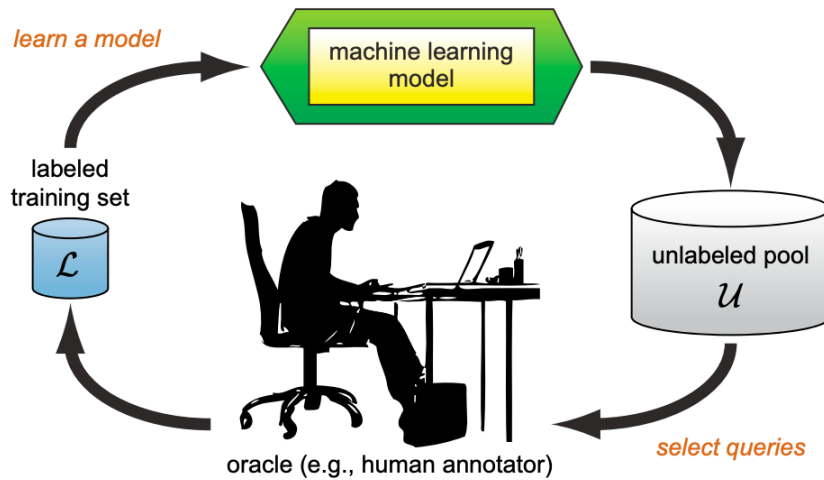


Fig. 2.6: The iterative process of active learning using a pool of data. | Source: [Set12]

2.3.4 Uncertainty Sampling

Uncertainty sampling is a widely used criterion for selecting examples in active learning and perhaps the simplest one [LG94]. Uncertainty sampling is a technique in which the learner inquires about the instances it is most uncertain about in its predictions and then presents them to the oracle for annotation [LC94]. Over the years, various forms of uncertainty have been employed; however, selecting a technique to quantify uncertainty is straightforward when using probabilistic measures. An example of uncertainty sampling can be seen in classification tasks, where a probabilistic classifier is used. In this case, instances are chosen where the probability of the predicted class does not exceed a predetermined threshold. This threshold is used to determine which instances are uncertain and require further labeling [LC94]. A more general uncertainty sampling strategy (and possibly the most popular) uses entropy [Sha48] as an uncertainty measure:

$$\mathcal{H}_x = - \sum_{i \in N} P(y_i|x) \log P(y_i|x).$$

where y_i spans across all possible labels of x , where x is a discrete random variable, which takes values in the range $[0,1]$. Entropy is a measure of information theory that quantifies the amount of data needed to be “encode” a distribution. It is often seen as a measure of uncertainty or impurity in machine learning. For binary classification, entropy-based sampling is equivalent to the margin and least confident strategies, which query the instance with a class posterior closest to 0, 5. Entropy

can be extended to probabilistic multilabel classifiers and probabilistic models for more complex structured instances, such as sequences and trees [SC08; Hwa04]. It is possible to use uncertainty sampling strategies with nonprobabilistic classifiers too. Lewis [LC94] was one of the first to investigate this using a decision tree classifier. Fujii et al. [Fuj+98] and Lindenbaum [LMR04] have applied similar approaches to active learning with nearest neighbor (also known as “memory-based” or “instance-based”) classifiers, where each neighbor casts a vote on the class label of x , the proportion of these votes representing the posterior label probability. Tong [TK02] also experimented with an uncertainty sampling strategy for support vector machines (SVMs) that involves querying the instance closest to the decision boundary.

2.3.5 Diversity Sampling

In diversity sampling, the learner aims to identify data points that represent the diversity present in the unannotated dataset. This technique can be used to discover the "unrecognized unknowns" of a model, which can reveal what the model is unaware of or does not comprehend. This is crucial because machine learning models often need to operate in a constantly changing environment, so it is important to identify and learn from diverse data. The diversity sampling approach can be implemented in various ways. Some renowned methods are model-based outliers, cluster-based sampling, representative sampling, and sampling for real-world diversity.

Model-based outliers are data points that do not conform to the patterns seen in most data, as determined by a statistical model. These outliers can be identified by fitting a model to the data and then measuring the distance of each data point from this model. For example, in a linear regression model, outliers would be points that are far away from the best fit line. A commonly used measure to detect these outliers is the standardized residual [ES10], which measures the number of standard deviations a data point is away from the predicted value. If the absolute value of the standardized residual is greater than a certain threshold (often 2 or 3), the data point is considered an outlier [Mon21].

Cluster-based sampling is an effective way to collect training data. Instead of randomly selecting samples, the data is divided into multiple clusters and an equal number of samples are taken from each cluster. The reason for this is simple. For example, if we have a dataset with many articles on different topics, random sampling for human review could lead to a lot of time being spent manually labeling similar articles in a particular category. However, by pre-clustering the data, these similar articles are likely to be grouped together, allowing us to label only a few examples from that cluster. This not only saves time but also helps to allocate resources to label data from other clusters. These other clusters may contain rare

types of articles that are important but too infrequent to be reliably captured by random sampling. Therefore, clustering not only makes the labeling process more efficient but also increases the diversity of the labeled dataset [Mon21].

2.3.6 Representativeness Sampling

Representative sampling is a sampling technique in which the sample drawn from the population represents the characteristics of the population as a whole. This method is used to make generalizations from the sample to the population. The goal of representative sampling is to obtain a sample that is as similar as possible to the population in terms of key characteristics. The representativeness of an example can be determined by the number of unlabeled examples that are similar to it [AA07]. Examples that are highly representative are less likely to be outliers, and their inclusion in the training set will have an impact on a great number of unlabeled examples. According to [AA07], if the examples were grouped together based on similarity, the centroids of the clusters would be the most representative examples. During their experiments, [She+04] employed cosine similarity and Dynamic Time Warping to measure the similarity between the examples.

System Design

3.1 Introduction

In this chapter, we give a comprehensive description of the theoretical design of OCaTS. We will also provide an outline of each element of the system and an explanation of how they interact with each other. Additionally, we will examine various decision-making criteria and introduce a new evaluation metric to assess the effectiveness of the system.

3.2 System Design

OCaTS uses the teacher-student architecture initially proposed in knowledge distillation (Sec 2.3.1). As illustrated in Figure 3.1, when a customer submits a new instance, the student model is asked to provide an answer and its level of certainty. For example, in a classification task, the student's confidence can be expressed as the entropy of its probabilities. The system then evaluates the student's certainty and verifies if the student was provided with representative data to handle the new instance. To do this, it measures the similarity of the new instance with a sample of previously trained data stored in the cache. If the criteria are met, the system will give the student response to the customer. If not, the more costly, yet more powerful teacher model will be asked to give its answer. The customer will be given the teacher's response, and the new instance with the teacher response will be kept in the cache for future training. Next, we will provide an overview of its component separately.

3.2.1 Cache

The cache is a storage system that stores the original training data (if provided) and new incoming customer instances that the student could not provide an answer for. If training data are provided, the cache is considered to have a *warm* start. It can take the form of a database that contains the data, a vector database that holds representations of the data, or even a text file that holds the records. However, if no training data are given, the cache is empty and filled with new queries; this

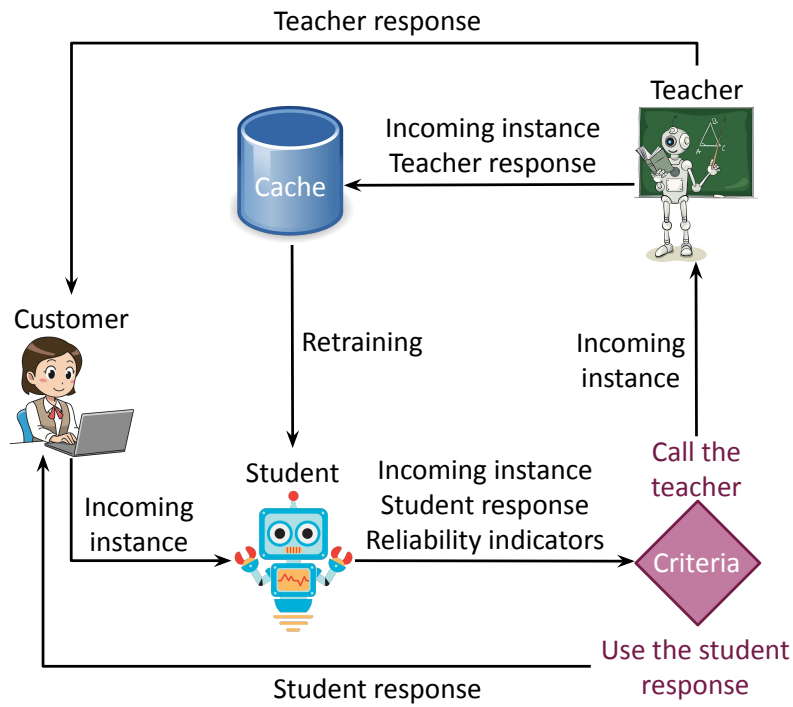


Fig. 3.1: OCATS architecture

is considered a *cold* start. When it is time for the framework to determine the representativeness of the training data for the new query, the cache is asked to provide the k most similar cached instances. This can be achieved by performing a similarity search on instances if the cache is a vector database or a traditional information retrieval search (e.g., BM25 [RZ09]) if the cache is a collection of documents.

3.2.2 Teacher

The teacher model is a large-scale model that has been trained with a great volume of task-specific data, or an LLM that can achieve great results in downstream tasks when given limited task-specific data as context and prompted correctly. The teacher is asked to provide its answer when the student cannot do so or when the data the teacher has been trained on are not very representative. The teacher's answers are always considered to be correct and are stored to further train the student. Due to its size, hosting the teacher can be costly, and if it is an LLM provided by an API, a payment per call is necessary. We would like to limit the number of calls to the teacher based on the financial capacity of the framework's user.

3.2.3 Student

The student model is a task-specific model on the SME side that has been trained with limited data, resulting in lower performance and capacity than the teacher model. However, it is lightweight and can be deployed and maintained at minimal cost. Whenever a new user instance is received, the student is asked to provide an answer along with its confidence. If the confidence of the model is sufficiently high, the response is considered to be of satisfactory quality. The purpose of OCATS is to maximize the use of the student model and only invoke the teacher when absolutely necessary, thereby enhancing the student's progress over time while limiting monetary expenses.

3.2.4 Criteria

In order to achieve the most beneficial balance between involving the teacher or depending on the student while keeping costs to a minimum, we must identify appropriate criteria. To identify these criteria, we take our cues from active learning. We evaluate the certainty of the student model and the representativeness of the training data in comparison to the new user instance. In a classification task, we can measure the confidence of the classifier by setting a maximum threshold for its entropy. This implies that even if the classifier is able to accurately predict the correct class but is uncertain, the larger model will carry out the classification. We can assess the representativeness of the training and previously answered data by calculating the mean distance of the new query from the k most similar cached instances and limiting the demand to less than a certain value. Therefore, these criteria can be seen as thresholds for various measures that can be both task-specific and data-specific. A more detailed discussion of the criteria used in the experiments will be provided in Section 4.3.2

3.2.5 Discounted evaluation metric

The purpose of the proposed architecture is to reduce the number of times the costly teacher model is called by storing past teacher responses and using them to train a cheaper student model on the SME side. This creates a trade-off between the operational cost of calling the teacher and the dissatisfaction of the end-users when the less precise student model is used instead. To measure how well the two factors are balanced, we introduce a *discounted* version $\hat{\phi}$ of any conventional evaluation measure ϕ (e.g., accuracy, F1), as follows:

$$\hat{\phi} = \phi - \lambda \cdot \frac{M}{N} = \phi - \lambda \cdot \rho, \quad (3.1)$$

where N is the number of incoming instances that have been processed (on which ϕ is measured), M is the number of calls made to the teacher while processing the N instances, $\rho = \frac{M}{N}$ shows for what percentage of incoming instances we call the teacher, and λ is a scalar specifying how intensively the conventional measure should be discounted.

Assume, for example, that the accuracy of the teacher-student combination is $\phi = 0.8$, but that this accuracy is achieved with $\rho = \frac{1}{3}$. If OCaTS users view the value ρ (which could be translated into a monthly cost) as costly as a one percent decrease in accuracy, then $\hat{\phi} = 0.7$, and Eq. 3.1 becomes $0.7 = 0.8 - \lambda \cdot \frac{1}{3}$, resulting in $\lambda = 0.3$. If OCaTS users consider the same ρ value to be more or less costly in terms of accuracy points, then larger or smaller λ values will be obtained. Eq. 3.1 can also be expressed as $\delta = \lambda \cdot \rho$, where $\delta = \phi - \hat{\phi}$ indicates how much ϕ is discounted to account for the cost of ρ . In this case, λ can be thought of as a currency exchange rate, showing how expensive ρ is in terms of δ (e.g., loss of accuracy in percentage points). In this thesis, we assume that the exchange rate λ is constant for all values of δ and ρ . In practice, it may be different for different ranges of δ and ρ , but this is left for future work.

Experimental Setup

4.1 Introduction

In this chapter, we explore different versions of the framework presented in the preceding chapter. We analyze the datasets and the various scenarios employed and provide a summary of the experiment conducted. Lastly, we assess our framework in terms of our objectives and review related research in the field.

4.2 Datasets

4.2.1 BANKING77 Dataset

The BANKING77 dataset [Cas+20] is an intent recognition dataset from the banking customer service domain. While there are various intent recognition datasets, they fail to capture the examples for each domain, and may not sufficiently capture the full complexity of each domain as encountered in real world applications. The single-domain focus of BANKING77 with a large number of intents makes it more challenging. Certain intent categories have some degree of overlap with each other, necessitating careful consideration when making decisions. BANKING77 includes 13,083 customer messages. The ground truth assigns to each message a single label (intent) of the 77 available. The dataset is divided into training (10,003 instances) and test (3,080) subsets. Figure 4.1 shows the distribution of labels for each subset.

Message	Intent
Can I track my card while it is in the process of delivery?	card_arrival
Are cards available in the EU?	country_support
What is the identity verification process?	verify_my_identity
I need to renew a card that's about to expire	card_about_to_expire

Tab. 4.1: Examples of the Banking77 dataset.

4.2.2 Large Movie Review Dataset

The Large Movie Review (LMR) dataset consists of 50,000 reviews from IMDB, with no more than 30 reviews per movie. The authors created the dataset to contain an equal number of positive and negative reviews, so that randomly guessing would yield 50% accuracy. The reviews are highly polarized, with negative reviews having a score of 4 or less out of 10, and positive reviews having a score of 7 or more out of 10. Neutral reviews are not included. The dataset is divided into 25,000 reviews for training and 25,000 for testing, with additional unlabeled data available. Each review is labeled as either negative (0) or positive (1), as shown in Figure 4.2.

4.2.3 Few-shot scenario

Assuming that a small business can only manage to create a limited number of training examples, for the BANKING77 dataset we use only 231 instances from the original training set, three per class, as a few-shot version of the training set. These examples and their corresponding labels are included in each prompt call to the teacher as demonstrations (in-context learning). We handpicked 231 examples taking care to avoid any confusion, such as instances that had different labels. Additionally, we generated a development set of 1,001 instances from the original training set to adjust the hyperparameters (see below). We employed a different approach to create the few-shot scenario for the LMR dataset, given the large number of examples. We randomly selected 5 instances per sentiment based on the average length (in tokens) of the reviews. We used the same method to generate the development and test sets of 1000 and 5000 instances, respectively.

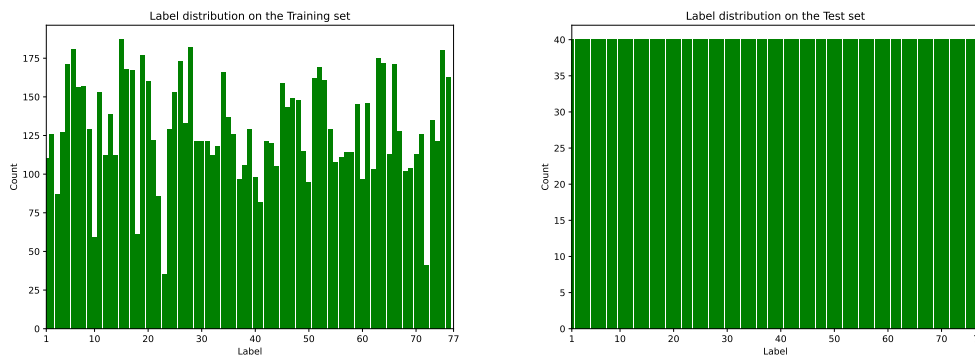


Fig. 4.1: Data distribution of the *original* BANKING77 dataset. Left figure shows the distribution of the labels in the training set while the right figure show the distribution in the test set.

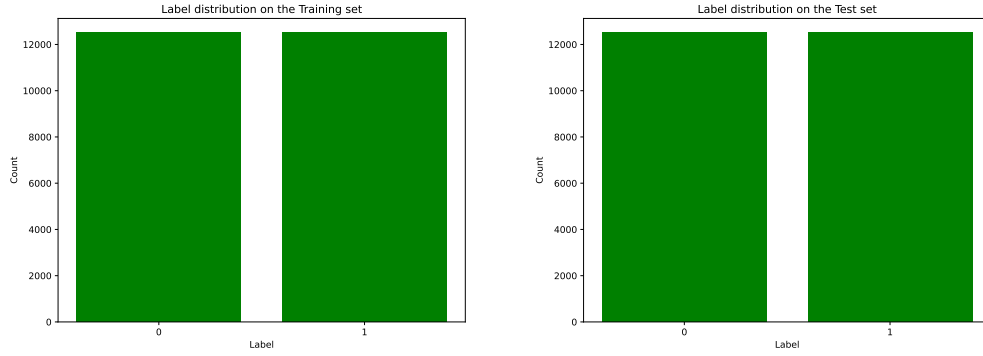


Fig. 4.2: Data distribution of the *original* LMR dataset. Left figure shows the distribution of the labels in the training set while the right figure show the distribution in the test set.

4.3 Experiments and System Evaluation

4.3.1 Online setting and evaluation measure

We use the BANKING77 and LMR test sets as incoming instances in an online setting and repeat each experiment with five random shuffles of the test set. Average scores of the shufflings are reported. Since both test sets are balanced (Figures 4.1 & 4.2), we set ϕ (Eq 3.1) to accuracy.

4.3.2 Criteria

For both tasks, we require two conditions to be met: one that captures the certainty of the classifier and one that captures the representativeness of the training (or previously streamed) data, in order for the student’s response to be accepted. If either of these conditions is not satisfied, the teacher will be invoked. The first condition is that the cosine distance between the vector representation of the incoming message and the weighted centroid vector \mathbf{c} of the k nearest neighbors must be less than a threshold t_c . These vector representations are generated with the Sentence Transformer variation of MPNet [RG19; Son+20].¹ The weighted centroid vector is given by the following equation:

$$\mathbf{c} = \sum_i \hat{w}_i \cdot \mathbf{v}_i$$

where

$$\hat{w}_i = \frac{w_i}{\sum_{j=1}^k w_j}$$

¹We used `all-mpnet-base-v2` as the sentence transformer.

and w_i is the weight assigned by distance weighting to the neighbor i -th, and \mathbf{v}_i is the vector representation (MPNet-based) of the neighbor. The weight of the neighbor i , denoted by w_i , is inversely proportional to the square of the cosine distance between the vectors of the incoming instance and the neighbors. This can be expressed as

$$w_i = \frac{1}{d_i^2}$$

This condition ensures that the incoming instance is close to the stored instances in an intuitive manner.

To define the second condition, let C be the set of all potential labels (classes). Let p_c be the probability of each $c \in C$ assigned by a probabilistic classifier $f(x; \theta_0)$ where x denotes the input of the classifier and θ_0 the classifier parameters. The entropy \mathcal{H} of the probabilities p_c of the labels are calculated as:

$$\mathcal{H} = - \sum_{c \in C} p_c \log p_c$$

The second criterion requires that the entropy \mathcal{H}_w is below a certain threshold $t_{\mathcal{H}}$. This implies that the classifier must be certain of its prediction.

4.3.3 Experiment on BANKING77 with a k -NN student and GPT-4 teacher

We chose the BANKING77 domain for our initial experiment because of its complexity and its ability to simulate a real world situation. Hereafter, we will provide more information about the other components of OCaTS.

Teacher: We employed GPT-4 [Ope23]² as the teacher for our few-shot learning tasks, which was the most advanced language model available at the time. The prompt we used included instructions, examples, and the instance to be classified. We provide GPT-4 with the following prompt to carry out the classification:

You are an expert assistant in the field of customer service. Your task is to help workers in the customer service department of a company. Your task is to classify the customer's question in order to help the customer service worker to answer the question. In order to help the worker you MUST respond with the number and the name of one of the following classes you know. In case you reply with something else, you will be penalized.

²We used gpt-4-0314

The classes are:

- activate_my_card
- age_limit
- ...

User: My new card is here, what's the process for activating it?

Assistant: activate_my_card

User: I am unable to activate my card, it won't let me.

Assistant: activate_my_card

User: Can you help me activate my card

Assistant: activate_my_card

...

User: What is the youngest age for an account?

Assistant: age_limit

User: What is the appropriate age for my child to be able to open an account?

Assistant: age_limit

User: How do I set up an account for my children?

Assistant: age_limit

...

Student: A distance-weighted k -NN classifier was employed as the student. Vector representations of the incoming instances were generated with a Sentence-Transformer variation of MPNet. The distance weighting of the k -NN classifier was performed by calculating the cosine distance between the incoming instance and the nearest neighbor i , and assigning the weight of the neighbor. The weights of the neighbors in each class were then summed and passed through a softmax function to create a probability distribution for each class. The incoming instance was classified into the class with the highest probability. No retraining was necessary for the k -NN student, since the cache was the same as the memory of the k -NN classifier. The cache (memory) was initialized with the 3-shot training examples of the classes (231 instances in total), as illustrated in Fig. 3.1.

Hyperparameter tuning: There are three hyperparameters: the number of neighbors k , and the thresholds t_c and t_h . We set k to 5, as it is a reasonable number given that there are 3 examples per class. For each of the four indicative λ values (0.05, 0.1, 0.2, 0.3), we use Bayesian optimization [SLA12] in the development set we created (Section 4.2.3) to find the optimal combination of the two thresholds that maximizes $\hat{\phi}$. We start by conducting a 10×10 grid search on the range of values of the two thresholds, and then use Optuna's [Aki+19] implementation of the Tree-Structured Parzen Estimator (TSPE) algorithm [Ber+11]. The contour maps and optimal values of the two thresholds for each λ value can be seen in

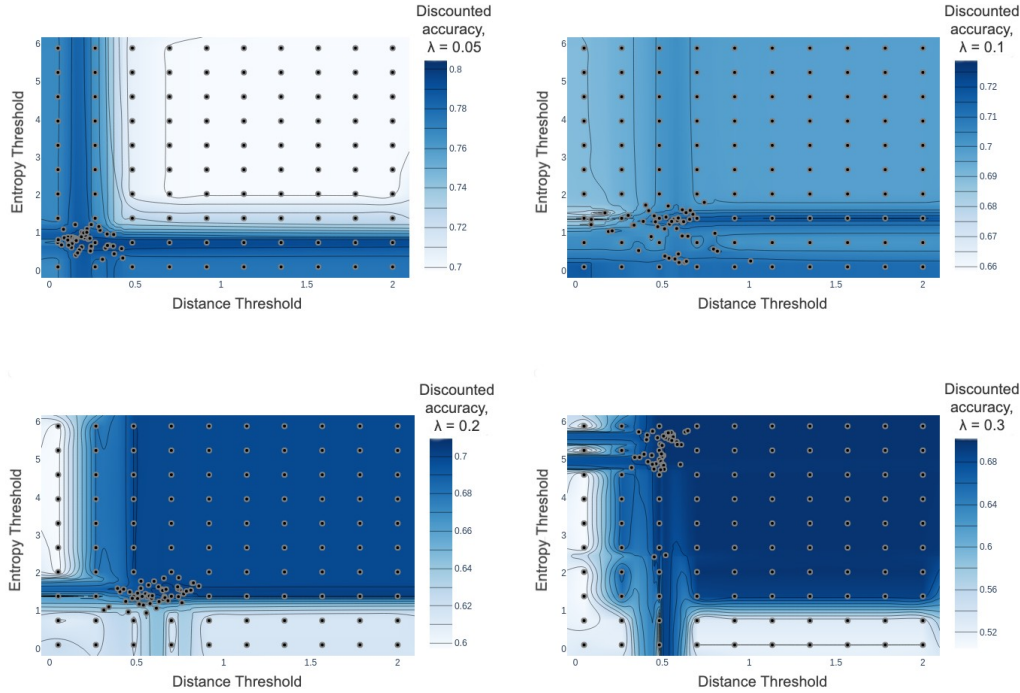


Fig. 4.3: Contour plots (discounted accuracy) obtained during threshold tuning, for various λ values.

Figure 4.3. This figure shows that there are multiple points close to the best point that maximize the discounted metric to a great extent. This implies that several threshold combinations can be considered optimal. Additionally, there are large regions with a wide range of values for a threshold that are almost equally effective in maximizing the discounted accuracy measure.

Results: We assess the performance of OCaTS for four different values of the discount factor λ , using the same incoming instances (the original BANKING77 test set) and the λ -specific tuned thresholds t_c , t_H . As seen in Fig. 4.4, OCaTS is successful in finding a balance between teacher calls and accuracy: on the left, we observe that as λ increases, fewer calls to the teacher are made; in the middle, we see the amount of accuracy sacrificed for this cost reduction. Specifically, for $\lambda = 0.05$, the accuracy of OCaTS is very close to the accuracy of the GPT-4 teacher, with a difference of 0.37 percentage points (83.05% compared to 82.68% for the entire test set), while the teacher is called for only one third of the incoming instances (1050 out of 3080). For larger values of λ , we observe the expected decrease in accuracy to achieve an even smaller number of calls to the teacher. Figure 4.4 (right) also shows that the discounted accuracy $\hat{\phi}$ of OCaTS (solid lines, one per λ value) is always higher than the discounted accuracy of always calling the GPT-4 teacher (dashed lines), as intended.

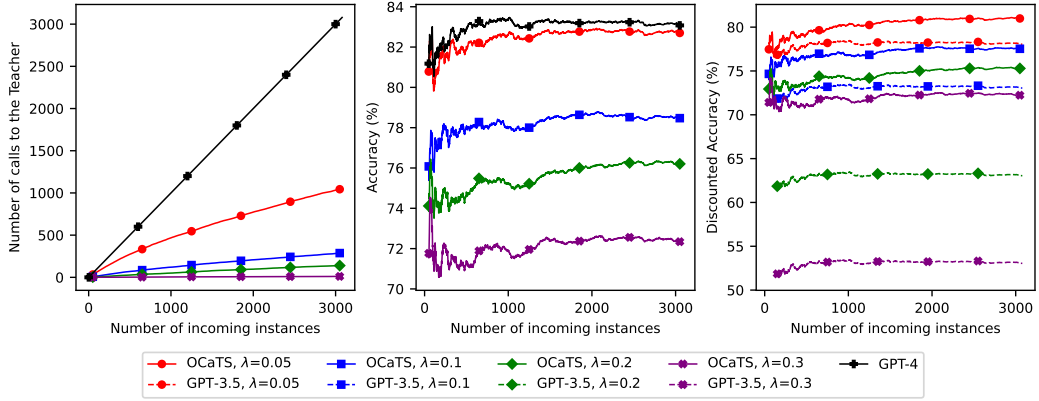


Fig. 4.4: The number of calls to the teacher (on the left), accuracy (in the middle), and discounted accuracy (on the right) are shown for various λ values when using a GPT-4 teacher and a k -NN student. The dashed lines indicate the discounted accuracy when GPT-4 is called for all incoming instances. The larger the λ , the more the SME prefers fewer calls, even though it may lead to increased user frustration. As expected, OCaTS has a better discounted accuracy than when the GPT-4 teacher is called for all incoming instances.

4.3.4 Experiment on BANKING77 with an MLP student and GPT-4 teacher

In order to further explore the performance-cost trade-off of OCaTS, we repeated the experiment as in Section 4.3.3, but changed the student.

Student: We considered two options for our model: a DistillBert classifier fine-tuned for text classification (referred to as DBert), and a custom Multi-Layer Perceptron (MLP) classifier built on top of frozen MPNet embeddings. To determine the best option, we evaluated them based on two factors. The first factor involved evaluating the overall performance of each model on the task. To gauge this, we examined the learning curves of both models, noting which showed faster and more substantial performance improvements. Before creating these curves, we identified the optimal architecture for each model through hyperparameter tuning. This tuning process was conducted on a few-shot dataset that we had generated (Section 4.2.3).

For DBert, we tuned parameters such as the learning rate, l2 regularization factor, and whether to freeze the body. In the case of MLP, we utilized one hidden layer, a tunable activation function, a dropout layer, and an output layer with 77 neurons (in accordance with the number of classes) and a softmax. The MPNet embeddings were incorporated without altering the model parameters, and we employed cross-entropy loss. We tuned the number of neurons for both the hidden and output layers, the activation function, and the dropout rate. The best hyperparameters for each model are presented in Table 4.2.

Once we determined the optimal architectures, we plotted the corresponding learning curves alongside those of the original weighted- k NN student used in a previous experiment (depicted in Figure 4.5). The learning curves indicated that the MLP outperformed DBert. Consequently, we chose to utilize the MLP as our student model. However, it is important to note that in this scenario, the chosen student requires retraining when new instances arrive. Therefore, we decided to retrain the MLP after every 100 calls to the teacher.

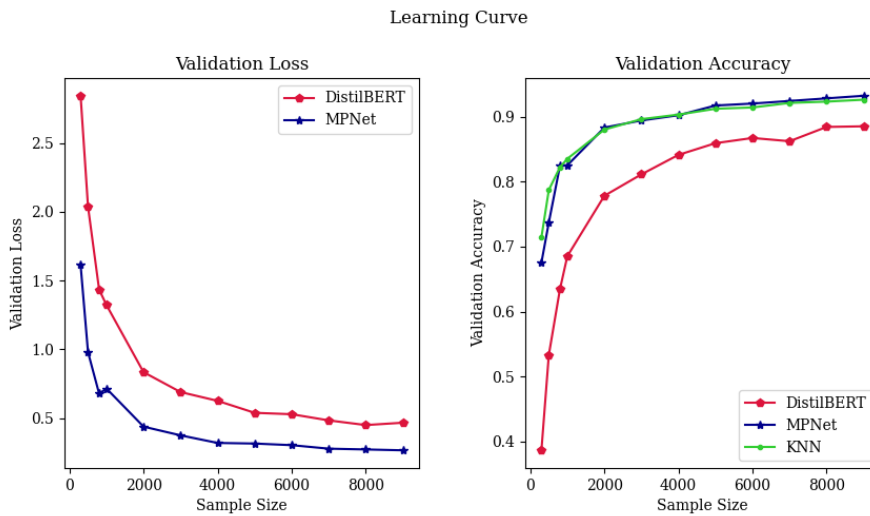


Fig. 4.5: The learning curves of the models demonstrate that MLP (blue line) is much more adept at assimilating the data and performs comparably to k -NN, while DBert’s performance is significantly weaker.

Parameters	DBert	MLP
Learning rate	$7 \cdot 10^{-5}$	$1.6 \cdot 10^{-5}$
L2 Norm	0.29	-
Hidden size	-	1024
Dropout	-	0.22
Activation	-	ReLU

Tab. 4.2: Parameters that were tuned for defining the architectures of DBert and MLP respectively.

Hyperparameter tuning: We repeat the same tuning process as before (Section 4.3.3) to pinpoint the best combination of two thresholds that maximizes $\hat{\phi}$ in the development set we constructed (Section 4.2.3). We begin by doing a 10×10 grid search across the range of values of the two thresholds, and then use Op-tuna’s Tree-Structured Parzen Estimator (TSPE) algorithm implementation [Aki+19; Ber+11]. The number of nearest neighbors is set to 5 while thresholds t_C and t_H are tuned. For each of the four indicative λ values (0.05, 0.1, 0.2, 0.3), we attempt to find the optimal combination of the two thresholds that maximizes $\hat{\phi}$. Figure 4.6 displays the contour maps of the tuning process.

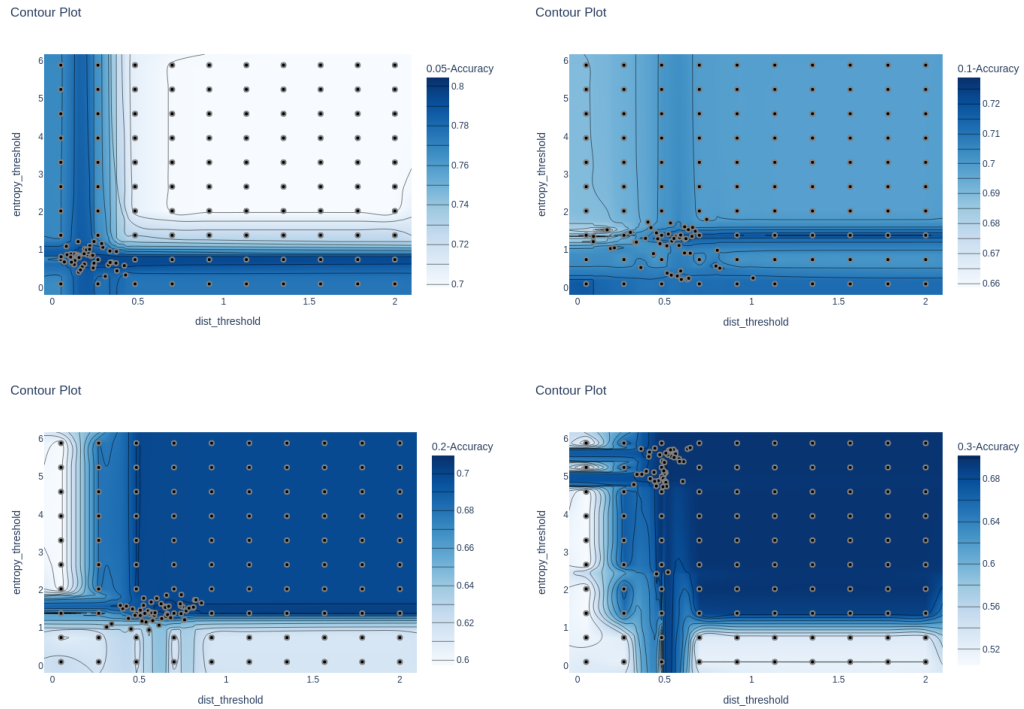


Fig. 4.6: Contour plots of discounted accuracy, obtained during threshold tuning, are presented for different λ values when using an MLP student and a GPT-4 teacher.

Results: We assessed the performance of OCaTS using the same method as in the preceding experiment in Section 4.3.3. As shown in Figure 4.7, OCaTS achieved, as in the previous experiment, a balance between trusting the student and prompting the teacher for an answer with only a slight decrease in accuracy. However, in this experiment, we noticed a greater performance decrease for $\lambda = 0.05$ compared to the results in Section 4.3.3, with a difference of 1.1 percentage points (83.05% compared to 81.95% for the entire test set) with a similar number of calls. Nevertheless, with the MLP student, OCaTS appears to perform better than with the k -NN when the cost of consulting the teacher is higher. This can be attributed to the fact that when the student is retrained, it adjusts better to the new data.

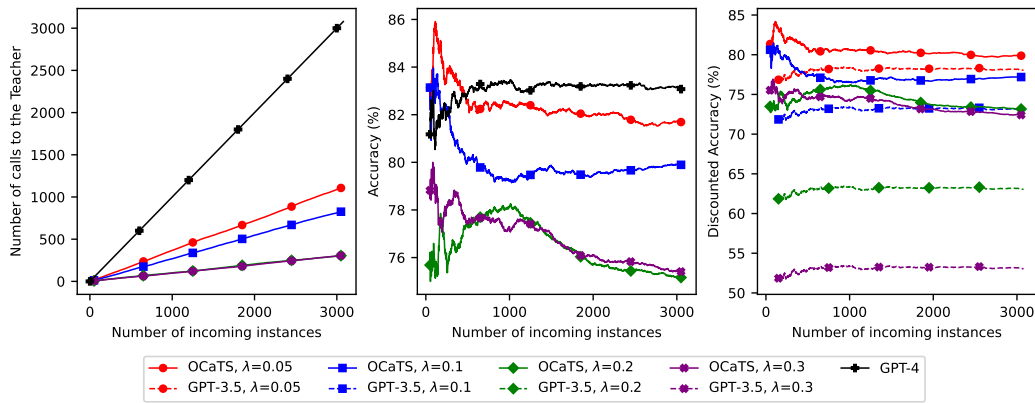


Fig. 4.7: The number of calls to the teacher (on the left), accuracy (in the middle), and discounted accuracy (on the right) are shown for various λ values when using a GPT-4 teacher and a MLP student.

4.3.5 Experiment on LMR with the k -NN student and GPT-3.5 teacher

We conducted a final experiment in which we employed a new dataset and instructor, while retaining the same k -NN student as in Section 4.3.3. This was done in response to the results of the first experiment and our desire to investigate less proficient students.

Teacher: We employed GPT-3.5-turbo [Ouy+22], one of the most advanced language models available, to act as a few-shot learning teacher. This model was used to classify a given instance, with instructions, examples, and the instance itself included in the prompt. The examples we used were the few-shot examples we sampled in Section 4.2.3. The following is the structure prompt used for the task:

Analyze the sentiment of the following reviews and classify them as either 'positive' or 'negative'

User: When I started watching this movie I saw the dude from Buffy, Xander, and figured ah how nice that he's still making a living acting in movies. Now a weird movie I can stand, given that it's a good dose of weird like for example David Lynch movies, twin peaks, lost highway etc. And you sort of have to be in the mood for one. This one however made me mockingly remember the crazy websites about there about conspiracy theory's that make absolutely no sense. I mean come on people Nazi's who conspire with America to make an unholy trinity of evil powers? I was surprised they didn't

mention the hollow earth in this movie with Hitler flying saucers and lizard people. Maybe if you had like 60 grams of heroine with this movie it would make some sort of sense, but seriously I don't condone drugs like I don't condone this movie. It should be burned, shredded and forgotten just so good ol' Xander might get another acting job. It wasn't his acting though, that was alright, but the script just didn't make any sense. Sorry.

Assistant: negative

User: As part of the celebration of the release of Casino Royale, this film with the new Bond starring in it was shown, from director Roger Michell (Notting Hill). I almost turned it off for being a bit boring, but I'm glad I stuck with it. Basically May (Anne Reid) is a single mother of Helen (Anna Wilson-Jones) who hardly sees anyone and has not had a boyfriend in years. Her daughter says that she might want to get married to her new boyfriend, Darren (Daniel Craig, of course). After knowing each other only a few days, May and Darren have a secret affair. And at her age, with a 30-something, and the new Bond?! Anyway, they obviously want to keep it a secret, but May has regrets and wonders if Helen will find out. When she does, Darren gets less hassle than May. In fact, Helen asks her permission to hit her.

Also starring Peter Vaughan as Toots, Danira Govich as Au Pair, Harry Michell as Harry, Rosie Michell as Rosie and Johnny English's Oliver Ford Davies as Bruce.

Very good!

Assistant: positive

...

Student: For the student, the same distance-weighted k -NN approach mentioned in Section 4.3.3 was used. No retraining was necessary for the k -NN student, as the cache was the same as the memory of the k -NN classifier. The cache (memory) was initialized with the 5-shot training examples of the classes (10 instances in total), as shown in Fig. 3.1.

Hyperparameter tuning: We followed the same tuning process as in prior experiments, but we did not do the initial grid search since the hyperparameter space was smaller. We used Bayesian optimization [SLA12] on the development set we created (Section 4.2) with Optuna's Tree-Structured Parzen Estimator (TSPE) algorithm

implementation [Aki+19; Ber+11]. The number of nearest neighbors was set to 5, while the thresholds t_c and t_H were tuned. We also set the maximum value for the entropy threshold t_H to 0.7 since there were only two possible classifications, positive or negative. Therefore, the maximum entropy, when using the natural logarithm, was 0.69. The contour maps of the tuning process are shown in Figure 4.8.

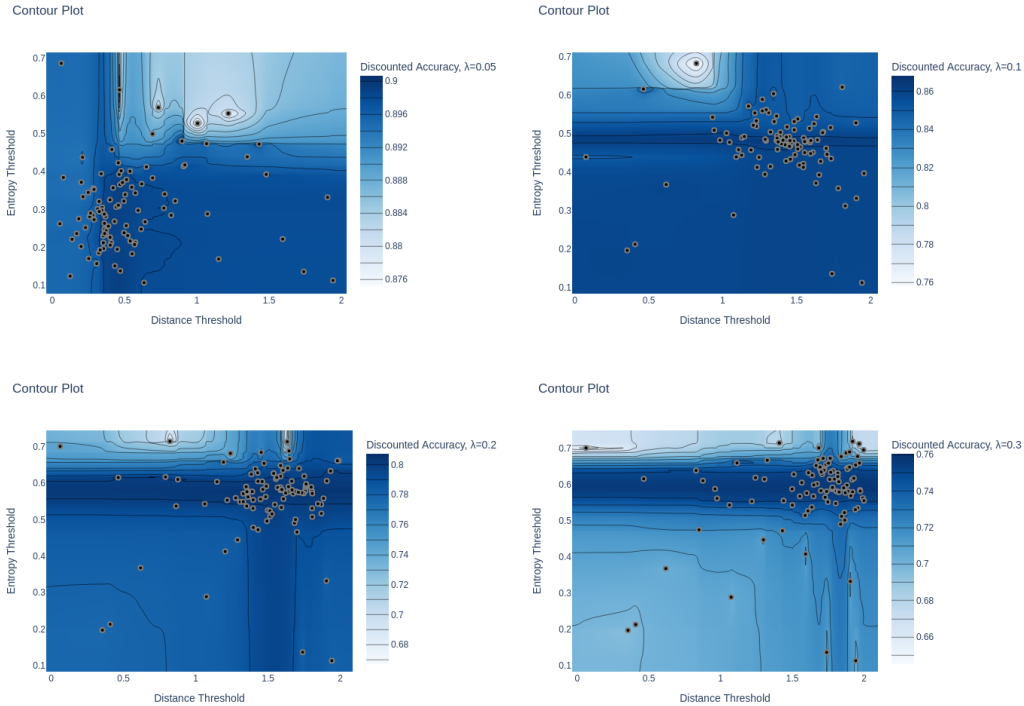


Fig. 4.8: Contour plots of discounted accuracy, obtained during threshold tuning, are presented for different λ values when using a k -NN student and a GPT3.5 teacher.

Results: To assess the performance of our system, we used the few-shot testing set outlined in Section 4.2.3. As shown in Figure 4.9, OCaTS manages to achieve a good equilibrium between utilizing the k -NN student and prompting the GPT-3.5 teacher to provide an answer with only a slight decrease in accuracy. We can see that in this experiment, OCaTS called the teacher more times than in the previous experiments, even for λ values greater than 0.05. This can be attributed to the fact that the k -NN student, initialized with the 10 samples ($2 \cdot 5$ for each sentiment), achieved an accuracy of 0.64%, while the GPT-3.5 teacher, provided with the same examples as the in-context demonstrators, achieved a remarkable accuracy of 94%. We also observed that for $\lambda = 0.05$ we obtained a high performance of 93% while calling the teacher for only 60% of the set, resulting in 40% cost reduction.

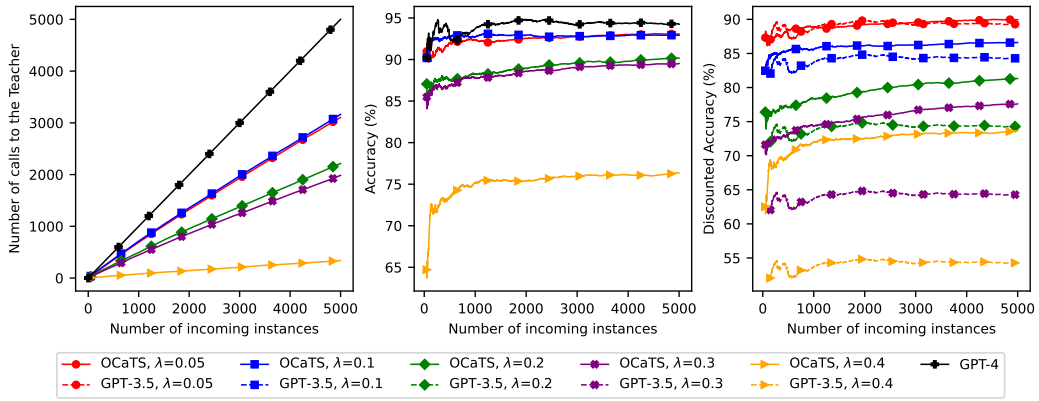


Fig. 4.9: The number of calls to the teacher (on the left), accuracy (in the middle), and discounted accuracy (on the right) are shown for various λ values when using a GPT-3.5 teacher and a k -NN student. The dashed lines indicate the discounted accuracy when GPT-4 is called for all incoming instances.

4.4 Related work on this topic

During the completion of this thesis, a recent study was conducted in this area by Ramir ez et al. [Ram+23]. They proposed the concept of neural caching, which uses active learning sampling techniques to answer questions with a smaller student model instead of a larger one available via API. The student model used was $T5_{base}$, a version of the T5 [Raf+19] model, which is a pre-trained encoder-decoder Transformer model belonging to the class of LLMs. LoRA adapter layers were added for parameter-efficient fine-tuning [Hu+22]. The teacher model used was `text-davinci-003`, a variation of Instruct-GPT [Ouy+22]. They evaluated their framework in text classification using four distinct data sets. The results of this work further support the assertion that active learning criteria can be used to have a local, cost-effective model to reduce the calls to an LLM. This work is similar to the work performed in this thesis, as both try to minimize calls to expensive LLMs. Our work is different in that we adjust the criteria for when to trust the student and when to invoke the teacher based on a significance factor λ . This factor assigns a specific importance to the cost of calling the teacher, with a higher value of λ leading to a more severe penalty for such calls. Furthermore, we do not limit ourselves to a fixed budget, as we are not aiming to train the student in an online active learning setting. Instead, our focus is on creating an ensemble that combines the teacher and the student, balancing the cost-performance trade-off while continually updating the student.

Zhu et al. [Zhu+23] recently studied two strategies that are analogous to the problem discussed in this thesis: (1) they implemented a caching system to store past queries and (2) created a model multiplexer to select the most suitable model from a set of models to process each new query. When a query or prompt is received,

they first check if it is stored in the cache. If not, a model multiplexer is used to decide which model should be used first, based on the cost of each model. This cost considers a range of factors ranging from floating point operations (FLOPS) to API calls. The primary distinction between this work and the one discussed in this thesis is that we train a local, inexpensive student to reduce the number of API calls required for an expensive LLM, while their primary focus is on achieving lower latency both online and offline.

Liu et al. [Liu+23b] recently proposed a method called online speculative decoding to accelerate the inference of language models. The concept is to use the extra computing power in an LLM serving cluster to continually update multiple draft models with user query data. Since LLM inference is limited by memory, the extra computing power in a typical language model serving cluster can be used to retrain the draft models without additional cost. Since the query distribution of an LLM service is relatively simple, retraining on the query distribution allows the draft model to more accurately predict the output of the target model, especially for data from the query distribution. As the draft model is updated online, it is adjusted to the query distribution in real time, reducing the impact of distribution shifts. They create a prototype for online speculative decoding, employing online knowledge distillation, and assess its performance using both synthetic and actual query data across various well-known LLM. This work is similar to that of Liu et al., yet there are some key differences. We do not use a cluster of LLMs, as our focus is on training a much smaller and cheaper student. Furthermore, we do not train our student on all the observed previous queries, instead we use active learning techniques to determine which data the student is unable to handle and we train on the responses the teacher provides for these queries. Our main goal and contribution is to reduce the total number of calls to LLMs available through the API while still maintaining performance. On the other hand, Liu et al. emphasize reducing the overall latency while still maintaining performance.

Conclusion and Future work

5.1 Conclusion

In this thesis, we explored how small and medium-sized businesses (SMEs) can reduce their operational expenses (OpEx) when using third-party API-provided commercial Large Language Models (LLMs) by utilizing previously streamed instances. We proposed OCaTS, an Online Cost-aware Teacher-Student framework, which helps SMEs cut down operational costs by caching the answers of commercial LLMs and training local students while taking into consideration the certainty of the student and the representativeness of previously streamed instances. We also presented a novel discounted version of common evaluation metrics to measure the performance-user frustration trade-off, which penalizes the overall performance based on the number of calls to the commercial LLM. We tested OCaTS in three experiments with GPT3.5 and GPT-4 as teachers and a k -NN and an MLP as students on two different tasks: Sentiment Analysis and Intent Recognition. We used the BANKING77 and Large Movie Reviews datasets and showed that OCaTS successfully captures the performance-cost trade-off and manages to reduce or increase the calls to the teacher depending on the significance of the cost the teacher has. Finally, we found that the total calls to the teacher can be reduced by one third with only a slight decrease in performance.

5.2 Future work

We are aiming to carry out experiments in the future that involve different scenarios, datasets and tasks for LLMs, such as question answering, summarization, and text generation. We also want to experiment with an empty cache with no training or development data. Additionally, we plan to investigate scenarios where the classes of the problem change dynamically. We also plan to dive deeper into discounted metrics by providing adaptive policies for λ to allow more expenditure in the short term and be more selective in the long run. Additionally, we intend to include more financial metrics, such as the cost of retraining and hosting the local student, in the discounted versions of the evaluation measures and explore more complex

strategies for when the student should refrain from responding using techniques such as Reinforcement Learning and Game Theory.

5.3 Limitations

This thesis proposed a framework (OCaTS) to help small businesses save money by incorporating commercial LLMs into their solutions. We used 231 (three per class) which were chosen to be diverse and representative of the classes, for the Banking77 dataset, and 10 (5 per sentiment), for the LMR dataset, to train the students. To adjust t_c and t_H , we used a development set taken from the training data, which is not always available in practice, but was used for the sake of analysis. Our analysis can be used by SMEs as a starting point to reduce the number of trials needed to find suitable values for t_c and t_H . Furthermore, $\hat{\phi}$ only takes into account the cost of calling the teacher (ρ) and the performance drop, but OCaTS can be extended to include the student cost and other financial metrics with different weights. We did not compare with existing caching libraries, such as GPTCache ¹, as they are less flexible than OCaTS and can only be used with a limited number of teacher-student settings.

5.4 Ethics statement

Constantly querying LLMs to solve everyday tasks is not only costly; it also has a large energy footprint. Our framework aims to alleviate both phenomena. However, our study required a significant amount of resources. Nevertheless, we believe that by making the framework and the analysis publicly available, we can pave the way towards reducing the resources required by SMEs to handle their day-to-day tasks in the long run.

¹<https://github.com/zilliztech/GPTCache>

Bibliography

- [AA07] Shilpa Arora and Sachin Agarwal. “Active learning for natural language processing”. In: *Language Technologies Institute School of Computer Science Carnegie Mellon University 2* (2007).
- [AI08] Alexandr Andoni and Piotr Indyk. “Near-Optimal Hashing Algorithms for Approximate Nearest Neighbor in High Dimensions”. In: *Commun. ACM* 51.1 (Jan. 2008), pp. 117–122.
- [Aki+19] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. “Optuna: A Next-generation Hyperparameter Optimization Framework”. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2019.
- [Alo+19] Md. Zahangir Alom, Tarek Taha, Chris Yakopcic, et al. “A State-of-the-Art Survey on Deep Learning Theory and Architectures”. In: *Electronics* 8 (Mar. 2019), p. 292.
- [Ary+98] Sunil Arya, David M. Mount, Nathan S. Netanyahu, Ruth Silverman, and Angela Y. Wu. “An Optimal Algorithm for Approximate Nearest Neighbor Searching Fixed Dimensions”. In: *J. ACM* 45.6 (Nov. 1998), pp. 891–923.
- [BCN06] Cristian Bucilun, Rich Caruana, and Alexandru Niculescu-Mizil. “Model Compression”. In: *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD ’06. Philadelphia, PA, USA: Association for Computing Machinery, 2006, pp. 535–541.
- [Ber+11] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. “Algorithms for Hyper-Parameter Optimization”. In: *Advances in Neural Information Processing Systems*. Ed. by J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger. Vol. 24. Curran Associates, Inc., 2011.
- [BM98] Avrim Blum and Tom Mitchell. “Combining Labeled and Unlabeled Data with Co-Training”. In: *Proceedings of the Eleventh Annual Conference on Computational Learning Theory*. COLT’ 98. Madison, Wisconsin, USA: Association for Computing Machinery, 1998, pp. 92–100.

- [Bro+20] Tom B. Brown, Benjamin Mann, Nick Ryder, et al. “Language Models are Few-Shot Learners”. In: *ArXiv abs/2005.14165* (2020).
- [BSM21] Meghana Moorthy Bhat, Alessandro Sordoni, and Subhabrata Mukherjee. “Self-training with few-shot rationalization”. In: *Proceedings of the 2021 conference on empirical methods in natural language processing*. 2021, pp. 10702–10712.
- [Cas+20] Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. “Efficient Intent Detection with Dual Sentence Encoders”. In: *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*. Online: Association for Computational Linguistics, July 2020, pp. 38–45.
- [CGJ96] David A Cohn, Zoubin Ghahramani, and Michael I Jordan. “Active learning with statistical models”. In: *Journal of artificial intelligence research* 4 (1996), pp. 129–145.
- [CH67] Thomas M. Cover and Peter E. Hart. “Nearest neighbor pattern classification”. In: *IEEE Trans. Inf. Theory* 13 (1967), pp. 21–27.
- [Das11] Sanjoy Dasgupta. “Two faces of active learning”. In: *Theoretical Computer Science* 412.19 (2011). Algorithmic Learning Theory (ALT 2009), pp. 1767–1781.
- [Dev+19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *ArXiv abs/1810.04805* (2019).
- [Du+21] Jingfei Du, Edouard Grave, Beliz Gunel, et al. “Self-training Improves Pre-training for Natural Language Understanding”. In: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Ed. by Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, et al. Online: Association for Computational Linguistics, June 2021, pp. 5408–5418.
- [Dud76] Sahibsingh A. Dudani. “The Distance-Weighted k-Nearest-Neighbor Rule”. In: *IEEE Transactions on Systems, Man, and Cybernetics SMC-6.4* (1976), pp. 325–327.
- [ES10] Brian S Everitt and Anders Skrondal. *The Cambridge dictionary of statistics*. 2010.
- [Fuj+98] Atsushi Fujii, Kentaro Inui, Takenobu Tokunaga, and Hozumi Tanaka. “Selective Sampling for Example-based Word Sense Disambiguation”. In: *Computational Linguistics* 24.4 (1998), pp. 573–597.
- [Gou+21] Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. “Knowledge Distillation: A Survey”. In: *International Journal of Computer Vision* 129.6 (Mar. 2021), pp. 1789–1819.
- [HS04] Klaus Hechenbichler and Klaus Peter Schliep. “Weighted k-Nearest-Neighbor Techniques and Ordinal Classification”. In: 2004.

- [Hu+22] Edward J Hu, yelong shen, Phillip Wallis, et al. “LoRA: Low-Rank Adaptation of Large Language Models”. In: *International Conference on Learning Representations*. 2022.
- [HVD15a] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. *Distilling the Knowledge in a Neural Network*. 2015.
- [HVD15b] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. “Distilling the Knowledge in a Neural Network”. In: *ArXiv abs/1503.02531* (2015).
- [Hwa04] Rebecca Hwa. “Sample Selection for Statistical Parsing”. In: *Computational Linguistics* 30.3 (2004), pp. 253–276.
- [JDJ17] Jeff Johnson, Matthijs Douze, and Hervé Jégou. *Billion-scale similarity search with GPUs*. 2017. eprint: arXiv:1702.08734.
- [LBH15] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep Learning”. In: *Nature* 521 (2015), pp. 436–444.
- [LC94] David D. Lewis and Jason Catlett. “Heterogenous Uncertainty Sampling for Supervised Learning”. In: *Proceedings of the Eleventh International Conference on International Conference on Machine Learning*. ICML’94. New Brunswick, NJ, USA: Morgan Kaufmann Publishers Inc., 1994, pp. 148–156.
- [LG94] David D. Lewis and William A. Gale. “A Sequential Algorithm for Training Text Classifiers”. In: *CoRR abs/cmp-lg/9407020* (1994). arXiv: cmp-1g/9407020.
- [Li+21a] Shiyang Li, Semih Yavuz, Wenhua Chen, and Xifeng Yan. “Task-adaptive Pre-training and Self-training are Complementary for Natural Language Understanding”. In: *Findings of the Association for Computational Linguistics: EMNLP 2021*. Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 1006–1015.
- [Li+21b] Shiyang Li, Semih Yavuz, Wenhua Chen, and Xifeng Yan. “Task-adaptive Pre-training and Self-training are Complementary for Natural Language Understanding”. In: *Findings of the Association for Computational Linguistics: EMNLP 2021*. Ed. by Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih. Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 1006–1015.
- [Liu+23a] Nelson F. Liu, Kevin Lin, John Hewitt, et al. *Lost in the Middle: How Language Models Use Long Contexts*. 2023. arXiv: 2307.03172 [cs.CL].
- [Liu+23b] Xiaoxuan Liu, Lanxiang Hu, Peter Bailis, et al. *Online Speculative Decoding*. 2023. arXiv: 2310.07177 [cs.AI].
- [LMR04] Michael Lindenbaum, Shaul Markovitch, and Dmitry Rusakov. “Selective Sampling for Nearest Neighbor Classifiers”. In: *Machine Learning* 54.2 (Feb. 2004), pp. 125–152.

- [Mar+23] Katerina Margatina, Timo Schick, Nikolaos Aletras, and Jane Dwivedi-Yu. *Active Learning Principles for In-Context Learning with Large Language Models*. 2023. arXiv: 2305.14264 [cs.CL].
- [Mi+21a] Fei Mi, Wanhao Zhou, Feng Cai, et al. “Self-training Improves Pre-training for Few-shot Learning in Task-oriented Dialog Systems”. In: *Conference on Empirical Methods in Natural Language Processing*. 2021.
- [Mi+21b] Fei Mi, Wanhao Zhou, Lingjing Kong, et al. “Self-training Improves Pre-training for Few-shot Learning in Task-oriented Dialog Systems”. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 1887–1898.
- [Mon21] Robert Monarch. *Human-in-the-Loop Machine Learning*. Manning Publications, 2021.
- [Niu+20] Yilin Niu, Fangkai Jiao, Mantong Zhou, et al. “A Self-Training Method for Machine Reading Comprehension with Soft Evidence Extraction”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Ed. by Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault. Online: Association for Computational Linguistics, July 2020, pp. 3916–3927.
- [Ope23] OpenAI. “GPT-4 Technical Report”. In: *ArXiv abs/2303.08774* (2023).
- [Ouy+22] Long Ouyang, Jeffrey Wu, Xu Jiang, et al. “Training language models to follow instructions with human feedback”. In: *Advances in Neural Information Processing Systems*. Ed. by S. Koyejo, S. Mohamed, A. Agarwal, et al. Vol. 35. Curran Associates, Inc., 2022, pp. 27730–27744.
- [Rad+18] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. “Improving language understanding by generative pre-training”. In: (2018).
- [Raf+19] Colin Raffel, Noam Shazeer, Adam Roberts, et al. “Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer”. In: *CoRR abs/1910.10683* (2019). arXiv: 1910.10683.
- [Ram+23] Guillem Ramírez, Matthias Lindemann, Alexandra Birch, and Ivan Titov. *Cache & Distil: Optimising API Calls to Large Language Models*. 2023. arXiv: 2310.13561 [cs.CL].
- [RG19] Nils Reimers and Iryna Gurevych. “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks”. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 3982–3992.

- [RHW86] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. *Learning representations by back-propagating errors*. en. Oct. 1986.
- [RZ09] Stephen Robertson and Hugo Zaragoza. “The Probabilistic Relevance Framework: BM25 and Beyond”. In: *Found. Trends Inf. Retr.* 3.4 (Apr. 2009), pp. 333–389.
- [SC08] Burr Settles and Mark Craven. “An Analysis of Active Learning Strategies for Sequence Labeling Tasks”. In: *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*. Honolulu, Hawaii: Association for Computational Linguistics, Oct. 2008, pp. 1070–1079.
- [Set12] Burr Settles. *Active Learning*. Morgan & Claypool Publishers, 2012.
- [Sha48] C. E. Shannon. “A mathematical theory of communication”. In: *The Bell System Technical Journal* 27.3 (1948), pp. 379–423.
- [She+04] Dan Shen, Jie Zhang, Jian Su, Guodong Zhou, and Chew-Lim Tan. “Multi-Criteria-based Active Learning for Named Entity Recognition”. In: *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)*. Barcelona, Spain, July 2004, pp. 589–596.
- [SJ21] Thakur Ashutosh Suman and Abhinav Jain. “AStarTwice at SemEval-2021 Task 5: Toxic Span Detection Using RoBERTa-CRF, Domain Specific Pre-Training and Self-Training”. In: *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*. Ed. by Alexis Palmer, Nathan Schneider, Natalie Schluter, et al. Online: Association for Computational Linguistics, Aug. 2021, pp. 875–880.
- [SLA12] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. “Practical bayesian optimization of machine learning algorithms”. In: *Advances in neural information processing systems* 25 (2012).
- [Son+20] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. “MPNet: Masked and Permuted Pre-training for Language Understanding”. In: *NeurIPS 2020*. ACM. Sept. 2020.
- [TK02] Simon Tong and Daphne Koller. “Support Vector Machine Active Learning with Applications to Text Classification”. In: *J. Mach. Learn. Res.* 2 (Mar. 2002), pp. 45–66.
- [Tou+23a] Hugo Touvron, Louis Martin, Kevin Stone, et al. *Llama 2: Open Foundation and Fine-Tuned Chat Models*. 2023. arXiv: 2307.09288 [cs.CL].
- [Tou+23b] Hugo Touvron, Louis Martin, Kevin Stone, et al. *Llama 2: Open Foundation and Fine-Tuned Chat Models*. 2023. arXiv: 2307.09288 [cs.CL].
- [Vas+17] Ashish Vaswani, Noam Shazeer, Niki Parmar, et al. “Attention Is All You Need”. In: *CoRR abs/1706.03762* (2017). arXiv: 1706.03762.

- [Yan+23] Chuanguang Yang, Xinqiang Yu, Zhulin An, and Yongjun Xu. *Categories of Response-Based, Feature-Based, and Relation-Based Knowledge Distillation*. 2023. arXiv: 2306.10687 [cs.CV].
- [Yar95] David Yarowsky. “Unsupervised Word Sense Disambiguation Rivaling Supervised Methods”. In: *Proceedings of the 33rd Annual Meeting on Association for Computational Linguistics*. ACL '95. Cambridge, Massachusetts: Association for Computational Linguistics, 1995, pp. 189–196.
- [Zha+19] L. Zhang, J. Song, A. Gao, et al. “Be Your Own Teacher: Improve the Performance of Convolutional Neural Networks via Self Distillation”. In: *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. Los Alamitos, CA, USA: IEEE Computer Society, Nov. 2019, pp. 3712–3721.
- [Zhu+23] Banghua Zhu, Ying Sheng, Lianmin Zheng, et al. *On Optimal Caching and Model Multiplexing for Large Model Inference*. 2023. arXiv: 2306.02003 [cs.LG].

List of Acronyms

OCaTS Online Cost-aware Teacher-Student

NLP Natural Language Processing

SME Small and Medium Enterprise

OpEx Operational Expenses

LLM Large Language Model

LMR Large Movie Reviews

List of Figures

2.1	The Transformer’s Encoder-Decoder architecture is illustrated on the left and right sides. The components of the encoder are displayed on the left, while the components of the decoder are shown on the right. Source:[Vas+17]	6
2.2	Bert’s encoder-only architecture and GPT-2’s decoder-only architecture. Source: http://jalammr.github.io/illustrated-gpt2/	7
2.3	The value of k affects how the k -NN algorithm works. In this example, when k is set to 3, the red point is classified as Class B, but when k is increased to 7, the same point is classified as Class A.	9
2.4	The teacher-student framework for knowledge distillation	11
2.5	Different knowledge distillation training schemes	12
2.6	The iterative process of active learning using a pool of data. Source: [Set12]	14
3.1	OCaTS architecture	18
4.1	Data distribution of the <i>original</i> BANKING77 dataset. Left figure shows the distribution of the labels in the training set while the right figure show the distribution in the test set.	22
4.2	Data distribution of the <i>original</i> LMR dataset. Left figure shows the distribution of the labels in the training set while the right figure show the distribution in the test set.	23
4.3	Contour plots (discounted accuracy) obtained during threshold tuning, for various λ values.	26
4.4	The number of calls to the teacher (on the left), accuracy (in the middle), and discounted accuracy (on the right) are shown for various λ values when using a GPT-4 teacher and a k -NN student. The dashed lines indicate the discounted accuracy when GPT-4 is called for all incoming instances. The larger the λ , the more the SME prefers fewer calls, even though it may lead to increased user frustration. As expected, OCaTS has a better discounted accuracy than when the GPT-4 teacher is called for all incoming instances.	27
4.5	The learning curves of the models demonstrate that MLP (blue line) is much more adept at assimilating the data and performs comparably to k -NN, while DBert’s performance is significantly weaker.	28

4.6	Contour plots of discounted accuracy, obtained during threshold tuning, are presented for different λ values when using an MLP student and a GPT-4 teacher.	29
4.7	The number of calls to the teacher (on the left), accuracy (in the middle), and discounted accuracy (on the right) are shown for various λ values when using a GPT-4 teacher and a MLP student.	30
4.8	Contour plots of discounted accuracy, obtained during threshold tuning, are presented for different λ values when using a k -NN student and a GPT3.5 teacher.	32
4.9	The number of calls to the teacher (on the left), accuracy (in the middle), and discounted accuracy (on the right) are shown for various λ values when using a GPT-3.5 teacher and a k -NN student. The dashed lines indicate the discounted accuracy when GPT-4 is called for all incoming instances.	33

List of Tables

4.1	Examples of the Banking77 dataset.	21
4.2	Parameters that were tuned for defining the architectures of DBert and MLP respectively.	28

List of Algorithms

1	Typical active learning procedure	13
---	---	----

