

# “Call My Big Sibling (CMBS)” – A Confidence-Based Strategy Leveraging Instance Selection to Combine Small and Large Language Models for Cost-Effective Text Classification

Anonymous ACL submission

## Abstract

Transformers have achieved cutting-edge results, with Large Language Models (LLMs) considered SOTA in many NLP tasks. However, the literature has not yet fully demonstrated that LLMs are always superior to first-generation Transformers (a.k.a. Small Language Models (SLMs)) in all NLP tasks and scenarios. This study compares four SLMs (BERT, RoBERTa, Qwen, BART) with four open LLMs (LLaMA 3.1, Mistral, Falcon, DeepSeek) across 9 sentiment analysis and 4 topic classification datasets, totaling over 1000 results. Findings indicate that open LLMs can moderately outperform or tie with SLMs in all tested datasets, though only when fine-tuned, at a very high computational cost. To address the cost-effectiveness trade-off, we propose “Call My Big Sibling” (CMBS)<sup>1</sup>, a novel confidence-based framework that efficiently integrates *calibrated SLMs with open LLMs using advanced instance selection techniques*. CMBS assigns high-confidence predictions to the cheaper SLM, while low-confidence cases are directed to LLMs in zero-shot, in-context, or partially-tuned (leveraging state-of-the-art **instance selection** approaches) modes, optimizing cost-effectiveness. Experiments show that CMBS significantly outperforms SLMs and achieves *LLM-level effectiveness at a fraction of the cost*, offering a practical alternative for cost-sensitive NLP applications.

## 1 Introduction

Automatic text classification (ATC) is essential in diverse contexts, ranging from organizing large data volumes to personalizing user experiences. ATC has experienced a huge revolution with the advent of semantically enriched Transformer models (Devlin et al., 2019) that have achieved state-of-the-art performance (de Andrade et al., 2023; Cunha et al., 2023a; Zanutto et al., 2021).

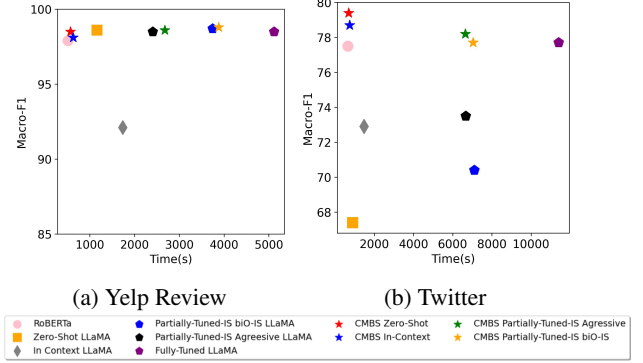


Figure 1: Total Time (seconds) and Macro-F1 in RoBERTa, version of LLaMA and CMBS. All CMBS proposals outperform the other baselines, being much cheaper.

More recently, Large Language Models (LLMs) emerged, built on top of the first generation of Transformers (a.k.a., small language models – SLMs). Studies have identified LLMs as the current SOTA for several NLP tasks (Liang et al., 2023). Although the literature reports LLMs’ superiority for tasks such as summarization and translation, for others, such as sentiment analysis and topic classification, it is not yet clear whether LLMs’ complexity and size (e.g., in terms of a number of parameters) translate into statistical and mainly *practical* gains. In fact, several studies point to the SLM RoBERTa as a very strong sentiment classifier (Cunha et al., 2023b) ranking prominently on leaderboards such as the GLUE benchmark<sup>2</sup>.

Depending on the type of training (or its absence), LLM strategies can be categorized into four groups: *zero-shot*, *in-context*, *partially-*, and *fully-tuned*. In zero-shot, the model performs tasks without specific training. In in-context, a few examples are provided in the prompt for learning. Partially-tuned models use a significant part of the training set, while fully-tuned ones use the entire set for better optimization. As shown in (Table 1), open LLMs surpass SLMs mainly in the

<sup>1</sup>Code available at <https://github.com/Anonymous>

<sup>2</sup><https://gluebenchmark.com/leaderboard/>

066 **fully-tuned** scenario, albeit at a **very high cost**,  
067 highlighting the need for cost-effective alternatives.

068 Accordingly, our first research question is **RQ1**:  
069 “Are (open) LLMs more effective than SLMs  
070 in sentiment and topic classification?” Recent  
071 work (Fields et al., 2024) shows no consensus  
072 on whether LLMs always excel in ATC tasks.  
073 To investigate, we compared four SLMs (BERT,  
074 RoBERTa, Qwen<sup>3</sup>, BART) and four open LLMs  
075 (LLaMA 3.1 8B, Mistral 7B, Falcon 7B, DeepSeek  
076 8B) on 9 sentiment and 4 topic datasets, including  
077 two post-LLM releases (IMDB2024, RottenT2024)  
078 to reduce data contamination (Liang et al., 2023).  
079 We focus on open LLMs, as closed ones (e.g., Chat-  
080 GPT) lack transparency<sup>4</sup>. Our statistically rigorous  
081 experiments<sup>5</sup> show open LLMs can outperform  
082 SLMs, achieving up to 8.3% (avg. 3.1%) effective-  
083 ness gains, though mostly in fully-tuned mode.

084 Given the (much) higher computational costs  
085 associated with fully fine-tuning open LLMs (the  
086 most effective approach), a natural question we  
087 posited is **RQ2**: “How does the computational cost  
088 of using open LLMs for ATC compare to SLMs’  
089 cost?”. To answer this question, we conducted  
090 a thorough analysis of our experimental results,  
091 considering zero-shot, in-context, partially-tuned,  
092 and fully-tuned strategies, to assess the trade-offs  
093 between effectiveness and costs in terms of  
094 computational cost<sup>6</sup> to train and employ the  
095 models and their impact on carbon emission. We  
096 found that LLMs are orders of magnitude more  
097 costly to fully fine-tune when compared to SLMs  
098 – fully fine-tuned LLMs are up to 1700% more ex-  
099 pensive than SLMs. As current LLMs produce just  
100 moderate gains over SLMs and only through highly  
101 costly full fine-tuning processes, depending on the  
102 scenario, the benefits may not be worth the costs.

103 All this leads to our final research question **RQ3**:  
104 “Is it possible to perform a combination of SLM and  
105 (open) LLMs to achieve a better effectiveness/cost  
106 trade-off than using either SLM or LLM alone?”  
107 To answer this question, we proposed a novel  
108 confidence-based strategy called “Call My Big Sib-  
109 ling” (**CMBS**), which smoothly combines SLM  
110 and (open) LLMs based on *calibrated confidences*.

<sup>3</sup>Among the multiple Qwen versions, we use the smallest one, with 0.5B parameters, which we consider as an SLM due its effectiveness comparable to the other SLMs used in this paper and its fast response time.

<sup>4</sup>Closed LLMs are irreproducible (Gao et al., 2024).

<sup>5</sup>Please see Section 4 for a full description of our experimental setup.

<sup>6</sup>Including both, time to train the model and predict test.

111 In CMBS, we rely on *fully fine-tuned SLMs*,  
112 which have already attained effectiveness and  
113 efficiency and are calibrated<sup>7</sup> for ATC tasks. We  
114 then use the classification confidence to determine  
115 whether the LLM should classify a low-confidence  
116 document. In other words, the fully-tuned SLM  
117 classifies high-certainty documents (i.e., with high  
118 certainty calibrated scores), while low-confidence  
119 documents are sent to the *zero-shot, in-context or*  
120 *partially-tuned versions of the LLMs* for ATC. In  
121 the particular case of the partially-tuned LLMs,  
122 our CBMS solution introduces an innovation  
123 by exploiting recent SOTA **instance selection**  
124 methods (Cunha et al., 2025) to reduce the cost of  
125 tuning while maintaining effectiveness. Such com-  
126 bination with a cheaper LLM version (compared to  
127 the fully-tuned LLM) brings potential effectiveness  
128 gains to the SLM and is on par with fully-tuned  
129 LLM, being a very attractive, cost-effective option  
130 in basically all tested cases, as we shall see.

131 Our results show that for sentiment classification,  
132 combining an SLM with a zero-shot LLM (**CMBS**  
133 **Zero-shot**) improves effectiveness at the lowest  
134 cost, proving its practicality. Figure 1 illustrates the  
135 effectiveness (Macro-F1) and efficiency (Time(s))  
136 of our approach versus baselines in two datasets,  
137 with our proposals marked by stars. In Figure 1a,  
138 CMBS Zero-Shot matches Fully-Tuned LLaMA’s  
139 effectiveness at a fraction of the cost. Similarly, in  
140 Figure 1b, all CMBS methods outperform baselines  
141 while remaining much cheaper. Indeed, our experi-  
142 ments reveal that CBMS Zero-shot outperforms the  
143 SLM in 8 out of 9 sentiment datasets, tying with the  
144 remaining one, with an increase in computational  
145 cost over SLMs of only 8%. Moreover, compared  
146 to fully-tuned LLaMA, CBMS Zero-Shot delivers  
147 comparable effectiveness at a significantly lower  
148 cost. In 4 of the 13 datasets, CMBS Zero-Shot ties  
149 with the fully-tuned LLM, with minimal losses (on  
150 average, just 2%) in the other datasets, at  $\frac{1}{10}$  of the  
151 cost. CMBS Partially-Tuned-IS (Aggressive), on  
152 turn, ties with fully-tuned LLaMA in *all* sentiment  
153 datasets at half of the cost.

154 For topic classification with several classes (up  
155 to 11) and uneven distributions, CBMS zero-shot  
156 and **CMBS In-Context** struggle with effectiveness.  
157 Only when combined with **CMBS Partially-**  
158 **Tuned-IS** can it surpass the SLM. Among 4  
159 datasets, CMBS Partially-Tuned-IS (Aggressive)

<sup>7</sup>The confidence of the SLM’s softmax function is highly calibrated as we shall discuss.

outperforms RoBERTa in 2 and ties in the others. Against fully-tuned LLaMA, it achieves statistical parity in 3 datasets (with just a 2% loss in the 4th) while operating at half the computational cost.

In sum, the main contributions of this paper are:

- A comprehensive comparative evaluation of SLMs and (open) LLMs regarding cost/effectiveness trade-offs.
- The proposal of “Call My Big Sibling” (CMBS), a novel confidence-based strategy to combine calibrated SLMs and zero-shot or partially-tuned with *instance selection* LLMs, aimed at optimizing the effectiveness-cost trade-off.
- A thorough evaluation of our proposals considering 13 distinct datasets, in 2 tasks<sup>8</sup>: sentiment (binary) and topic (multi-class) classification tasks, 4 SLMs and zero-shot, in-context, partially-tuned and fully-tuned versions of 4 open LLMs.

## 2 Related Work

LLMs’ computational costs have led to numerous studies highlighting their financial and environmental impacts. (Strubell et al., 2019) illustrates the substantial financial costs propelled by the continuous need for investment in specialized hardware to manage progressive LLMs. This trend not only limits access to these models but also escalates energy consumption, affecting the environment by increasing carbon dioxide (CO<sub>2</sub>) emissions.

Among LLMs, proprietary and closed-source ones, such as *GPT*, operate as black boxes. This opacity poses challenges in comprehending their training and internal structures, obstructing reproducibility. Utilizing such LLMs often entails transmitting data through web platforms or APIs, a sensitive issues when privacy is a concern. As a result, numerous studies advocate restricting scientific evaluations to locally-run, open-source LLMs (Spirling, 2023).

(Xu et al., 2024) combines SLMs with LLMs aiming to improve effectiveness. In that work, computational costs are not evaluated, and LLMs are used to classify the entire test set, unlike ours, which only uses them to classify a subset of hard instances. On the other hand, we select only low-confidence documents to be forwarded to the LLM, a strategy that greatly reduces computational costs as the LLM is significantly more expensive. Moreover, a closed LLM is employed via an API,

which provides no control over the computational structure or the model architecture. Finally, only a single sentiment dataset is used in the experiments.

Liang et al., 2023 investigate various LLMs across multiple tasks, prompts, metrics, and datasets. Like Liang et al., we investigate the effectiveness-cost tradeoff of LLMs. Unlike their study, which focuses on the breadth of evaluation with several domains (including only one sentiment dataset), our work is depth-oriented into the specific task of sentiment and topic classification, covering multiple datasets with diverse characteristics and domains. Moreover, although Liang et al. evaluates several models, they do not compare them with an SLM such as RoBERTa, considered SOTA in sentiment and topic classification (Bai et al., 2023a; Cunha et al., 2021a, 2020; França et al., 2024; Belém et al., 2024). Last, but not least, they do not provide any solution for the effectiveness-cost trade-off problem, which *we do!*.

Yue et al., 2024 proposed a hybrid approach using GPT-3.5 (weak) and GPT-4 (strong), where the weak model is the first to respond, and a decision model checks response consistency across multiple prompts, calling the strong model if needed. This potentially improves effectiveness but does not reduce model costs – indeed, only financial costs based on the number of OpenAI API calls are considered in their work. In contrast, our approach considers both financial and computational costs (time) and uses SLM confidence, not response consistency, to decide when to invoke a cost-optimized LLM.

Finally, it is worth noticing that the recent success and impact, beyond the academic world, of DeepSeek (DeepSeek-AI et al., 2025), which matched or improved the effectiveness of state-of-the-art pre-trained LLMs while reducing computational demands, show how important the issue of trade-off effectiveness vs. cost is to the community. We investigate DeepSeek as one of the LLMs for our CBMS solution.

## 3 The CBMS Solution

One of the main contributions of our work is the proposal of a novel strategy to combine simpler, more efficient, but perhaps less effective SLMs with potentially more effective but costly LLMs, aiming to promote effectiveness while minimizing computational costs. Our solution, “Call-My-Big-Sibling” (CMBS), metaphorically

<sup>8</sup>An additional application of the CBMS on a different NLP task from the GLUE benchmark is discussed in Appendix J.



conjuges up the image of a small (but smart) child who, in a challenging situation, seeks help from a bigger sibling. CBMS pursues the best trade-off between effectiveness and costs with a confidence-based pipeline of Language Models.

CMBS seamlessly integrates SLMs and (open) LLMs by leveraging **instance selection** and **calibrated confidences**. In this framework, we first employ *fully-tuned SLMs models*<sup>9</sup>, which are already highly effective in some classification tasks (and faster to tune compared to LLMs). In our solution, (test) documents classified below a certain confidence threshold (a parameter) by the SLM are sent to an open LLM to be classified. We have 5 versions of LLMs: Zero-shot, In-Context, Partially-Tuned-IS biO-IS and Aggressive, and Fully-Tuned.

In particular, in the case of partially-tuned LLMs, we have two options. Partially-Tuned-IS biO-IS leverages the state-of-the-art just-released instance selection method biO-IS (Cunha et al., 2025) to reduce the training set size, and thus the computational cost while trying to maintain effectiveness. In biO-IS, the “optimal” training set reduction is variable and determined by its own algorithm (Cunha et al., 2025). Partially-tuned-IS Aggressive always applies a randomly stratified reduction of 50% in the training set despite potential effectiveness losses in the LLM. We use this percentage drawing on recent work in Instance Selection (Cunha et al., 2023a) that determined this is the maximum reduction rate that can still assure good efficiency while producing minimal effectiveness losses. The procedure is illustrated in Figure 2.

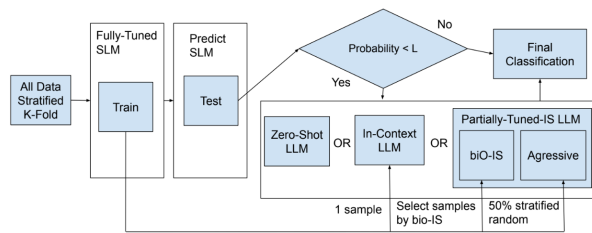


Figure 2: CBMS FlowChart.

For CBMS to properly work, we have to trust the probability outputs, or, in other words, the probabilities need to be calibrated<sup>10</sup>. Wolfe et al., 2017 argue that RoBERTa’s softmax function provides calibrated probabilities as it is a generalization of logistic regression. To demonstrate this, Table 13 (Ap-

pendix G) presents the Brier score used to measure model calibration, in datasets used in our experiments, by two transformers: BERT and RoBERTa. This score is calculated based on the model probabilities and actual labels. The score ranges from 0 to 1, with values closer to 1 indicating a better alignment between probabilistic predictions and actual outcomes. As we can observe, the table reinforces that RoBERTa is a very calibrated model (Brier score > 0.8), being as calibrated as well-calibrated classifiers such as Logistic Regression and Random Forests that obtain similar Brier scores for some of these benchmark datasets (Cunha et al., 2025). In addition to Table 13, Figure 5 presents the calibration curve for the BERT (red) and RoBERTa (blue) models. The dashed line along the main diagonal represents ideal calibration. As shown in Figure 5, RoBERTa is consistently closer to the ideal calibration across all datasets, further supporting the results presented in Table 13.

We select a document set for which the classifier is least confident about its classification (Probability <  $L$ ) to send to an LLM for final prediction. Due to computational costs, we employ either the zero-shot, in-context, or partially-tuned strategies for this LLM. Finally, our final prediction set is built using the following procedure: 1) we evaluate the probability the model provides and compare it with the threshold parameter; 2) we decide whether the prediction will be made using an SLM or an LLM (zero-shot, in-context or partially-tuned).

In this proposal, the choice of confidence threshold  $L$  is essential to evaluate the documents that will be sent to the LLM. To illustrate this point, Figure 3 presents the effectiveness based on prediction confidence for the SST2 dataset. On the Y-axis, we have RoBERTa’s effectiveness, and on the X-axis, we have confidence. We can observe the more confident RoBERTa is, the more effective its predictions are. The figure highlights the importance of selecting an appropriate confidence threshold, showing that it is more advantageous to forward low-confidence documents, as high-confident ones are classified with high accuracy by the cheaper SLM.

## 4 Experimental Methodology and Setup

### 4.1 Datasets

Our study draws on **thirteen** datasets developed for sentiment analysis and topic classification. The sentiment analysis datasets include **Finance**, **IMDB**, **PangMovie**, **SemEval17**, **SST**, **SST2** and

<sup>9</sup>Tuned with the full training data.

<sup>10</sup>A calibrated classifier has a strong correlation between class prediction probabilities and frequency of correctly predicted instances belonging to each probability range.

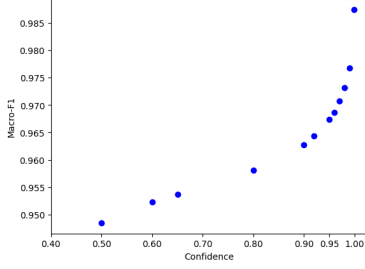


Figure 3: RoBERTa’s Macro-F1 vs Confidence for SST2.

**Yelp Review** (Yelp2L), while prompt fication we used **ACM**, **DBLP**, **Twitter** and **Webkb**. With the significant amount of data used in building LLMs, several authors express concerns about contamination in evaluation data. To minimize this issue, we collected and curated two datasets with data post-LLMs release (**RottenT2024** from Jan-Nov 2024 and **IMDB2024** from Jan-May 2024), ensuring no contamination in the training of these LLMs. Further details about the datasets in Appendix B: domain, number of documents, density, and skewness (class imbalance). Our benchmark covers a wide variety of heterogeneous scenarios.

## 4.2 Prompt Template

We evaluated four open LLMs—Falcon 7B, Mistral 7B, DeepSeek 8B, and LLaMA 3.1 8B—using the prompt template from Liang et al. (2023), who found that the most effective format includes: (i) task description, (ii) examples with expected responses, and (iii) the text to be evaluated. We adapted this for sentiment (Table 5) and topic classification (Table 6, Appendix A). The prompt provides instructions, class examples, and the text to evaluate, with the LLM generating the class as the “next word.” The In-Context LLM template (Table 7) modifies this by replacing a generic example with the closest training document, determined via cosine similarity using RoBERTa embeddings.

## 4.3 Zero-shot, In-context, Partially-tuning or Fully-tuning for Text Classification

Applying SLM or LLM pre-trained models to ATC can be done within our CBMS solution through five strategies: zero-shot, in-context, partially-tuned-IS biO-IS, partially-tuned-IS Aggressive, and fully-tuned. **Zero-shot** strategy predicts text classes without using training examples or performing model fine-tuning. In an **In-context** approach, the model relies on a prompt containing the nearest neighbors of the evaluated example inserted into it to provide context for making predictions without adjusting its weights. In the **partially-tuned (IS)** strategies,

a portion of labeled data is employed to adjust the model weights, simulating a scenario of data scarcity. As described in Section 3, we have two options for this strategy: using biO-IS instance selection for training set reduction, with variable reduction rates, depending on the dataset, and Aggressive, in which we fixed 50% of the training partition data for model training, selecting these instances in a stratified random manner. An evaluation using different fixed training data sizes was conducted in Appendix E, further justifying the 50% rate choice.

Lastly, the **fully-tuned** strategy utilizes all available labeled data in the model’s training partition to maximize model adjustment for the task and data domain. While this strategy typically achieves better effectiveness, it has a very high computational cost. In our paper, the fully-tuned strategy is used as *baseline* to compare with our proposals.

We only employ the SLMs fully tuned, which is essential for effectiveness (de Andrade et al., 2023). Fully tuning SLMs involves fine-tuning the SLM’s text representation (CLS token) and a fully connected layer that performs class prediction, utilizing all available training samples.

## 4.4 Method-Specific Parameter Tuning

All data is divided using stratified 5-fold cross-validation, a widely accepted technique in model evaluation. This method enhances the robustness and reliability of the model by splitting the dataset into five parts: three for training, validation, and testing. In each of the five iterations, the roles of the partitions alternate between training, validation, and testing, ensuring that the class distribution is preserved in the test partition. The validation set is crucial for parameter tuning, as detailed below.

For SLMs, we adopted Cunha et al. (2023b)’s hyper-parameterization, fixing the learning rate in  $2 \times 10^{-5}$ , the batch size with 64 documents, adjusted the model for five epochs and set the maximum size of each document to 256 tokens. We adopted the following parameters for the LLM models: all LLMs use 4-bit quantization, with QLoRA and PEFT enabling fine-tuning on reasonably equipped machines. For LLaMA, we used 1024 maximum tokens, a learning rate of  $2 \times 10^{-4}$ , and a temperature equal to 0.6. All other parameters were set at their default values. For fully-tuning processes, which are more costly due to the model’s weight adjustment (backpropagation), we had to reduce the maximum number of tokens to 256. We performed training for three epochs.

We introduce a confidence threshold parameter: if the SLM’s confidence is below this threshold, predictions are forwarded to the LLM, the "Big Sibling." This ensures that the complex LLM takes over only when the SLM is not confident. Using the validation set, we vary this parameter to optimize Macro-F1 without increasing cost. Table 12 (Appendix F) shows the selected threshold for sample datasets, the percentage of forwarded instances, and LLM/SLM effectiveness. For example, in SST2, documents with confidence below 0.9 go to the LLM; otherwise, the SLM classifies them. Higher thresholds mean more documents are sent to the LLM. Notably, the optimal threshold (around 0.9) is consistent across most datasets, and the LLM outperforms the SLM in these difficult cases, supporting CBMS gains.

#### 4.5 Metrics and Experimental Protocol

We evaluated SLMs and (open) LLMs regarding the effectiveness/cost tradeoff. All models were assessed on identical hardware configuration: a 4-core processor, 32GB of system memory, and an Nvidia Tesla P100 GPU. Classification effectiveness is assessed using Macro-F1 due to imbalance in several datasets. To ensure statistical validity of the results and demonstrate model generality, models were evaluated using the test set from a 5-fold stratified cross-validation methodology and a t-test with 95% confidence with Bonferroni correction to account for multiple comparisons.

To assess cost-effectiveness tradeoff, we evaluate each method’s total model-building time, including learning and class prediction for the test set. For CMBS Zero-Shot, this includes tuning the SLM, predicting part of the test set with the SLM, and a smaller portion with the LLM. CMBS in-context also accounts for finding the  $k$  nearest neighbors, while CMBS partially-tuned includes training the LLM using biO-IS or 50% stratified random of the training for Aggressive.

### 5 Experimental Results and Analyses

#### 5.1 SLMs vs. LLMs - Effectiveness

To address RQ1, we first evaluated several popular open-source LLMs, including Falcon 7B, Mistral 7B, DeepSeek 8B, and LLaMA 3.1 8B. We began by comparing the performance of these four LLMs on sentiment and topic tasks in the Zero-Shot setting. Table 9 in Appendix C presents the Macro-F1 scores, highlighting the best results in bold, includ-

ing statistical ties. LLaMA 3.1 8B consistently achieves the best results (statistically) across most datasets in both sentiment and topic tasks. Due to the high computational cost of fully tuning LLMs, we selected LLaMA 3.1 8B for all subsequent tests.

Regarding SLMs, we did a similar experiment and compared four widely used SLMs – BERT, BART, RoBERTa and QWen. Results in Table 10 in Appendix D show that, among SLMs, RoBERTa achieves the highest effectiveness (or ties for it) in all cases, confirming findings reported in the literature (Cunha et al., 2023b; Bai et al., 2023b).

Still answering RQ1, we compare RoBERTa with five LLaMA versions in Table 1—zero-shot, in-context, partially-tuned (biO-IS and Aggressive), and fully-tuned—using the full training set. Zero-shot LLaMA 3.1 ties or underperforms compared to RoBERTa in sentiment datasets and performs much worse in topic classification. Similarly, In-Context LLaMA does not outperform RoBERTa in sentiment classification. Since Zero-Shot LLaMA lacked training labels, while RoBERTa was fully tuned, RoBERTa achieved better effectiveness in most cases. The performance gap is even larger in topic classification due to more classes and higher complexity of the ATC task.

Only partially-tuned-IS (in both versions, biO-IS and Aggressive, for most datasets) and fully-tuned LLaMA can outperform RoBERTa, with some advantage for fully-tuned, especially for topic classification. However, in some datasets, such as Finance and Yelp2L (sentiment), and Twitter (topic), fully-tuned and RoBERTa are statistically tied, mainly due to larger confidence intervals (meaning higher variability) for the RoBERTa results. In some other datasets, RoBERTa’s and fully-tuned LLaMA’s effectiveness are also close.

The instance selection versions of the LLM are also very competitive with the fully-tuned version in most datasets, but at a much cheaper cost (between 30%-50%). When comparing both instance selection alternatives, we can see that biO-IS preserves effectiveness in more datasets, which is consistent with (Cunha et al., 2025), probably due to its lower training set reduction rate (around 40% on average when compared to Aggressive (always 50%). These results further motivate us to combine SLMs and LLMs with our proposed CMBS pipeline for the sake of optimizing the effectiveness-cost trade-off. This trade-off is the core of our subsequent analyses.



Dataset	RoBERTa	Zero-Shot LLaMA	In-Context LLaMA	Partially-Tuned-IS biO-IS LLaMA	Partially-Tuned-IS Aggressive LLaMA	Fully-Tuned LLaMA
Finance	<b>98.1±1.9</b>	95.4±1.2	<b>98.6±1.8</b>	88.6±1.3	<b>98.6±0.1</b>	<b>98.7±1.6</b>
Imdb	93±0.5	93±0.3	78.9±1.2	<b>95.7±0.3</b>	<b>95.8±0.2</b>	<b>95.9±0.4</b>
PangMovie	88.7±0.9	88.8±0.9	89.9±0.7	<b>93.5±0.4</b>	<b>93.1±0.4</b>	<b>93.7±0.5</b>
SemEval17	91.2±0.7	89.7±0.6	90.1±0.7	<b>92.7±0.6</b>	<b>92.7±0.6</b>	<b>93.5±0.3</b>
Sst	87.3±1	87.9±0.7	88.5±1	<b>90.7±0.9</b>	<b>90.9±0.8</b>	<b>91.1±1</b>
Sst2	94.6±0.2	91.4±0.4	93.5±0.4	<b>95.8±0.3</b>	<b>95.7±0.2</b>	<b>96±0.1</b>
Yelp2L	<b>97.9±0.5</b>	<b>98.6±0.3</b>	92.1±1	<b>98.7±0.3</b>	<b>98.5±0.6</b>	<b>98.5±0.5</b>
IMDB2024	97.6±1	96.5±1	93.9±1	<b>98.4±0.7</b>	<b>98.6±0.7</b>	<b>98.7±0.7</b>
RottenT2024	93.7±1.1	95.2±1.4	95.3±1	<b>96.3±0.4</b>	<b>96.6±0.7</b>	<b>96.7±0.4</b>
ACM	70.7±1.5	35.6±1.1	50.5±1.6	<b>74.6±2.7</b>	72.4±1.6	<b>76.6±2.1</b>
DBLP	81.9±0.7	53.7±0.8	53.2±1	<b>86.7±1</b>	85.9±0.8	<b>87.8±0.7</b>
Twitter	<b>77.5±2.7</b>	67.4±2.7	72.9±1.6	70.4±2.3	73.5±3.1	<b>77.7±2.5</b>
Webkb	<b>82.3±2.6</b>	41.9±1.5	64±1.8	83.4±1	82.4±2.1	<b>86±1.3</b>

Table 1: Average Macro-F1 and 95% confidence interval for SLMs and versions Llama 3.1 8B. Best results (including statistical ties) are marked in **bold**.

Dataset	RoBERTa	CMBS Zero-Shot	CMBS In-Context	CMBS Partially-Tuned-IS biO-IS	CMBS Partially-Tuned-IS Aggressive	Fully-Tuned LLaMA
Finance	<b>98.1±1.9</b>	<b>98±2.1</b>	<b>98.2±1.7</b>	<b>97.7±1.6</b>	<b>98.3±1.3</b>	<b>98.7±1.6</b>
Imdb	93±0.5	94±0.6	92.5±0.6	<b>95.7±0.3</b>	<b>95.8±0.2</b>	<b>95.9±0.4</b>
PangMovie	88.7±0.9	90.2±0.9	89.9±0.8	<b>93.5±0.4</b>	<b>93.1±0.3</b>	<b>93.7±0.5</b>
SemEval17	91.2±0.7	92.2±0.6	92±0.5	<b>92.9±0.5</b>	<b>92.9±0.6</b>	<b>93.5±0.3</b>
Sst	87.3±1	89±0.6	88.5±1.2	<b>90.8±0.9</b>	<b>90.9±0.9</b>	<b>91.1±1</b>
Sst2	94.6±0.2	95.1±0.2	94.8±0.3	<b>95.8±0.2</b>	<b>95.7±0.2</b>	<b>96±0.1</b>
Yelp2L	97.9±0.5	<b>98.5±0.2</b>	<b>98.1±0.2</b>	<b>98.8±0.2</b>	<b>98.6±0.5</b>	<b>98.5±0.5</b>
IMDB2024	97.6±1	<b>98.2±0.9</b>	97.3±1.2	<b>98.5±0.6</b>	<b>98.7±0.8</b>	<b>98.7±0.7</b>
RottenT2024	93.7±1.1	95.6±1	95.7±0.7	<b>96±0.7</b>	<b>96.3±0.7</b>	<b>96.7±0.4</b>
ACM	70.7±1.5	70.5±1.2	70.6±1.2	<b>74.7±2.7</b>	<b>73.3±2.4</b>	<b>76.6±2.1</b>
DBLP	81.9±0.7	81.9±0.6	82±1.6	<b>86.7±0.9</b>	86±0.8	<b>87.8±0.7</b>
Twitter	<b>77.5±2.7</b>	<b>79.4±2.7</b>	<b>78.7±2.5</b>	<b>77.7±2.2</b>	<b>78.2±1.8</b>	<b>77.7±2.5</b>
Webkb	<b>82.3±2.6</b>	82.1±2.3	<b>82.2±2.7</b>	<b>83.6±1.1</b>	<b>83.8±2.5</b>	<b>86±1.3</b>

Table 2: Average Macro-F1 and 95% confidence intervals for RoBERTa, versions of CMBS and Fully-Tuned LLaMA. Best results (including statistical ties) are marked in **bold**.

## 5.2 SLMs vs. LLMs - Computational Cost

Table 3 presents total time (in seconds) required to obtain final predictions for each solution. The Table shows that RoBERTa’s time is the shortest, followed by LLM Zero-Shot, which is around 76% more expensive than the SLM, on average. LLM In-context, in turn, is 176% slower than RoBERTa and 56% costlier than LLM Zero-Shot.

In Partially-Tuned-IS, cost increases significantly due to weight adjustment performed via LLM backpropagation. Aggressive version is always around 9%-35% faster than biO-IS due to its larger training set reduction, which impacts the fine-tuning time. Of course, fully-tuned LLM is the most expensive solution, which is 1700% more expensive than RoBERTa. With an average improvement of 3.3% across all datasets (peaking at 8.3% in ACM), it is not clear whether such improvements justify the significant computational cost increases. Depending on the application scenario, costly solutions may not be ideal or even feasible. To address this, our proposed solution reduce costs associated with using LLMs while preserving effectiveness gains.

Dataset	RoBERTa	Zero-Shot LLaMA	In-Context LLaMA	Partially-Tuned-IS biO-IS LLaMA	Partially-Tuned-IS Aggressive LLaMA	Fully-Tuned LLaMA
Finance	79	103	123	514	484	896
Imdb	2615	6295	11548	33554	25176	39257
PangMovie	934	1200	1490	7169	5892	10921
SemEval17	2416	3160	4251	19541	15154	28087
Sst	1027	1230	1562	7913	6544	11791
Sst2	5817	7800	10936	47813	37435	65428
Yelp2L	510	1161	1736	3743	2407	5116
IMDB2024	681	1623	2538	9015	5822	12304
RottenT2024	789	983	1708	5766	4393	8130
ACM	2665	3163	7896	18539	16877	28207
DBLP	4140	8113	17311	47501	27564	139250
Twitter	651	892	1478	7105	6664	11406
Webkb	910	2877	3274	12547	10150	26021

Table 3: Average Total Time for RoBERTa and versions of Llama3.1-8B.

Dataset	RoBERTa	CMBS Zero-Shot	CMBS In-Context	CMBS Partially-Tuned-IS biO-IS	CMBS Partially-Tuned-IS Aggressive	Fully-Tuned LLM
Finance	79	84	89	542	515	896
Imdb	2615	2930	3245	32813	25273	39257
PangMovie	934	994	1054	7386	6237	10921
SemEval17	2416	2574	2732	20004	16055	28087
Sst	1027	1089	1150	8149	6917	11791
Sst2	5817	6207	6597	48848	39508	65428
Yelp2L	510	568	626	3879	2676	5116
IMDB2024	681	762	844	8795	5921	12304
RottenT2024	789	838	887	5978	4743	8130
ACM	2665	2823	2981	19349	17854	28207
DBLP	4140	4546	4951	46891	28948	139250
Twitter	651	695	740	7045	6648	11406
Webkb	910	1054	1198	12202	10045	26021

Table 4: Average Total Time for RoBERTa and versions of CMBS and Fully-Tuned LLaMA.

## 5.3 RQ3: Evaluating the CBMS Approach

Focusing now on our proposals, we assess the four CBMS implementations: CMBS Zero-Shot, In-Context, Partially-Tuned-IS (bio-IS) and Partially-Tuned-IS (Aggressive). Starting with sentiment classification, Table 2 presents results for RoBERTa, each CMBS version, and Fully-Tuned LLaMA. CMBS Zero-Shot outperforms RoBERTa in 8 out of 9 sentiment datasets, tied only in Finance. These gains come with a small increase in computational cost over SLMs of only 8%. Moreover, in 4 of the 9 datasets, **CMBS Zero-Shot** ties with fully-tuned LLM, with minimal losses in others (on average, just 2% less effective). These excellent effectiveness results come at 10% of the fully-tuned cost, as demonstrated in Table 4, which presents total time results for all alternatives. Moreover, both **CMBS partially-tuned-IS** versions tie with fully-tuned LLaMA in *all* sentiment datasets at 30%-50% of the fully-tuned cost.

For topic classification, with more categories (up to 11) and uneven distributions, CBMS Zero-Shot and **CMBS In-Context** struggle with effectiveness. Significant gains over SLMs occur only with **CMBS Partially-Tuned-IS** versions. Among four topic datasets, both CBMS Partially-Tuned-IS

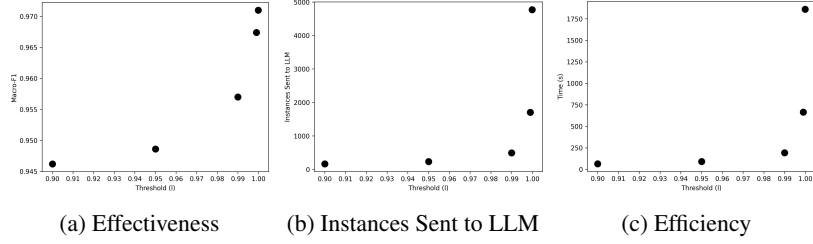


Figure 4: Effectiveness, Size of the Test Set Sent to LLM and Efficiency for IMDB dataset.

(Aggressive) and (bio-IS) outperform RoBERTa in two datasets, tying with the other two. Both also surpass partially-tuned LLaMA in all cases, with up to 6.4% gains in Twitter. Compared to fully-tuned LLaMA, CMBS Partially-Tuned-IS (bio-IS) achieves statistical equivalence in all datasets and (Aggressive) in three, having just a small deficit of around 2% in the fourth (DBLP). (Bio-IS) cuts computational costs by 40% on average while Aggressive achieves 50% of cost reduction.

If choosing between the two partially-tuned-IS CBMS versions, we would prefer Aggressive. Both show similar effectiveness in sentiment and topic datasets, but Aggressive offers a better effectiveness-cost tradeoff. Its main cost reduction comes from randomly cutting the training set by 50%. Surprisingly, its effectiveness matches state-of-the-art bio-IS when used within CBMS. However, remember that partially-tuned LLaMA (Aggressive), used within CBMS, only needs to handle the hardest cases. Outside CBMS, Bio-IS performs better (Table 1), but within our solution, Aggressive achieves the best effectiveness-cost balance.

Summarizing, for sentiments, best effectiveness tradeoff is achieved by CBMS Zero-Shot. If effectiveness is mandatory, the choice is CBMS Partially-Tuned-IS (Aggressive), which ties with LLaMA Fine-tuned at half the cost. For topics, the choice is also CBMS Partially-Tuned-IS (Aggressive), which ties with LLaMA fine-tuned in 3 out of four datasets, losing minimally (by 2%) in the fourth, being twice more efficient.

Finally, we calculated CO<sub>2</sub> emissions for obtaining final model predictions using Lacoste et al. (2019)’s methodology. Results in Table 14 (Appendix H) show emissions are significantly higher for LLMs, by orders of magnitude. Financial costs, analyzed in Appendix I, lead to similar conclusions.

#### 5.4 Confidence Threshold Sensitivity Analysis

We analyze confidence threshold impact. Figure 4 shows IMDB results, where CMBS Partially-tuned-IS (Aggressive) achieved a strong cost-benefit

tradeoff—outperforming RoBERTa and matching Fully-tuned LLaMA, at half the cost. Figures 4a, 4b, and 4c depict effectiveness gains, instances sent to LLM, and cost increases. Trends are similar across all graphs, despite different metrics. With appropriate thresholds, effectiveness can improve, though at a higher cost.

#### 5.5 CBMS Applied to Other NLP Tasks-CoLA GLUE

To demonstrate the potential of our solution for other NLP tasks, Appendix J shows the results of applying CBMS-Partially-Tuned-IS (Aggressive) to CoLA (Corpus of Linguistic Acceptability), a GLUE task assessing a model’s ability to classify English sentences as grammatically acceptable or unacceptable. As seen in Table 16, our solution outperforms RoBERTa and is statistically equivalent to Fully-Tuned at half of its computational cost.

### 6 Conclusion

We proposed Call-My-Big-Sibling (CMBS), an ATC solution that optimizes the effectiveness-cost trade-off by combining efficient, calibrated SLMs with more effective but costlier LLMs. CMBS uses LLMs only when the SLM is not confident, minimizing costs through zero-shot inference or instance selection when tuning is needed. Experiments on 13 sentiment and topic classification datasets confirmed its strong effectiveness-cost balance. CMBS Zero-Shot outperformed SLMs in 8 of 9 sentiment datasets with minimal overhead, while CMBS Partially-Tuned-IS matched fully tuned LLMs in sentiment tasks at half the cost. For topics, it surpassed partially tuned LLaMA and achieved near-equal effectiveness to fully tuned LLMs at twice the efficiency. Such efficiency gains are crucial in real-world scenarios. Future work includes extending CMBS to tasks like hate speech, irony detection, summarization, and Q&A, refining instance selection to reduce training needs, and exploring additional LLMs, configurations, and alternative prediction methodologies.



## 7 Limitations

Despite relevant contributions, our study has some limitations. Our current work covers only two classification tasks, which we have pursued to evaluate in depth. In this study, we used 13 datasets, 9 on sentiment analysis and 4 on topic classification, all with distinct characteristics. Additionally, we covered a single GLUE task on grammatical correctness to demonstrate generalization.

In spite of such limitations, there are only few studies - most of them covered in our related work - that analyze the complex issue of balancing effectiveness and the associated computational costs of NLP tasks in the Language Models context, particularly through the combination of LMs, as proposed in our paper. Indeed, this is a challenging trade-off to address: improving effectiveness often leads to higher costs. The recent success and real-world impact, including financial, of DeepSeek (DeepSeek-AI et al., 2025), which matched or surpassed the effectiveness of state-of-the-art pre-trained LLMs while reducing computational demands, highlights the importance of this topic (effectiveness vs. cost trade-off) to the community.

Although we have explored two (three if we consider GLUE CoLA) classification tasks - topic classification and sentiment analysis (and grammatical correctness) - we emphasize that ATC remains an open challenge in several NLP applications, such as hate speech detection, misinformation identification, and sentiment analysis—all of which can be framed as classification problems. Papers on these topics are recurrent in major ACL conferences, which attests to the relevance of this problem. Moreover, many classification tasks are imbalanced, posing significant challenges, particularly in multi-class classification. For instance, in the ACM dataset, which contains documents distributed across 11 imbalanced classes, the highest reported effectiveness is 76.6 Macro-F1, highlighting that this task remains challenging and far from being successfully solved.

We believe our proposal is comprehensive, as it evaluates four SLM methods, four LLMs - including four versions of LLaMA (zero-shot, in-context, partially tuned, and fully tuned) - and three versions of our solution across 13 different datasets. This extensive effort comprises more than 1000 experimental results (13 datasets x 16 methods x 5 runs). Nonetheless, more SLMs and LLMs, alone and in combination, will be tested in the future.

We focused our evaluation on open LLMs for the sake of reproducibility in subsequent research using our method. Hence we have refrained from using proprietary and closed-source LLMs, such as GPT, which operate as black boxes. This opacity poses challenges in understanding their training methodologies or internal structures, thereby obstructing reproducibility in research reliant on these models.

LLMs have been made available for different purposes. Some of these LLMs have high execution costs, such as Falcon 180B (Penedo et al., 2023), which requires an expensive infrastructure to use it. In our work, we limited our study to the best evaluated LLMs in the Hugging Face platform<sup>11</sup>, with around 7 billion parameters, which have a reasonable structure allowing us to evaluate zero-shot, in-context, partially and fully-tuned versions of our solutions.

Regarding, probability estimation, we employ a straightforward heuristic that exclusively considers the token with the maximum predicted probability. While this approach offers simplicity, we recognize its limitations and intend to explore alternative methodologies that incorporate a broader spectrum of high-probability tokens, thereby transcending the constraints of a singular, most likely prediction.

Finally, our work focused on applying our proposals with four open LLMs—Falcon 7B, Mistral 7B, DeepSeek 8B, and LLaMA 3.1 8B. However, new LLMs emerge all the time and it is hard to test all of them. In any case, we intend to explore new LLMs in future work. Given that these new LLMs tend to be increasingly more complex and costly, optimizing the cost-effectiveness of our combined approach using SLMs and LLMs will remain not only relevant but also an even more compelling objective.

## References

- Dimosthenis Antypas, Asahi Ushio, Jose Camacho-Collados, Leonardo Neves, Vitor Silva, and Francesco Barbieri. 2022. Twitter Topic Classification. In *Proceedings of the 29th International Conference on Computational Linguistics*, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Jun Bai, Xiaofeng Zhang, Chen Li, Hanhua Hong, Xi Xu, Chenghua Lin, and Wenge Rong. 2023a. [How to determine the most powerful pre-trained language model without brute force fine-tuning? an empirical](#)

<sup>11</sup><https://huggingface.co/models>



885	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris	<a href="#">treebank</a> . In <i>Proceedings of the 2013 Conference on</i>	944
886	Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian	<i>Empirical Methods in Natural Language Processing,</i>	945
887	Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Ku-	<i>EMNLP 2013, 18-21 October 2013, Grand Hyatt</i>	946
888	mar, Benjamin Newman, Binhang Yuan, Bobby Yan,	<i>Seattle, Seattle, Washington, USA, A meeting of SIG-</i>	947
889	Ce Zhang, Christian Alexander Cosgrove, Christo-	<i>DAT, a Special Interest Group of the ACL</i> , pages	948
890	pher D Manning, Christopher Re, Diana Acosta-	1631–1642. ACL.	949
891	Navas, Drew Arad Hudson, Eric Zelikman, Esin		
892	Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren,	Arthur Spirling. 2023. Why open-source generative	950
893	Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel	ai models are an ethical way forward for science.	951
894	Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun,	<i>Nature</i> , 616(7957):413–413.	952
895	Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar		
896	Khattab, Peter Henderson, Qian Huang, Ryan An-	Emma Strubell, Ananya Ganesh, and Andrew McCal-	953
897	drew Chi, Sang Michael Xie, Shibani Santurkar,	lum. 2019. Energy and policy considerations for deep	954
898	Surya Ganguli, Tatsunori Hashimoto, Thomas Icard,	learning in nlp. <i>arXiv preprint arXiv:1906.02243</i> .	955
899	Tianyi Zhang, Vishrav Chaudhary, William Wang,		
900	Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Ko-	Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang,	956
901	reeda. 2023. <a href="#">Holistic evaluation of language models</a> .	and Zhong Su. 2008. <a href="#">Arnetminer: Extraction and</a>	957
902	<i>Transactions on Machine Learning Research</i> . Fea-	<a href="#">mining of academic social networks</a> . KDD '08, page	958
903	tured Certification, Expert Certification.	990–998, New York, NY, USA. Association for Com-	959
		puting Machinery.	960
904	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,		
905	Dan Huang, Andrew Y. Ng, and Christopher Potts.	Felipe Viegas, Sergio Canuto, Washington Cunha, Celso	961
906	2011. <a href="#">Learning word vectors for sentiment analysis</a> .	França, Claudio Valiense, Leonardo Rocha, and Mar-	962
907	In <i>The 49th Annual Meeting of the Association for</i>	cos André Gonçalves. 2023. <a href="#">Clusent – combining se-</a>	963
908	<i>Computational Linguistics: Human Language Tech-</i>	<a href="#">mantic expansion and de-noising for dataset-oriented</a>	964
909	<i>nologies, Proceedings of the Conference, 19-24 June,</i>	<a href="#">sentiment analysis of short texts</a> . In <i>Proceedings of</i>	965
910	<i>2011, Portland, Oregon, USA</i> , pages 142–150. The	<i>the 29th Brazilian Symposium on Multimedia and the</i>	966
911	Association for Computer Linguistics.	<i>Web</i> , WebMedia '23, page 110–118, New York, NY,	967
		USA. Association for Computing Machinery.	968
912	Pekka Malo, Ankur Sinha, Pekka J. Korhonen, Jyrki		
913	Wallenius, and Pyry Takala. 2014. <a href="#">Good debt or bad</a>	J. Wolfe, X. Jin, T. Bahr, and N. Holzer. 2017. <a href="#">Ap-</a>	969
914	<a href="#">debt: Detecting semantic orientations in economic</a>	<a href="#">plication of softmax regression and its validation</a>	970
915	<a href="#">texts</a> . <i>J. Assoc. Inf. Sci. Technol.</i> , 65(4):782–796.	<a href="#">for spectral-based land cover mapping</a> . <i>The Inter-</i>	971
		<i>national Archives of the Photogrammetry, Remote</i>	972
916	Luiz Felipe Mendes, Marcos Gonçalves, Washington	<i>Sensing and Spatial Information Sciences</i> , XLII-	973
917	Cunha, Leonardo Rocha, Thierson Couto-Rosa, and	1/W1:455–459.	974
918	Wellington Martins. 2020. <a href="#">"keep it simple, lazy" –</a>		
919	<a href="#">metalyze: A new metastrategy for lazy text classifica-</a>	Canwen Xu, Yichong Xu, Shuohang Wang, Yang Liu,	975
920	<a href="#">tion</a> . In <i>Proceedings of the 29th ACM International</i>	Chenguang Zhu, and Julian McAuley. 2024. <a href="#">Small</a>	976
921	<i>Conference on Information &amp; Knowledge Manage-</i>	<a href="#">models are valuable plug-ins for large language mod-</a>	977
922	<i>ment</i> , CIKM '20, page 1125–1134, New York, NY,	<a href="#">els</a> . In <i>Findings of the Association for Computational</i>	978
923	USA. Association for Computing Machinery.	<i>Linguistics: ACL 2024</i> , pages 283–294, Bangkok,	979
		Thailand. Association for Computational Linguistics.	980
924	Bo Pang and Lillian Lee. 2005. <a href="#">Seeing stars: Exploit-</a>		
925	<a href="#">ing class relationships for sentiment categorization</a>	Murong Yue, Jie Zhao, Min Zhang, Liang Du, and Ziyu	981
926	<a href="#">with respect to rating scales</a> . In <i>Proceedings of the</i>	Yao. 2024. <a href="#">Large language model cascades with</a>	982
927	<i>43rd Annual Meeting of the Association for Compu-</i>	<a href="#">mixture of thought representations for cost-efficient</a>	983
928	<i>tational Linguistics (ACL'05)</i> , pages 115–124, Ann	<a href="#">reasoning</a> . In <i>The Twelfth International Conference</i>	984
929	Arbor, Michigan. Association for Computational Lin-	<i>on Learning Representations</i> .	985
930	guistics.		
931	Guilherme Penedo, Quentin Malartic, Daniel Hesslow,		
932	Ruxandra Cojocaru, Alessandro Cappelli, Hamza	Bruna Stella Zanotto, Ana Paula Beck da Silva Etges,	986
933	Alobeidli, Baptiste Pannier, Ebtesam Almazrouei,	Avner Dal Bosco, Eduardo Gabriel Cortes, Renata	987
934	and Julien Launay. 2023. <a href="#">The refinedweb dataset</a>	Ruschel, Ana Claudia De Souza, Claudio MV An-	988
935	<a href="#">for falcon LLM: outperforming curated corpora with</a>	drade, Felipe Viegas, Sergio Canuto, Washington	989
936	<a href="#">web data, and web data only</a> . <i>CoRR</i> , abs/2306.01116.	Luiz, et al. 2021. Stroke outcome measurements	990
		from electronic medical records: cross-sectional	991
937	Sara Rosenthal, Noura Farra, and Preslav Nakov. 2019.	study on the effectiveness of neural and nonneural	992
938	<a href="#">Semeval-2017 task 4: Sentiment analysis in twitter</a> .	classifiers. <i>JMIR Medical Informatics</i> , 9(11):e29120.	993
939	<i>CoRR</i> , abs/1912.00741.		
940	Richard Socher, Alex Perelygin, Jean Wu, Jason		
941	Chuang, Christopher D. Manning, Andrew Y. Ng,		
942	and Christopher Potts. 2013. <a href="#">Recursive deep mod-</a>		
943	<a href="#">els for semantic compositionality over a sentiment</a>		



## A Prompt Templates

Below are examples of the structure of the prompts we used for our experiments’ Zero-Shot and In-Context versions of the LLMs. They all draw on (Liang et al., 2023). Table 5 provides the prompt used for sentiment classification, while Table 6 presents the prompt for topic classification. Both prompts include the tag [Evaluate Text], which represents the (test) text to be classified, and the tag [Response from LLM], which contains the model’s output. If the model’s output does not match any of the given alternatives (due to hallucination), we predict the majority class from the training set.

We employ a straightforward heuristic that exclusively considers the token with the maximum predicted probability. While this approach offers simplicity, we recognize its limitations and intend to explore alternative methodologies that incorporate a broader spectrum of high-probability tokens, thereby transcending the constraints of a singular, most likely prediction.

The third example of prompt, shown in Table 7, is tailored for in-context learning. For the evaluated test document, “I spent a day at a 5-star hotel, which was amazing.” the most similar example from the training set included in the prompt was “5-star hotels have many food options.”. A vector representation is generated for each evaluated example using the fully-tuned RoBERTa as an encoder. By comparing the vector of the evaluated (test) document with the vectors of the training set documents, we identify the most similar document based on the cosine similarity between the vectors and use it as a training example in the prompt.

---

Classify the sentiment in the text exclusively as positive or negative:  
Input: I love you.  
Reference:  
A. Positive  
B. Negative  
Answer: A  
Input: The product is bad.  
Reference:  
A. Positive  
B. Negative  
Answer: B  
Input: {**Evaluate Text**}  
Reference:  
A. Positive  
B. Negative  
Answer: {**Response from LLM**}

---

Table 5: Prompt template for sentiment classification.

---

Classify the topic of the text exclusively with one of the references:  
Input: Messi scored a goal against France.  
Reference:  
A. Pop culture  
B. Sports or gaming  
C. Daily life  
D. Science or technology  
E. Business or entrepreneurs  
F. Arts or culture  
Answer: B  
Input: {**Evaluate Text**}  
Reference:  
A. Pop culture  
B. Sports or gaming  
C. Daily life  
D. Science or technology  
E. Business or entrepreneurs  
F. Arts or culture  
Answer: {**Response from LLM**}

---

Table 6: Prompt template for topic classification.

---

Classify the sentiment in the text exclusively as positive or negative:  
Input: 5-star hotels have many food options.  
Reference:  
A. Positive  
B. Negative  
Answer: A  
Input: I spent a day at a 5-star hotel, which was amazing.  
Reference:  
A. Positive  
B. Negative  
Answer: {**Response from LLM**}

---

Table 7: Prompt template for sentiment classification for In-Context Llama and CMBS In-Context.

## B Datasets

Our study draws on **thirteen** datasets developed for topic and sentiment classification. Our choice was strategically purposeful due to the effort to perform an in-depth analysis of this task. The datasets include **Finance** (Malo et al., 2014) focusing on economic news, **IMDB** (Maas et al., 2011)<sup>12</sup> compiling movie reviews as well as **PangMovie** (Pang and Lee, 2005) including Rotten Tomatoes<sup>13</sup> data, **SemEval17** (Rosenthal et al., 2019) containing Twitter texts used in a significant text classification challenge, and the Stanford Sentiment Treebank (SST) (Socher et al., 2013) and **SST2** (Socher et al., 2013), where sentiment classification relies on a *treebank*, a corpus with sentiment labels and labeled parse trees. **Yelp Review** is a subset of Yelp data widely used in sentiment classification studies

<sup>12</sup><https://www.imdb.com/>

<sup>13</sup><https://www.rottentomatoes.com/>

(Canuto et al., 2016; Viegas et al., 2023; Mendes et al., 2020). **IMDB2024** and **RottenT2024** were collected to avoid data contamination by LLM. For topic classification, we have ACM Digital Library (Cunha et al., 2021b), DBLP (Tang et al., 2008), Twitter Topic (Antypas et al., 2022) and WebKB (Craven et al., 1998).

	Dataset	Domain	Idl	Avg Words	Classes	Minor Class	Major Class
Sentiment	Finance	Finance	873	24.88	2	303	570
	IMDB	Movie	24904	234	2	12432	12472
	PangMovie	Movie	10662	21.02	2	5331	5331
	SemEval17	Twitter	27413	19.85	2	7745	19668
	Sst	Movie	11841	19.18	2	5905	5936
	Sst2	Movie	66973	10.45	2	29643	37330
	Yelp2L	Place	4995	131.8	2	2495	2500
	IMDB2024	Movie	6572	163.02	2	2057	4515
	RottenT2024	Movie	7948	46.13	2	3315	4633
	Acem	Article	24897	63.52	11	63	6562
Topic	Dblp	Article	38128	141.43	10	1414	9746
	Twitter	Twitter	6997	28.68	6	152	2738
	Webkb	Pages	8199	208.81	7	137	3705

Table 8: Datasets Statistics.

As detailed in Table 8, we can observe an ample diversity in many aspects of these datasets: domain, number of documents (Idl), density (the average number of words per document), etc.

## C Evaluating LLMs

We evaluate four LLMs in Zero-shot mode: Falcon 7B, Mistral 7B, LLaMA 3.1-8B and DeepSeek 8B. Table 9 presents the results regarding Macro-F1, with the best outcomes highlighted in bold. As observed, LLaMA is consistently the best performer, either alone or tied with Mistral, across **all** datasets, except WebKB.

Dataset	Falcon 7B	Mistral 7B	Llama 3.1 8B	DeepSeek 8B
Finance	46.7±4.8	<b>94.3±1.9</b>	<b>95.4±1.2</b>	<b>95.6±2.1</b>
Imdb	68.4±0.7	68.4±0.7	<b>93±0.3</b>	84.2±0.5
PangMovie	43.6±0.5	82.3±0.9	<b>88.8±0.9</b>	82.1±0.5
SemEval17	54.4±0.6	81±0.9	<b>89.7±0.6</b>	87.2±0.9
Sst	47±1.2	82±0.8	<b>87.9±0.7</b>	82.3±0.8
Sst2	38.6±0.1	86.2±0.5	<b>91.4±0.4</b>	84.4±0.6
Yelp2L	79.9±1.3	96.2±0.9	<b>98.6±0.3</b>	96.4±0.6
IMDB2024	78.4±0.8	94.9±0.9	<b>96.5±1</b>	94.6±0.8
RottenT2024	65.8±1.3	93.8±1.2	<b>95.3±1</b>	92.1±1.0
ACM	2.6±0.2	18.2±0.9	<b>35.6±1.1</b>	14.2±0.4
DBLP	3.1±0.2	50.2±0.6	<b>53.7±0.8</b>	39.2±0.5
Twitter	13±0.3	<b>62.2±2.1</b>	<b>63.5±1.7</b>	<b>63.6±1.2</b>
Webkb	3.8±0.3	<b>42.1±0.6</b>	37±2.1	24.1±3.5

Table 9: Effectiveness in Macro-F1 for sentiment and topic classification tasks with the LLMs in Zero-shot version, Falcon 7B, Mistral 7B, and Llama 3.1 8B. Best results (including statistical ties) marked in **bold**.

## D Evaluating SLMs

We evaluate four fully-tuned SLMs: BART, BERT, Qwen 0.5B and RoBERTa. Table 10 presents the results regarding Macro-F1, with the best outcomes

highlighted in bold. As observed, RoBERTa is consistently the best performer, alone or tied with another SLM, across **all** datasets, with no exception.

Dataset	BERT	BART	Qwen 0.5B	RoBERTa
Finance	94.1±3.8	<b>97±1.7</b>	67.8±14.2	<b>98.1±1.9</b>
Imdb	91.7±0.4	<b>92.8±0.4</b>	81.6±1.3	<b>93±0.5</b>
PangMovie	87.5±0.7	<b>88.4±1</b>	72.5±1.8	<b>88.7±0.9</b>
SemEval17	90.3±0.3	<b>91±0.4</b>	79.9±0.7	<b>91.2±0.7</b>
Sst	86.1±0.4	<b>87.7±1.1</b>	69.5±8.7	<b>87.3±1</b>
Sst2	<b>94.8±0.1</b>	94.2±0.3	87.9±0.7	<b>94.6±0.2</b>
Yelp2L	96.8±0.4	<b>97.7±0.2</b>	89.4±3.7	<b>97.9±0.5</b>
IMDB2024	96.6±0.5	<b>97.5±0.6</b>	90.2±2.2	<b>97.6±1</b>
RottenT2024	92.5±1	<b>93.5±0.5</b>	84.6±0.8	<b>93.7±1.1</b>
ACM	<b>69.8±1.8</b>	68±2.8	46.4±13.7	<b>70.7±1.5</b>
DBLP	<b>82.1±0.9</b>	<b>81.9±0.6</b>	68.5±8.4	<b>81.9±0.7</b>
Twitter	<b>76.6±4.4</b>	<b>76.9±3.3</b>	39.1±5.7	<b>77.5±2.7</b>
Webkb	<b>80.8±3.8</b>	<b>81.7±3.5</b>	45.9±13.6	<b>82.3±2.6</b>

Table 10: Average Macro-F1 and 95% confidence interval for SLMs . Best results (including statistical ties) marked in **bold**.

## E Evaluating Aggressive Reduction in the LLMs with Different Reduction Levels

As mentioned, fine-tuning is essential for LLM effectiveness. Here, we illustrate the impact of training data size on the LLM effectiveness using the validation set and a sample of two datasets. The pattern of results is basically the same in all other datasets we experimented with.

Table 11 presents the effectiveness results when utilizing 30%, 50%, and 70% of the training data in Twitter and WebKB, two topic datasets in which CMBS performs very well. As we can see in the Table, 30% of training generally is not enough for achieving reasonable effectiveness, while the improvements of using 70% are either marginal or incur in higher costs.

As discussed in Section 5, the CMBS Partially-Tuned-IS (Aggressive) version we employed in our experiments uses 50% of the training data, randomly selected in a stratified manner, based on results of instance selection experiments (Cunha et al., 2023a). In all datasets, such a choice produced the best tradeoff between effectiveness and computational cost.

Dataset	Portion Train	Macro-F1
Twitter	30	66.2
Twitter	50	71.9
Twitter	70	76.1
Webkb	30	76.7
Webkb	50	83.4
Webkb	70	85.2

Table 11: Evaluate amount training LLM.

## F Evaluating Threshold L

We evaluate the impact of the parameter L, which determines the number of documents sent to the LLM. The higher the value of L, the more documents fall below the threshold, increasing the number of documents forwarded to the LLM. Table 12 presents this evaluation for a sample of four datasets, showing the dataset name, the percentage of instances sent to the LLM relative to the total of test instances, and the effectiveness of both the SLM and LLM on this subset of instances. It is interesting to notice that the choice of L that maximizes the cost-effectiveness threshold (around 0.9) is similar in all datasets, and that the LLM effectiveness in these hard-to-classify instances is better than the SLM, which justifies the CBMS Partially-Tuned-IS Aggressive gains.

Dataset	Percentage of Instances	SLM Macro-F1	LLM Macro-F1	Threshold (L)
Sst	27.0	0.65	0.76	0.9
Sst2	25.3	0.82	0.85	0.9
IMDB2024	7.8	0.72	0.87	0.9
Webkb	13.9	0.56	0.67	0.9
Twitter	13.9	0.51	0.53	0.9

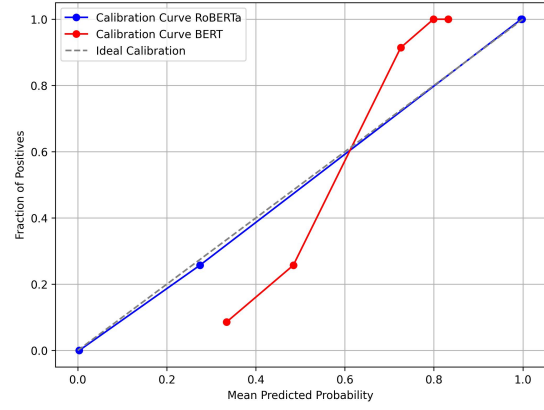
Table 12: Evaluation Threshold L.

## G Calibration Model

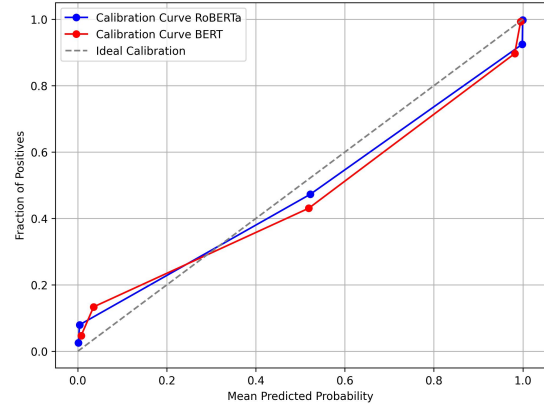
We analyze the behavior of the Transformer calibrations using Brier Score (BS) (BRIER, 1950), a scoring rule applied to measure the accuracy of probabilistic predictions. Brier (BRIER, 1950) defines  $BS = \frac{1}{n} \sum_{i=1}^n \sum_{c=1}^C (P(Y = y_c | x_i) - o_{ci})^2$ , where  $o_{ci}$  is the one-hot vector with a value 1 in the index of the true class of  $x_i$ , and 0 otherwise. BS ranges from 0 (worst) to 1 (best) – the closer to one, the better in achieving more calibrated probability estimations. Table 13 presents the obtained Brier Score averaged across the datasets by applying BERT and RoBERTa classifiers on all datasets considered in our experiments. Meanwhile, Figure 5 presents the calibration curves for BERT and RoBERTa, reinforcing that RoBERTa is a better-calibrated model compared to BERT.

## H CO<sub>2</sub> emissions

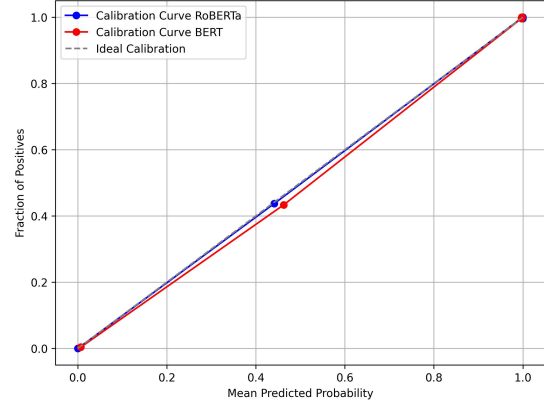
We calculated the CO<sub>2</sub> emissions associated with the execution of the model using the methodology developed by Lacoste et al. (2019). It is possible to associate the value of emission 0.14 kg of CO<sub>2</sub>eq per hour with a machine of similar structure to



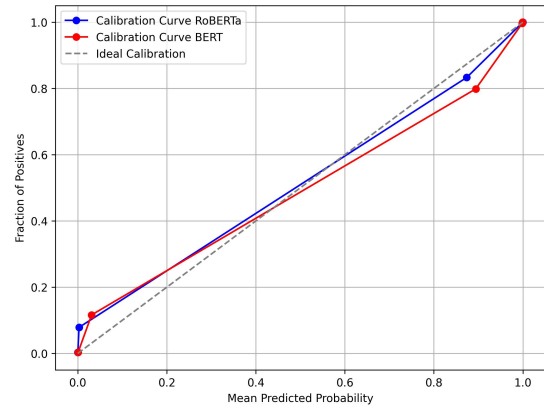
(a) Finance



(b) Pang Movie



(c) IMDB2024



(d) RottenT2024

Figure 5: Calibration curve for BERT and RoBERTa.



Dataset	BERT	RoBERTa
Finance	0.784	0.989
IMDB	0.854	0.873
PangMovie	0.793	0.804
SemEval17	0.879	0.887
Sst	0.783	0.792
Sst2	0.919	0.912
Yelp2L	0.968	0.968
IMDB2024	0.944	0.959
RottenT2024	0.898	0.911
ACM	0.691	0.662
DBLP	0.765	0.747
Twitter	0.810	0.828
Webkb	0.779	0.80

Table 13: Brier score for BERT and RoBERTa.

the one used in our experiments<sup>14</sup>. The emission values are presented in Table 14.

## I Financial Cost

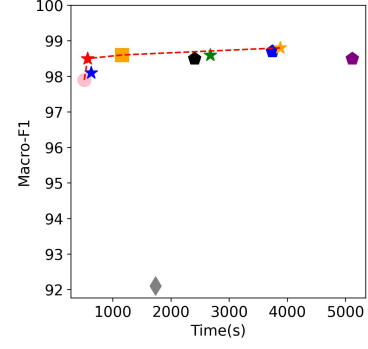
Following (Griggs et al., 2024), Table 15 presents the financial cost in dollars for executing the main methods discussed in this paper. As a reference, we used the hourly price of a setup similar to the one used in this research<sup>15</sup>, offered by a large cloud company, which currently costs \$0.752 per hour.

## J GLUE CoLA

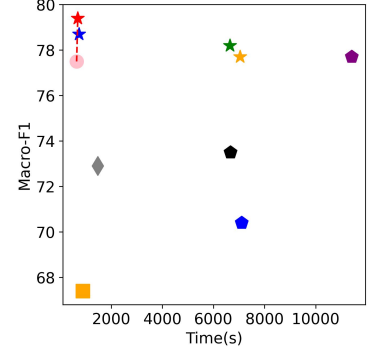
To demonstrate the potential of our solutions for other NLP tasks, we present the application of the best version of CBMS (CBMS Partially-Tuned (Aggressive) to CoLA (Corpus of Linguistic Acceptability), a GLUE task that evaluates a model’s ability to predict the grammaticality of an English sentence by identifying it as either grammatically acceptable or unacceptable. For example, the sentence “The cat sat on the mat” is grammatically correct, whereas the sentence “On the mat sat cat the” is grammatically incorrect. Table 16 presents the results. As it can be observed, our solution outperforms RoBERTa and is statistically equivalent to Fully-Tuned at approximately half the computational cost.

## K Pareto-optimal

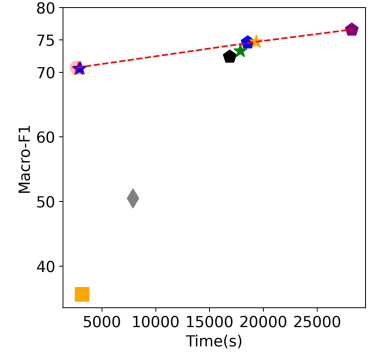
Figure 6 presents the graph of the effectiveness-efficiency trade-off across all methods evaluated in this study. The Pareto frontier is represented by the dashed red line, with all points overlapping this line being Pareto-optimal. We can observe that CMBS is the most frequent method on the Pareto frontier.



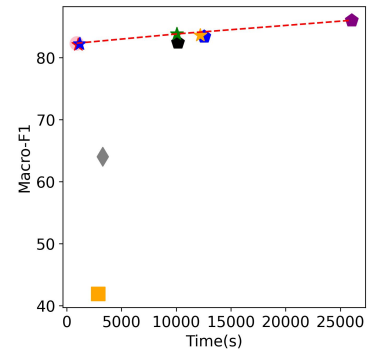
(a) Yelp2L



(b) Twitter



(c) ACM



(d) Webkb

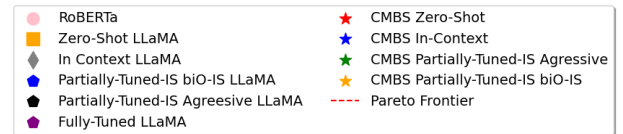


Figure 6: Total time (seconds) and Macro-F1 for RoBERTa, different versions of LLaMA, and CMBS. The Pareto frontier is represented by a dashed red line, with all points overlapping this line being Pareto-optimal.

<sup>14</sup><https://mlco2.github.io/impact/#co2eq>

<sup>15</sup>[aws.amazon.com/ec2/instance-types/g4/](https://aws.amazon.com/ec2/instance-types/g4/)

Dataset	RoBERTa	Zero-Shot LLaMA	In-Context LLaMA	Partially-Tuned-IS biO-IS LLaMA	Partially-Tuned-IS Agressive LLaMA	CMBS Zero-Shot	CMBS In-Context	CMBS Partially-Tuned-IS biO-IS	CMBS Partially-Tuned-IS Agressive	Fully-Tuned LLaMA
Finance	0.02	0.02	0.02	0.1	0.09	0.02	0.02	0.11	0.1	0.17
Imdb	0.51	1.22	2.25	6.52	4.9	0.57	0.63	6.38	4.91	7.63
PangMovie	0.18	0.23	0.29	1.39	1.15	0.19	0.21	1.44	1.21	2.12
SemEval17	0.47	0.61	0.83	3.8	2.95	0.5	0.53	3.89	3.12	5.46
Sst	0.2	0.24	0.3	1.54	1.27	0.21	0.22	1.58	1.34	2.29
Sst2	1.13	1.52	2.13	9.3	7.28	1.21	1.28	9.5	7.68	12.72
Yelp2L	0.1	0.23	0.34	0.73	0.47	0.11	0.12	0.75	0.52	0.99
IMDB2024	0.13	0.32	0.49	1.75	1.13	0.15	0.16	1.71	1.15	2.39
RottenT2024	0.15	0.19	0.33	1.12	0.85	0.16	0.17	1.16	0.92	1.58
ACM	0.52	0.62	1.54	3.6	3.28	0.55	0.58	3.76	3.47	5.48
DBLP	0.81	1.58	3.37	9.24	5.36	0.88	0.96	9.12	5.63	27.08
Twitter	0.13	0.17	0.29	1.38	1.3	0.14	0.14	1.37	1.29	2.22
Webkb	0.18	0.56	0.64	2.44	1.97	0.2	0.23	2.37	1.95	5.06

Table 14: Emission CO<sub>2</sub>. Calculation based on the work of [Lacoste et al. \(2019\)](#).

Dataset	RoBERTa	Zero-Shot LLaMA	In-Context LLaMA	Partially-Tuned-IS biO-IS LLaMA	Partially-Tuned-IS Agressive LLaMA	CMBS Zero-Shot	CMBS In-Context	CMBS Partially-Tuned-IS biO-IS	CMBS Partially-Tuned-IS Agressive	Fully-Tuned LLaMA
Finance	0.08	0.11	0.13	0.54	0.51	0.09	0.09	0.57	0.54	0.94
Imdb	2.73	6.57	12.06	35.05	26.29	3.06	3.39	34.27	26.4	41
PangMovie	0.98	1.25	1.56	7.49	6.15	1.04	1.1	7.71	6.51	11.41
SemEval17	2.52	3.3	4.44	20.41	15.83	2.69	2.85	20.89	16.77	29.33
Sst	1.07	1.28	1.63	8.26	6.84	1.14	1.2	8.51	7.22	12.31
Sst2	6.08	8.15	11.42	49.94	39.1	6.48	6.89	51.02	41.26	68.34
Yelp2L	0.53	1.21	1.81	3.91	2.51	0.59	0.65	4.05	2.8	5.34
IMDB2024	0.71	1.7	2.65	9.42	6.08	0.8	0.88	9.19	6.18	12.85
RottenT2024	0.82	1.03	1.78	6.02	4.59	0.88	0.93	6.24	4.95	8.49
ACM	2.78	3.3	8.25	19.36	17.63	2.95	3.11	20.21	18.65	29.46
DBLP	4.32	8.47	18.08	49.61	28.79	4.75	5.17	48.97	30.23	145.44
Twitter	0.68	0.93	1.54	7.42	6.96	0.73	0.77	7.36	6.94	11.91
Webkb	0.95	3.01	3.42	13.1	10.6	1.1	1.25	12.74	10.49	27.18

Table 15: Finance Cost in dollars (\$) for RoBERTa, Zero-Shot LLaMA, In-Context LLaMA, Partially-Tuned LLaMA, CMBS Zero-Shot, CMBS In-Context, CMBS Partially-Tuned, and Fully-Tuned LLaMA.

Dataset	RoBERTa	CMBS Partially-Tuned-IS Agressive	Fully-Tuned LLM
GlueCola	80.4±1.8	<b>83.4±2.1</b>	<b>84.5±0.4</b>

Table 16: Average Macro-F1 and 95% confidence interval RoBERTa, CMBS Partially-Tuned-IS Agressive and Fully-Tuned LLaMA. Best results (including statistical ties) marked in **bold**.