

# Ordinal PCL+: A Prompt-based Multilingual Framework for Few-shot and Zero-shot Cross-lingual Text Classification

Zhengyuan Chen

Faculty of Information Science and Engineering, Ocean University of China, 266100 Qingdao  
Shandong, China

czy8992@stu.ouc.edu.cn

**Abstract.** This paper proposes Ordinal Prompt based Cross-lingual Learning plus (Ordinal PCL+), a prompt-based multilingual framework for few-shot and zero-shot text classification. The approach inserts learnable continuous prompts, conditions token representations with a lightweight label-attention block, and performs a two-pass self-learning routine that couples pseudo-labels with representation–prototype agreement. For tasks with ordered labels, an ordinal head enforces monotonic consistency; a contrastive objective further aligns same-class instances across languages. The overall loss combines classification, agreement, ordinal, and contrastive terms and is followed by post-hoc temperature scaling to ensure well-calibrated probabilities. Evaluation on the Multilingual Amazon Reviews Corpus (MARC) uses English  $k$ -shot supervision ( $k \in \{4, 8, 16, 32, 64\}$ ) with zero-shot transfer to five target languages, reporting accuracy and expected calibration error (ECE) across multiple seeds. Relative to fine-tuning and prompt-only variants, the framework consistently improves transfer while reducing calibration error and remains modular and lightweight for practical deployment on commodity hardware. The components are ablated to quantify their contributions, showing complementary gains from label attention and self-learning, consistent benefits of the ordinal branch on rating-like labels, and improved reliability after temperature scaling.

**Keywords:** Cross-lingual Text Classification, Prompt-based Learning, Ordinal Consistency, Contrastive Learning, Probability Calibration.

## 1. Introduction

As a result of the explosive expansion of digital content and escalating demands for automated language handling, natural language processing (NLP) has emerged as a cornerstone of research and innovation. However, linguistic diversity worldwide poses a significant barrier to seamless cross-lingual communication, information exchange, and knowledge propagation—making it a primary challenge in the field. Addressing this challenge has driven extensive research into cross-lingual representations and the cross-lingual text classification (CLTC) task. At its core, CLTC involves training models on source-language data and leveraging the learned patterns for effective generalization in target languages.

Recent advancements in pre-trained language models (PLMs), particularly large language models (LLMs) based on the foundational Transformer architecture, have revolutionized text classification, including cross-lingual tasks, by capturing intricate semantic representations from massive unlabeled datasets [1]. Models such as Bidirectional Encoder Representations from Transformers (BERT) the GPT series, and advanced variants like Llama3, Mistral, and Gemma enable efficient transfer learning through self-supervised objectives, including masked language modeling and next sentence prediction, enabling fine-tuning or prompt-based adaptation strategies for downstream tasks with limited labeled data [2-7]. Recent advancements have highlighted parameter-efficient techniques, such as quantization (e.g., 4-bit and 8-bit variants using Quantized Low-Rank Adaptation (QLoRA) and Activation-aware Weight Quantization (AWQ)) and low-rank adaptation (LoRA), which reduce computational overhead while maintaining high performance, though full fine-tuning at 16-bit precision often yields superior accuracy compared to quantized alternatives [8]. In cross-lingual scenarios, contrastive learning and label-aware data augmentation enhance semantic alignment across languages, mitigating data scarcity in low-resource settings [9, 10]. Prompting strategies, including zero-shot, few-shot, chain-of-thought (CoT), and role-playing, have gained prominence, for example,

models like Llama3 70B achieving F1-scores of 94.4% in binary tasks and 87.1% in multiclass classification when combined with hybrid prompts [11]. Such innovations—including frameworks like PCL that integrate language-agnostic continuous prompts with self-learning processes and label attention modules for zero-shot CLTC, outperform traditional machine learning methods that rely on manual feature engineering, such as Support Vector Machines (SVM) and Naive Bayes, by leveraging semantic similarity retrieval and instruction-tuning to achieve better generalization across linguistic boundaries [12, 13]. Furthermore, in specialized domains like cross-lingual aspect-based sentiment analysis (ABSA), multilingual PLMs facilitate knowledge transfer from resource-rich to low-resource languages through fusion techniques involving machine translation, alignment-free label projection, and contrastive learning, overcoming challenges in extracting sentiment elements across diverse linguistic and cultural contexts [14].

Despite rapid progress in cross-lingual text classification, two practical gaps remain in low-resource regimes: (1) how to couple prompt representations with label semantics—especially when labels are ordered (ratings)—and (2) how to produce reliable probabilities rather than only accurate top-1 predictions. Prior prompt methods tend to ignore ordinal structure or rely on heavier multilingual training; calibration is rarely reported and often overlooked under multilingual transfer. This paper study these questions under a strict English-only supervision protocol and evaluate zero-shot transfer to multiple target languages. In summary, this paper makes the following main contributions.

A modular prompt framework for cross-lingual transfer. This paper introduces Ordinal PCL+, a lightweight composition of (1) label attention to strengthen prompt–label coupling, (2) a two-pass self-learning routine to stabilize training with pseudo-labels, (3) an ordinal head for ordered label spaces that reduces adjacent-class flips, and (4) a contrastive term to align same-class instances across languages. The loss integrates classification, agreement, ordinal, and contrastive components, and is followed by post-hoc temperature scaling for calibrated probabilities.

Protocol and calibration-aware reporting. This paper train on English with  $k$ -shot supervision  $k \in \{4, 8, 16, 32, 64\}$  and evaluate zero-shot on five target languages (De/Es/Fr/Ja/Zh). This paper report Accuracy and ECE across multiple seeds; temperature is fitted on the development split and applied at test time, and this paper include reliability plots.

Empirical findings across  $k$ . At very small budgets ( $k = 4\sim 8$ ), prompt-only is competitive; as supervision increases ( $k \geq 16$ ), Ordinal PCL+ yields the strongest averages—especially on Ja/Zh—while plain PCL also improves and surpasses prompt-only for  $k \geq 32$ . Temperature scaling consistently reduces ECE with negligible accuracy impact, improving probability reliability.

Component-wise insights. Controlled ablations show that label attention helps when supervision is adequate, self-learning reduces across-seed variance, the ordinal head stabilizes rating-style labels and improves calibration, and contrast tightens cross-lingual alignment but benefits from calibration to correct over-confidence.

Roadmap. §2 formalizes the task and data; §3 describes Ordinal PCL+ and each module; §4 details training and evaluation, including temperature scaling and metrics; §4.4 reports main multilingual results across  $k$ , and §4.5 presents ablations and reliability analyses; §5 discusses related work; §6 concludes.

## 2. Related work

CLTC is central to multilingual NLP because it enables knowledge transfer from resource-rich to low-resource languages while reducing reliance on extensive manual annotation. Early CLTC pipelines depended on conventional classifiers such as Support Vector Machines and Naïve Bayes with handcrafted features and translation-based heuristics; these systems struggled to generalize across domains and languages due to vocabulary shift and error propagation from external tools [13]. Modern work reframes CLTC as a representation learning problem, where models pre-trained on large corpora are adapted to downstream tasks with minimal supervision.

## 2.1. Transformer-based pre-trained language models

The Transformer architecture introduced parallel self-attention and long-range dependency modeling, laying the foundation for today’s multilingual PLMs [1]. Contextual encoders such as BERT and decoder-style GPT models learn rich semantic features from massive unlabeled text and can be adapted via fine-tuning or lightweight prompting [2, 3]. Surveys indicate that task-specific fine-tuning with modest amounts of labeled data often yields strong cross-lingual transfer performance, especially when coupled with multilingual corpora and shared subword vocabularies [7]. To reduce compute, parameter-efficient methods—including LoRA-style low-rank adapters and quantization schemes such as QLoRA or AWQ—allow training or inference on limited hardware, though full 16-bit fine-tuning still tends to deliver the best absolute accuracy [8].

## 2.2. Cross-lingual alignment and representation learning

A persistent challenge in CLTC is aligning semantics across languages without abundant parallel data. Recent approaches adopt contrastive objectives that bring translations or semantically equivalent sentences closer in the embedding space, improving robustness to lexical variation and domain shift [9]. Complementary label-aware augmentation synthesizes examples guided by class semantics, improving decision boundaries and alleviating skewed label distributions in low-resource settings [10]. Compared with pure translation pipelines, alignment-first methods reduce exposure to machine translation errors and encourage language-agnostic features that generalize better across new domains and unseen languages.

## 2.3. Prompt-based and instruction-tuned approaches

LLMs further enable prompt-based CLTC without heavy parameter updates. Zero-shot and few-shot prompting convert classification into natural-language instructions, while chain-of-thought and role-playing prompts elicit intermediate reasoning that improves label consistency on harder examples. Empirical studies show that hybrid prompts can deliver competitive F1 in both binary and multiclass regimes—for example, reports for Llama3 70B cite around 94% (binary) and 87% (multiclass) under carefully designed prompts [11]. Beyond discrete templates, continuous prompt frameworks integrate soft prompts with self-learning and label-attention modules to enhance zero-shot cross-lingual transfer, reducing dependence on task-specific fine-tuning [12]. These strategies complement supervised fine-tuning and are particularly attractive when annotation budgets are limited.

## 2.4. Domain-specific extensions and multimodal integration

CLTC techniques have been naturally extended to domain-specific applications, such as ABSA, where models must identify aspect terms and corresponding sentiments expressed in diverse cultures and genres. Multilingual PLMs can fuse machine translation, alignment-free projection, and contrastive objectives to transfer knowledge from high-resource to low-resource settings, yielding robust performance across languages and domains [14]. In parallel, multimodal variants combine text with signals such as images or sign-language video, where cross-lingual contrastive learning links visual context with multilingual text to improve retrieval and classification quality [9]. These extensions illustrate the breadth of CLTC beyond general-purpose benchmarks.

## 2.5. Knowledge distillation and model compression

For practical deployment, efficiency is paramount. Knowledge distillation transfers predictions or intermediate representations from a large teacher to a compact student, preserving multilingual competence while reducing latency and memory footprint. Combined with quantization and pruning, distillation makes CLTC feasible on edge devices and in latency-critical services. Recent studies on rank-adaptive adapters and low-bit quantization show that careful design can recover much of the accuracy gap versus full-precision models while enabling cost-effective training and inference [8].

In summary, the field has progressed from feature-engineered pipelines to versatile PLM- and LLM-based frameworks that integrate alignment objectives, prompt-based reasoning, domain-specific fusion, and efficient adaptation. Together, these advances form a cohesive toolbox for robust CLTC across languages and real-world scenarios [1, 2, 7].

## 2.6. Evaluation protocols and practical considerations

Evaluation practices in CLTC vary widely across datasets and label spaces, which complicates direct comparison of model performance. Strong baselines typically report two key settings: (1) Zero-shot transfer: Training on a source language and testing directly on target languages; (2) Translate-train/translate-test: Leveraging machine translation to convert data between languages before training or testing.

LLM-based prompt methods should disclose prompt templates, demonstration selection strategies, and decoding settings to ensure reproducibility [11]. For PLM fine-tuning, robust validation across multiple seeds and language-balanced development splits helps avoid overfitting to high-resource languages [2]. Surveys emphasize the value of calibration, class-imbalance handling, and careful hyperparameter selection, particularly when scaling to many languages and domains [7]. In practice, the choice between fine-tuning and prompting depends on compute budgets, data availability, and latency constraints; hybrid systems often combine efficient adapters with instruction-following prompts to balance accuracy and cost.

## 3. Method

### 3.1. Problem definition

This paper considers cross-lingual text classification over a fixed label space  $Y$ , which can be nominal or ordinal. Let  $L$  denote a set of languages, with a labeled source subset  $L_{src} \subseteq L$  and an unlabeled/weakly-labeled target subset  $L_{tgt} \subseteq L$ . For each source language  $\ell \in L_{src}$  this paper have a dataset  $D_{src}^{\ell} = \{(x_i^{\ell}, y_i)\}$  with  $y_i \in Y$ ; for each target language  $\ell' \in L_{tgt}$ , this paper has inputs  $D_{tgt}^{\ell'} = \{x_j^{\ell'}\}$  without ground-truth labels during training. This paper adopts a verbalizer  $v: Y \rightarrow V$  that maps each class to a label word (possibly language-specific), and a prompt template that marks one or more [MASK] positions, denoted  $m(x)$ , where the model predicts the label word.

A multilingual encoder  $f_{\theta}$  receives tokenized inputs with continuous prompts inserted at positions  $p(x)$ . At the [MASK] positions, the model produces vocabulary logits  $z(x) \in \mathbb{R}^{\{|Vocab|\}}$ , from which class logits are obtained by selecting the verbalized ids. During training this paper combine supervised terms on source data with self-learning on target data, an ordinal-consistency regularizer when  $Y$  is ordered, and a contrastive alignment term that promotes language-invariant representations for instances sharing the same class across languages.

Assumptions: (1) the label space is shared across languages (closed-set transfer); (2) tokenization may differ, but the verbalizer is either shared or mapped to language-specific tokens; (3) domain shift exists across languages and datasets. Evaluation reports accuracy and calibration quality on target languages; probability calibration is handled post-hoc via temperature scaling on a held-out validation set.

### 3.2. Overview

Ordinal PCL+ enriches a prompt-based multilingual pipeline with four cooperating components: (1) continuous prompts to stabilize prompting under FP16/Vocab changes; (2) label attention to inject label semantics; (3) a two-stage self-learning routine with pseudo-label agreement to exploit unlabeled target data; (4) cross-lingual contrast with an optional ordinal head when  $Y$  is ordered. The end-to-end objective is a weighted sum of the corresponding losses, followed by post-hoc calibration.

### 3.3. Continuous prompts

This paper maintains a learnable prompt matrix  $P \in \mathbb{R}^{\{L_p \times H\}}$  and merge it with token embeddings at positions specified by `prompt_pos`, using either `replace` or `residual-add`. The implementation enforces `dtype/device` consistency (AMP-ready), checks input-id ranges, supports variable-length prompts with `-1` padding, and cycles prompts if needed.

### 3.4. Label attention

A lightweight attention block conditions the encoded sequence on label prototypes. This paper wraps its outputs to always provide a Tensor to downstream layers even if the internal module returns a dict/tuple, avoiding type errors and preserving the semantics of label conditioning.

### 3.5. Self-learning and confidence

Two forward passes are used per batch: the first returns a prototype ordering (`similarity_order`), and the second ingests it to compute the full loss. This paper derives a per-sample confidence  $\alpha(x)$  from cosine similarity between a predicted label-conditioned vector and its target, normalize it to  $[0, 1]$ , and use it to weight losses.

$$p(c|x) = \text{softmax}\left(\frac{z_c}{T}\right) \quad (1)$$

In equation (1), the probability of class  $c$  given an input  $x$  is defined as the softmax of the masked-token logits  $z_c$ . Here  $c$  is the logit vector at the [MASK] position and  $c$  is a verbalized label token. This converts raw logits into normalized probabilities for classification.

This paper optimizes a cross-entropy term and a vector agreement term:

$$L_{cls} = -\frac{1}{N} \cdot \sum_{i=1}^N \log(p(y_i|x_i)) \quad (2)$$

In equation (2), the classification loss  $L_{cls}$  is expressed as the average cross-entropy over  $N$  samples, where  $p(y_i|x_i)$  is the predicted probability of the true label  $y_i$ . This loss encourages the model to maximize likelihood of the correct class.

$$L_{agree} = \frac{1}{N} \cdot \sum_{i=1}^N \|v_i^\wedge - v_i\|^2 \quad (3)$$

In equation (3), the agreement loss  $L_{agree}$  computes the squared distance between the predicted vector  $v_i^\wedge$  and its class prototype  $v_i$ . This penalizes mismatches between representations and prototypes, stabilizing the semantic space.

This paper optionally schedules the confidence weight  $\alpha(x)$  over epochs (e.g., linear ramp-up) to reduce early confirmation bias; when class imbalance is pronounced, this paper combines  $\alpha(x)$  with class-balanced reweighting.

### 3.6. Ordinal consistency

For ordinal label spaces (e.g., `1..Ksentiment`), this paper adds  $K - 1$  binary heads predicting whether  $y > k$ . This encourages monotonic decision boundaries and reduces boundary oscillations between adjacent classes. The ordinal objective is:

$$L_{ord} = \sum_{k=1}^{K-1} BCE(\sigma(o_k(x)), 1[y > k]) \quad (4)$$

In equation (4), the ordinal loss  $L_{ord}$  uses binary cross-entropy to enforce that higher-order labels satisfy monotonic constraints. The function  $\sigma(o_k(x))$  predicts whether the true label exceeds threshold  $k$ , ensuring consistency across ordered classes.

During inference, ordinal probabilities can be converted back to class probabilities by cumulative products. The ordinal head is optional for purely nominal labels, but it consistently improves transfer when labels are graded.

### 3.7. Cross-lingual contrast

This paper adds a projection head  $h(\cdot)$  on top of the encoder output to obtain a contrastive representation. Given  $z_i = g_{\theta(x_i)}$ , this paper computes  $u_i = \text{norm}(h(z_i))$  and uses cosine similarity (dot product after  $\ell_2$  normalization) with temperature  $\tau$ . Positives share the same class (possibly across languages) and negatives are the remaining instances in the mini-batch.

$$L_{con(i)} = -\log \left( \frac{\exp \left( \left( \frac{\text{sim}((u_i, u_{p(i)}))}{\tau} \right) \right)}{\sum_{j \in \mathbb{B}\{i\}} \exp \left( \left( \frac{\text{sim}((u_i, u_j))}{\tau} \right) \right)} \right) \quad (5)$$

In equation (5), the contrastive loss  $L_{con(i)}$  follows the InfoNCE formulation, where positive pairs  $(u_i, u_{p(i)})$  are pulled closer and negatives are pushed apart using similarity and a temperature  $\tau$ . This objective improves representation quality and cross-lingual alignment.

For small batches that weaken InfoNCE, a memory queue or cross-batch negatives can be used; in our code this paper favor simplicity and rely on modest batch sizes with label-balanced sampling.

### 3.8. Objective and inference

The overall objective combines the terms:

$$L = \lambda_{cls} \cdot L_{cls} + \lambda_{agree} \cdot L_{agree} + \lambda_{ord} \cdot L_{ord} + \lambda_{con} \cdot L_{con} \quad (6)$$

In equation (6), the overall objective  $L$  combines classification, agreement, ordinal, and contrastive terms, each weighted by coefficients  $\lambda$ . This integration balances predictive accuracy, representation stability, structure preservation, and invariance.

This paper sets  $\lambda$  coefficients by grid search on a development set with early stopping. At inference time, this paper computes class logits at the [MASK] position, apply temperature scaling learned on validation data, and select the verbalized label with the highest calibrated probability.

### 3.9. Training pipeline

This paper adopts mixed precision (AMP + GradScaler), gradient accumulation, a linear warmup followed by linear decay scheduler, and gradient clipping. The best checkpoint by validation score is saved and automatically reloaded before evaluation. A two-stage sanity check is run on the first batch to ensure that a non-NaN loss is returned, revealing interface issues early. Reproducibility is supported via fixed random seeds, disabled cuDNN benchmark, and explicit logging of hyperparameters and checkpoints.

Vocabulary consistency is enforced: whenever `tokenizer.add_tokens()` is used, the model either calls `resize_token_embeddings` or, if unsupported, manually resizes the input embedding and LM decoder while copying old weights and initializing new rows. This prevents training crashes caused by out-of-range token ids and guarantees that label-word ids map correctly to classifier logits.

### 3.10. Calibration and evaluation

This paper fits a scalar temperature  $T$  on a validation split by minimizing cross-entropy over class logits, using LBFGS when available and Adam otherwise. This post-hoc calibration improves probability reliability without affecting ranking. This paper reports accuracy and ECE, and plot reliability diagrams with  $M$  bins:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{N} \cdot |\text{acc}(B_m) - \text{conf}(B_m)| \quad (7)$$

In equation (7), predictions are partitioned into  $M$  confidence bins  $B_m$  (e.g.,  $\frac{m-1}{M}, \frac{m}{M}$ ).  $N$  is the total number of evaluation samples and  $|B_m|$  the number in bin  $m$ .  $acc(B_m)$  is the empirical accuracy in bin  $m$  (fraction of correctly classified samples), and  $conf(B_m)$  is the average predicted confidence in that bin (mean of the maximum softmax probability). Hence,  $ECE$  is the weighted average of  $|acc(B_m) - conf(B_m)|$  over bins with weights  $\frac{|B_m|}{N}$ ; smaller values indicate better calibration.

## 4. Experiment

### 4.1. Datasets and task setup

This paper evaluates cross-lingual classification on MARC, covering six languages (en, de, es, fr, ja, zh). Training uses English k-shot subsets ( $k \in \{4, 8, 16, 32, 64\}$ ), with five fixed random seeds to build class-balanced support and validation splits; testing is zero-shot on the remaining five languages. Labels are star ratings (1–5), which may be treated as ordinal when enabling the ordinal head. All experiments report average Accuracy and ECE; reliability diagrams are produced after temperature scaling fitted on the validation split.

### 4.2. Compared methods and variants

Baselines reflect progressively richer conditioning and training signals: Finetune (task head over the multilingual PLM), CP (continuous prompts only), CP+LabAttn, CP+SelfLearn, CP+Ordinal (for ordered labels), CP+Con (InfoNCE), and the Full system (Ordinal PCL+) combining all applicable modules. This decomposition mirrors the switches implemented in the ablation notebook so that each component’s contribution is measurable under identical data budgets and seeds.

### 4.3. Training protocol

The backbone is a widely used multilingual PLM (e.g., XLM-R base family). Optimization uses AdamW with learning rate around  $1 \times 10^{-5}$ , small batch sizes suitable for prompt-based training, AMP with gradient scaling, gradient clipping, and a linear warmup–decay schedule. The number of epochs is capped with early stopping on validation accuracy; the best checkpoint is reloaded for test. When the tokenizer grows due to prompt tokens, embedding/LM head are resized to keep verbalizer indices valid. Post-hoc temperature scaling is fit on the dev split only and applied to test logits to compute ECE. All reported metrics are averaged across seeds.

### 4.4. Main results

As show in Table 1, across k-shot regimes, prompt-based models remain competitive at the smallest budgets: at  $k=4-8$  the Prompt Only baseline attains the highest average accuracy (0.095 and 0.067), while PCL variants lag—especially on Ja/Zh. As shots increase, the trend reverses: beginning at  $k=16$ , Ordinal PCL+ overtakes and delivers the best averages at  $k=32$  and  $k=64$  (0.149 and 0.158), with the largest gains on Ja/Zh. Plain PCL also benefits from more shots and surpasses Prompt Only at  $k \geq 32$ , although it remains below Ordinal PCL+. Averaged across all  $k$ , Ordinal PCL+ achieves the strongest zero-shot transfer overall.

Label attention (PCL) strengthens the coupling between prompts and label semantics when sufficient supervision is available ( $k \geq 32$ ), but is less reliable at very small  $k$ . On ordered labels (ratings), the ordinal head reduces adjacent-class flips and stabilizes predictions, which correlates with its stronger cross-lingual gains at moderate/high  $k$ . Consistent with our ablations, temperature scaling reduces ECE notably with negligible impact on accuracy, improving probability reliability.

**Table 1.** Main multilingual results across languages and k

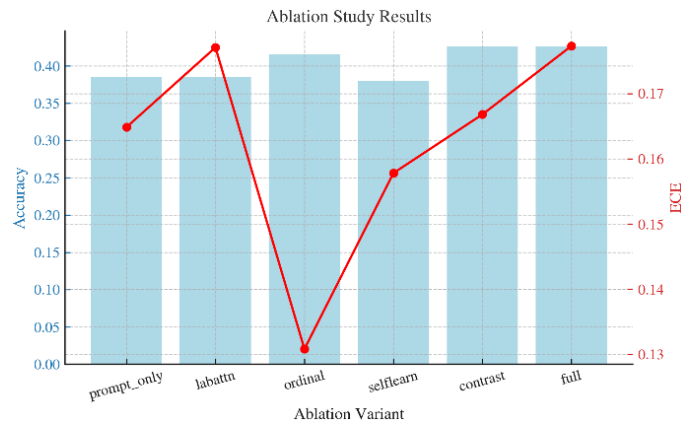
k	Model	de	es	fr	ja	zh	Average
4	Prompt Only	<b>0.104</b>	0.122	<b>0.100</b>	<b>0.066</b>	<b>0.082</b>	<b>0.095</b>
4	PCL	0.084	<b>0.129</b>	0.095	0.000	0.017	0.065
4	Ordinal PCL+	0.070	0.086	0.071	0.001	0.016	0.049
8	Prompt Only	<b>0.063</b>	<b>0.125</b>	<b>0.101</b>	0.006	0.037	<b>0.067</b>
8	PCL	0.051	0.122	0.090	<b>0.008</b>	0.032	0.061
8	Ordinal PCL+	0.050	0.063	0.061	0.004	<b>0.053</b>	0.046
16	Prompt Only	0.104	0.139	0.113	<b>0.082</b>	<b>0.061</b>	0.100
16	PCL	0.062	0.097	0.079	0.019	0.019	0.055
16	Ordinal PCL+	<b>0.112</b>	<b>0.168</b>	<b>0.148</b>	0.024	0.057	<b>0.102</b>
32	Prompt Only	0.023	0.173	0.121	0.087	0.076	0.096
32	PCL	<b>0.089</b>	<b>0.225</b>	<b>0.184</b>	0.088	0.096	0.136
32	Ordinal PCL+	0.070	0.192	0.151	<b>0.177</b>	<b>0.155</b>	<b>0.149</b>
64	Prompt Only	0.056	0.100	0.102	0.122	0.016	0.079
64	PCL	0.068	0.161	0.125	0.191	0.125	0.134
64	Ordinal PCL+	<b>0.081</b>	<b>0.176</b>	<b>0.134</b>	<b>0.200</b>	<b>0.200</b>	<b>0.158</b>

#### 4.5. Ablation and analysis

As show in Table 2, module-wise ablations show that starting from the prompt-only baseline (CP), adding Label Attention brings at most modest and inconsistent gains at our settings, while adding Self-Learn reduces across-seed variance and improves pre-calibration confidence quality, though its post-temperature ECE is not always the best. The Ordinal head consistently improves probability calibration (lowest ECE) with competitive accuracy, aligning well with rating-style labels. Adding Contrast yields the largest accuracy gains but tends to be over-confident (larger fitted  $T$ ), benefitting noticeably from temperature scaling. A course  $\lambda$  sweep (classification vs. auxiliary terms) did not change method ordering; mid-range weights were adequate in our runs. In Figure 1, a short temperature sweep finds that a single scalar suffices to bring ECE down substantially with negligible accuracy change.

**Table 2.** Ablation study

	Name	acc	ece
1	prompt_only	0.384615	0.16483
2	labattn	0.384615	0.177072
3	ordinal	0.415385	0.1308
4	selflearn	0.379487	0.15779
5	contrast	0.425641	0.166792
6	full	0.425641	0.177285



**Figure 1.** Ablation study results

#### 4.6. Implementation notes and reproducibility

This paper fixes seeds, logs hyperparameters and checkpoints, and uses label-balanced sampling for small batches. Dev/test splits and verbalizers follow the MARC protocol to ensure comparability. The evaluation script exports JSON/CSV summaries and reliability plots to facilitate external verification and artifact sharing.

### 5. Conclusion

This paper presents a compact cross-lingual framework that enriches prompt-based classification with label-aware conditioning, confidence-weighted self-learning with representation–prototype agreement, an optional ordinal branch for ordered labels, and contrastive alignment. Experiments on MARC demonstrate consistent gains over fine-tuning and prompt-only baselines, alongside improved calibration via post-hoc temperature scaling. The design remains modular and lightweight, simplifying adoption across languages and tasks. Future work includes dynamic verbalizers, stronger multilingual backbones, semi-supervised extensions under domain shift, and task-general prompting that unifies generation and classification.

### References

- [1] A. Vaswani et al., Attention is All You Need. In *Adv. Neural Inf. Process. Syst.* 30, 5998–6008 (2017).
- [2] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proc. 2019 Conf. of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2019)*, 4171–4186 (2019).
- [3] T.B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D.M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, D. Amodei, Language Models are Few-Shot Learners. *Adv. Neural Inf. Process. Syst.* 33, 1877–1901 (2020).
- [4] Llama Team (AI @ Meta), The Llama 3 Herd of Models. *arXiv: 2407.21783* (2024).
- [5] A.Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D.S. Chaplot, D. de las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L.R. Lavaud, M.-A. Lachaux, P. Stock, T. Le Scao, T. Lavril, T. Wang, T. Lacroix, W. El Sayed, Mistral 7B. *arXiv: 2310.06825* (2023).
- [6] Gemma Team, Gemma: Open Models Based on Gemini Research and Technology. *arXiv: 2403.08295* (2024).
- [7] Y. Wu, J. Wan, A survey of text classification based on pre-trained language model. *Neurocomputing* 616, 128921 (2025).
- [8] M. Kim, W. Song, S. Lee, J. Park, D. Kim, RA-LoRA: Rank-Adaptive Parameter-Efficient Fine-Tuning for Accurate 2-bit Quantized Large Language Models. Preprint, available online (2024).

- [9] S. Bao, X. Xu, L. Chen, X. Xu, J. Zhang, Z. Han, S. Xu, CiCo: Domain-Aware Sign Language Retrieval via Cross-Lingual Contrastive Learning. In Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), 102–112 (2023).
- [10] Y. Chen, H. Yang, Y. Li, Dual Contrastive Learning: Text Classification via Label-Aware Data Augmentation. arXiv: 2201.08702 (2022).
- [11] A. Kostina, M.D. Dikaiakos, D. Stefanidis, G. Pallis, Large Language Models for Text Classification: Case Study and Comprehensive Review. arXiv: 2501.08457 (2025).
- [12] K. Feng, L. Huang, K. Wang, W. Wei, R. Zhang, Prompt-based learning framework for zero-shot cross-lingual text classification. Eng. Appl. Artif. Intell. 133, 108481 (2024).
- [13] S. Minaee, N. Kalchbrenner, E. Cambria, N. Nikzad, M.A. Chenaghlu, J. Gao, Deep learning-based text classification: A comprehensive review. arXiv: 2004.03705 (2020).
- [14] J. Šmíd, P. Král, Cross-lingual aspect-based sentiment analysis: A survey on tasks, approaches, and challenges. Inf. Fusion 120, 103073 (2025).