STRUCTURING SEMANTIC EMBEDDINGS FOR PRINCIPLE EVALUATION: A KERNEL-GUIDED CONTRASTIVE LEARNING APPROACH

Anonymous authorsPaper under double-blind review

ABSTRACT

Evaluating principle adherence in high-dimensional text embeddings is challenging because principle-specific signals are often entangled with general semantic content. Our kernel-guided contrastive learning framework learns to disentangle these signals by projecting embeddings into a structured subspace. In this space, each principle is centered on a learnable **prototype kernel**—an optimized vector that embodies its core meaning—while a jointly learned **semantic basis** preserves context. A novel **offset penalty**, a loss term designed to create structure, then enforces a margin around each prototype. This ensures that even semantically similar principles are clearly separated while capturing their inherent contextual variability. Experiments show our optimized embeddings significantly improve performance on principle classification and ordinal regression, outperforming few-shot Large Language Models and demonstrating the value of specialized representations for reliable principle evaluation.

1 Introduction

Ensuring that the generated text adheres to human-defined principles—such as fairness, honesty, and safety—is a critical challenge for outputs from powerful language models (Weidinger et al., 2021; Bommasani et al., 2021; Hendrycks et al., 2023). Most research on AI alignment has focused on controlling models during the text generation process, aiming to make their outputs inherently safe or helpful (Christiano et al., 2017; Ouyang et al., 2022; Bai et al., 2022a). However, a separate and less-explored problem is how to reliably evaluate a text's adherence to principles after it has already been generated, a task known as post-hoc evaluation (Gehman et al., 2020; Rae et al., 2022). To tackle this challenge at scale, the dominant approach is to first encode the text into high-dimensional embeddings, which are designed to capture its rich semantic meaning, and then evaluate these embeddings against structured principle representations to determine whether the text adheres to the specified principles. The core difficulty, however, begins with these very embeddings. Standard general-purpose embeddings capture rich context and broad semantics, yet they are not explicitly structured or optimized to isolate the specific, often subtle, features indicative of principle adherence from general linguistic content (Devlin et al., 2019; Liu et al., 2022). This leads to embeddings where signals for nuanced principles, whose manifestations are frequently context-dependent (e.g., subtle bias, implied intent), are entangled with unrelated information (Bolukbasi et al., 2016; Caliskan et al., 2017). This entanglement poses a fundamental bottleneck to reliable post-hoc evaluation, with serious consequences in sensitive domains like automated content moderation and medical diagnosis (Birhane et al., 2021; Riegel et al., 2021; Esteva et al., 2021).

Addressing this representational bottleneck is therefore the logical first step towards bridging the evaluation gap. We therefore propose a novel **kernel-guided contrastive learning** framework that actively remodels the representation space. Specifically, our method first defines a learnable, central **prototype** for each principle, and then uses a novel **contrastive objective** to sculpt the space around these prototypes—pulling similar texts closer while enforcing a clear margin from dissimilar ones. The resulting structured embeddings can then be used as high-quality input features for any standard downstream classifier (e.g., an SVM or a simple regressor) to perform the final, reliable evaluation. To validate this two-stage approach, we test it on measurable proxy tasks—such as fine-grained

emotion and toxicity detection—that embody the same core challenge of disentangling subtle signals from a noisy background.

Our key contributions, aimed at enabling reliable principle alignment evaluation, are as follows: (1) A novel conceptual framework for principle alignment evaluation, utilizing prototype vectors as structured principle representations and a semantic basis to capture contextual dependencies. (2) A neural principle architecture that materializes this framework, using an attention mechanism to map embeddings to a structured subspace defined by learnable prototype kernels. (3) A kernel-guided contrastive learning objective featuring a novel offset penalty, specifically designed to organize this subspace by imposing a structured geometry around the prototypes for fine-grained principle separation. (4) Extensive experimental validation demonstrating the effectiveness of our optimized embeddings on various downstream principle evaluation tasks, including principle classification (GoEmotions), ordinal regression (Amazon Reviews), and classification in a sensitive domain (Toxic Comment Classification Challenge). Experiments show significant performance improvements over raw embeddings and superior results compared to few-shot Large Language Models in these specific evaluation contexts, validating our central thesis that optimizing representation geometry is a crucial and effective step towards reliable principle evaluation.

2 METHODOLOGY

Our framework is designed to restructure fixed, pre-computed text embeddings by projecting them into a principle-aligned subspace. This approach is distinct from end-to-end finetuning; instead of updating a large encoder's weights, we learn a separate, lightweight transformation that disentangles principle-specific features from general contextual information. The core idea is that learnable prototype kernels can serve as explicit anchors for principles within this subspace, much like basis vectors defining a coordinate system. This allows varied, context-dependent manifestations of a principle to be organized around a shared, abstract representation. Our framework realizes this by training a dedicated neural architecture with a novel, geometry-aware contrastive objective.

2.1 Framework Overview

Our framework's central goal is to map a *fixed*, *high-dimensional* text embedding $\mathbf{X}_i \in \mathbb{R}^D$ into a low-dimensional, structured representation $\mathbf{e}_i \in \mathbb{R}^d$ (where $d \ll D$). This new representation is designed to make principle-specific features, which are entangled in the original space, readily discernible for downstream evaluators. Our framework consists of two core components:

A Neural Architecture (f_{θ}) . We introduce a dedicated architecture, the *neural principle extractor*, which performs the mapping from \mathbf{X}_i to \mathbf{e}_i . It uses an attention mechanism to project the input onto a subspace defined by a set of learnable prototypes $\{\mathbf{c}_k\}$, which represent the core meaning of each principle. The learnable nature of these prototypes is crucial, as their optimal positions are data-dependent (Details in Section 2.2).

A Geometry-Aware Learning Objective (\mathcal{L}_{total}). This is a composite loss function that trains the neural architecture. It combines a supervised contrastive loss with our novel prototype offset penalty to explicitly organize the subspace, ensuring that the learned representations \mathbf{e}_i are structured around their corresponding principle prototypes \mathbf{c}_k (Details in Section 2.3).

2.2 NEURAL PRINCIPLE EXTRACTOR

The neural principle extractor, f_{θ} , is designed to untangle principle-specific features from general context. It processes the input embedding X_i through two parallel streams that are ultimately fused.

The first stream produces the **Semantic Basis** (S_i) by processing the input through a **shared Multi-Layer Perceptron** (**MLP**). This vector's role is to capture the general topic or context of the text. The second stream produces the **Prototype Mapping** (m_i) via an **attention mechanism**, which computes a weighted combination of a set of learnable **prototype kernels** { c_k }. These kernels act as abstract, optimizable anchors for each principle in the learned subspace (initialization strategies are detailed in Appendix A.1).

 For instance, given "This movie was a crushing disappointment," the Semantic Basis captures 'a movie review,' while the Prototype Mapping captures 'disappointment.' The final representation \mathbf{e}_i is a weighted fusion of these two components, with their relative contribution controlled by a learnable parameter α . This makes the final representation both principle-aware and contextually grounded. Further architectural details are available in Appendix A.2.

2.3 Kernel-guided Contrastive Learning

The neural principle extractor is trained by minimizing a composite loss function, \mathcal{L}_{total} , which we designed to progressively sculpt the geometry of the principle subspace. Our design philosophy is to combine standard learning procedures with novel, structure-enforcing objectives. We first use a **Contrastive Loss** to pull embeddings of the same principle into coarse clusters, and then introduce our novel **Offset Loss** to refine the local geometry around each prototype kernel by enforcing explicit separation margins. Finally, auxiliary terms including an **Orthogonality Loss** and a **Classification Loss** help ensure feature disentanglement and stabilize the training process. The total loss is a weighted combination of these components:

Contrastive Loss ($\mathcal{L}_{contrastive}$). This term encourages samples from the same principle to be closer in the learned embedding space while pushing apart samples from different principles. Based on the InfoNCE loss (Oord et al., 2018), our implementation incorporates adaptations like hard negative mining and dynamic class weights to enhance discrimination and handle class imbalance. The core form is:

$$\mathcal{L}_{\text{contrastive}} = \text{AdaptedInfoNCE}(\mathbf{e}_i, \{\mathbf{e}_k\}_{k \neq i}; \tau) \tag{1}$$

where e_i 's are the enhanced features and τ is the temperature.

Offset Loss ($\mathcal{L}_{\text{offset}}$). While the contrastive loss encourages general class clustering, it does not explicitly control the shape of the clusters or the guaranteed separation between them. To impose a more precise geometric structure, we introduce the novel Offset Loss. This term acts as a geometric regularizer, controlling the position of a sample's kernel mapping \mathbf{m}_i relative to the prototype kernels. It is composed of two penalties that work in tandem: an *Intra-Class Penalty* to allow for variation within a principle's cluster, and an *Inter-Class Penalty* to enforce separation between different principle clusters.

The Intra-Class Penalty introduces a "safe radius" (δ_{intra}) around each prototype, only penalizing samples that fall outside this radius. This encourages compactness while preserving diversity. Here, \mathbf{c}_{y_i} denotes the learnable prototype kernel corresponding to the true class label y_i . Its formulation is:

$$P_{\text{intra},i} = \max(0, ||\mathbf{m}_i - \mathbf{c}_{y_i}|| - \delta_{\text{intra}})^2$$
(2)

The Inter-Class Penalty enforces a clear margin (δ_{inter}) between clusters by ensuring that each sample is closer to its true prototype than to any incorrect one. Its formulation is:

$$P_{\text{inter},i} = \max(0, ||\mathbf{m}_i - \mathbf{c}_{y_i}|| - \min_{k \neq y_i} ||\mathbf{m}_i - \mathbf{c}_k|| + \delta_{\text{inter}})^2$$
(3)

The total Offset Loss is a weighted sum of these two penalties:

$$\mathcal{L}_{\text{offset}} = \frac{1}{B} \sum_{i=1}^{B} w_{y_i} (\lambda_{\text{inclass}} P_{\text{intra},i} + \lambda_{\text{crossclass}} P_{\text{inter},i})$$
 (4)

where B is the batch size, w_{y_i} is the class weight for label y_i , and λ_{inclass} , $\lambda_{\text{crossclass}}$ are hyperparameters controlling the relative contributions of the intra-class and inter-class penalties.

Orthogonality Loss ($\mathcal{L}_{\text{orthogonality}}$). To encourage the semantic basis \mathbf{s}_i and kernel mapping \mathbf{m}_i to capture largely distinct information while allowing necessary interaction, this orthogonality loss promotes their 'soft' orthogonality. It is based on the cosine similarity between the two vectors, penalized when exceeding a dynamic margin $\delta_{\text{orthogonal}}$:

$$\mathcal{L}_{\text{orthogonality}} = \frac{1}{B} \sum_{i=1}^{B} w_{y_i} \cdot \max(0, |\cos(\mathbf{s}_i, \mathbf{m}_i)| - \delta_{\text{orthogonal}})$$
 (5)

where w_{y_i} is a class weight.

Classification Loss ($\mathcal{L}_{classification}$) As an auxiliary objective, a standard classification loss is applied directly to the model's attention scores, which are interpreted as logits for class membership. This loss penalizes the model if it fails to assign a high score to the correct principle for a given input. Through backpropagation, this penalty signal directly trains the transformations that produce the query and key vectors, forcing the attention mechanism to learn the alignment between inputs and their corresponding principle prototypes. To mitigate class imbalance and focus on harder examples, we employ Focal Loss (Lin et al., 2017):

$$\mathcal{L}_{\text{classification}} = -\sum_{i} w_{y_i} (1 - p_{y_i})^{\gamma} \log p_{y_i}$$
 (6)

where p_{y_i} is the model's predicted probability for the true class y_i (derived from softmax over the attention scores), γ is the focusing parameter, and w_{y_i} is a class weight.

Magnitude Loss ($\mathcal{L}_{\text{magnitude}}$). Specifically for ordinal regression tasks, this loss enforces the natural ordering of labels by encouraging the magnitude of the kernel mapping $||\mathbf{m}_i||$ to be proportional to the intensity of the ordinal label $I(y_i)$. This helps the learned subspace reflect the graded nature of ordinal values. The loss is:

$$\mathcal{L}_{\text{magnitude}} = \frac{1}{B} \sum_{i=1}^{B} (||\mathbf{m}_i|| - \lambda_{\text{mag_scale}} I(y_i) \cdot ||\mathbf{c}_{y_i}||)^2$$
(7)

where $\lambda_{\text{mag_scale}}$ is a scaling factor and $I(y_i)$ maps the label to a numerical intensity (e.g., 1-5 for star ratings).

Total Loss (\mathcal{L}_{total}). The model is trained end-to-end by minimizing the total loss, a weighted sum of the above components:

$$\mathcal{L}_{total} = \lambda_{contrastive} \cdot \mathcal{L}_{contrastive} + \lambda_{offset} \cdot \mathcal{L}_{offset} + \lambda_{class} \cdot \mathcal{L}_{classification} + \lambda_{orth} \cdot \mathcal{L}_{orthogonality} (+\lambda_{mag} \cdot \mathcal{L}_{magnitude})$$
(8)

The weights (λ values, including λ_{mag} for ordinal regression) are key hyperparameters determined through optimization, such as Bayesian optimization. A detailed analysis of the computational complexity can be found in Appendix A.4.

3 EXPERIMENT

This section evaluates the performance of our kernel-guided contrastive learning framework in enhancing principle evaluation in text embeddings. We detail our experimental setup in Section 3.1 and present results on three distinct datasets representing different principle evaluation tasks: GoEmotions, Amazon Reviews, and the Toxic Comment Classification Challenge. Our experiments demonstrate that embeddings optimized by our framework significantly improve downstream evaluation performance compared to using raw embeddings.

3.1 EXPERIMENTAL SETUP

All experiments utilize text embeddings (dimension D=1024) generated by the <code>jina-embeddings-v3</code> (Jina AI, 2024) model, which has demonstrated strong performance in semantic similarity tasks relevant to principle alignment.

GoEmotions Dataset. GoEmotions (Demszky et al., 2020) is a large-scale corpus of Reddit comments annotated with 27 fine-grained emotion categories. For our experiments, we focus on a challenging subset of five principles (Disappointment, Sadness, Disapproval, Gratitude, Approval), including clusters of semantically similar emotions, to test the method's ability to distinguish subtle differences. This selection represents common positive and negative social emotions and allows rigorous evaluation of fine-grained discriminative capacity.

Amazon Reviews Dataset. The Amazon Reviews dataset (Ni et al., 2019) comprises user reviews and corresponding 1-5 star ratings for products on Amazon. These star ratings serve as indicators of sentiment intensity and represent ordinal values. The dataset provides a platform for validating our approach through sentiment classification (treating ratings as distinct classes) and ordinal regression

tasks (treating ratings as ordered values). A subset of this dataset was sampled and preprocessed, preserving the original, unbalanced distribution of star ratings.

Toxic Comment Classification Challenge. This dataset (Jigsaw & Kaggle, 2018) is critical for evaluating principle alignment in a sensitive domain - toxicity detection. It presents a highly unbalanced binary classification task (toxic vs. non-toxic), where we extract samples labeled 'nontoxic' (all toxicity labels are 0) and samples with only the 'toxic' label as 1. This forms an extremely unbalanced dataset (approximately 1:25 toxic vs. non-toxic labels in the test set, reflecting the realworld imbalance where toxic comments are much rarer). The training set is balanced to approximately 1:3 using negative oversampling and positive undersampling. The task's difficulty is exacerbated by

228 229 230

225

226

227

231

232 233 234

235

240

251

258

259

267

268

269

the less clearly defined nature of 'non-toxic' and 'toxic' principles.

Detailed data distributions are in Appendix C.3. Our principle extractor maps embeddings to d=64dimensions, chosen based on preliminary studies comparing different dimensions (Appendix A.5, showing d=64 provides an optimal balance between performance and efficiency). Performance metrics are reported as Mean \pm Standard Deviation over 10-fold cross-validation.

3.2 Performance Evaluation on Downstream Tasks

GoEmotions Dataset: Classification. Our evaluation on the GoEmotions dataset focuses on fine-grained principle classification. We train standard classifiers (SVM, Random Forest, Logistic Regression, XGBoost, Transformer) on both raw and optimized embeddings of the GoEmotions five-principle subset. Overall performance averaged over principles is summarized in Table 1. Detailed per-principle performance highlighting improvements on challenging principles is presented in Appendix B.1.

As summarized in Table 1, our optimized embeddings yield consistent and statistically significant improvements across all classifiers compared to raw embeddings. While XGBoost and Transformer also show improvements, the relative gains are most substantial for models that performed less strongly with raw embeddings, indicating our method particularly benefits classifiers struggling with the original representation. This highlights the effectiveness of our framework in creating a more discriminative representation space for principles, especially improving the performance floor.

Table 1: Overall (Avg. Principle) Performance on GoEmotions Five-Principle Set (Mean \pm Std. Dev.)

Metric	Emb. Type	SVM	RF	LR	XGBoost	Transformer
Precision	Raw Emb. Opt. Emb.	0.748 ± 0.049 0.787 ± 0.035	0.733 ± 0.059 0.789 ± 0.036	0.737 ± 0.031 0.791 ± 0.033	0.747 ± 0.020 0.773 ± 0.028	0.785 ± 0.031 0.787 ± 0.029
Recall		0.721 ± 0.045 0.764 ± 0.034				
F1		0.729 ± 0.046 0.770 ± 0.033				

Amazon Reviews Dataset: Ordinal Regression. On the Amazon Reviews dataset, we assess the utility of our optimized embeddings for principle evaluation, particularly in the context of ordinal tasks. The 1-5 star ratings in this dataset naturally represent ordered values, making it suitable for ordinal regression tasks, which capture the intensity of sentiment as an ordinal principle. We also evaluate performance on the associated *classification* tasks (treating ratings as distinct categories), with detailed results presented in Appendix B.3. Our framework is designed to enhance representations for both scenarios, but we focus the main text discussion on ordinal regression as it directly leverages the magnitude learning objective.

Table 2 summarizes overall ordinal regression performance across different regressors. Optimized embeddings consistently improve overall metrics (MSE, RMSE, R²) compared to raw embeddings. Detailed per-class MSE, showing significant reductions for most star ratings, is provided in Appendix B.2. This demonstrates improved ability to capture the graded nuances of sentiment as an ordinal principle.

Table 2: Overall Ordinal Regression Performance on Amazon Reviews (Mean \pm Std. Dev.)

Metric	Emb. Type	SVM	RF	LR	XGBoost	Transformer
MSE	Raw Emb.	0.668 ± 0.135	0.506 ± 0.120	0.635 ± 0.097	0.546 ± 0.163	0.602 ± 0.173
	Opt. Emb.	0.365 ± 0.158	0.392 ± 0.143	0.394 ± 0.149	0.377 ± 0.097	0.359 ± 0.086
RMSE	Raw Emb.	0.813 ± 0.083	0.706 ± 0.083	0.795 ± 0.059	0.731 ± 0.111	0.768 ± 0.110
	Opt. Emb.	0.593 ± 0.119	0.617 ± 0.107	0.618 ± 0.112	0.609 ± 0.080	0.595 ± 0.071
\mathbb{R}^2		0.604 ± 0.086 0.785 ± 0.089				

Toxic Comment Classification Challenge. For principle alignment evaluation in a sensitive domain, we assess our framework's performance on the Toxic Comment Classification Challenge dataset.

Table 3 summarizes the key results. Optimized embeddings consistently yield statistically significant improvements in Average F1 and Minority F1 across all classifiers compared to raw embeddings.

Table 3: Performance on Toxic Comment Classification Challenge (Mean \pm Std. Dev.)

Metric	Emb. Type	SVM	RF	LR	XGBoost	Transformer
Avg. F1				0.897 ± 0.004 0.936 ± 0.003		
Minority F1	Raw Emb. Opt. Emb.	0.497 ± 0.025 0.507 ± 0.024	0.405 ± 0.044 0.537 ± 0.035	0.396 ± 0.018 0.493 ± 0.023	0.433 ± 0.024 0.518 ± 0.028	0.574 ± 0.027 0.589 ± 0.023

3.3 COMPARISON WITH FEW-SHOT LARGE LANGUAGE MODELS

We compare our method's performance with that of few-shot prompted Large Language Models (LLMs), where the models directly perform the evaluation tasks based on prompting, serving as a generalist baseline. We evaluated LLama3.3_70b_Q4_K (Meta AI Team, 2025), and additionally tested deepseek-chat-v3-0324 (DeepSeek Team, 2025) and gemini-2.5-pro-exp-03-25 (Google DeepMind Team, 2025) via API calls.

Table 4 summarizes the comparison. Our method, using optimized embeddings with Transformer, consistently outperforms few-shot LLMs on these principle alignment tasks. This highlights the advantage of task-specific representation refinement for precise principle alignment compared to general-purpose LLM prompting. Detailed LLM evaluation results are in Appendix D.2.

Table 4: Performance Comparison with Few-shot Large Language Models

Dataset	Metric	LLama3.3	DeepSeek-chat-v3	Gemini-2.5-pro	Opt. Emb. + Transformer
GoEmotions Amazon Reviews		0.67 0.60	0.70 0.45	0.70 0.56	0.77 0.36
Toxic Comment	Avg. F1	0.91	0.89	0.91	0.96

3.4 ABLATION STUDY

To understand the contribution of each component of our kernel-guided contrastive learning framework to the observed performance improvements, we conducted an ablation on the GoEmotions dataset study by selectively removing or isolating each loss term during training. Table 5 summarizes the performance across different configurations, showing F1 scores per principle and the average.

The results in Table 5 demonstrate that the combination of all three loss components yields the highest overall performance (0.78 average F1) and generally the best per-principle scores on this dataset. The **Classification Loss** is essential for basic discriminability, but insufficient on its own without

Table 5: Ablation study on the GoEmotions dataset (F1 score Mean). Disappt.-Disappointment, Sad.-Sadness, Disapprv.-Disapproval, Grat.-Gratitude, Apprv.-Approval.

Configuration	Disappt.	Sad.	Disapprv.	Grat.	Apprv.	Average
Only Contrastive Loss	0.37	0.75	0.72	0.92	0.74	0.75
Only Offset Loss	0.40	0.75	0.72	0.94	0.77	0.77
Only Classification Loss	0.26	0.58	0.61	0.85	0.61	0.64
Without Contrastive Loss	0.42	0.77	0.72	0.94	0.77	0.77
Without Offset Loss	0.44	0.71	0.73	0.93	0.76	0.77
Without Classification Loss	0.36	0.74	0.70	0.94	0.77	0.76
Raw Embeddings Optimized Embeddings (Full Model)	0.36 0.49	0.64 0.77	0.66 0.72	0.93 0.94	0.72 0.76	0.72 0.78

explicit structure guidance. The **Contrastive Loss** and **Offset Losses** are crucial for improving the structural separation in the learned space, as shown by their better performance compared to using only classification loss. While removing either structural loss slightly reduces the average F1, their *combined effect* alongside the classification objective is necessary for optimal performance, particularly evident for the challenging 'Disappointment' principle where the full model achieves the highest F1 (0.49). Consistent with the main results, all configurations involving the learned losses significantly outperform using raw embeddings (0.72 on average), validating the overall approach to optimize the embedding space. This study, conducted on the GoEmotions dataset, demonstrates that integrating these complementary loss components is key to learning a balanced and robust principle-aware representation for effective evaluation.

3.5 EMBEDDING SPACE ANALYSIS

To understand the impact of our kernel-guided contrastive learning framework on the structure of the embedding space, we perform both qualitative visualization and quantitative geometric analysis.

Qualitative Visualization using t-SNE. We visualize the embedding spaces before and after optimization using t-SNE (Van der Maaten & Hinton, 2008) to qualitatively assess the separation of different principles or ratings.

Figure 1 compares the raw and optimized embedding spaces. Raw embeddings show significant overlap, especially for semantically similar principles or adjacent ratings. In contrast, optimized embeddings exhibit much clearer separability, forming distinct clusters for each principle/rating. For ordinal regression, the optimized clusters also show a clear ordered arrangement, consistent with the magnitude loss objective.

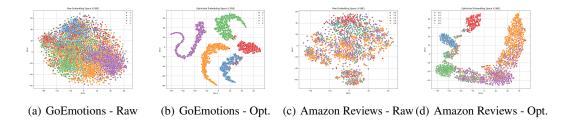


Figure 1: Comparison of embedding spaces using t-SNE. (a) and (b) show raw and optimized embeddings for GoEmotions (classification). (c) and (d) show raw and optimized embeddings for Amazon Reviews (ordinal regression). Optimized embeddings demonstrate clearer separability and structure.

Quantitative Geometric Analysis. This qualitative observation of improved structure and separation in the embedding space is quantitatively supported by metrics analyzing cluster separability and quality. These metrics include the ratio between Within-class Variance and Between-class Variance

(Fisher, 1936), Silhouette Score (Rousseeuw, 1987), Class Overlap (Dom, 2001), and Fisher's Discriminant Ratio (Fisher, 1936). Table 6 summarizes the average results and standard deviations for each metric on the GoEmotions dataset embeddings.

Table 6: Quantitative Geometric Quality Metrics (Mean \pm Std. Dev.)

Metric	Emb. Type	Mean \pm Std. Dev.	Interpretation	Improvement (%)
Within/Between Ratio	Raw	8.76 ± 0.53	Lower is better	94.35%
	Optimized	0.50 ± 0.06		
Silhouette Score	Raw	0.018 ± 0.004	Higher is better	833.55%
	Optimized	0.164 ± 0.029		
Class Overlap	Raw	0.563 ± 0.017	Lower is better	34.12%
	Optimized	0.371 ± 0.038		
Fisher's Discriminant Ratio	Raw	0.165 ± 0.008	Higher is better	1008.87%
	Optimized	1.825 ± 0.236		

The results in Table 6 demonstrate a consistent and significant improvement across all measured geometric quality metrics after applying our kernel-guided contrastive learning optimization. This quantitative evidence strongly supports the visual observations from the t-SNE plots regarding enhanced class separability and improved cluster quality in the optimized embedding space. The improvements are statistically significant (as indicated by the separation of mean values relative to standard deviations), confirming that our optimization effectively structures the embedding space to enhance principle separability and cluster quality.

3.6 BEYOND PERFORMANCE: ADVANTAGES OF OPTIMIZED EMBEDDINGS FOR PRINCIPLE EVALUATION

Beyond the significant performance improvements on downstream evaluation tasks (3.2), our kernel-guided contrastive learning framework yields optimized text embeddings that offer several key advantages for principle evaluation. (1) It provides a *reusable intermediate representation* specifically structured for principle evaluation, enabling a modular pipeline across different tasks. This also facilitates *simplified downstream modeling*, allowing simple classifiers to achieve strong performance comparable to complex models on raw embeddings, significantly reducing model selection and tuning efforts. (2) Using these lower-dimensional embeddings enables *enhanced computational efficiency* for downstream evaluation. Training times for models like XGBoost were reduced by up to 96.5% compared to raw embeddings, making the evaluation process more practical. (3) The method yields a *structured subspace* for principle features (3.5), enhancing their discernibility and potential interpretability. These advantages highlight the utility of our approach for building robust and efficient principle evaluation systems beyond just metric gains.

4 RELATED WORK

Principle Alignment in Text. Prior efforts related to principle alignment often focus on constraining language models during text generation (e.g., RLHF (Christiano et al., 2017; Ouyang et al., 2022), DPO, Constitutional AI (Bai et al., 2022c)). Evaluating the principle alignment of already generated text or explicitly structuring embeddings for this evaluation purpose remains less explored. Our framework focuses on this evaluation gap by learning a discernible representation space.

Semantic Contextualization and Subspace Learning. Text embeddings capture context but may conflate subtle principle distinctions (Devlin et al., 2019; Liu et al., 2022). Subspace learning or task-specific projections aim to extract relevant features (Guo & Mackey, 2022; Wei et al., 2022; Gu & Roth, 2023). However, general-purpose methods like UMAP (McInnes et al., 2018) or Spectral Embedding (Von Luxburg, 2007), while useful for dimensionality reduction or structure visualization, are not optimized to specifically disentangle predefined principle-specific features. Our neural principle extractor learns a structured subspace specifically for principle evaluation, uniquely combining principle-specific extraction with a semantic basis via attention to preserve context dependence.

Contrastive Learning. Contrastive learning enhances representation separability by pulling positives together and pushing negatives apart (Hadsell et al., 2006; Khosla et al., 2020). However, many supervised methods struggle with semantically similar principles because they contrast sample pairs, rather than explicitly modeling class structure. While Prototypical Contrastive Learning (PCL) (Li et al., 2021) addresses this by contrasting samples with prototypes computed from batch data, our approach differs by using learnable prototype kernels as explicit model parameters. We further impose a structured geometry around these kernels using a novel offset penalty for more fine-grained separation.

Kernel-Based Methods and Prototype Learning. Kernel methods and prototype learning impose structure or learn representative points in embedding spaces (Schölkopf & Smola, 2002; Snell et al., 2017; Yang et al., 2016). Kernel offset constraints can refine separation (Wilson et al., 2017; Zhang et al., 2021). Our work integrates these ideas into a contrastive framework. Unlike methods that use prototypes for inference-time classification (Snell et al., 2017) or as batch-computed centers (Li et al., 2021), our learnable kernels serve as foundational parameters that actively define the target geometry of the representation space.

Metric Learning. Metric learning aims to learn embedding spaces where distances reflect desired relationships (Yao et al., 2021). Objectives like triplet or N-pair loss are common, relying on relative distances between samples (Weinberger & Saul, 2009; Sohn, 2016). Our kernel-guided framework is a form of supervised deep metric learning, but differentiates itself by using explicit learnable prototype kernels as anchors and an offset penalty to structure the space around them, addressing the challenge of distinguishing semantically close principles beyond generic distance constraints based on sample pairs.

Geometric Embeddings. Geometric embeddings represent structured data (e.g., knowledge graphs) as shapes for tasks like reasoning or querying (e.g., Query2Box (Ren et al., 2020) for KG querying). Our method similarly structures a representation space but applies learned point kernels to general text embeddings. Unlike geometric embeddings focusing on structured data relations, our objective is to enhance the distinguishability of principle-specific features in text via contrastive learning, using learned prototype kernels as anchors and attention to map text inputs to this space.

5 CONCLUSION AND FUTURE WORK

This paper addresses the challenge of principle evaluation, a task hindered by the entangled, unstructured nature of standard text embeddings. We introduce a **kernel-guided contrastive learning** framework that tackles this representational bottleneck directly. Our method remodels the embedding space by learning a structured subspace where each principle is anchored by a learnable prototype kernel. A novel offset penalty then enforces a clear geometric separation between these principles, resulting in representations with significantly improved discernibility. Our experiments demonstrate that these structured embeddings boost the performance of various downstream evaluators and outperform few-shot LLMs. These results validate our central thesis: that explicitly remodeling the geometry of representation space is a critical and effective step towards building more reliable principle evaluation systems.

Current limitations include reliance on supervised data for known principles and unverified performance in diverse domains/languages. Application risks like misuse or bias amplification also warrant consideration. Future work targets evaluation with reduced supervision or for unseen principles, and extending the scope to diverse domains/languages. We will also investigate ethical considerations including misuse and bias. Technical extensions include continuous regression, automated annotation strategies, and dynamic environments. Integrating into evaluation components within RLAIF pipelines (Bai et al., 2022b; Lee et al., 2023) could enable more robust feedback signals.

REFERENCES

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nisanth DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2208.03282*, 2022a.

- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a harmless language model with constitutional ai: An overview. *arXiv preprint arXiv:2212.08073*, 2022b.
 - Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a harmless language model with constitutional ai: An overview. *arXiv preprint arXiv:2212.08073*, 2022c.
 - Abeba Birhane, Pranav Kalluri, Dallas Card, William Agnew, Ravit Dotan, and Michelle Bao. The values encoded in machine learning research. *arXiv* preprint arXiv:2106.15597, 2021.
 - Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Tauman Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
 - Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
 - Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186, 2017.
 - Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Advances in neural information processing systems*, pp. 4299–4307, 2017.
 - DeepSeek Team. Deepseek v3: Scaling language model training with efficient architecture. *arXiv* preprint arXiv:2412.19437, 2025. URL https://arxiv.org/abs/2412.19437. Accessed: May 11, 2025.
 - Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Katherin Razavi, Anish Gupta, Jatin Singh, and Vered Shwartz. Goemotions: A dataset of fine-grained emotions. *arXiv preprint arXiv:2005.00547*, 2020.
 - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2019.
 - Byron Dom. An information-theoretic external cluster-validity measure. In *Proceedings of the 18th International Conference on Machine Learning (ICML)*, 2001.
 - Andre Esteva, Katherine Chou, Serena Yeung, Nicholas Codella, Anthony Cooper, Roxana Novoa, Richard Socher, William Mitchell, Martin Witkowski, et al. Deep learning-enabled medical computer vision. *NPJ digital medicine*, 4(1):1–9, 2021.
 - Ronald A Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7 (2):179–188, 1936.
 - Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 3356–3369, 2020.
 - Google DeepMind Team. Gemini 2.5 pro: Advancements in multimodal ai capabilities, 2025. URL https://blog.google/products/gemini/gemini-2-5-pro-updates/. Google Blog, Accessed: May 11, 2025.
 - Yuchen Gu and Dan Roth. Learning task-specific representation for novel words by retrieving related entries. *arXiv preprint arXiv:2304.06765*, 2023.
 - Yutong Guo and Lester Mackey. Dimension reduction for supervised learning with kernel embeddings. *The Journal of Machine Learning Research*, 23(1):10308–10371, 2022.
 - Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), 2:1735–1742, 2006.

- Dan Hendrycks, Mantas Mazeika, Andy Zou, Joseph Simon, Dawn Banneb, Joel Gehman, and Thomas Pell. An overview of catastrophic ai risks. *arXiv preprint arXiv:2306.12938*, 2023.
- Jigsaw and Kaggle. Toxic Comment Classification Challenge, 2018. URL https://www.kaggle.com/c/toxic-comment-classification-challenge. Accessed: [Date you accessed].
 - Jina AI. Jina Embeddings v3. https://jina.ai/embeddings/, 2024.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:18661–18673, 2020.
 - Kimin Lee, Lantao Yu, Xingyou Zhan, David Krueger, Pieter Abbeel, and Chelsea Finn. Rlaif: Scaling reinforcement learning from ai feedback. *arXiv preprint arXiv:2309.00267*, 2023.
 - Junnan Li, Pan Zhou, Caiming Xiong, and Steven C. H. Hoi. Prototypical contrastive learning for imbalanced unsupervised representation learning. In *International Conference on Learning Representations (ICLR)*, 2021.
 - Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.
 - Yuchen Liu, Lucas Chilton, Yitong Li, Lajanugen Logeswaran, Honglak Lee, and Dan Roth. A closer look at how fine-tuning changes bert's representations. *arXiv preprint arXiv:2206.01360*, 2022.
 - Leland McInnes, John Healy, Geoffrey Hinton, and Sam Roweis. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
 - Meta AI Team. Llama 3.3: A multilingual large language model for research and commercial use. Technical report, Meta AI, 2025. URL https://www.llama.com/docs/model-cards-and-prompt-formats/llama3_3/. Accessed: May 11, 2025.
 - Jianmo Ni, Jiacheng Li, and Julian McAuley. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. *IEEE Transactions on Knowledge and Data Engineering*, 33(4): 1812–1824, 2019.
 - Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
 - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35: 27730–27744, 2022.
 - Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Scott Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2022.
 - Xiang Ren, Taylor Berg-Kirkpatrick, Chris Dyer, and Noah A Smith. QUERY2BOX: Reasoning over Knowledge Graphs in Vector Space using Box Embeddings. In *International Conference on Learning Representations*, 2020.
 - Simon Riegel, David WeLTe, Stefan Geier, Halil Kilicoglu, Jens Gruendner, Hans-Ulrich Prokosch, Joerg Christoph, and Martin Sedlmayr. Towards a standard for the automated evaluation of clinical text. *Journal of biomedical informatics*, 116:103759, 2021.
 - Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53–65, 1987.
 - Bernhard Schölkopf and Alexander J Smola. Learning with kernels: Support vector machines, regularization, optimization, and beyond. *MIT Press*, 2002.

- Jake Snell, Kevin Swersky, and Richard S Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017.
 - Kihyuk Sohn. N-pair loss less than triplet loss: Large-scale deep metric learning via proxy neighborhood assignment. In *European Conference on Computer Vision*, pp. 500–516. Springer, 2016.
 - Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.
 - Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
 - Tingwu Wei, Yixiang Zhang, Yali Li, and Qiang Wang. Learning to solve complex tasks in low-dimensional environments. *arXiv preprint arXiv:2208.03412*, 2022.
 - Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Megan Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*, 2021.
 - Kilian Q Weinberger and Lawrence K Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, 10:207–244, 2009.
 - Andrew G Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P Xing. Kernel interpolation for scalable structured gaussian processes (kiss-gp). *International Conference on Machine Learning*, 2017.
 - Jie Yang, Jia Liu, Xiaofeng Wu, Chuan Yang, Keyu Li, and Jianlong Zhang. Joint embedding and clustering for heterogeneous information networks. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 10(4):1–26, 2016.
 - Jiayi Yao, Guiguang Ding, Jungong Han, Yi Wang, Xiaohui Zhang, Zhili Li, and Yu Xu. Deep metric learning: A survey. *Neurocomputing*, 461:131–153, 2021.
 - Yue Zhang, Linfeng Sun, and Hao Chen. Kernel offset constraints for high-dimensional representation learning. *ICLR*, 2021.

THE USE OF LARGE LANGUAGE MODELS (LLMS)

In preparing this work, we made limited and appropriate use of Large Language Models (LLMs) as follows:

- Writing aid and polishing: LLMs were used to assist in improving grammar, clarity, and style. The substantive content, ideas, and technical contributions remain the authors' own.
- Retrieval and discovery: LLMs were employed to support literature search and discovery (e.g., identifying related work). All cited references were verified by the authors.

A IMPLEMENTATION DETAILS

This appendix provides further details regarding the neural principle extractor's architecture, prototype kernel initialization strategies, specific training hyperparameters and loss function configurations, computational complexity, and justification for the principle subspace dimension, as referenced in the main paper.

A.1 PROTOTYPE KERNEL INITIALIZATION DETAILS

The K learnable prototype kernels $\mathbf{c}_k \in \mathbb{R}^d$ are initialized based on the task type to encourage structured learning.

For Classification Tasks: For classification tasks (GoEmotions 5-principle set, Amazon Reviews classification), the K prototype kernels are initialized randomly on the unit hypersphere in \mathbb{R}^d . To ensure distinct starting points and facilitate separation during training, we apply a procedure to guarantee a minimum pairwise Euclidean distance between any two initialized kernels. While not strictly enforcing orthogonality, this initial separation prevents kernels from collapsing onto the same point early in training. A target minimum distance of $\sqrt{2}$ (the Euclidean distance between orthogonal unit vectors) is aimed for during this initialization step.

For Ordinal Regression Tasks: For ordinal regression tasks (Amazon Reviews ordinal regression), the initialization incorporates the inherent order of the labels. The kernels are first initialized randomly on the unit hypersphere. Subsequently, the norm of the kernel \mathbf{c}_k corresponding to ordinal label k is scaled by a factor proportional to its numerical intensity I(k). For the 1-5 star ratings in Amazon Reviews, I(k) = k. The scaling factor used is $(1.0 + (k-1) \cdot \text{scale_multiplier})$, where scale_multiplier is a small constant (e.g., 0.1) to ensure that kernels corresponding to higher ratings have progressively larger initial magnitudes. This guides the magnitude loss and encourages the learned principle representations to exhibit an ordered structure.

A.2 NEURAL NETWORK ARCHITECTURE DETAILS

The neural principle extractor f_{θ} is implemented as a neural network that maps the input text embedding $\mathbf{X}_i \in \mathbb{R}^{1024}$ to a d=64 dimensional principle-aware representation \mathbf{e}_i . The architecture is composed of a shared Multi-Layer Perceptron (MLP) and an attention mechanism.

The shared MLP used to compute the semantic basis \mathbf{s}_i consists of two fully connected layers with LeakyReLU activation functions and Batch Normalization. Dropout is applied after each hidden layer for regularization. The layer dimensions are as follows:

- Input layer: $\mathbb{R}^{1024} \to \mathbb{R}^{512}$
- Hidden layer 1: $\mathbb{R}^{512} \to \mathbb{R}^{256}$ (followed by LeakyReLU, Batch Norm, Dropout)
- Hidden layer 2: $\mathbb{R}^{256} \to \mathbb{R}^d$ (followed by LeakyReLU, Batch Norm, Dropout), where d = 64. The output of this layer is the semantic basis \mathbf{s}_i .

The Dropout rate used throughout the MLP is 0.2.

The attention mechanism involves linear transformations of the input embedding and the prototype kernels to compute queries, keys, and values:

- Query projection: query_fc : $\mathbb{R}^{1024} \to \mathbb{R}^d$
- Key projection: key fc: $\mathbb{R}^d \to \mathbb{R}^d$
- Value projection: value fc: $\mathbb{R}^d \to \mathbb{R}^d$

These projected vectors are used in the scaled dot-product attention calculation as described in Section 2.3.

The learnable parameter α that weights the semantic basis and kernel mapping in the final fusion layer is a scalar variable initialized to 0.05.

A.3 TRAINING AND LOSS FUNCTION DETAILS

This appendix provides detailed information regarding the training procedure and the full mathematical formulations and specific configurations for each component of the kernel-guided contrastive learning objective, as referenced from the main paper.

The neural principle extractor is trained end-to-end using the AdamW optimizer. The initial learning rate is set to 1e-4, with a weight decay of 1e-5. A learning rate scheduler, such as Cosine Annealing or ReduceLROnPlateau, is employed to dynamically adjust the learning rate during training based on validation performance. Training is performed for a maximum of 100 epochs, with early stopping based on performance on a validation set to prevent overfitting. The batch size used throughout our experiments is 128.

To handle class imbalance, which is particularly pronounced in datasets like Amazon Reviews and the Toxic Comment Classification Challenge, dynamic class weights w_{y_i} are applied to the loss calculations. These weights are computed as the inverse frequency of each true class within the current training batch, providing stronger gradients for minority classes.

The training objective is to minimize a composite loss function \mathcal{L}_{total} , which is a weighted sum of several components (detailed below). The full mathematical formulations and specific parameters are provided for clarity and reproducibility.

Contrastive Loss ($\mathcal{L}_{\text{contrastive}}$): While the core idea is based on InfoNCE applied to enhanced features \mathbf{e}_i , our implementation incorporates specific adaptations for hard negative mining and handles class imbalance via dynamic weights w_{y_i} . The temperature parameter τ is set to 0.1. The specific adapted InfoNCE formulation used for $\mathcal{L}_{\text{contrastive}}$ is as follows, incorporating a term for hard negative samples:

$$\mathcal{L}_{\text{contrastive}} = \frac{1}{B} \sum_{i=1}^{B} w_{y_i} \left[-\log \frac{\exp(\text{sim}(\mathbf{e}_i, \mathbf{e}_{p_i})/\tau)}{\sum_{k=1, k \neq i}^{B} \exp(\text{sim}(\mathbf{e}_i, \mathbf{e}_k)/\tau)} - \lambda_{\text{hard}} \sum_{n \in \text{HardNegativeSet}(i)} \log(1 - \exp(\text{sim}(\mathbf{e}_i, \mathbf{e}_n)/\tau)) \right]$$
(9)

where \mathbf{e}_{p_i} is a positive sample for \mathbf{e}_i (another sample with the same true label in the batch), $\mathrm{sim}(\cdot,\cdot)$ denotes cosine similarity, $\tau=0.1$ is the temperature, w_{y_i} is the dynamic class weight, $\mathrm{HardNegativeSet}(i)$ is a subset of hard negative samples for \mathbf{e}_i identified based on criteria like high attention scores towards y_i or small embedding distance, and λ_{hard} is a weighting factor for the hard negative term (tuned during optimization, typically between 0 and 1).

Offset Loss ($\mathcal{L}_{\text{offset}}$): This novel term regulates the position of a sample's kernel mapping \mathbf{m}_i relative to its true principle prototype kernel \mathbf{c}_{y_i} and other kernels, enforcing proximity to the correct kernel while maintaining distance from incorrect ones with controlled margins. The full loss definition is:

$$\mathcal{L}_{\text{offset}} = \frac{1}{B} \sum_{i=1}^{B} w_{y_i} (\lambda_{\text{inclass}} P_{\text{intra},i} + \lambda_{\text{crossclass}} P_{\text{inter},i})$$
 (10)

where B is the batch size, w_{y_i} is the dynamic weight, and λ_{inclass} , $\lambda_{\text{crossclass}}$ are hyperparameters (tuned during optimization). The penalty terms $P_{\text{intra},i}$ and $P_{\text{inter},i}$ are defined based on Euclidean distances to kernels and margins δ_{intra} and δ_{inter} (tuned during optimization, typically within [0.1, 0.5] for δ values):

$$P_{\text{intra},i} = \max(0, \|\mathbf{m}_i - \mathbf{c}_{y_i}\| - \delta_{\text{intra}})^2$$
(11)

$$P_{\text{inter},i} = \max(0, \|\mathbf{m}_i - \mathbf{c}_{y_i}\| - \min_{k \neq y_i} \|\mathbf{m}_i - \mathbf{c}_k\| + \delta_{\text{inter}})^2$$
(12)

Orthogonality Loss ($\mathcal{L}_{\text{orthogonality}}$): This term promotes "soft" orthogonality between the semantic basis \mathbf{s}_i and kernel mapping \mathbf{m}_i to encourage separation of information types. It is based on the

absolute cosine similarity, penalized when exceeding a dynamic margin $\delta_{orthogonal}$:

$$\mathcal{L}_{\text{orthogonality}} = \frac{1}{B} \sum_{i=1}^{B} w_{y_i} \cdot \max(0, |\cos(\mathbf{s}_i, \mathbf{m}_i)| - \delta_{orthogonal})$$
(13)

where w_{y_i} is a class weight. The dynamic margin $\delta_{orthogonal}$ is annealed linearly from an initial value of 0.5 down to a final value of 0.05 over the course of training epochs.

Classification Loss ($\mathcal{L}_{\text{classification}}$): A standard classification loss is applied using the attention scores as logits. To handle class imbalance and focus on harder examples, we employ Focal Loss (Lin et al., 2017). The full formula is:

$$\mathcal{L}_{\text{classification}} = -\sum_{i} w_{y_i} (1 - p_{y_i})^{\gamma} \log p_{y_i}$$
 (14)

where p_{y_i} represents the predicted probability for the true class y_i , $\gamma = 2$ is the focusing parameter, and w_{y_i} is a class weight.

Magnitude Loss ($\mathcal{L}_{\text{magnitude}}$): Used specifically for ordinal regression tasks, this loss enforces the natural ordering by encouraging the magnitude of \mathbf{m}_i to be proportional to the intensity of $I(y_i)$. The full formula is:

$$\mathcal{L}_{\text{magnitude}} = \frac{1}{B} \sum_{i=1}^{B} w_{y_i} (||\mathbf{m}_i|| - \lambda_{mag_scale} I(y_i) \cdot ||\mathbf{c}_{y_i}||)^2$$
(15)

where w_{y_i} is a class weight, λ_{mag_scale} is a learnable scaling factor (initialized to 1.0), and $I(y_i)$ maps the ordinal label to a numerical intensity (e.g., $I(y_i) = y_i$ for 1-5 star ratings). This loss is applied only for ordinal regression tasks.

Total Loss (\mathcal{L}_{total}): The model is trained end-to-end by minimizing the total loss, which is a weighted combination of the above components:

$$\mathcal{L}_{total} = \lambda_{contrastive} \cdot \mathcal{L}_{contrastive} + \lambda_{offset} \cdot \mathcal{L}_{offset} + \lambda_{class} \cdot \mathcal{L}_{classification} + \lambda_{orth} \cdot \mathcal{L}_{orthogonality} (+\lambda_{mag} \cdot \mathcal{L}_{magnitude})$$
(16)

The weights (λ values, including λ_{mag} specifically for ordinal regression) are key hyperparameters that balance the contribution of each loss term. These weights, along with other hyperparameters like $\tau, \delta_{intra}, \delta_{inter}, \delta_{orthogonal}, \gamma, \lambda_{hard}$, and λ_{mag_scale} initialization, are determined through optimization, such as Bayesian optimization or extensive grid search on a validation set. Specific optimized lambda values used for each task/dataset are typically reported alongside the experimental results or in a dedicated section on hyperparameter tuning (e.g., Appendix C.2).

A.4 COMPUTATIONAL COMPLEXITY ANALYSIS

We analyze the computational complexity of our framework during training and inference.

Training Complexity: The primary computational cost during training arises from the forward and backward passes through the neural principle extractor and the calculation of the loss components over a batch of size B. The extractor involves:

- Shared MLP: A sequence of matrix multiplications. Given input dimension D=1024, output dimension d=64, and hidden dimensions $h_1=512, h_2=256$, the complexity is $O(D \cdot h_1 + h_1 \cdot h_2 + h_2 \cdot d)$ per sample.
- Attention Mechanism: Involves linear projections $(O(D \cdot d + d^2 \cdot K))$ for a batch of size B, where K is the number of principles), computing attention scores $(O(B \cdot K \cdot d))$, and weighted summation $(O(B \cdot K \cdot d))$.

The dominant part of the forward pass per batch is approximately $O(B \cdot (D \cdot h_1 + h_1 \cdot h_2 + h_2 \cdot d + K \cdot d))$. Loss calculations involve vector operations and distance calculations on the d-dimensional embeddings and K kernels:

• Contrastive Loss: $O(B^2 \cdot d)$ in the standard form, often optimized to $O(B^2)$ or $O(B \cdot P \cdot d)$ with P positives per sample.

- Offset Loss: Involves distances to K kernels, $O(B \cdot K \cdot d)$.
- Orthogonality, Classification, Magnitude Losses: $O(B \cdot d)$ or O(B).

 The overall training complexity per batch is dominated by the forward/backward passes and loss calculations, roughly $O(B \cdot (D \cdot h_{max} + K \cdot d) + B^2 \cdot d)$ in the worst case (for contrastive) or $O(B \cdot (D \cdot h_{max} + K \cdot d))$ with typical batch sizes and optimizations. This is comparable to other deep metric learning or contrastive learning frameworks.

Inference Complexity: Inference requires a single forward pass through the extractor. The complexity per sample is $O(D \cdot h_1 + h_1 \cdot h_2 + h_2 \cdot d + K \cdot d)$, which is linear with respect to D and K. This makes obtaining the optimized embedding efficient.

Downstream Efficiency Gains: A practical benefit is the reduced computational cost for downstream tasks operating on the d=64 dimensional embeddings compared to D=1024 raw embeddings. This reduction is significant for many standard classifiers and contributes to faster downstream training and inference times. As noted in Section 3.6, this led to substantial training time reductions for downstream models.

A.5 Justification for Principle Subspace Dimension (d = 64)

The choice of the principle subspace dimension d=64 for the output embeddings was guided by preliminary experiments. We evaluated model performance on a validation set using various output dimensions (e.g., 32, 128, 256). d=64 was found to provide a robust balance, offering significant dimensionality reduction from the input (1024 dimensions) while preserving sufficient information for effective principle discrimination in downstream tasks, yielding performance comparable to or better than higher dimensions with reduced computational cost for both model training and subsequent downstream task training/inference.

A.6 COMPUTE RESOURCES

All experiments, including the training of the Neural Principle Extractor and evaluation of downstream models, were conducted on a machine equipped with four NVIDIA RTX 4090 GPUs (24GB VRAM each) and 128GB of system RAM. The CPU used was an Intel(R) Xeon(R) Platinum 8336C CPU @ 2.30GHz, running on Ubuntu 24.04 LTS.

Training of the Neural Principle Extractor is computationally efficient. A full training run typically completed within 3 to 15 minutes on a single NVIDIA RTX 4090, depending on the dataset size and complexity. Using multiple GPUs can further reduce this time. Inference using the trained extractor to produce optimized embeddings is significantly faster, requiring only a single forward pass per sample. Evaluating downstream models on the optimized embeddings is also substantially more efficient than using raw embeddings, as discussed in Section 3.6 and detailed in Appendix A.4.

B DETAILED EXPERIMENTAL RESULTS

This appendix provides supplementary detailed results for the experiments presented in Section 3.

B.1 Goemotions Per-Principle Performance

This appendix provides detailed per-principle F1 performance for the GoEmotions dataset, complementing the overall results presented in Section 3.2. Table 7 shows the Mean \pm Standard Deviation F1 scores for each of the five selected emotion principles.

The improvements are most pronounced for semantically similar and initially challenging principles with lower initial F1 scores, such as Disappointment and Sadness. Conversely, for principles like Gratitude, which already achieved high F1 scores with raw embeddings, the relative improvement is more modest across most classifiers. These results demonstrate that our method is particularly effective at refining distinctions for principles that are difficult to classify using standard embedding techniques, raising the performance ceiling for challenging cases while maintaining strong performance on easier ones.

Table 7: Per-Principle F1 Performance on GoEmotions Five-Principle Set (Mean \pm Std. Dev.) in Appendix. Principles are abbreviated as Disappt., Sad., Disapprv., Grat., Apprv.

Principle	Emb. Type	SVM	RF	LR	XGBoost	Transformer
Disappt.	Raw Emb. Opt. Emb.		0.237 ± 0.202 0.410 ± 0.087		0.359 ± 0.144 0.439 ± 0.099	0.315 ± 0.073 0.386 ± 0.117
Sad.	Raw Emb. Opt. Emb.	0.643 ± 0.059 0.687 ± 0.048	0.728 ± 0.069 0.734 ± 0.031		0.711 ± 0.082 0.698 ± 0.031	
Disapprv.	Raw Emb. Opt. Emb.		0.691 ± 0.050 0.740 ± 0.064		0.703 ± 0.032 0.724 ± 0.082	
Grat.	Raw Emb. Opt. Emb.	0.925 ± 0.032 0.939 ± 0.021			0.905 ± 0.031 0.934 ± 0.032	0.915 ± 0.026 0.938 ± 0.029
Apprv.	Raw Emb. Opt. Emb.	0.732 ± 0.090 0.769 ± 0.053	0.707 ± 0.031 0.747 ± 0.078		0.730 ± 0.060 0.762 ± 0.067	0.716 ± 0.075 0.758 ± 0.053

B.2 AMAZON REVIEWS PER-RATING PERFORMANCE

This appendix provides detailed per-rating performance for the Amazon Reviews dataset, supplementing the summarized classification and ordinal regression results presented in Section 3.2.

Table 11 shows the F1 performance for each star rating (1-5 S) on the Amazon Reviews dataset using raw and optimized embeddings.

Table 8: Classification F1 Performance per Rating on Amazon Reviews (Mean \pm Std. Dev.) in Appendix

Ratings	Emb. Type	SVM	RF	LR	XGBoost	Transformer
1 - S			0.772 ± 0.166 0.874 ± 0.085		0.744 ± 0.235 0.894 ± 0.093	
2 - S	Raw Emb. Opt. Emb.		0.204 ± 0.213 0.708 ± 0.315		0.288 ± 0.221 0.760 ± 0.170	
3 - S	Raw Emb. Opt. Emb.		0.556 ± 0.176 0.697 ± 0.073		0.520 ± 0.141 0.657 ± 0.076	
4 - S	Raw Emb. Opt. Emb.		0.598 ± 0.114 0.639 ± 0.120			
5 - S	Raw Emb. Opt. Emb.		0.710 ± 0.099 0.766 ± 0.093		0.736 ± 0.069 0.741 ± 0.101	

Table 9 provides the per-rating Mean Squared Error (MSE) for the Amazon Reviews ordinal regression task.

B.3 AMAZON REVIEWS CLASSIFICATION RESULTS

This appendix section provides detailed classification performance results on the Amazon Reviews dataset, supplementing the main text discussion which focuses on ordinal regression. For this task, the 1-5 star ratings are treated as distinct discrete categories.

Table 10 summarizes the overall (average per rating) classification performance across different classifiers using both raw and optimized embeddings.

Table 11 presents the F1 performance for each individual star rating (1-5) using both raw and optimized embeddings. Optimized embeddings generally show improved performance across most individual ratings, particularly for the intermediate ratings (2, 3, 4 stars) which are often more challenging to distinguish.

Table 9: Ordinal Regression MSE Performance per Rating on Amazon Reviews (Mean \pm Std. Dev.) in Appendix

Metric	Emb. Type	SVM	RF	LR	XGBoost	Transformer
1-S MSE	Raw Emb. Opt. Emb.			0.567 ± 0.533 0.140 ± 0.254		
2-S MSE	Raw Emb. Opt. Emb.		1.415 ± 0.880 0.435 ± 0.547	1.250 ± 0.972 0.405 ± 0.334	1.465 ± 0.912 0.335 ± 0.338	
3-S MSE		0.726 ± 0.431 0.386 ± 0.260		0.804 ± 0.378 0.442 ± 0.367	0.712 ± 0.269 0.509 ± 0.361	
4-S MSE	Raw Emb. Opt. Emb.	0.653 ± 0.221 0.387 ± 0.171		0.567 ± 0.211 0.453 ± 0.195		0.653 ± 0.260 0.500 ± 0.189
5-S MSE	Raw Emb. Opt. Emb.	0.607 ± 0.348 0.427 ± 0.389		0.467 ± 0.163 0.400 ± 0.394	0.247 ± 0.095 0.307 ± 0.200	0.567 ± 0.438 0.313 ± 0.193

Table 10: Overall (Avg. Rating) Classification Performance on Amazon Reviews (Mean \pm Std. Dev.)

Metric	Emb. Type	SVM	RF	LR	XGBoost	Transformer
Precision	Raw Emb. Opt. Emb.	0.541 ± 0.038 0.728 ± 0.077	0.627 ± 0.093 0.726 ± 0.094	0.595 ± 0.058 0.716 ± 0.084	0.630 ± 0.073 0.718 ± 0.077	0.639 ± 0.067 0.725 ± 0.062
Recall		0.522 ± 0.043 0.721 ± 0.073				
Avg. F1		0.521 ± 0.041 0.717 ± 0.074				

C ADDITIONAL EXPERIMENTAL DETAILS

C.1 DETAILS ON USED ASSETS AND LICENSES

This appendix provides details on the licenses and terms of use for the external datasets, embedding models, and language models used in this research, as referenced from the main paper. Our use of these assets adheres to their respective licenses and terms.

Datasets.

- GoEmotions Dataset (Demszky et al., 2020): This dataset is released under the Creative Commons Attribution-ShareAlike 4.0 International License (CC BY-SA 4.0). Available at https://github.com/google-research/goemotions.
- Amazon Reviews Dataset (Ni et al., 2019): This dataset is provided for research purposes. Its use is subject to the terms specified by the data providers (e.g., Stanford/UCSD). Researchers should refer to the original source for specific usage guidelines. Available via the cited research project website.
- Toxic Comment Classification Challenge: This dataset, originally hosted on Kaggle (Jigsaw & Kaggle, 2018), is made available under the CCO 1.0 Universal Public Domain Dedication. Available at https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge.

Embedding Model.

• Jina Embeddings v3 (Jina AI, 2024): The embeddings used were generated by the jina-embeddings-v3 model. Jina AI models are typically licensed under the Apache 2.0 License. Researchers should consult the official Jina AI model documentation or Hugging Face model card for the most precise license information and terms of use.

Table 11: Classification F1 Performance per Rating on Amazon Reviews (Mean \pm Std. Dev.) in Appendix

Ratings	Emb. Type	SVM	RF	LR	XGBoost	Transformer
1 - S	Raw Emb. Opt. Emb.	0.712 ± 0.219 0.869 ± 0.112		0.713 ± 0.208 0.875 ± 0.084	0.744 ± 0.235 0.894 ± 0.093	0.731 ± 0.209 0.888 ± 0.090
2 - S	Raw Emb. Opt. Emb.	0.277 ± 0.118 0.691 ± 0.187		0.297 ± 0.168 0.667 ± 0.176	0.288 ± 0.221 0.760 ± 0.170	0.432 ± 0.163 0.711 ± 0.141
3 - S	Raw Emb. Opt. Emb.			0.503 ± 0.081 0.662 ± 0.112	0.520 ± 0.141 0.657 ± 0.076	0.584 ± 0.085 0.696 ± 0.143
4 - S	Raw Emb. Opt. Emb.	0.478 ± 0.071 0.650 ± 0.094	0.598 ± 0.114 0.639 ± 0.120	0.565 ± 0.096 0.637 ± 0.129	0.613 ± 0.078 0.622 ± 0.113	0.558 ± 0.123 0.620 ± 0.105
5 - S	Raw Emb. Opt. Emb.	0.614 ± 0.074 0.764 ± 0.112		0.676 ± 0.087 0.764 ± 0.115	0.736 ± 0.069 0.741 ± 0.101	0.724 ± 0.103 0.768 ± 0.082

Large Language Models (for Comparison).

- LLama 3.3 (Meta AI Team, 2025): The Llama 3 family of models is available under the Llama 3 Community License. Use of the quantized version (LLama 3.3_70b_Q4_K) adheres to the terms of this license.
- DeepSeek-Chat-v3 (DeepSeek Team, 2025): Used via API. Use is subject to DeepSeek AI's API Terms of Service.
- Gemini 2.5 Pro (Google DeepMind Team, 2025): Used via API. Use is subject to Google's API Terms of Service (e.g., Google AI or Google Cloud terms).

C.2 HYPERPARAMETERS

C.3 DATA DISTRIBUTION

Hyperparameter Name

input_dim

num_classes

hidden_dims

Description

shared MLP.

Table 12: List of Hyperparameters

Neural Principle Extractor

Dimension of the input feature vector.

Number of classes in the classification task.

Dimensions of the hidden layers in the

Current Value

(Variable)

(Variable)

[512, 256]

		STATE OF THE STATE	
4	output_dim	Dimension of the output feature vector from	64
		the extractor.	
5	kernel_margin	Minimum distance between the initialized	1.414
		prototype kernels.	
6	alpha	Weighting factor for combining semantic	0.3 (trainable)
		basis and principle-specific features.	
		Kernel-guided Contrastive Learnin	ng
7	temperature	Temperature coefficient for adjusting the	0.1
	1	similarity scaling.	
8	class_weights	Weights for each class in the loss function.	(Computed by sklearn.utils.class_weight)
9	lambda_contrastive	Weight for the contrastive loss.	(Optimized using Bayesian optimization)
10	lambda_offset	Weight for the kernel offset loss.	(Optimized using Bayesian optimization)
11	lambda_classification	Weight for the classification loss.	(Optimized using Bayesian optimization)
12	lambda_orthogonality	Weight for the orthogonality loss.	(Optimized using Bayesian optimization)
13	lambda_magnitude	Weight for the magnitude loss.	(Optimized using Bayesian optimization)
14	offset_delta	Tolerance radius for in-class offset.	0.63
15	offset_margin	Additional penalty term for cross-class sep-	2
		aration.	
16	contrastive_k	Number of hard negative samples used in	10
		the contrastive loss.	
17	focal_alpha	Weights for each class in the classification	(Computed by class_weights)
		loss.	_
18	focal_gamma	Controls the influence of easy and hard sam-	2
		ples in the classification loss.	
19	orthogonality_margin	Controls the strength of orthogonality in the	0.5
		optimization.	

Table 13: Data Distribution for GoEmotions, Amazon Reviews, and Toxic Comment Classification Challenge Datasets

Dataset	Label	Train Count	Validation Count	Test Count
GoEmotions Dataset	Disappointment (0)	709	91	88
	Sadness (1)	817	84	102
	Disapproval (2)	1200	212	195
	Gratitude (3)	1200	261	260
	Approval (4)	1200	258	236
Amazon Reviews Dataset	1 star (0)	249	53	54
	2 stars (1)	198	43	42
	3 stars (2)	424	91	91
	4 stars (3)	700	150	150
	5 stars (4)	700	150	150
Toxic Comments Dataset	Toxic (1)	11064	580	870
	Non-toxic (0)	33192	14493	21739

EVALUATION LOG

D.1 Embedding Log

GoEmotions

Model: SVM_RAW

AUC: Mean=0.9198, Std=0.0172

```
Precision: Mean=0.7483, Std=0.0485
1081
         Recall: Mean=0.7207, Std=0.0454
1082
         F1 Score (Minority Class 1): Mean=0.3867, Std=0.0899
1083
         Overall F1 Score: Mean=0.7286, Std=0.0457
1084
         F1 Score (Class 0): Mean=0.3867, Std=0.0899
         F1 Score (Class 1): Mean=0.6434, Std=0.0586
1085
         F1 Score (Class 2): Mean=0.6626, Std=0.0738
1086
         F1 Score (Class 3): Mean=0.9251, Std=0.0322
1087
         F1 Score (Class 4): Mean=0.7315, Std=0.0899
1088
1089
      Model: RF_RAW
1090
         AUC: Mean=0.9064, Std=0.0260
1091
         Precision: Mean=0.7334, Std=0.0592
1092
         Recall: Mean=0.7366, Std=0.0305
1093
         F1 Score (Minority Class 1): Mean=0.2372, Std=0.2015
1094
         Overall F1 Score: Mean=0.7218, Std=0.0354
1095
         F1 Score (Class 0): Mean=0.2372, Std=0.2015
         F1 Score (Class 1): Mean=0.7275, Std=0.0691
1096
         F1 Score (Class 2): Mean=0.6905, Std=0.0498
1097
         F1 Score (Class 3): Mean=0.9202, Std=0.0244
1098
      2025-05-09 11:58:19,124 - INFO - F1 Score (Class 4): Mean
1099
         =0.7067, Std=0.0311
1100
      2025-05-09 11:58:19,124 - INFO -
1101
      Model: LR_RAW
1102
      2025-05-09 11:58:19,124 - INFO -
                                          AUC: Mean=0.9222, Std=0.0163
1103
      2025-05-09 11:58:19,124 - INFO -
                                          Precision: Mean=0.7367, Std
1104
         =0.0310
1105
      2025-05-09 11:58:19,124 - INFO -
                                          Recall: Mean=0.7219, Std=0.0311
1106
      2025-05-09 11:58:19,124 - INFO -
                                          F1 Score (Minority Class 1):
         Mean=0.3753, Std=0.0537
1107
         Overall F1 Score: Mean=0.7261, Std=0.0306
1108
         F1 Score (Class 0): Mean=0.3753, Std=0.0537
1109
         F1 Score (Class 1): Mean=0.6715, Std=0.0811
1110
         F1 Score (Class 2): Mean=0.6523, Std=0.0703
1111
         F1 Score (Class 3): Mean=0.9214, Std=0.0201
1112
         F1 Score (Class 4): Mean=0.7270, Std=0.0609
1113
1114
     Model: XGB_RAW
1115
         AUC: Mean=0.9266, Std=0.0204
1116
         Precision: Mean=0.7470, Std=0.0203
1117
         Recall: Mean=0.7412, Std=0.0144
         F1 Score (Minority Class 1): Mean=0.3591, Std=0.1436
1118
         Overall F1 Score: Mean=0.7365, Std=0.0179
1119
         F1 Score (Class 0): Mean=0.3591, Std=0.1436
1120
         F1 Score (Class 1): Mean=0.7106, Std=0.0818
1121
         F1 Score (Class 2): Mean=0.7034, Std=0.0315
1122
         F1 Score (Class 3): Mean=0.9054, Std=0.0307
1123
         F1 Score (Class 4): Mean=0.7298, Std=0.0600
1124
1125
     Model: TRANSFORMER RAW
1126
         AUC: Mean=0.9344, Std=0.0146
1127
         Precision: Mean=0.7715, Std=0.0358
1128
         Recall: Mean=0.7639, Std=0.0347
1129
         F1 Score (Minority Class 1): Mean=0.4150, Std=0.1180
         Overall F1 Score: Mean=0.7637, Std=0.0355
1130
         F1 Score (Class 0): Mean=0.3154, Std=0.0728
1131
         F1 Score (Class 1): Mean=0.7107, Std=0.0872
1132
         F1 Score (Class 2): Mean=0.6774, Std=0.0545
1133
         F1 Score (Class 3): Mean=0.9154, Std=0.0261
```

```
1134
         F1 Score (Class 4): Mean=0.7162, Std=0.0747
1135
1136
      Model: SVM_OPT
1137
         AUC: Mean=0.9360, Std=0.0143
1138
         Precision: Mean=0.7868, Std=0.0350
         Recall: Mean=0.7639, Std=0.0343
1139
         F1 Score (Minority Class 1): Mean=0.4818, Std=0.1023
1140
         Overall F1 Score: Mean=0.7698, Std=0.0330
1141
         F1 Score (Class 0): Mean=0.4818, Std=0.1023
1142
         F1 Score (Class 1): Mean=0.6870, Std=0.0483
1143
         F1 Score (Class 2): Mean=0.7200, Std=0.0737
1144
         F1 Score (Class 3): Mean=0.9389, Std=0.0209
1145
         F1 Score (Class 4): Mean=0.7689, Std=0.0531
1146
1147
      Model: RF_OPT
1148
         AUC: Mean=0.9279, Std=0.0154
1149
         Precision: Mean=0.7758, Std=0.0345
         Recall: Mean=0.7685, Std=0.0385
1150
         F1 Score (Minority Class 1): Mean=0.4104, Std=0.0866
1151
         Overall F1 Score: Mean=0.7669, Std=0.0363
1152
         F1 Score (Class 0): Mean=0.4104, Std=0.0866
1153
         F1 Score (Class 1): Mean=0.7338, Std=0.0312
1154
         F1 Score (Class 2): Mean=0.7403, Std=0.0635
1155
         F1 Score (Class 3): Mean=0.9395, Std=0.0283
1156
         F1 Score (Class 4): Mean=0.7468, Std=0.0775
1157
1158
      Model: LR_OPT
1159
         AUC: Mean=0.9359, Std=0.0126
1160
         Precision: Mean=0.7913, Std=0.0334
         Recall: Mean=0.7718, Std=0.0328
1161
         F1 Score (Minority Class 1): Mean=0.4794, Std=0.1059
1162
         Overall F1 Score: Mean=0.7763, Std=0.0317
1163
         F1 Score (Class 0): Mean=0.4794, Std=0.1059
1164
         F1 Score (Class 1): Mean=0.7207, Std=0.0544
1165
         F1 Score (Class 2): Mean=0.7279, Std=0.0634
1166
         F1 Score (Class 3): Mean=0.9411, Std=0.0236
1167
         F1 Score (Class 4): Mean=0.7707, Std=0.0547
1168
1169
      Model: XGB_OPT
1170
         AUC: Mean=0.9310, Std=0.0135
1171
         Precision: Mean=0.7728, Std=0.0364
         Recall: Mean=0.7651, Std=0.0391
1172
         F1 Score (Minority Class 1): Mean=0.4386, Std=0.0989
1173
         Overall F1 Score: Mean=0.7643, Std=0.0358
1174
         F1 Score (Class 0): Mean=0.4386, Std=0.0989
1175
         F1 Score (Class 1): Mean=0.6978, Std=0.0313
1176
         F1 Score (Class 2): Mean=0.7238, Std=0.0816
1177
         F1 Score (Class 3): Mean=0.9339, Std=0.0315
1178
         F1 Score (Class 4): Mean=0.7616, Std=0.0670
1179
1180
     Model: TRANSFORMER OPT
1181
         AUC: Mean=0.9298, Std=0.0130
1182
         Precision: Mean=0.7797, Std=0.0307
        Recall: Mean=0.7685, Std=0.0306
1183
        F1 Score (Minority Class 1): Mean=0.3870, Std=0.0734
1184
        Overall F1 Score: Mean=0.7701, Std=0.0302
1185
        F1 Score (Class 0): Mean=0.3859, Std=0.1174
1186
        F1 Score (Class 1): Mean=0.7137, Std=0.0339
1187
        F1 Score (Class 2): Mean=0.7333, Std=0.0593
```

```
1188
1189
1189
1190
1191
1192
F1 Score (Class 3): Mean=0.9375, Std=0.0291
F1 Score (Class 4): Mean=0.7583, Std=0.0527
Comparison script finished.
```

Amazon Reviews

```
1194
1195
      Model: SVM RAW
1196
         AUC: Mean=0.8540, Std=0.0200
         Precision: Mean=0.5409, Std=0.0383
1197
         Recall: Mean=0.5215, Std=0.0434
1198
         F1 Score (Minority Class 1): Mean=0.7120, Std=0.2189
1199
         Overall F1 Score: Mean=0.5209, Std=0.0409
1200
         F1 Score (Class 0): Mean=0.7120, Std=0.2189
1201
         F1 Score (Class 1): Mean=0.2771, Std=0.1176
1202
         F1 Score (Class 2): Mean=0.4332, Std=0.1575
1203
         F1 Score (Class 3): Mean=0.4781, Std=0.0709
1204
         F1 Score (Class 4): Mean=0.6141, Std=0.0736
1205
         MSE: Mean=0.6675, Std=0.1352
         RMSE: Mean=0.8127, Std=0.0834
1206
         R^2: Mean=0.6037, Std=0.0855
1207
         MSE (Class 1-Star): Mean=0.4833, Std=0.5309
1208
         MSE (Class 2-Star): Mean=1.0800, Std=0.8667
1209
         MSE (Class 3-Star): Mean=0.7256, Std=0.4307
1210
         MSE (Class 4-Star): Mean=0.6533, Std=0.2207
1211
         MSE (Class 5-Star): Mean=0.6067, Std=0.3483
1212
1213
      Model: RF RAW
1214
         AUC: Mean=0.8734, Std=0.0288
1215
         Precision: Mean=0.6267, Std=0.0929
1216
         Recall: Mean=0.6279, Std=0.0831
         F1 Score (Minority Class 1): Mean=0.7721, Std=0.1663
1217
         Overall F1 Score: Mean=0.6091, Std=0.0820
1218
         F1 Score (Class 0): Mean=0.7721, Std=0.1663
1219
         F1 Score (Class 1): Mean=0.2038, Std=0.2129
1220
         F1 Score (Class 2): Mean=0.5562, Std=0.1758
1221
         F1 Score (Class 3): Mean=0.5983, Std=0.1135
1222
         F1 Score (Class 4): Mean=0.7095, Std=0.0985
1223
         MSE: Mean=0.5059, Std=0.1198
1224
         RMSE: Mean=0.7064, Std=0.0831
1225
         R^2: Mean=0.6998, Std=0.0739
1226
         MSE (Class 1-Star): Mean=0.7500, Std=1.2071
1227
         MSE (Class 2-Star): Mean=1.4150, Std=0.8804
1228
         MSE (Class 3-Star): Mean=0.6233, Std=0.2354
         MSE (Class 4-Star): Mean=0.3200, Std=0.1543
1229
         MSE (Class 5-Star): Mean=0.2867, Std=0.1492
1230
1231
      Model: LR_RAW
1232
         AUC: Mean=0.8498, Std=0.0205
1233
         Precision: Mean=0.5952, Std=0.0577
1234
         Recall: Mean=0.5830, Std=0.0550
1235
         F1 Score (Minority Class 1): Mean=0.7133, Std=0.2077
1236
         Overall F1 Score: Mean=0.5816, Std=0.0519
1237
         F1 Score (Class 0): Mean=0.7133, Std=0.2077
1238
         F1 Score (Class 1): Mean=0.2971, Std=0.1682
         F1 Score (Class 2): Mean=0.5032, Std=0.0808
1239
         F1 Score (Class 3): Mean=0.5649, Std=0.0964
1240
         F1 Score (Class 4): Mean=0.6762, Std=0.0871
1241
         MSE: Mean=0.6349, Std=0.0965
```

```
1242
         RMSE: Mean=0.7946, Std=0.0592
1243
         R^2: Mean=0.6236, Std=0.0601
1244
         MSE (Class 1-Star): Mean=0.5667, Std=0.5325
1245
         MSE (Class 2-Star): Mean=1.2500, Std=0.9724
1246
         MSE (Class 3-Star): Mean=0.8044, Std=0.3784
         MSE (Class 4-Star): Mean=0.5667, Std=0.2113
1247
         MSE (Class 5-Star): Mean=0.4667, Std=0.1633
1248
1249
      Model: XGB_RAW
1250
         AUC: Mean=0.8862, Std=0.0250
1251
         Precision: Mean=0.6300, Std=0.0733
1252
         Recall: Mean=0.6344, Std=0.0599
1253
         F1 Score (Minority Class 1): Mean=0.7444, Std=0.2354
1254
         Overall F1 Score: Mean=0.6192, Std=0.0591
1255
         F1 Score (Class 0): Mean=0.7444, Std=0.2354
1256
         F1 Score (Class 1): Mean=0.2876, Std=0.2212
1257
         F1 Score (Class 2): Mean=0.5197, Std=0.1405
         F1 Score (Class 3): Mean=0.6128, Std=0.0784
1258
         F1 Score (Class 4): Mean=0.7362, Std=0.0692
1259
         MSE: Mean=0.5462, Std=0.1627
1260
         RMSE: Mean=0.7307, Std=0.1110
1261
         R^2: Mean=0.6768, Std=0.0949
1262
         MSE (Class 1-Star): Mean=0.9567, Std=1.1450
1263
         MSE (Class 2-Star): Mean=1.4650, Std=0.9116
1264
         MSE (Class 3-Star): Mean=0.7122, Std=0.2693
1265
         MSE (Class 4-Star): Mean=0.3400, Std=0.1873
1266
         MSE (Class 5-Star): Mean=0.2467, Std=0.0945
1267
1268
      Model: TRANSFORMER_RAW
         AUC: Mean=0.8916, Std=0.0340
1269
         Precision: Mean=0.6389, Std=0.0667
1270
         Recall: Mean=0.6279, Std=0.0640
1271
         F1 Score (Minority Class 1): Mean=0.7306, Std=0.2088
1272
         Overall F1 Score: Mean=0.6216, Std=0.0605
1273
         F1 Score (Class 0): Mean=0.7306, Std=0.2088
1274
         F1 Score (Class 1): Mean=0.4323, Std=0.1628
1275
         F1 Score (Class 2): Mean=0.5835, Std=0.0853
1276
         F1 Score (Class 3): Mean=0.5580, Std=0.1231
1277
         F1 Score (Class 4): Mean=0.7237, Std=0.1026
1278
         MSE: Mean=0.6020, Std=0.1734
1279
         RMSE: Mean=0.7680, Std=0.1101
         R^2: Mean=0.6434, Std=0.1027
1280
         MSE (Class 1-Star): Mean=0.3500, Std=0.4493
1281
         MSE (Class 2-Star): Mean=1.0250, Std=1.0724
1282
         MSE (Class 3-Star): Mean=0.5389, Std=0.2358
1283
         MSE (Class 4-Star): Mean=0.6533, Std=0.2596
1284
         MSE (Class 5-Star): Mean=0.5667, Std=0.4384
1285
1286
      Model: SVM OPT
1287
         AUC: Mean=0.9268, Std=0.0203
1288
         Precision: Mean=0.7281, Std=0.0773
1289
         Recall: Mean=0.7208, Std=0.0731
1290
         F1 Score (Minority Class 1): Mean=0.8693, Std=0.1118
         Overall F1 Score: Mean=0.7170, Std=0.0738
1291
         F1 Score (Class 0): Mean=0.8693, Std=0.1118
1292
         F1 Score (Class 1): Mean=0.6912, Std=0.1865
1293
         F1 Score (Class 2): Mean=0.6690, Std=0.1056
1294
         F1 Score (Class 3): Mean=0.6495, Std=0.0939
1295
         F1 Score (Class 4): Mean=0.7642, Std=0.1115
```

```
1296
         MSE: Mean=0.3653, Std=0.1581
1297
         RMSE: Mean=0.5925, Std=0.1192
1298
         R^2: Mean=0.7847, Std=0.0885
1299
         MSE (Class 1-Star): Mean=0.1767, Std=0.2599
1300
         MSE (Class 2-Star): Mean=0.2900, Std=0.3859
         MSE (Class 3-Star): Mean=0.3856, Std=0.2603
1301
         MSE (Class 4-Star): Mean=0.3867, Std=0.1707
1302
         MSE (Class 5-Star): Mean=0.4267, Std=0.3890
1303
1304
      Model: RF OPT
1305
         AUC: Mean=0.9314, Std=0.0239
1306
         Precision: Mean=0.7261, Std=0.0939
1307
         Recall: Mean=0.7290, Std=0.0797
1308
         F1 Score (Minority Class 1): Mean=0.8735, Std=0.0851
1309
         Overall F1 Score: Mean=0.7212, Std=0.0849
1310
         F1 Score (Class 0): Mean=0.8735, Std=0.0851
1311
         F1 Score (Class 1): Mean=0.7075, Std=0.3149
         F1 Score (Class 2): Mean=0.6966, Std=0.0731
1312
         F1 Score (Class 3): Mean=0.6394, Std=0.1204
1313
         F1 Score (Class 4): Mean=0.7655, Std=0.0934
1314
         MSE: Mean=0.3915, Std=0.1433
1315
         RMSE: Mean=0.6165, Std=0.1073
1316
         R^2: Mean=0.7695, Std=0.0795
1317
         MSE (Class 1-Star): Mean=0.3600, Std=0.5426
1318
         MSE (Class 2-Star): Mean=0.4350, Std=0.5473
1319
         MSE (Class 3-Star): Mean=0.4422, Std=0.3120
1320
         MSE (Class 4-Star): Mean=0.4333, Std=0.1961
1321
         MSE (Class 5-Star): Mean=0.3333, Std=0.3651
1322
1323
      Model: LR_OPT
         AUC: Mean=0.9227, Std=0.0224
1324
         Precision: Mean=0.7160, Std=0.0842
1325
         Recall: Mean=0.7146, Std=0.0781
1326
         F1 Score (Minority Class 1): Mean=0.8752, Std=0.0844
1327
         Overall F1 Score: Mean=0.7099, Std=0.0813
1328
         F1 Score (Class 0): Mean=0.8752, Std=0.0844
1329
         F1 Score (Class 1): Mean=0.6672, Std=0.1763
1330
         F1 Score (Class 2): Mean=0.6619, Std=0.1115
1331
         F1 Score (Class 3): Mean=0.6366, Std=0.1294
1332
         F1 Score (Class 4): Mean=0.7643, Std=0.1152
1333
         MSE: Mean=0.3938, Std=0.1491
         RMSE: Mean=0.6175, Std=0.1120
1334
         R^2: Mean=0.7684, Std=0.0827
1335
         MSE (Class 1-Star): Mean=0.1400, Std=0.2538
1336
         MSE (Class 2-Star): Mean=0.4050, Std=0.3343
1337
         MSE (Class 3-Star): Mean=0.4422, Std=0.3666
1338
         MSE (Class 4-Star): Mean=0.4533, Std=0.1950
1339
         MSE (Class 5-Star): Mean=0.4000, Std=0.3944
1340
1341
      Model: XGB OPT
1342
         AUC: Mean=0.9316, Std=0.0247
1343
         Precision: Mean=0.7184, Std=0.0771
1344
         Recall: Mean=0.7125, Std=0.0687
1345
         F1 Score (Minority Class 1): Mean=0.8944, Std=0.0927
         Overall F1 Score: Mean=0.7078, Std=0.0688
1346
         F1 Score (Class 0): Mean=0.8944, Std=0.0927
1347
         F1 Score (Class 1): Mean=0.7598, Std=0.1698
1348
         F1 Score (Class 2): Mean=0.6569, Std=0.0760
1349
         F1 Score (Class 3): Mean=0.6215, Std=0.1125
```

```
1350
         F1 Score (Class 4): Mean=0.7411, Std=0.1006
1351
         MSE: Mean=0.3774, Std=0.0968
1352
         RMSE: Mean=0.6091, Std=0.0798
1353
         R^2: Mean=0.7768, Std=0.0552
1354
         MSE (Class 1-Star): Mean=0.4200, Std=0.5546
         MSE (Class 2-Star): Mean=0.3350, Std=0.3377
1355
         MSE (Class 3-Star): Mean=0.5089, Std=0.3609
1356
         MSE (Class 4-Star): Mean=0.3800, Std=0.1607
1357
         MSE (Class 5-Star): Mean=0.3067, Std=0.2004
1358
1359
      Model: TRANSFORMER_OPT
1360
         AUC: Mean=0.9191, Std=0.0259
1361
         Precision: Mean=0.7253, Std=0.0616
1362
         Recall: Mean=0.7227, Std=0.0564
1363
         F1 Score (Minority Class 1): Mean=0.8876, Std=0.0895
1364
         Overall F1 Score: Mean=0.7175, Std=0.0577
1365
         F1 Score (Class 0): Mean=0.8876, Std=0.0895
         F1 Score (Class 1): Mean=0.7110, Std=0.1411
1366
         F1 Score (Class 2): Mean=0.6959, Std=0.1430
1367
         F1 Score (Class 3): Mean=0.6196, Std=0.1046
1368
         F1 Score (Class 4): Mean=0.7680, Std=0.0816
1369
         MSE: Mean=0.3592, Std=0.0856
1370
         RMSE: Mean=0.5951, Std=0.0714
1371
         R^2: Mean=0.7880, Std=0.0473
1372
         MSE (Class 1-Star): Mean=0.1567, Std=0.2495
1373
         MSE (Class 2-Star): Mean=0.2650, Std=0.3627
1374
         MSE (Class 3-Star): Mean=0.3756, Std=0.2606
1375
         MSE (Class 4-Star): Mean=0.5000, Std=0.1892
1376
         MSE (Class 5-Star): Mean=0.3133, Std=0.1933
1377
      Comparison script finished.
1378
```

Toxic

```
1381
      Model: SVM_RAW
1382
         AUC: Mean=0.9577, Std=0.0081
1383
         F1 Score (Minority Class 1): Mean=0.4974, Std=0.0252
1384
         Overall F1 Score: Mean=0.9320, Std=0.0041
1385
1386
      Model: SVM_OPT
1387
         AUC: Mean=0.8548, Std=0.0253
1388
         F1 Score (Minority Class 1): Mean=0.5069, Std=0.0241
1389
         Overall F1 Score: Mean=0.9375, Std=0.0036
1390
     Model: RF_RAW
1391
         AUC: Mean=0.9222, Std=0.0116
1392
         F1 Score (Minority Class 1): Mean=0.4051, Std=0.0440
1393
         Overall F1 Score: Mean=0.9487, Std=0.0032
1394
1395
      Model: RF_OPT
1396
         AUC: Mean=0.9000, Std=0.0156
1397
         F1 Score (Minority Class 1): Mean=0.5370, Std=0.0350
1398
         Overall F1 Score: Mean=0.9488, Std=0.0041
1399
1400
      Model: LR_RAW
         AUC: Mean=0.9544, Std=0.0080
1401
         F1 Score (Minority Class 1): Mean=0.3958, Std=0.0177
1402
         Overall F1 Score: Mean=0.8969, Std=0.0044
1403
```

```
1404
     Model: LR_OPT
1405
         AUC: Mean=0.9430, Std=0.0122
1406
         F1 Score (Minority Class 1): Mean=0.4931, Std=0.0233
1407
         Overall F1 Score: Mean=0.9357, Std=0.0034
1408
      Model: XGB_RAW
1409
         AUC: Mean=0.9374, Std=0.0102
1410
         F1 Score (Minority Class 1): Mean=0.4326, Std=0.0239
1411
         Overall F1 Score: Mean=0.9182, Std=0.0042
1412
1413
      Model: XGB_OPT
1414
         AUC: Mean=0.9442, Std=0.0095
1415
         F1 Score (Minority Class 1): Mean=0.5184, Std=0.0280
1416
         Overall F1 Score: Mean=0.9427, Std=0.0041
1417
1418
     Model: TRANSFORMER_RAW
1419
         AUC: Mean=0.9536, Std=0.0086
         F1 Score (Minority Class 1): Mean=0.5743, Std=0.0267
1420
         Overall F1 Score: Mean=0.9561, Std=0.0040
1421
1422
     Model: TRANSFORMER_OPT
1423
         AUC: Mean=0.9508, Std=0.0112
1424
         F1 Score (Minority Class 1): Mean=0.5885, Std=0.0229
1425
         Overall F1 Score: Mean=0.9585, Std=0.0035
1426
1427
      Cross-validation with all models completed.
1428
```

D.2 LLM Log

1429 1430

```
1432
      LLM few shot prompting on GoEmotions
1433
      You are an emotion classifier. I will provide you with a text, and
1434
          you need to determine the main emotion expressed in the text
1435
         and output the corresponding index of that emotion.
1436
1437
      Here are some examples:
1438
1439
      Text: "I hope Dallas gets close and loses in heartbreak."
1440
      Emotion: 0. disappointment
1441
1442
      Text: "It could be worse. We could get [NAME] or that Philly
1443
         traitor [NAME] back."
1444
      Emotion: 0. disappointment
1445
      Text: "To be fair, the world (especially politics) has been kind
1446
         of a shit show since 2016"
1447
      Emotion: 0. disappointment
1448
1449
      Text: "Don't let your high expectations of government disappoint
1450
         you."
1451
      Emotion: 0. disappointment
1452
1453
      Text: "I think it was, it was do scary, i honestly never wanna do
1454
         that stuff again"
1455
      Emotion: 0. disappointment
1456
      Text: "You sound upset."
1457
      Emotion: 1. sadness
```

```
1458
1459
      Text: "I'm so sorry. Read about you getting kicked out at home.
1460
         That must be devastating"
1461
      Emotion: 1. sadness
1462
      Text: "I used to do the same exact thing! Now I love the fat on my
1463
          steak and watching them cut it off in japanese restaurants
1464
         makes me sad."
1465
      Emotion: 1. sadness
1466
1467
      Text: "I miss my Shrek Ghrok :("
1468
      Emotion: 1. sadness
1469
1470
      Text: "For some reason I wanted to be a bartender when I was 8.
1471
         After hearing this story, makes me feel I missed out. "
1472
      Emotion: 1. sadness
1473
      Text: "it is very clearly incorrect, and it's apparent that your
1474
         descent into reactionary politics has made you fully
1475
         delusional. "
1476
      Emotion: 2. disapproval
1477
1478
      Text: "This is why I could never work in construction"
1479
      Emotion: 2. disapproval
1480
1481
      Text: "I don't like [NAME] but I'd never wish this fate to any
1482
         parent."
1483
      Emotion: 2. disapproval
1484
      Text: "i say no or in da club if they come to damce with me i walk
1485
          away"
1486
      Emotion: 2. disapproval
1487
1488
      Text: "Not only that, the "improved controls" aren't a thing at
1489
         all. I'm not re-buying Blood Money for nicer graphics. "
1490
      Emotion: 2. disapproval
1491
1492
      Text: "Not gonna lie. Sucked out a few, but am really trying to
1493
         analyze my play afterwards. Thanks!"
1494
      Emotion: 3. gratitude
1495
1496
      Text: "ok thanks I'll give it a read and try to fact check"
      Emotion: 3. gratitude
1497
1498
      Text: "Thanks for the recommendations!"
1499
      Emotion: 3. gratitude
1500
1501
      Text: "Nice, I didn't know that! Thanks for the info."
1502
      Emotion: 3. gratitude
1503
1504
      Text: "Really glad you were there for her. I wish you both the
1505
         best."
1506
      Emotion: 3. gratitude
1507
      Text: "Ah, fair enough."
1508
      Emotion: 4. approval
1509
1510
      Text: "Needed that. A [NAME] ridiculous play is a game changer"
1511
      Emotion: 4. approval
```

```
1512
1513
      Text: "it's horrid :/"
      Emotion: 4. approval
1515
1516
      Text: "I agree with this statement"
1517
      Emotion: 4. approval
1518
      Text: "Came on her face and she told you? Wow disrespect. Either
1519
          she didn't know how you really felt or didn't care. Sorry
1520
         about your luck"
1521
      Emotion: 4. approval
1522
1523
      Emotion labels and their corresponding indices are as follows:
1524
1525
      0. disappointment
1526
      1. sadness
1527
      2. disapproval
1528
      3. gratitude
      4. approval
1529
1530
      Please output only the index corresponding to the emotion, without
1531
           any other content.
1532
1533
      Here is the text to be classified:
1534
1535
      I didn't know that, thank you for teaching me something today!
1536
```

LLM outputs on GoEmotions

1537

```
1539
      Processing text 1/881...
1540
        Text: I'm really sorry about your situation : ( Although I love
1541
            the names Sapphira, Cirilla, and Scarlett!
        Original output: 1
1542
        Predicted emotion index (0-4, -1 \text{ for invalid}): 1
1543
        Mapped predicted emotion index: 25
1544
        Actual emotion index: 25
1545
        Prediction correct
1546
1547
      Processing text 2/881...
1548
        Text: I didn't know that, thank you for teaching me something
1549
            today!
1550
        Original output: 3
1551
        Predicted emotion index (0-4, -1 \text{ for invalid}): 3
        Mapped predicted emotion index: 15
1552
        Actual emotion index: 15
1553
        Prediction correct
1554
1555
      Processing text 3/881...
1556
        Text: Thank you for asking questions and recognizing that there
1557
            may be things that you don't know or understand about police
1558
             tactics. Seriously. Thank you.
1559
        Original output: 3
1560
        Predicted emotion index (0-4, -1 \text{ for invalid}): 3
1561
        Mapped predicted emotion index: 15
1562
        Actual emotion index: 15
1563
        Prediction correct
1564
      Processing text 4/881...
1565
        Text: You're welcome
```

```
Original output: 3
1567
        Predicted emotion index (0-4, -1 \text{ for invalid}): 3
        Mapped predicted emotion index: 15
1569
        Actual emotion index: 15
1570
        Prediction correct
1571
      Processing text 5/881...
1572
        Text: 100%! Congrats on your job too!
1573
        Original output: 4
1574
        Predicted emotion index (0-4, -1 \text{ for invalid}): 4
1575
        Mapped predicted emotion index: 4
1576
        Actual emotion index: 15
1577
        Prediction incorrect
1578
1579
      Processing text 6/881...
1580
        Text: Girlfriend weak as well, that jump was pathetic.
1581
        Original output: 2
        Predicted emotion index (0-4, -1 \text{ for invalid}): 2
1582
        Mapped predicted emotion index: 10
1583
        Actual emotion index: 25
1584
        Prediction incorrect
1586
      . . . . . . . . . . . .
1587
```

LLM few shot prompting on Amazon Reviews

```
1589
1590
      You are a product rating classifier. I will provide you with a
1591
         customer review text, and you need to determine the product
1592
         rating (number of stars) that the customer provided in the
1593
         review, and output the corresponding integer from 1-5.
1594
1595
      Here are some examples:
1596
      Text: "Lesson learned, next time I am going to research this a bit
1597
          more and find foam that will actually LINE UP. This foam is
1598
         the worst for lining up the pieces to create a seamless studio
1599
          look. You can clearly see every cut out when hung up. Will
         definitely not be buying again. Next time I am going to look
1601
         into pyramid foam, which is similar in price, and will
1602
         actually look seamless by the time I am done hanging it on the
1603
          wall.
1604
1605
      Event when some of the egg foam pieces lined up between panels,
         they would eventually get off from one another How is that
1606
         even possible when they started off being lined up??!!??
1607
         Bummer.
1608
1609
      The description should state "foam will not line up when paneling
1610
         together." Do not purchase if you want to think the foam will
1611
         line up!"
1612
      Rating: 1.0
1613
1614
      Text: "I am in pro audio & video for 30 years. I recently bought
1615
         new Bose L1 Model II speakers and needed patch cables. The
1616
         reviews on these are great so I bought 6 cables, some
1617
         different lengths or ends. After 6 months I noticed one
         speaker was lower in volume and then eventually starting
1618
         cutting off. I tested the cables and 3 of the 6 cables have a
1619
         short between the RING & SLEEVE conductors. They tested fine
```

when I received them new. I am now replacing all 6 cables with another brand.

I liked the cables because they seemed sturdy and well made but I can no longer count on them in a professional environment.

If you own and of the Monoprice cables, I recommend that you test them with simple meter for shorts between the conductors.

I take very good care of my cables because I need them to work and I wish these would have worked but they have failed me. Keep testing these cables"

Rating: 1.0

Text: "When I pay top dollar for a Rode product I expect it to work in a professional environment and this product has been a disappointment. Yes, it does isolate the handling noise of a boom pole quite well but the rubber bands do not support the weight of my Rode NTG-3 mic. The mic sags in the bands and tends to fall out of the mount in use. Good thing the mic cable was still holding the mic or my mic would have crashed to the ground several times. Plus, you look like a fool when you flub the take because your mic slid out of the mount. I even wrapped a rubber band around the mic base to "catch" on the mount's rubber and the mic still slid out of the mount. This mount sort of works on a static boom but is pretty much worthless on a hand held boom pole. A nightmare in the field" Rating: 1.0

Text: "Based on the positive reviews I ordered this, but had to return it right away. I've owned another shure wireless mic system with a wireless beta 58 mic and had good experiences, but it was a different design than this. I've also used other wireless mic systems from other manufacturers. From the reviews I thought this would be a good economical way to get another shure system with two mic's. When it arrived I hooked it up for a test and right away it was amplifying every single movement of my hand, and very loudly. Maybe if I leave it in a mic stand I could do away with that, but this is a wireless mic so my intention is to be holding it so I can move around the stage or room. Any slight movement of my hand on any part of the body of the mic is heard throughout the room. It's not acceptable. Picks up noise from handling mic body" Rating: 1.0

Text: "I'm a long time musician, and a long time user of Ernie Ball Strings (on electric guitars). Unfortunately, I can't endorse these strings, even though they sound great and the light gauge saves my fingers. The problem is that they break, in my humble opinion, excessively easily. I put a set on and broke a g string within 2 days of moderate use. It snapped near the saddle, so I didn't think much of it. I replaced the set, and within a few more days I noticed that the wrapping on the g string was broken and becoming unraveled near my third fret. I bought 7 sets, so I still have 5 more sets to go through, but honestly I'll probably go back to elixirs when these are gone, since they seem to last longer.

update

I purchased these strings on 12/19. It' now 1/18, and I've broken 3 strings from 3 different packs.

At this point, I really hate these strings. Break easily"

Rating: 1.0

 Text: "My hopes for this pad was that it would be soft a sqishy on
 my shoulder. The gel is pretty dense, so it's kind of hard
 and squishy and very heavy! I found that it just added more
 weight. I'm not using it. I would look elswhere for a
 different product. Not as good as I expected"

Rating: 2.0

Text: "Works wonderfully. I don't think a snare stand will ever fit completely into any of the pouches if using the large one for the hihat stand, but having it stick out isn't a big deal. My one complaint is the idiotic placement of the shoulder strap loops — one is placed on the bottom and the other at the same spot on the top so that when lifting the bag from the strap it just rolls one way or the other along an axis through the center of the bag. It's very unstable and the handle in addition to the shoulder strap pretty much has to be used when carrying the bag.

Update after a year or two: this piece of gear is shredded on the inside of the pouches. It's extremely frustrating taking the gear in and out. The inside of the pouches is made of some sort of soft felt, which sounds great until you realize that there are no drum stands that have perfectly smooth shapes. Rubber feet, hat-stand spikes, pegs, screws, anything that sticks out can and will get caught on the felt and before long it turns into a web that instead of allowing your hardware to slide in and out makes it a struggle to pack and unpack. A struggle filled with cursing and people waiting to get on the dance floor you just played on because the DJ has set up already during the time it took you to wrestle all your stuff in to the shredded pouches.

Another thing I have found after many, many gigs and a few tours with this thing is that it is simply not conducive to setting up in a confined area. In order to get everything out you have to unroll it and take up a large footprint. So I end up taking up a ton of space while setting up whereas a regular bag or case would take up significantly less space. Very frustrating on a tiny stage or even on a big stage with a large band. Great piece of gear, could be better"

Rating: 2.0

Text: "Bought a brand new Behringer PMP1680S powered mixer and wanted a couple of new quality cables to go along with. Treated these cables like they were made of glass. In other word, very carefully. Hooked em up and I wasn't getting any sound out of the mains. After fooling around for 45 minutes thinking it was something I was doing wrong we gave up and just used the one channel. When I got home, I unscrewed the speakon side only to find that the wires literally fell out of the casing. I assume they were originally attached but the slightest tug must have made them fall out. Yea! China strikes again. I screwed it down and now it appears to be ok. I emailed Pyle Pro and it's been a week ago and ya know what they said......Nothing. THEY NEVER RESPONDED. SPEAKON/SPEAKOFF

Rating: 2.0

1728 1729 Text: "No instructions on using the nut slotting depth benchmark. 1730 All of the files seem to be about the same width, which makes 1731 the smaller string cuts way to wide. The smaller strings 1732 slots need to be much narrower. Files too thick and no instructions" 1733 Rating: 2.0 1734 1735 Text: "I will make this brief. I should have listened to other 1736 reviewers. I had to return it immediately because of the 1737 following. It has a major design flaw. 1738 1739 It is heavy and will not support itself unless the piece that 1740 attaches it to the mic stand has something immediately 1741 underneath it to support its weight, otherwise it will slide 1742 down and the screw will scratch your mic stand badly. 1743 1744 The part in my On Stage mic stand that had this criteria did not fit inside the quiklok piece which is made of metal so it can' 1745 t be modified in any way. So to me it was useless. Doesn't fit 1746 standard mic stands! (major design flaw)" 1747 Rating: 2.0 1748 1749 Text: "We have home parties alot..... 1750 1751 I had a cheap 400 watter that did the trick using Froggy's Swamp 1752 Juice, but it eventually died. No doubt from using too thick 1753 of fog juice. BUT, I absolutely swear by Froggy's. If your 1754 machine can't run it, get a better machine because Froggy's does an amazing job filling the area with fog. 1755 1756 Anyhow, this 1200 watter pushing Froggy's Backwood Bay juice cut 1757 at 2-3 parts juice to 1 part distilled water is really 1758 something else. I hit the button once and the entire party 1759 area is filled - about 900 square feet. I have to turn off 1760 the smoke detectors or it sets them off. If I hit the button 1761 2 or 3 times within 30 minutes, you honestly cannot see 10 1762 feet in front of you. 1763 1764 I tend to hold the button down until it cycles off. Wait about 15 1765 minutes and do it again, but note that the machine is ready 1766 again in about 5. Then maybe once or twice more over the next couple of hours I cycle it again. Then I unplug the machine 1767 and just let the fog hang in the air and eventually dissipate 1768 after maybe 4 hours from that first hit. 1769 1770 Put the right fog juice in THIS machine and it is the bomb! Fills 1771 the room with fog that has great hang time and, cutting it 1772 with distilled water, it really is cost effective over the 1773 cheaper stuff. The right machine and the right juice, cut 1774 down = 5 stars all the way! 1775 1776

Bad things (sorta):

1777

1778

1779

- 1. Had to move the fogger about 20 feet from the action as the fog shoots out really far with FORCE.
- Must turn off the breaker and unplug the smoke detectors. They will surely trip.

3. Placement is a bit difficult given the size and weight of the machine, but especially the 20+ foot stream of forceful fog juice that comes shooting out the machine.

Lastly, I don't care about a timer. Since I only need to hit the button 2 or 3 times total over a 4 hour period, a timer isn't warranted.

UPDATE AUGUST 2013: HAD ONE DIE WITHIN A MONTH. HAD TO DO A RETURN. REPLACEMENT ONE GOING STRONG MONTHS LATER. WOW! Amazing fogger - QC issue"

Rating: 3.0

Text: "PROS:

- * Decent quality strings that tune quickly and hold a tuning well
- * Comfy to play on low EAD end-- which is how silk & steel is supposed to be
- * Affordable price

CONS:

- * Seems to my ear a bit more "tinny" in sound... not quite as deep , vibrant and mellow as standard strings. This surprised me as these are supposed to have a "softer, more mellow" tone (which is why I bought these)
- \star The trebles of course seem no different than any treble of similar measure.
- * Since silk & steel traditionally doesn't wear as long or as well as phosphor bronze, aside from the price I see no advantage to this type of string.
- SPECIAL NOTE: I like how the strings come two-to-a-paper. The E is packed with g, A with b, D with e. At first glance one thinks "that's cheap" however, it's smart marketing for several reasons:
- It reduces cost of production so allows to keep retail costs as low as possible (it does cost something to make and print those sleeves)
- 2) It saves trees. One would not believe how much wood is saved by simple, small conservation steps. I remember a report from Celestial Tea which stated that annually they save something like 1 million trees by not putting paper tags on their teabags. One never thinks about tiny things like that, but if we're going to save forests we need to reduce the use of paper however we can.
- So 5 stars on conservation, but a "they're okay" on the strings themselves. They're by no means a "bad" string; rather good in fact. They just don't live up to the "silk & steel" claims I read before deciding to try out this type of string. I'm probably going to be back with phosphor bronze before long. The low price however might seem attractive for students up front, but the low wear-time in the long run might not prove satisfactory, as from what I've heard these strings need to be replaced about twice as often as standard strings.
- They're worth trying to see if you like the sound better. These are just my personal observations. Overall a good string; they just didn't strike me as "a reason to switch to silk &

1836 steel" on a permanent basis Good strings but not a reason to 1837 switch" 1838 Rating: 3.0 1839 1840 Text: "Sounds ok but the position to rotate the dials is hard for it to be mounted above dslr camera. Plus xlr is not locking. 1841 Meaning accidentally you can pull the cable off and you would 1842 lose recording. 1843 1844 I also did a sound test. 1845 1846 Mics: Sennheiser MKE600 on line 1, Audio Technica AT8035 on line 1847 2, both are set on mic stand and about 2.5' away from the 1848 computer monitor that is playing a feature length film of mine 1849 from 2004 that was ADRed. The setup is in our current decent 1850 ADR room, so it's quiet enough. 1851 1852 Both Tascam and Zoom were providing phantom power to the mics 1853 I was at the other room with the door closed, without doing much 1854 of setting adjustments on either devices (assuming an indie 1855 filmmaker like myself would not have time to go and carefully 1856 set each as needed). I have only set recording to 48khz/24bits 1857 1858 1859 On Tascam I have the gain set in mid and the dial at 3pm line (1860 both channels) 1861 On Zoom, I have the dial set both on 5.1 (it's hard to get exact 1862 as I was having seeing issue - old age). 1863 Volume on both are set high for headphone. 1864 1865 so as you can see, the setting is not ideal but assuming that I'm 1866 going to be running to set a shoot at a location and with the 1867 chaos I may encounter when it's a one man crew or two men crew 1868 set... 1869 1870 btw, mounting the Zoom right above the camera is not a good thing 1871 cause you can't see the dials or the screen (when tripod is 1872 set at 5'6" tall). You will need a cage or additional hot 1873 shoes support to mount the Zoom and the mic around the camera (DSLR in this case), where as the Tascam can be mounted below 1874 so with limited crew members, it is a fairly practical setup 1875 than the Zoom. both limiters were set off and such. 1876 1877 I made sure the audio from the monitor is of the similar lines to 1878 gage the recorded info. Tascam picks up better audio than the 1879 Zoom did (maybe the settings were not correct). I do like the 1880 packaging of the Zoom plus how it saves the audio into two 1881 mono files unlike Tascam saving it the two tracks into a 1882 stereo file with 2 channels. 1883 1884 Placing onto Premiere CS6, as it is, audio audible but quiet, not hearing much of differences, then I added normalization to 1885 both to -12db, Zoom sounded good but with a lot of noise... (1886 some friends of mine told me that the H4N had noise as well, 1887 so not sure if this is the same issue or new problem). 1888

Also: I did notice that for practicality wise, using the Tascam 1891 for a one man crew shoot would work best while the Zoom H6 is 1892 not feasible to access or read the screen when it is mounted 1893 above the camera. (Unless you use one of the magic arm). 1894 Quick disclaimer: I'm not an audio engineering, but just an indie 1895 filmmaker that does a lot of stuff around Ohio and have 1896 learned to always have good audio with video (or worst case to 1897 do ADR... NOT). So I try to invest in decent mics and make 1898 sure get good audio that helps my visual creation... I did 1899 these tests on practicality and real-life environment since 1900 not every shoot will have a perfect sound, perfect room, 1901 perfect environment, etc... Look nice but not practical" 1902 Rating: 3.0 1903 1904 Text: "I have to edit my once five star review because twice in a 1905 row I recieved a tuner in the mail, and after less than one week the screen goes totally dim. I Can't imagine this is the 1906 battery since that is supposed to last a year and I only used 1907 the tuner 4-5 times for less than a minute each time. It's so 1908 dim I can barley see the reading when I squint. It's only \$10 1909 just a pain to wait two weeks for amazon to give a refund. I' 1910 ll be buying another one and trying it again! Good tuner, 1911 unlucky with two lemons I guess. Two were dead after 3 days, 1912 third times the charme!" 1913 Rating: 3.0 1914 1915 Text: "I purchased this six-pack of colored cords several years 1916 ago. All of the cords worked fine when I purchased them, but over the years, I've had to resolder the connections on almost 1917 all of the cords. The problem is due to the conductors 1918 gradually sliding further and further out of the insulation 1919 until there's so much slack, they start coming into contact 1920 with each other at the pins and shorting out the audio signal. 1921 If you're handy with a soldering iron, a relatively easy 1922 resoldering job will fix the problem. I was able to extend the 1923 lifespan of all my cords. 1924 1925 I originally bought these colored mic cords for live use so that I 1926 could quickly trace a performer's mic to the mixer. They 1927 worked great for that, but they also posed a unique problem -the cords can really jump out in photographs! I never would 1928 have realized it, but then I saw pictures people had taken of 1929 our band. The multicolored spaghetti is a bit of an eyesore in 1930 my opinion, especially in flash photos against a dark stage. 1931 (FYI, I didn't dock the product any stars for this, because 1932 the cords are obviously intended to stand out.) 1933 1934 For the reason above, I went back to using all black cords on 1935 stage and moved these GLS into my recording studio. In the 1936 studio, multiple-colored mic cords are a real asset. In the 1937 years since, I've had to make a few additional repairs, but 1938 all six cords are still functional. Very handy, great for 1939 studio use, but longevity is questionable" Rating: 3.0 1940

26

Text: "i was initially curious to know whether these strings were

marketing gimmick of selling old wine in a new bottle. many

really customized for flamenco or if it was all just a

1941

1942

1944 companies just package their products under a different name 1945 or packaging to sell it as if its new. i just decided to go 1946 for it & see for myself. 1947 yes, its true that these strings are different from standard 1948 daddario proarte. while proartes are delicate & rich in harmonics, these strings are like a slap in the face, they've 1949 a lot of bite to cut thru all kinds of noise.u don't hear any 1950 sustaining harmonic trails when u strike a note, they just 1951 strike really sharp. i changed from the proartes to the 1952 flamenco set & there was this spike in the midrange & treble 1953 frequencies. this helps the rhtyhm to really cut through. 1954 especially my picado runs sound almost dirty, like someone 1955 belting the guitar in a frenzy. the sustain is also a bit less 1956 ,so rasguedos don't get muddy, they get short & clear. 1957 so its well suited for flamenco, dont buy it for any classical 1958 playing. the black trebles also look kinda cool. price at \$9 1959 may be a bit steep, but if u r particular about having a sharp biting flamenco sound, this is it. flamenco strings for 1960 flamenco players" 1961 Rating: 4.0 1962 1963 Text: "This little ukulele does the job and doesn't have any 1964 quality problems. It arrived only 3 days after I ordered it, 1965 and in perfect shape. 1966 The only annoying thing about the particular one I received is 1967 that the volumne of the C-string dwarfs all of the other 1968 strings. It's practically all you can hear when you strum (1969 with the open C) while tuned up to pitch. I'm guessing that's 1970 just a fluke of this particular one, like a certain tone might 1971 resonate louder than any other in a particular room. There are better-made ukuleles out there that come with gig-bags 1972 and tuners for only 10-20 dollars more, but they don't have 1973 the goofy pineapple print or cutaway in the headstock. :-) 1974 Not knowing a thing about these prior to getting this one, I'd opt 1975 for a concert sized version next time. Good item for the 1976 price" 1977 Rating: 4.0 1978 1979 Text: "<a data-hook="product-link-linked" class="a-link-normal" 1980 href="/Cordoba-22T-CE-Tenor-Cutaway-Ukulele/dp/B00JPN1XEK/ref= 1981 cm_cr_arp_d_rvw_txt?ie=UTF8">Cordoba 22T-CE Tenor Cutaway UkuleleI gave this Crdoba 4 stars because the set up was 1982 horrible. Saddle as leaning forward so strings would not 1983 intonate. With some repair work and adjustments I was able to 1984 fix it and changing strings took care of the intonation 1985 problem. Now it plays in tune. Euke has plenty of volume and 1986 sustain plugged in my acoustic amp or unplugged. Workmanship 1987 is fine on the rosewood back and sides and solid spruce top. 1988 No fret buzzes. 1989 1990 I bought this instrument used from Amazon with the understanding 1991 that it had no dings, scratches etc. It arrived in a single 1992 cardboard box with no padding and the top of the box was not 1993 taped or secured in anyway. Just amazing it arrived with only a small dent in the lower bout on the back. Given the way it 1994 was packaged it could have been destroyed. I thought about 1995 sending it back but I like the instrument and can live with

workmanship. They just need to do a better job with their set

the ding. Crdoba makes fine instruments with supeb

1996

1998 ups. I am surprised Amazon would ship an instrument in a 1999 flimsy box with little support. There was also no paperwork 2000 with the instrument. ... this Crdoba 4 stars because the set 2001 up was horrible. Saddle as leaning forward so strings would 2002 not ..." Rating: 4.0 2003 2004 Text: "This is the most inexpensive keyboard stand that I have 2005 found. But for the price, you get an item that does what it 2006 should. It is obviously not meant to hold a heavy or large 2007 keyboard, but a 61 key instrument like the kind that Casio or 2008 Yamaha makes so many of will be fine with this. 2009 2010 If you need a more sturdy unit because you are using a keyboard 2011 with more weight to it, then really you should invest in a 2012 more sturdy and expensive model. I noticed that people are 2013 saying that the screws are not included. I at first thought 2014 that this was the case as well as they were not in the bag with the wrench and instructions. For the benefit of anyone 2015 who might be confused like I was, the allen wrench is strapped 2016 to the leg and the screws themselves are in the tops of the ' 2017 X^{\prime} section. You need to remove the screws and then put the 2018 legs on and mount the base legs. Why did they not put all the 2019 items in the plastic bag I do not know, but now you know where 2020 to look for them. 2021 2022 Inexpensive stand. Easy to set up. Works fine for light weight 2023 keyboards. 2024 I am pleased with it. cheap stand that does the job" 2025 Rating: 4.0 2026 Text: "I have these in three guitars after trying out some 2027 2028 2029

undersaddle and sound hole pickups. Some peope miss the "quack " ,electric guitar type sound from these other types. I do not . I also like that the bridge and saddle do not have to be modified , nor is any battery required. I run it through a

Baggs Para-acoustic preamp or my Fishman Aura. It really has to be installed by a professional repairman to be sure the sensors are located correctly.

Product was delivered safely and promptly. Most Natural Sound for the Money"

Rating: 4.0

2030

2031

2032

2033

2034

2035

2036

2037

2038

2039

2040

2041

2042 2043

2044

2045

2046

2047

2048

2049

2050

2051

Text: "Where do I begin??? For starters, I have been playing for almost 20 years and have been doing repairs on solid body electrics and acoustic guitars for 5 years now commercially. Not to sound arrogant - but I KNOW WHAT I'M TALKING ABOUT PEOPLE....

The craftsmanship for this quitar is excellent. The construction is solid feeling and the paint finish is handsomely done. bridge is glued straight and even onto the body with no exposed seams. The nut is precisely cut for the strings - they sit there nice and snug. On cheaply done 12 strings, you see that the nut is not cut very evenly. This guitar feature a NUBONE brand nut & compensated saddle & that makes for great intonation on this guitar and your strings will stay in tune much longer.

```
2052
      The tuners are standard sealed tuners and do well holding the
2053
         guitar in tune. I normally like to replace the tuners
2054
         immediately on my acoustics with "GROVER" brand tuners but
2055
         after playing this guitar for a week straight - no less than 3
2056
          hours a day (seriously folks), I've only had to tune the
2057
         guitar twice.
2058
      The guitar arrived with the action set perfectly. The strings
2059
         were nice and low from the start and still had enough
2060
         clearance for hard strumming with NO BUZZING. The strings
2061
         didn't need changing either since the guitar has D'Addario XL
2062
         strapped on - great bright/jangly sounding highs with nice
2063
         clear lows!!!
2064
2065
      The body has a spruce top with mahogany sides & back. This guitar
2066
         sounds great. I compared it to my friends 18 year old FENDER
2067
         brand 12 string and it sounds better than hers!
                                                           (I told her
2068
         it sounds just as good as her guitar so as not to hurt her
         feelings. But in reality this guitar beats hers....LOL)
2069
2070
      The electronics are not top of the line - obviously. BUT... The
2071
         wiring is neat and the sound when amplified is pretty darn
2072
         good. Naturally, you need to reduce how much gain you use if
2073
         you play really loud on your amp. But that's just the way it
2074
         is when dealing with acoustic guitars anyway. I might in the
2075
         future change it out for a PIEZO system but then again - I
2076
         might not. It holds up well as is. The instrument control
2077
         board also includes a nice handy tuner which always helps & it
2078
          does the job accurately.
      Kudos to the manufacturer for including an output jack for direct
2079
         line connection to a soundboard/mixer besides the standard
2080
         output jack you would plug into your amp. That's real
2081
         versatility and very helpful. You usually don't find such
2082
         versatility except on more expensive instruments!!! Very nice
2083
         , indeed.
2084
2085
        On the control panel for the electronics, the buttons are a
2086
           little stiff when you press them. The knobs are a little too
2087
            small to grab and a little too tight when you turn them.
           You have to have a moderately "soft touch" when pushing the
2089
           buttons and a little more strength when using the knobs.
2090
           Personally speaking, I tend to leave the "tone tweaking"
           alone on the actual guitars I own and always use my amp or
2091
           pedals for tone shaping alteration anyway. So i won't be
2092
           using the buttons or knobs too much anyway.
2093
2094
      Lastly, I wish it came with a pick guard. Oh well, I just ordered
2095
          a nice one with a fancy hummingbird design on it.
2096
      Problem solved.
2097
2098
      What's the bottom line, people???
2099
2100
      BUY THIS GUITAR !!!!!!!!
2101
      You'd be a fool not to. You get so much for under $200.
2102
          heck are u waiting for???
2103
      I'm seriously gonna get a second one. GREAT INSTRUMENT FOR THE
2104
         PRICE !!!"
2105
      Rating: 5.0
```

2122

2123

2124 2125

2126

2127

2128

2129

2130

2131

2132

2133 2134

2135

2136

2137

2138

2139

2140

2141 2142

2143

2144

2145

2147

2148

2149 2150

2151

2152

2153

2154

2155

2156

2157

2158

2159

2106 2107 Text: "My application of this is to plug three guitars in to my 2108 pedalboard (see pics). It is noiseless, and shows no tone-2109 sucking qualities at all. It runs well off a 9-volt or a 2110 Dunlop DC Brick, although that is just for the function of the lights. It also works with no power at all (but the lights 2111 come in handy to know which selection is ON). As good as I 2112 hoped!" 2113 Rating: 5.0 2114 2115 Text: "Being the owner of a couple different Heil microphones (the 2116 Heil Goldline Pro and the Heil Heritage) which never fail to 2117 2118 see what it was capable of as well. I purchased the Heil PR 2119 2120 my ham radio, but also for use in VOIP applications, and 2121

- impress, I decided it was time to step up to the Heil PR-40 to -40 to serve a few different purposes, mainly as a new mic for finally to use for recording narration for videos. It suits all of these purposes just fine, but of course I knew this already since the same source is used for all those applications, that being my voice.
- The mic itself is everything I would expect of a Heil microphone. It build ruggedly, and has a beautiful, flawless finish. The PR-40 comes in a padded leatherette carrying case that also holds the included mic clamp. The clamp has an adapter that screws in to allow the clamp to be used on different sized stands and boom arms. Also included was a Heil Sound decal. I have to say I'm disappointed the Heil mics don't come in the wooden presentation boxes they used to, but I am quite happy with the padded case they use now too.
- The XLR jack on the mic is a little tight, it took a lot of effort to get the Neutrik XLR connector on my mic cable to lock into place. Removing the o-ring from the Neutrik connector on the mic cable allows the cable to lock into the PR-40 with no effort, but I prefer the slight amount of compression the oring provides to help keep dust out as well as prevent any rattling.
- The thing about Heil mics is they pretty much occupy their own audio space in the world of mics. Nothing else comes close to sounding like a Heil mic, their timbre is unmistakable, but in a good way. A lot of people have trouble getting used to Heil mics because they're used to older design mics that need a lot of EQ to make them sound good. With Heil mics, they don't require an obscene amount of EQ to make them sound great, they pretty much sound great with the EQ set flat.
- What I really like about the Heil PR-40 is it has a slightly scooped mid-range that takes the nasal "honk" and stuffiness out of my voice. I've never used a mic that sounded so broadcast-ready right out of box, and much of that is helped by the extended low-frequency response (when compared to most other dynamic mics) that picks up more of the deep chest resonance of the person talking or singing into the mic. Of course, this also makes the mic more prone to picking up low frequency rumble, but installing the mic in the optional accessory shock mount takes care of that problem without having to EQ out the low frequencies, which EQ really should

2160 always be left as a last-ditch bandaid only anyways, as you're 2161 essentially EQ'ing out an entire octave of the human voice. 2162 2163 Proximity effect is really sweet on this mic. With the PR-40 2164 anything more than 6 inches away from the mouth starts 2165 sounding way too thin, with the sweet spot being around 2-4 inches from the grill. However, you can literally get your 2166 lips right up against the grill and the proximity effect is 2167 still controlled to the point that the bass never becomes too 2168 thick or muddy, the mic just takes on an extremely warm, 2169 intimate sound that is absolutely spectacular. 2170 2171 Even though the mic comes with an internal sorbothane shock 2172 mounted capsule, the mic is still fairly prone to picking up 2173 low frequency through the mic stand or boom arm. I highly 2174 recommend getting the accessory shock mount to go with this 2175 mic, I guarantee it will make a world of difference. 2176 What can I say? I'm more impressed by the mic than I thought I 2177 would be. I can gladly say I have absolutely no buyer's 2178 remorse whatsoever. This one is definitely a keeper. 2179 Remarkable Quality & Performance" 2180 Rating: 5.0 2181 2182 Text: "This is the best guitar humidifier on the market. What 2183 makes it the best? 2184 2185 1. The synthetic sponge that holds so much more water then a 2186 typical sponge and it does NOT drip. 2. Single one time purchase with nothing else to buy (like humidi 2187 paks) The Best!" 2188 Rating: 5.0 2189 2190 Text: "I've been playing the harmonica for over 43 years and this 2191 harmonica, made in Germany, exceeds my expectations. I can 2192 bend single notes for melodic playing and with the usual 2193 vibrato over a wide variety of musical genres. Many music 2194 stores sell this same harmonica for \$20 more than what I paid 2195 for it on Amazon. It pays to comparative shop. I've been 2196 playing the harmonica for over 43 years and ..." 2197 Rating: 5.0 2198 Please output only the integer from 1-5 corresponding to the 2199 rating, without any other content. 2200 2201 Here is the text to be classified: 2202 2203 I use several Behringer products (amps & pedals). I got this one 2204 last week and really hated the sound. The OD was WAY too harsh 2205 starting at level 1! 2206 2207 Lots of settings and tones, though. I changed out the tube for an 2208 Electro Harmonix 12AY7 and it sounds MUCH better. Don't expect "true" tube sound. It is, after all, only 1 preamp tube. It 2209 does give a "tube-like" OD sound, though. Close, but not exact 2210 . Good enough and built solid. 2211 2212

I was all ready to ship it back, but I'm now keeping it. Changed

the tube, now I like it.}

2214 2215 2216 **LLM outputs on Amazon Reviews** 2217 2218 Processing text 1/500... 2219 Text: My earlier review was for the Jr2- I don't know why it was 2220 posted for the Jrl! Well, anyway this is a nice guitar for 2221 the money. There is some initial buzzing in the beginning-2222 but not anymore. It has great tone- it's clear and bright. I like that. Great sustain. Great travel guitar for adults. 2223 comes w/ a gig bag. Jrl is a great guitar- buy one NOW! Good 2224 ! For Adults and travel. Jr1 2225 Original output: 5 2226 Predicted rating (1-5, -1 for invalid): 5 2227 Actual rating: 4 2228 Prediction incorrect 2229 2230 Processing text 2/500... 2231 Text: Don't waste your time with these cables. I bought 2 of 2232 them. One of them was already broken and the other broke 2233 after a month or so. 2234 Don't make the same mistake I did: spend a bit more money (even 2235 just \$10 more) and you can get cables that are 1000 times more 2236 reliable than these. To be honest I'm pretty shocked at how 2237 many good reviews these are getting. If I had only bought 1 2238 and it turned out not so good, I might've given them the 2239 benefit of the doubt, but the fact that both cables were 2240 complete duds makes it pretty evident that this is just a bad 2241 product. Bad cables, don't bother with these 2242 Original output: 1 2243 Predicted rating (1-5, -1 for invalid): 1Actual rating: 1 2244 Prediction correct 2245 2246 Processing text 3/500... 2247 Text: I really want to love it...but it's hard to part with \$400 2248 bucks for so little. I think the price on this should be 2249 more in the 250 range. Roland and Mogami...your paying for 2250 the name. That being said it is a great little portable amp 2251 . If you travel a lot it may be worth the investment. The 2252 looper is fine, but it only gives you a 40 sec loop. I like 2253 my boss looper better as I can store many full size songs. 2254 It's not as loud (even plugged in) as my little pignose, but the anti feedback works very well with my t-5 and 2255 acoustic taylors. 2256 I really do recommend this amp, I'm just still reeling from 2257 sticker shock. Nice amp...overpriced 2258 Original output: 4 2259 Predicted rating (1-5, -1 for invalid): 42260 Actual rating: 4 2261 Prediction correct 2262 2263 Processing text 4/500... 2264 Text: I owned the SN-1 tuner and loved it. However, the part 2265 that holds the stem in place broke (note: Don't carry any of these Snark tuners in your pocket!). I replaced it with 2266 the SN-8 because it is supposed to be the better model for 2267 not much more cost. It works okay, but I like the SN-1

```
2268
           better. The display on the SN-1 shows much finer gradations
2269
            of pitch. The SN-8 has much wider bars and does not
2270
           display a steady reading. The SN-8 has been harder to use
2271
           than any of my multiple previous tuners. I have used the SN
           -8 (and the SN-1 before it) on an upright bass, guitar,
2272
           mandolin & banjo and it picks up the pitch in any frequency
2273
           range. Overall, this tuner seems decent for the price.
2274
           Because of the display, I am considering going back to the
2275
           SN-1. As with any tuner, this will get you close, but you
2276
           still need to use your ears for exact fine tuning. Decent
2277
           Tuner For the Price - Like the SN-1 Better
2278
        Original output: 3
2279
        Predicted rating (1-5, -1 \text{ for invalid}): 3
2280
        Actual rating: 3
        Prediction correct
2282
2283
      Processing text 5/500...
        Text: I've always wanted a Gibson Les Paul, but not being a
2284
           professional I wasn't about to spend thousands for one.
2285
           was on the verge of going for a comparable Epiphone model,
2286
           then the OE20TS caught my eye. Reviews of the favorable
2287
           variety swayed me to go this way, as well as the beautiful
2288
           Tobacco SB finish. The price didn't hurt either.
2289
           quitar arrived packed well. No nicks, scratches, or dents,
2290
           and I received the one I ordered (seems to have been a
2291
           problem for some.) 2 for 2 so far. The instrument looks
2292
           wonderful. Nice finish, seems to be put together well, just
2293
            as I had hoped. Then I played it. It was obvious that it
2294
           needed quite a bit of adjustment. There was a ton of fret
2295
           buzz and the intonation was way off. After making
           adjustments to the truss rod and saddles I was ready to go.
2296
            I have to admit, this thing plays great. Fantastic tone
2297
           and great sustain. The action feels sharp and makes it a
2298
           pleasure to play. The trade off is you get a great guitar
2299
           for a great price, it may just need some setup. Don't be
2300
           scared off by this if you are a beginner. There are several
2301
            videos on youtube that show how to make these adjustments,
2302
           and they're really not very difficult. My guitar came with
2303
           a cable, a hex wrench (for truss adjustment), and a warranty
            card. Overall I'm very pleased with this purchase. Good
2305
           guitar becomes great with a little TLC.
2306
        Original output: 4
        Predicted rating (1-5, -1 \text{ for invalid}): 4
2307
        Actual rating: 5
2308
        Prediction incorrect
2309
2310
      Processing text 6/500...
2311
        Text: I am a brand new uke player as of Christmas. Never played
2312
            anything before in my life. That said, I have been keeping
2313
           my instrument up high for two reasons. 1.) have needed to
2314
           see the strings while I learn. 2.) Because I have had to
2315
           hold my instrument with the inside of my right elbow. It has
2316
            also been difficulty for me to focus on the finger/fret
           movement because I have also had to hold the uke with my
2317
           thumb. It also made my shoulder really q tight. Between the
2318
           shoulders and the thumb and the elbow working so hard, I
2319
           could not relax. I just got my strap today. I tried the
2320
           method where the strap goes around the neck and hooks into
2321
           the hole of the uke. That was a bit better tan nothing but
```

2322

2323

2325

2326

2327

2328

2329

2330

2331

2332

2333

2334

2335

2336

2337

2338

2339

2340

2341

2342

2343

2344

2345

2346

2347

2348

2349

2350

2351

23522353

235423552356

2357 2358

2359

2360

2361

23622363

2364

23652366

2367

2368

2369

2370

2371

2372

2373

2374

2375

the neck of the uke kept flopping around. So I tried the method that wraps around the shoulder. That was MUCH better. I LOVE the fact that the second way can hold the neck so it won't flop around and my thumb does not have to actually hold the instrument up. HOWEVER, I am a 45 yr adult adult and there is simply not enough strap left over to begin to wrap around my arm. I already know that if it was longer I could relax my shoulders and drop it down which is probably proper form. I don't want to pick up bad habits while I am teaching myself. So I am ordering ANOTHER strap just like it so I can lengthen it. I realize I could prob go to a fabric store and purchase the same type of strap. But I also want the plastic connectors. I intend to modifiy this. It is a GREAT idea... but does not work for me as manufactured. I hope to come up with a new way that does not go under my armpit because that was a little uncomfortable also when I moved the strap to the arm end just to see how it would feel (even though then it was WAY too short to reach the neck when I did that. This might be alright for a small child, or an adult that wants to use is and put it in the hole. But if you need the neck of the uke held so it does not flop around, this won't work for you... It is cheap enough It is inexpensive enough that it seems it would easy enough for the company to add another \$5-10 and add another foot or so of strap. I am going to end up spending that anyway. Would be nice to come from the manufacturer ready to use and not have to modify it. I will write another review after I see what I do. That was a bit better tan nothing but the neck of the uke kept ... Original output: 3 Predicted rating (1-5, -1 for invalid): 3Actual rating: 4 Prediction incorrect

LLM few shot prompting on Toxic

```
You are a toxicity classifier. I will provide you with a customer
   comment, and you need to determine if the comment is toxic.
Output 1 if the comment is toxic.
Output 0 if the comment is not toxic.
Here are some examples:
Comment: "Wikipedia: Future events
Weighing in on your conversation above; the examples I like to use
    are celestial events and the 2013 inauguration. Lunar and
   solar eclipses are predicted well past the end of this century
   , and the schedule is reliably sourced; there are many pages
   of these, and the possibility of at least one typographical
   error only means they should be checked afainst other
   references. The 2013 inauguation, meanwhile, is years away and
    the inauguree would be speculative, although references to an
    approved budget or scheduling of events would prompt the
   creation of its page, though it may occur before the election,
    and before the inauguree is known, if properly sourced. Wasn'
```

```
2376
         t positive how to link to the policy, but made it the subject
2377
         header. 75.203.136.226"
2378
      Toxicity: 0
2379
2380
      Comment: "LOL Nawlin, you are so predictable it's funny.
2381
      I knew you would spot my edits within a few hours and quickly
2382
         censor them because you are a Commie scum who hates the truth
2383
         . "
2384
      Toxicity: 1
2385
2386
      Comment: "On the discography it is labeled as both an EP (which it
2387
          is) and a studio album (which it also is). While it is
2388
         listed as a different color than the other studio albums, it
2389
         is under the list of Studio albums and also included in the
2390
         total count of studio albums."
2391
      Toxicity: 0
2392
      Comment: "There are two sources already provided for the statement
2393
          that he seems intent on removing, and the issue has been
2394
         discussed in the talk section of the article. The table lists
2395
         National Championships, not NCAA-recognized National
2396
         Championships. Stating that just because the NCAA doesn't yet
2397
         have a sanctioned system in place for recognizing a national
2398
         champion in a sport doesn't make it not a sport. His
2399
         assertions are laughably absurd and he has been reported for
2400
         his violation of 3RR and appears to be using sockpuppets to
2401
         wage an edit war in an article that has already seen too many.
2402
          Please don't indulge this kind of disruptive user."
      Toxicity: 0
2403
2404
      Comment: "Picking on opinions you do not like is obscene
2405
2406
      I followed the discussion. The proposal to userfy all userboxes
2407
         did not succeed, and in the meantime there is nothing ruling
2408
         out political userboxes, not to speak of deleting them from
2409
         the user space. Singling out someone you do not agree with
2410
         like Nicholas Turnbull does is just obscene, and as long as you
2411
          do not delete all the userboxes - which I would be perfectly
2412
         ok with - stop picking on mine."
2413
      Toxicity: 0
2414
      Comment: "Longhair, go take a quick eze. Spread yr misery
2415
         elsewhere."
2416
      Toxicity: 1
2417
2418
      Comment: "Relevant to a mature understanding of this topic is
2419
         Keyfabe - a term from the world of professional wrestling, but
2420
          which applies in a wider context. Individual cases require
2421
         thoughtful judgment, but one thing we should be clear on: not
2422
         everything in tabloids is true. A fair amount of it is staged
2423
          PR fluff. Another portion of it is simply bad reporting that
2424
          the stars don't complain about because it is harmless. There
2425
          are often good reasons to take it all with a grain of salt."
      Toxicity: 0
2426
2427
      Comment: "a region of the celestial sphere close to the ecliptic."
2428
      Toxicity: 0
2429
```

```
2430
      Comment: ""
2431
      Not true! The actual issue was that you deleted ALL the links I
2432
         added for him. That is the recap. The two other links worked
2433
         fine, but you deleted them as well in your hurry to smash
2434
         everything someone else did. In other words, you are the one
         who jumped to conclusions. Now you're trying to claim you only
2435
          deleted the one that was a problem. Even your edit summary is
2436
          wrong. The man's resignation was presented the same day I was
2437
          adding the links, unknown to me, so it seems the Parliament
2438
         link was being moved from current to former which caused
2439
         pointer problems when I clicked on it. Things happen. It only
2440
         needed to be fixed. I would have thought an Admin would be
2441
         capable of figuring that out. Nor could you figure out how to
2442
         leave a message on my Talk page so I could figure it out. You
2443
         just couldn't wait to smash everything. I only came back to
2444
         the article because I was going through the non-Cabinet people
2445
          on the list. Did you fix that list, once you realized this
         was a former MP? No, you did not. That would have been ""work
2446
         "", requiring ""thought"" and ""effort"". Same as you mis-
2447
         corrected the hat note on the other John Carter MP. Based on
2448
         your arrogance, I thought it likely you were an Admin. What a
2449
         surprise, I was right. And Jimbo wonders why the numbers of
2450
         actual contributors are going straight downhill. I have NO
2451
         intention of continuing to contribute to Wikipedia because you
2452
          obviously would prefer to do everything yourself. You're not
2453
         at all welcoming, helpful, polite, assuming good faith, or
2454
         anything else I was led to believe is part of the Wikipedia
2455
         ethos. So I leave you to it. You didn't even have the grace to
          apologize for deleting the other links, you just tried to
2456
2457
         blame me in this audit trail. Go lecture someone else, because
          I'm out of here, and you can be sure I won't be encouraging
2458
         anyone else to participate in what could have been a good
2459
         project. Your attitude is horrendous, and I'm sure I'm not the
2460
          first one you've chased out of here. Well, keep it up and you
2461
         'll soon have the whole thing to yourself. Enjoy. Go brag to
2462
         your friends that you got rid of another contributor. I'm sure
2463
          you're very proud of yourself.
2464
      Toxicity: 0
2465
2466
      Comment: "your conclusion (that Ryu isthe next Akuma or whatever)"
2467
      Toxicity: 0
2468
      Comment: ""
2469
2470
      I have not edited that section and I have not yet reviewed it.
2471
         However, it really does not matter what editors think, what
2472
         matters is what experts think and statements made supported by
2473
          reliable references. If you have them, use them. IF you don't
2474
          , then it is just bantering opinions, which is not fruitful. -
2475
         StormRider ""
2476
      Toxicity: 0
2477
2478
      Comment: ""
      What do you mean ""continue""? I did one edit after, which was for
2479
          something completely unrelated. (talk | ) ""
2480
      Toxicity: 0
2481
2482
      Comment: ""{| class=""wikitable""
2483
     DAMN RIGHT MARK!!! SERVES YOU RIGHT, DONT DO BAD HTINGS!
```

```
2484
      | } " "
2485
      Toxicity: 1
2486
2487
      Please output only the integer 0 or 1 corresponding to the
2488
         toxicity, without any other content.
2489
      You should use your reasoning skills to make your own decisions
         without consulting your memory or other external sources.
2490
2491
      Here is the comment to be classified:
2492
2493
      Comment: "Good Work
2494
2495
      Thank you for your contribution of File: Coastal City School Bus
2496
         crop.JPG to Green Lantern (film). It is nice to see editors
2497
          such as yourself taking the initiative and going out to
2498
         provide us with free content."
2499
      Toxicity:
2500
```

LLM outputs on Toxic

2501

```
2503
          Data loaded successfully. Initial shape: (150722, 3)
2504
      Shape after removing empty comments: (150722, 3)
2505
      Few-shot samples created: 12
2506
      Indices used for few-shot: [6435, 13708, 20264, 22988, 34089,
2507
         38490, 39356, 44010, 48934, 64535, 72318, 104250]
2508
2509
      Target test set: 500 non-toxic, 20 toxic.
2510
2511
     --- DEBUG: Data available for test set (after few-shot exclusion)
2512
2513
     Shape of df_test_candidates: (150710, 3)
2514
     Value counts in df_test_candidates['binary_label']:
     binary label
2515
          143336
2516
            7374
2517
     Name: count, dtype: int64
2518
2519
      --- DEBUG: Test set candidate splits ---
     Number of toxic candidates for test set: 7374
2521
     Number of non-toxic candidates for test set: 143336
2522
2523
      --- DEBUG: Test samples BEFORE shuffle ---
2524
     Total items in test samples (before shuffle): 520
     First 25 labels before shuffle: [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
2525
         1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0]
2526
     Last 25 labels before shuffle: [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
2527
         0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
2528
     Counts before shuffle: 0s=500, 1s=20
2529
2530
      Test set size (after shuffle): 520
2531
     Test set composition (true_labels overall): 20 toxic (label 1),
2532
         500 non-toxic (label 0).
2533
     First 50 true_labels after shuffle: [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
2534
          0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
2535
         Processing text 1/520...
2536
2537
        Comment: Notability of Ambridge country club
```

```
2538
      A tag has been placed on Ambridge country club, requesting that
2539
2540
        Original LLM output: 0
2541
        Predicted label (0/1, -1 \text{ for invalid}): 0
2542
        Actual label: 0
2543
      Processing text 2/520...
2544
        Comment: "
2545
2546
      Pakistani Language (Urdu)
2547
      Urdu in Persian (Parsi) means camp and ''Urdu'' language was the
2548
          langua...
2549
        Original LLM output: 1
2550
        Predicted label (0/1, -1 \text{ for invalid}): 1
2551
        Actual label: 0
2552
2553
      Processing text 3/520...
2554
        Comment: "::::::What I find a bit strange about the various
            guidelines and FAQ pages about categories is tha...
2555
        Original LLM output: 0
2556
        Predicted label (0/1, -1 \text{ for invalid}): 0
2557
        Actual label: 0
2558
2559
      Processing text 4/520...
2560
        Comment: Re: What is it with me and categories (very polite title
2561
            )
2562
      He's a Dark Lord. That's the same thing....
2563
        Original LLM output: 0
2564
        Predicted label (0/1, -1 \text{ for invalid}): 0
        Actual label: 0
2565
2566
      Processing text 5/520...
2567
        Comment: Plagerism
2568
      This entry is almost entirely taken from []. It is copyrighted.
2569
          Please fix the immediat...
2570
        Original LLM output: 0
2571
        Predicted label (0/1, -1 \text{ for invalid}): 0
2572
        Actual label: 0
2573
2574
      Processing text 6/520...
2575
        Comment: "
2576
        ""The use of the bombs""?
2577
2578
      Could the subsection titled ""The use of the bombs"" be renamed ""
2579
          Во...
2580
        Original LLM output: 0
2581
        Predicted label (0/1, -1 \text{ for invalid}): 0
2582
        Actual label: 0
2583
2584
      Processing text 7/520...
2585
        Comment: "
2586
2587
      Hi Smokefoot: Thanks for your comments. It became apparent that I
          could do further edits only aft...
2588
        Original LLM output: 0
2589
        Predicted label (0/1, -1 \text{ for invalid}): 0
2590
        Actual label: 0
2591
```

```
2592
      Processing text 8/520...
2593
         Comment: Photographs
2594
      A couple of photographs, at least, exist of Dilwar. I don't know
2595
          the copyright position ...
2596
        Original LLM output: 0
        Predicted label (0/1, -1 \text{ for invalid}): 0
2597
        Actual label: 0
2598
2599
      Processing text 9/520...
2600
         Comment: 1. search news about steam valve 2. add to steam
2601
            article whilst ignoring usefulness of content.
2602
2603
      NO....
2604
        Original LLM output: 0
2605
        Predicted label (0/1, -1 \text{ for invalid}): 0
2606
        Actual label: 0
2607
      Processing text 10/520...
2608
        Comment: "
2609
2610
       TVMediaInsights
2611
2612
      Recent discoveries show that people who work for this website tend
2613
           to be posti...
2614
        Original LLM output: 0
2615
        Predicted label (0/1, -1 \text{ for invalid}): 0
2616
        Actual label: 0
2617
2618
      . . . . . . . . . . . .
2619
```

- Writing aid and polishing: LLMs were used to assist in improving grammar, clarity, and style. The substantive content, ideas, and technical contributions remain the authors' own.
- **Retrieval and discovery:** LLMs were employed to support literature search and discovery (e.g., identifying related work). All cited references were verified by the authors.