### ProtoLens: Advancing Prototype Learning for Fine-Grained Interpretability in Text Classification

Anonymous ARR submission

#### Abstract

In this work, we propose ProtoLens, a novel prototype-based model that provides fine-grained, sub-sentence level interpretability for text classification. ProtoLens uses a Prototype-aware Span Extraction module to identify relevant text spans associated with learned prototypes and a Prototype Alignment mechanism to ensure prototypes are semantically meaningful throughout training. By aligning the prototype embeddings with humanunderstandable examples, ProtoLens provides interpretable predictions while maintaining competitive accuracy. Extensive experiments demonstrate that ProtoLens outperforms both prototype-based and non-interpretable baselines on multiple text classification benchmarks. Code and data are available at https://anonymous.4open. science/r/ProtoLens-CEOB/.

#### 1 Introduction

001

007

010

011

012

024

Deep neural networks (DNNs) have achieved remarkable success in various natural language processing tasks, including text classification (Kowsari et al., 2019), sentiment analysis (Medhat et al., 2014), and question answering (Allam and Haggag, 2012). However, their black-box nature presents significant challenges for interpretability, limiting their use in high-stakes applications where transparency, user trust, and accountability are paramount (Castelvecchi, 2016; Rudin, 2019). While post-hoc explanation methods attempt to address this (Jacovi et al., 2018; Mishra et al., 2017), they often lack faithfulness and consistency in explaining predictions (Rudin, 2019). In contrast, inherently interpretable models guarantee transparency, facilitating understanding and trust in model outputs (Molnar, 2020).

Among various approaches aimed at enhancing model interpretability, prototype-based methods have emerged as a prominent line of research.



Figure 1: Interpretable classification by ProtoLens.

041

043

045

047

051

053

054

059

060

061

062

063

064

065

066

067

068

069

These methods enable models to generate predictions by comparing inputs to prototypical examples, akin to human reasoning that relies on analogies to familiar cases. While prototype-based approaches have been extensively explored in computer vision (Dong and Xing, 2018; Sumbul et al., 2019; Zhang et al., 2023; Ming et al., 2019a; Gautam et al., 2022; Arik and Pfister, 2020; Willard et al., 2024; Nauta et al., 2023; Ma et al., 2024; Nauta et al., 2021; Xue et al., 2022), their application in natural language processing (NLP) is a relatively new area, with only a few works (Hong et al., 2023; Ming et al., 2019b; Sourati et al., 2023; Arik and Pfister, 2020) emerging in recent years. These models provide an intuitive form of interpretability, facilitating an understanding of predictions through direct reference to interpretable examples. For instance, in a movie review classification task, a prototype might represent a review like "This movie was amazing, with stunning visuals and a gripping storyline," which the model uses to classify new reviews with similar sentiments. The model explains its classification of a new review by highlighting its similarity to this prototypical example.

Despite the potential of prototype-based models for enhancing interpretability, existing approaches encounter significant limitations in text-based applications (Hong et al., 2023; Ming et al., 2019b). Typically, these models **define prototypes at the** 

instance/sentence level, which often lacks the granularity needed for effective interpretability 071 in complex or lengthy texts. For example, in a 072 movie review like "The visuals were stunning, but the plot was too predictable", an instance/sentencelevel prototype might only capture the general sentiment of the review, missing the nuance that the visuals were positive, while the plot had negative aspects. This coarse granularity makes it challenging to provide insightful interpretations when different sentiments or nuances coexist within a single text. In contrast, more fine-grained prototype modeling, such as sub-sentence level, is crucial for delivering detailed interpretative insights, allowing the model to explain specific aspects of the text, like "stunning 084 visuals" or "predictable plot".

> To address this challenge, a novel prototypebased model ProtoLens is designed for finergrained interpretability. ProtoLens builds on the concept of prototypical learning but extends it in key ways that make it better suited for handling the complexities inherent in textual data. The general reasoning process of ProtoLens is illustrated by the example in Figure 1: ProtoLens leverages three prototypes related to "emotion", "performance", and "script", and extracts prototype-specific text spans (sub-sentence level) from the input. Based on extracted spans, Prototype 1 and 2 are activated and thus positive prediction is derived.

087

094

098

100

101

102

103

104

106

108

109

110

111

112

113

There are two core modules in ProtoLens. First, for a specific prototype, the **Prototype-aware Span Extraction** module employs a Dirichlet Process Gaussian Mixture Model (DPGMM) (Görür and Edward Rasmussen, 2010; Rasmussen, 1999) to extract relevant text spans in inputs for model prediction and interpretation. This module enables sub-sentence extraction and offers a more accurate and finer-grained extraction of text spans for certain prototypes. Second, we devise the **Prototype Alignment** mechanism, which adaptively aligns the learned prototype embeddings with representative data samples throughout training. By this, we ensure that learned prototypes are semantically reasonable and effective for interpretation.

114Extensive experiments demonstrate that Pro-115toLens not only outperforms competitive baselines116on multiple text classification benchmarks but also117provides more intuitive and user-friendly explana-118tions for its predictions.

#### 2 Related Work

**Post-hoc Explanations.** Several post-hoc methods interpret DNN models by analyzing gradients or neuron activations, such as Integrated Gradients (Sayres et al., 2019; Qi et al., 2019), DeepLift (Li et al., 2021), and NeuroX (Nalls et al., 2015). Tsang et al. (2018) proposed a hierarchical method to capture interaction effects, later adapted by Jin et al. (2019) for text classification. In sentiment analysis, contextual decomposition (Murdoch et al., 2018) identifies sentiment words and their contributions. Attention-based models, such as Bahdanau (2014); Rocktäschel et al. (2015), analyze attention weights, though Jain and Wallace (2019) question their explanatory power. 119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

Prototype-based Deep Neural Networks. Prototype-based deep neural networks enhance interpretability by using prototypes as intuitive references for decision-making, a concept rooted in traditional models (Sørgaard, 1991; Fikes and Kehler, 1985; Kim et al., 2014). While prototypebased reasoning has been extensively developed in CV, with methods like ProtoPNet (Chen et al., 2019) for image classification and ProtoVAE (Gautam et al., 2022) introducing diverse and interpretable prototypes, its application in NLP is a relatively new area. Early works such as ProSeNet (Ming et al., 2019b) adapted prototype-based reasoning for text classification, followed by ProtoAttend (Arik and Pfister, 2020), which employed attention mechanisms for dynamic prototype selection. Recently, ProtoryNet (Hong et al., 2023) introduced prototype trajectory modeling to improve interpretability across domains. Despite these advances, prototype-based approaches in NLP remain underexplored, making our work a significant step forward in this emerging field.

Unlike previous methods, our approach directly embeds interpretability at the sub-sentence level, providing more granular insights than word- or sentence-level methods.

#### 3 Method

To deliver inherently interpretable predictions at a fine-grained level, we introduce **ProtoLens**, a prototype-based interpretable neural network. ProtoLens is designed to overcome two primary challenges: (C1) How to effectively extract text spans associated with a given prototype to provide interpretable predictions? and (C2) How to ensure



Figure 2: Model Structure. ProtoLens integrates Prototype-aware Span Extraction (via a GMM) and an interpretable classifier. The GMM models the similarity distribution between prototypes and text spans, identifying relevant spans. The classifier aggregates prototype contributions to predict outputs and provide interpretable explanations.

learned prototypes are semantically reasonable and effective for interpretation? To address C1, we propose a Prototype-Aware Span Extraction module, which extracts most relevant text spans for prototypes by a Dirichlet Process Gaussian Mixture Model. To address C2, we design a Prototype Alignment mechanism to adaptively align prototype embeddings to representative data samples through training. The overall model architecture is illustrated in Figure 2.

#### 3.1 Overall Structure

168

170

171

172

173

174

175

177

178

179

180

184

Given a corpus of textual data  $\mathcal{D} = \{(x_i, y_i)\},\$ where i = 1, ..., N, each instance  $x_i$  is associated with a label  $y_i \in \mathcal{Y}$ , our model processes the text through a text encoder, such as BERT (Devlin et al., 2019),  $\psi : \mathcal{X} \to \mathbb{R}^d$ , where  $\mathcal{X}$  represents the space of inputs and d is determined by the encoder.

For a text instance x, it is first inputted to an 185 Prototype-aware Span Extraction module, con-186 taining a set of trainable prototypes  $\mathcal{P} = \{\mathbf{p}_k \in$  $\mathbb{R}^d$ :  $k = 1, \ldots, K$ }, where each prototype is represented by a learnable embedding, and the hyperparameter K is the number of prototypes specified. The model will deliver classifications by comparing the input to these prototypes. For each prototype k, 192 we identify a relevant text span  $x^k \subseteq x$ , which rep-193 resents a sub-sentence capturing the most relevant 194 portion of x associated with that prototype. We then use an encoder  $\psi$  to compute an embedding 196

for each extracted span  $x^k$ :

$$z^k = \psi(x^k), \tag{1}$$

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

The similarity between  $\mathbf{z}^k$  and prototype  $\mathbf{p}_k$  is then computed as  $s^k = \text{RMSNorm}(cos(\mathbf{z}^k, \mathbf{p}_k))$ . The final prediction is computed via an interpretable model f applied to the similarity vector across all prototypes  $\mathbf{s} = [s^1, s^2, \dots, s^K]$ :  $\hat{y} = f(\mathbf{s})$ , where  $\mathbf{s}$  captures the proximity to all prototypes, serving as features for the final prediction; and f can be any interpretable models, such as decision tree or logistic regression. In this paper, we adopt the logistic regression as f.

Model Interpretation. The interpretability of ProtoLens is two-fold. First, it employs prototypes aligned with real-world text sentences to represent human-understandable concepts, assigning weights that reveal their presence and importance in predictions, ensuring intrinsic interpretability. Second, it extracts input spans most relevant to the activated prototypes, allowing users to intuitively compare these spans with the corresponding prototypes for fine-grained interpretability. These prototypes serve as features for an interpretable classifier, such as logistic regression, which provides an additional layer of transparency. Logistic regression assigns interpretable coefficients to each prototype, offering clear insights into how each prototype contributes to the final prediction. As illustrated in Figure 1, ProtoLens highlights spans from the in-

3

274

275

276

put text that relevant to prototypes. For example, spans like "powerful emotions" and "script was well-crafted" align with Prototype 1 and Prototype 3, respectively, contributing positively to the prediction. In contrast, Prototype 2, "The script is dull and uninspired", is not activated and thus has no contribution to the prediction.

#### 3.2 Prototype-aware Span Extraction

226

227

228

236

237

238

239

241

242

243

247

248

252

254

255

264

265

266

267

To extract the most relevant spans of the input text x for each prototype, we divide the input x into n-grams  $x = (c_t)_{t=1}^T$ , where  $c_t$  denotes the t-th n-gram, T is the total number of n-grams, and n is a hyperparameter. A text span is composed of consecutive n-grams. The text encoder processes each part  $c_t \in x$  to produce an embedding  $\mathbf{e}_t = \psi(c_t) \in \mathbb{R}^d$ . The similarity  $m_{t,k}$  between the part embedding  $\mathbf{e}_t$  and the prototype embedding  $\mathbf{p}_k$  is then measured using cosine similarity:  $m_{t,k} = cos(\mathbf{e}_t, \mathbf{p}_k)$ . The intermediate output of the module is the similarity vector between each text input and prototype k, denoted as  $\mathbf{m}_k = (m_{t,k})_{t=1}^T$ .

## 3.2.1 Similarity Distribution Modeling by DPGMM

Identifying the most relevant text spans that align with a prototype is a challenging task due to the inherent complexity and variability of patterns in natural language. The primary aim of employing "fine-grained prototypes" is to extract text spans of flexible lengths, rather than relying on rigid instance/sentence-level, or fixed-size windows.

To address this challenge, we use a Dirichlet Process Gaussian Mixture Model (DPGMM) (Görür and Edward Rasmussen, 2010; Rasmussen, 1999), which represents the relevance between prototypes and text spans as a probability distribution. By modeling similarity distributions in  $m_k$  with Gaussian components, DPGMM provides an effective framework for dynamically identifying highsimilarity regions in the input text, thereby facilitating the extraction of flexible and relevant text spans. DPGMM approximates  $m_k$  using up to G Gaussian components:

$$p(\mathbf{m}_k) = \sum_{g=1}^{G} \pi_g \cdot \mathcal{N}(\mathbf{m}_k \mid \mu_g, \sigma_g), \quad (2)$$

269 where  $\pi_g$  is the mixture weight, and  $\mathcal{N}(\mathbf{m}_k \mid \mu_g, \sigma_g)$  is the Gaussian distribution with mean  $\mu_g$ 270 and standard deviation  $\sigma_g$ . We deploy a neural 271 network based method to learn these parameters following existing works (Viroli and McLachlan, 2019; Bishop, 1994). Specifically, we first learn a hidden representation  $\mathbf{h} = MLP(\mathbf{m}_k)$  and compute these parameters as:

$$\boldsymbol{\mu} = \operatorname{sigmoid}(\mathbf{W}_{\mu}\mathbf{h} + \mathbf{b}_{\mu}) \times T, \qquad (3)$$

$$\boldsymbol{\sigma} = \exp(\mathbf{W}_{\sigma}\mathbf{h} + \mathbf{b}_{\sigma}), \qquad (4) \qquad 2$$

$$\boldsymbol{\nu} = \operatorname{sigmoid}(\mathbf{W}_{\pi}\mathbf{h} + \mathbf{b}_{\pi}), \qquad (5)$$

$$\pi_g = \nu_g \prod_{\ell=1}^{g-1} (1 - \nu_\ell), \quad g = 1, \dots, G.$$
 (6)

Here,  $\mu$  and  $\sigma$  are the parameters for the Gaussian components, while  $\pi$  is determined using the Stick-Breaking Process (Ren et al., 2011), allowing for an adaptive number of components. Detailed explanations can be found in the Appendix A.

#### 3.2.2 Span Extraction

a = 1

To extract a span that focuses on the most relevant area of the text, we select the Gaussian component with the highest mixture weight  $\pi_g = \max(\pi)$ , characterized by  $(\mu_g, \sigma_g)$ . Then,  $\mu_g$  serves as the center of the span, while  $\sigma_g$  defines its length. The span is thus given by:  $x^k = x[\mu_g - \sigma_g, \mu_g + \sigma_g]$ .

#### 3.3 Prototype Alignment

To ensure interpretable classifications, the learned prototypes must be semantically meaningful. However, these prototypes are numerical embeddings that are not inherently interpretable by human users. Therefore, we introduce a prototype alignment mechanism that maps each prototype to real-world training text sentences throughout the learning process.

**Representative Candidates.** We begin by encoding all sentences in the training instances (an instance can contain multiple sentences) into embeddings. In the embedding space, we apply the k-means to cluster sentences. The top 50 sentences closest to each cluster center obtained from k-means serve as representative examples of each cluster, making them suitable candidates for aligning prototypes.

**Prototype Alignment.** In Figure 3, we depict the prototype alignment process in ProtoLens. At one epoch during training, for each prototype with its current learned embedding  $\mathbf{p}_k$ , the top 3 most similar candidate sentences (green circles) from

351

352

353

354

355

356

357

358

360

361

362

364

365

366

367

369

370

371

372

373

374

375

376

378

379

380

381

383

385

386

387

389

390



Figure 3: Prototype Alignment.

the representative candidates are selected. These candidates are averaged to form a representative embedding  $\mathbf{c}_k$  (purple cross), which encapsulates the meaning from actual training data. The prototype is then updated towards  $\mathbf{c}_k$  (orange arrow), resulting in an updated prototype  $\mathbf{p}'_k$  (yellow star).

318

319

321

322

324

327

329

332

333

334

335

336

341

342

Specifically,  $\mathbf{p}_k$  is updated towards  $\mathbf{c}_k$  controlled by a weight factor  $\omega_k$ :

$$\omega_k = \operatorname{sigmoid}(\gamma \cdot (d_k - \tau)), \tag{7}$$

where  $d_k$  represents the Euclidean distance between  $\mathbf{p}_k$  and  $\mathbf{c}_k$ ,  $\tau$  is the movement threshold and  $\gamma$  controls the smoothness of the transition.

The updated prototype  $\mathbf{p}'_k$  is derived as a weighted combination of  $\mathbf{p}_k$  and the movement towards  $\mathbf{c}_k$ :

$$\mathbf{p}_{k}' = \omega_{k} \cdot (\mathbf{p}_{k} + \tau \cdot \mathbf{u}_{k}) + (1 - \omega_{k}) \cdot \mathbf{c}_{k}, \quad (8)$$

where  $\mathbf{u}_k$  is the unit vector pointing from  $\mathbf{p}_k$  to  $\mathbf{c}_k$ , defined as:

$$\mathbf{u}_k = \frac{\mathbf{c}_k - \mathbf{p}_k}{d_k + \epsilon},\tag{9}$$

with  $\epsilon$  being a small value to prevent division by zero. If  $\mathbf{p}_k$  is far from  $\mathbf{c}_k$  (i.e.,  $d_k \geq \tau$ ),  $\mathbf{p}_k$  will move a distance of  $\tau$  toward  $\mathbf{c}_k$ . Conversely, if  $d_k \leq \tau$ ,  $\mathbf{p}_k$  is directly aligned with  $\mathbf{c}_k$ . This process ensures that the prototypes shift toward semantically meaningful regions without abrupt changes.

#### 3.4 Learning Objectives

The learning objectives of the proposed model consist of three key components that contribute to both prediction accuracy and the interpretability of the learned representations.

#### 347 3.4.1 GMM Loss

To approximate complex similarity distributions between text samples and prototypes, we employ a Negative Log-Likelihood (NLL) loss for GMM jointly trained with the model, which is given by:

$$\mathcal{L}_{\text{NLL}} = -\log(\sum_{m=1}^{M} \pi_m \cdot \mathcal{N}(\tilde{s} \mid \mu_m, \sigma_m) + \epsilon),$$
(10)

where  $\pi_m$ ,  $\mu_m$ , and  $\sigma_m$  are the mixture weights, means, and standard deviations of the *m*-th Gaussian component, respectively, and  $\epsilon$  is a small constant added for numerical stability.

The overall loss for the GMM is defined as:

$$\mathcal{L}_{\text{GMM}} = \mathbf{E}[\mathcal{L}_{\text{NLL}}] + \mathcal{L}_{\text{L1}}, \qquad (11)$$

where an  $L_1$  regularization term is introduced to promote sparsity in the mixture weights:  $\mathcal{L}_{L1} = \lambda \sum_{m=1}^{M} |\pi_m|$ , where  $\lambda$  controls the regularization strength. This sparsity encourages the model to focus on a few significant Gaussian components.  $\lambda$ is set to  $1e^{-3}$  for all experiments.

#### 3.4.2 Diversity Loss

To encourage the model to learn high-quality and diverse prototypes, we introduce a **Diversity Loss** based on cosine distance:

$$\mathcal{L}_{\text{div}} = \sum_{i \neq j} (1 - \cos(\mathbf{p}_i, \mathbf{p}_j)).$$
(12)

Maximizing this diversity loss enhances generalization and interpretability by maintaining a diverse set of prototypes.

#### 3.4.3 Overall Objective

The final objective function for the proposed model is a weighted combination of the aforementioned loss components:

$$\mathcal{L} = \text{CrossEntropy}(y, \hat{y}) + \alpha \mathcal{L}_{\text{GMM}} - \beta \mathcal{L}_{\text{div}}, (13)$$

where y represents the true labels,  $\hat{y}$  denotes the prediction,  $\alpha$  and  $\beta$  are hyperparameters that control the balance between accuracy, Gaussian mixture modeling, and prototype diversity.  $\alpha$  and  $\beta$  is set to  $1e^{-1}$  and  $1e^{-3}$  for all experiments, respectively.

#### **4** Experiments

In this section, we conduct comprehensive experiments to evaluate the proposed model and answer the following research questions: **RQ1**: How does ProtoLens perform in terms of classification accuracy compared to state-of-the-art (SOTA) baselines? **RQ2**: What is the quality of the model interpretations? **RQ3**: What are the effects of the proposed Prototype Alignment mechanism and Diversity loss on ProtoLens? RQ4: What are the impacts of different hyperparameters on ProtoLens?

#### 4.1 Experimental Setup

396 397

400

426

Datasets. We evaluate ProtoLens on seven diverse text classification datasets spanning single-label, multi-label, and domain-specific classification tasks: IMDB, Yelp, Amazon, Hotel, Steam, DBPedia, and Consumer Complaint. Details are provided in Appendix B.

Reproducibility. The ProtoLens model was imple-401 mented using PyTorch. For training, the prototype 402 number K is selected from  $\{10, 20, 40, 50, 100\}$ . 403 The learning rate is selected from  $\{1e - 4, 1e -$ 404 5, 5e-5, with a decay of 10% every 10 epochs. 405 We used the AdamW optimizer (Loshchilov, 2017) 406 with a batch size of 16 for 25 epochs and 407 the n-gram size is selected from  $\{1, 3, 5, 7, 9\}$ . 408 The experiments were conducted on an NVIDIA 409 A100 80GB GPU. Code and data are available 410 at https://anonymous.4open.science/ 411 r/ProtoLens-CE0B/. 412

**Baselines.** We compare ProtoLens against a range 413 of baselines, encompassing both interpretable and 414 non-interpretable models. The interpretable base-415 lines include ProSeNet (Ming et al., 2019b) and 416 ProtoryNet (Hong et al., 2023), both are SOTA 417 prototype-based methods that provide insights into 418 their predictions via learned prototypical repre-419 sentations. Additionally, we include a zero-shot 420 Llama-3-8b (Touvron et al., 2023), MPNet (Song 421 et al., 2020a) and a Bag-of-Words model (Zhang 422 et al., 2010) using TF-IDF representations and Lo-423 gistic Regression for interpretable classification 424 (Hosmer Jr et al., 2013). 425

4.2 Prediction Accuracy (RQ1)

We evaluate the accuracy of ProtoLens against 427 several competitive baselines, including both 428 prototype-based and non-prototype-based methods. 429 The results are presented in Table 1. ProtoLens 430 consistently achieves the highest scores, outper-431 forming the baselines in all cases. The consis-432 433 tently higher performance of ProtoLens demonstrates its effectiveness and robustness across di-434 verse domains, highlighting its superiority in lever-435 aging fine-grained interpretability without sacrific-436 ing predictive accuracy. 437

	Prototype Aligned Interpretation				
e 0	top-3 representative candidates	Contribution to Positive Class			
Prototyp	<ol> <li>He, too does an excellent job in this movie.</li> <li>I have a lot of respect for his acting after viewing his performance in this movie.</li> <li>I was deeply impressed with the character he played.</li> </ol>	0.843			
Prototype 1	<ol> <li>This is supposed to be a horror film, but it's lacking in that area and isn't the least bit scary.</li> <li>I happen to to be a horror movie fan, but this film was just so poor, words fail me.</li> <li>Don't waste your time - even the tried and true horroring trigue classics fail in this movie.</li> </ol>	-0.809			
Prototype 2	<ol> <li>This show gave great laughs in premieres, and it still does during re-runs.</li> <li>When we started watching this series on cable, I had no idea how addictive it would be.</li> <li>it was actually a pretty funny show.</li> </ol>	0.723			
Prototype 3	<ol> <li>It lacks substance and style!</li> <li>It's silly, not thoughtful, and boring.</li> <li>It's talky, illogical, slow, and ultimately very boring.</li> </ol>	-0.956			
Prototype 4	<ol> <li>This movie was poorly acted, poorly filmed, poorly written, overall horribly executed.</li> <li>The film itself is poorly constructed and acted.</li> <li>The plot is slashed to bits and the acting is horrible.</li> </ol>	-0.854			

Figure 4: Sampled aligned interpretation of prototypes and their top-3 activated prototypes with aligned text sentences from the training set, with sentiment scores. Each prototype captures a distinct concept, and the aligned sentences provide interpretable explanations linked to sentiment contributions.

#### 4.3 Model Interpretations (RQ2)

ProtoLens offers two-fold interpretability. First, it uses prototypes aligned with training sentences to represent concepts with weights, revealing their presence and importance in predictions for intrinsic interpretability. Second, it extracts input spans most relevant to activated prototypes, allowing users to intuitively compare spans with prototypes for fine-grained interpretability. 438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

#### 4.3.1 Prototype Interpretation

In this section, we present an example of ProtoLens trained on the IMDB dataset with K = 10 prototypes. Figure 4 showcases five randomly selected prototypes along with their aligned sentence interpretations. These prototypes span a wide range of concepts, including acting, horror elements, humor, storyline, and film execution.

What stands out is that ProtoLens achieves high accuracy while relying on concise and interpretable prototypes, often represented by short sentences. This allows for rapid and straightforward comprehension of the model's reasoning process. Each prototype captures key characteristics of the corresponding text, providing insightful interpretations for various aspects of the movie, such as acting quality, humor, or poor execution. This feature enhances both the model's interpretability and us-

Table 1: Performance of ProtoLens in comparison with baselines.

Model	IMDB	Amazon	Yelp	Hotel	Steam	DBPedia	Consumer
Llama-3-8b	0.813	0.767	0.787	0.787	0.667	0.768	0.807
MPNet	0.846	0.899	0.950	0.961	0.913	0.991	0.933
Bag-of-words	0.877	0.830	0.908	0.905	0.844	0.993	0.930
ProSeNet	0.863	0.875	0.932	0.930	0.834	0.984	0.878
ProtoryNet	0.871	0.890	0.941	0.949	0.876	0.991	0.927
ProtoLens	0.903*	0.937*	0.962*	0.963*	0.931*	0.995*	0.945*

#### **Positive Class Text Instance**

465

466

467

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

Cosimo (Luis Guzmán) is told in prison about a perfect heist.Since he's behind bars and can't do it himself he has to leave it to his girl Rosalind (Patricia Clarkson).Soon there are five guys organizing the crime- five guys with very little brain capacity.Brothers Anthony and Joe Russo are the directors of Welcome to Collinwood (2002).It's a crime comedy that's often very funny.You can't help but laughing when everything goes wrong with these guys.There are some great actors playing these characters.William H.Macy plays Riley.Isaiah Washington is Leon.Sam Rockwell is Pero.Michael Jeter is Toto.Andy Davoli is Basil.Gabrielle Union plays his love interest Michelle.Jennifer Esposito plays Pero's love interest Carmela.George Clooney (also producer) plays Jerzy, the tattooed guy in a wheelchair. This is a highly entertaining flick.I certainly recommend it.

Top-3 Activated	Prototype 0	Prototype 2	Prototype 5	
Prototype				
	1. This is a very entertaining film with lots of comedy and plenty of laughs.	1. This film had some very funny moments.	<ol> <li>Some very good character actors in this fine film.</li> </ol>	
Aligned	2. I thought the whole film was decent and interesting.	2. Some of the comedy parts are really funny.	2. The acting by all of these actors is very good.	
	<ol><li>Overall, this is a fun film &amp; I highly recommend it.</li></ol>	3. A couple of the scenes are funny.	3. The actors deliver solid enough performances.	
Extracted Span	"highly entertaining flick"	"crime comedy that's often very funny"	"some great actors playing these characters"	
Similarity	0.708	0.549	0.730	
Contribution to	0.985	0.247	0.931	
Positive Class				
Prediction	Positive			

Figure 5: Case study of a positive class text instance.ProtoLens identifies relevant prototypes (e.g., "highly entertaining flick") and aligns them with specific spans in the input text. Extracted spans, similarity scores, and sentiment weights show how each prototype contributes to the positive prediction.

ability, as users can easily relate the prototypes to human-understandable concepts, making the predictions more transparent. Further examples and an in-depth analysis of prototype interpretations can be found in Appendix C.

#### 4.3.2 Classification Interpretation

When conducting classification on a text sample,
ProtoLens extracts the most relevant span from the sample for all prototypes. Similarities between spans and prototypes are then calculated to determine which concepts are activated for the sample. Last, interpretable classification is delivered based on the similarities. We present a positive example in Figure 5 and a negative example in Figure 6, both from the IMDB dataset.

As shown in Figure 5, the top three prototypes with the highest similarity scores significantly influence the classification. Prototype 0 captures the concept of a "highly entertaining flick" (similarity score 0.708, sentiment weight 0.985), Prototype 2 reflects humor with the span "crime comedy that's often very funny" (score 0.549, weight 0.247), and Prototype 5 highlights good acting with "some great actors playing these characters" (score 0.730, weight 0.931). These prototypes, focusing on entertainment, comedy, and acting, lead the model to correctly predict a "Positive" sentiment. In contrast, Figure 6 shows a negative example. Table 2: Performance of ProtoLens with different ablation settings on various datasets.

Dataset	ProtoLens	w/o Diversity	w/o Alignment
IMDB	0.903	0.882	0.886
Amazon	0.937	0.926	0.927
Yelp	0.962	0.931	0.943
Hotel	0.963	0.947	0.953
Steam	0.931	0.917	0.923

493

494

495

496

497

498

499

500

501

502

503

504

506

507

508

509

510

511

The text activates prototype 4, reflecting dissatisfaction with special effects, as captured in the span "problems with this film: 1 cheap special effects," with a similarity score of 0.657 and a sentiment weight of -0.956. Prototype 7 reflects frustration with the movie, highlighted by the span "ended up watching it the whole 2 hours," scoring 0.676 with a weight of -0.809. Prototype 9 captures disappointment with the lack of character development, aligned with the span "there was no character development," with a similarity score of 0.664 and weight of -0.756. These prototypes highlight negative aspects of the movie, leading the model to correctly predict the sentiment as "Negative". Further examples and an in-depth analysis of classification interpretations are shown in Appendix C.

#### 4.4 Ablation Study (RQ3)

To demonstrate the effectiveness of the Prototype Alignment and Diversity Loss, we compare Pro-

#### **Negative Class Text Instance**

I was sitting at home and flipping channels when I ran across what potentially sounded like an interesting film. I like Destruction type movies and decided to watch it. I don't know why but I ended up watching it the whole 2 hours. We have seen this type of movie I don't know how many times. Back in 1998 - 2000 there were dozen of films that dealt with global destruction of some sort. The best one on my list so far is Deep Impact which was more believable than this one. Here are my problems with this film: 1 cheap special effects, like something out of the old computer. 2 no background information or explanation on weather patterns. If you are going to make a movie about weather, at least have some decency to entertain the viewer with technical details. 3 How come only 2 or 3 people figure out that the storm is converging on Chicago... no more experts left in the field? 4 where are some interesting characters? I truly don't care for anyone except maybe the pregnant woman I felt that there was no character development. 5 no thought provoking moment what so ever and factually incorrect theme. And this is only the first part of the film. bet the conclusion will show us few destruction scenes and a search and rescue operation just like it has been done many times before. And judging by the special effects in the first part of the movie, I can only imagine what we are to expect. Of course, at the end, the main characters will survive and life will go on... how origina

	Top-3 Activated Prototype	Prototype 4	Prototype 7	Prototype 9
		1. But instead the Special Effects are poorly done.	1. I had to force myself to sit through it.	1. I just found it incoherent, tasteless, and boring.
	Aligned	<ol><li>The special effects are unconvincing.</li></ol>	2. I actually forced myself to watch the rest of it hoping it would get better.	2. It was too plain boring, uninteresting and unnecessary.
	Interpretation	<ol><li>Very poor and disorienting camera work and editing.</li></ol>	3. In fact, I stopped watching it halfway through, which is something I rarely do	3. It was too slow, too predictable, and not moving enough.
	Extracted Span	"problems with this film: 1 cheap special effects, like something out of the old computer"	"ended up watching it the whole 2 hours"	"there was no character development"
	Similarity	0.657	0.676	0.664
1	Contribution to Positive Class	-0.956	-0.809	-0.756
	Prediction		Negative	1

Figure 6: Case study of a negative class text instance. ProtoLens identifies relevant prototypes (e.g., "there was no character development") and aligns them with specific spans in the input text. Extracted spans, similarity scores, and sentiment weights show how each prototype contributes to the negative prediction.



Figure 7: Performance of ProtoLens in comparison with different number of prototypes. Performance improves with more prototypes, peaking at an optimal K (e.g., 40 for IMDB), before stabilizing or slightly decreasing.

toLens trained with and without these components. Prototype Alignment ensures that prototypes main-513 tain their semantic faithfulness. The Diversity Loss encourages prototypes to be distinct, reducing redundancy in representation. The results, shown in Table 2, indicate that both the Prototype Alignment and Diversity Loss are essential for maintaining ProtoLens's high performance and interpretability, as their removal leads to significant declines in accuracy across datasets. A detailed analysis is provided in Appendix E.

#### 4.5 Hyperparameter (RQ4)

512

514

515

516

517

518

520

521

522

525

527

529

531

We explored the impact of varying the number of prototypes K and n-gram sizes on ProtoLens's performance, identifying dataset-specific optimal values that balance model complexity and classification accuracy. In conclusion, the optimal number of prototypes K varies by dataset, with K = 50performing best for Amazon and Yelp, K = 40 for IMDB, and K = 20 for Hotel, while an n-gram



Figure 8: Accuracy of ProtoLens across IMDB, Amazon, and Hotel datasets as n-gram size varies. Larger n-grams improve contextual representation, but performance plateaus or slightly decreases beyond n=5, indicating a tradeoff between context and generalizability.

size of 5 consistently yields the best results across all datasets, balancing complexity and performance. A detailed analysis is provided in Appendix F.

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

#### 5 Conclusion

In this paper, we present ProtoLens, a prototypebased model offering fine-grained, sub-sentence level interpretability for text classification. we introduce a Prototype-aware Span Extraction module with a Prototype Alignment mechanism to ensure prototypes remain semantically meaningful and aligned with human-understandable examples. Extensive experiments across multiple benchmarks show that ProtoLens outperforms both prototypebased and non-interpretable baselines in accuracy while providing more intuitive and detailed explanations.

#### 6 Limitations

548

551

552

553

554

558

559

563

564

565

570

571

572

575

577

579

582

585

586

587

591

593

594

597

While ProtoLens offers significant advancements in interpretability through prototype-based reasoning and fine-grained sub-sentence level analysis, there are several limitations to consider. First, the quality of the learned prototypes heavily depends on the training data. If the data contains inherent biases, these biases may be reflected in the prototypes, potentially leading to biased predictions or explanations. This limitation underscores the importance of careful data curation and ongoing monitoring of the model's outputs to mitigate bias.

Second, ProtoLens currently focuses on text classification tasks and has not yet been evaluated on more complex natural language processing (NLP) tasks such as machine translation or summarization. Adapting ProtoLens to these tasks may require significant architectural changes to maintain interpretability without compromising performance.

Additionally, while we include results using large language models (LLMs) in a zero-shot setting, we have not yet explored their capabilities in fine-tuning, or in-context learning scenarios. A thorough comparison of ProtoLens across these settings with LLMs could provide deeper insights into its robustness, scalability, and utility in diverse tasks.

Future work could address these limitations by developing methods to automatically detect and mitigate biases, adapting ProtoLens to more complex tasks, conducting comprehensive comparisons across LLM learning settings, and improving the efficiency and usability of the learned interpretations.

#### 7 Ethics

We have carefully considered the ethical implications of our work. ProtoLens is designed to enhance interpretability in deep neural networks, particularly for text classification tasks. By providing more transparent and intuitive explanations, ProtoLens aims to improve trust and accountability in AI systems, which is crucial in high-stakes applications such as healthcare, legal, and financial domains.

We are committed to ensuring that the use of ProtoLens is aligned with ethical standards, promoting transparency and fairness in decision-making processes. However, as with all AI models, there is a potential risk of misuse or bias amplification if the model is trained on biased data. To mitigate this, we emphasize the importance of careful data curation and ongoing monitoring of model outputs to identify and address any unintended biases. We encourage users of ProtoLens to conduct thorough bias audits and ensure that the model is applied in a fair and responsible manner. 598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

645

646

Furthermore, the datasets used in our experiments, including IMDB, Yelp, Amazon, Hotel, and Steam reviews, are publicly available and widely used in the research community. We have ensured that no personally identifiable information (PII) is present in the data, and that our use of these datasets complies with relevant ethical guidelines.

In conclusion, we believe that ProtoLens contributes positively to the field of interpretable AI by improving transparency and user understanding. We acknowledge the importance of continuously evaluating and mitigating potential risks to ensure that AI systems remain fair, accountable, and ethical in their applications.

#### References

Ali Mohamed Nabil Allam and Mohamed Hassan Hag-	619
gag. 2012. The question answering systems: A sur-	620
vey. <i>International Journal of Research and Reviews</i>	621
<i>in Information Sciences (IJRRIS)</i> , 2(3).	622
Sercan O Arik and Tomas Pfister. 2020. Protoattend:	623
Attention-based prototypical learning. <i>Journal of</i>	624
<i>Machine Learning Research</i> , 21(210):1–35.	625
Dzmitry Bahdanau. 2014. Neural machine translation	626
by jointly learning to align and translate. <i>arXiv</i>	627
<i>preprint arXiv:1409.0473</i> .	628
Christopher M Bishop. 1994. Mixture density networks.	629
Davide Castelvecchi. 2016. Can we open the black box of ai? <i>Nature News</i> , 538(7623):20.	630 631
Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett,	632
Cynthia Rudin, and Jonathan K Su. 2019. This looks	633
like that: deep learning for interpretable image recog-	634
nition. <i>Advances in neural information processing</i>	635
<i>systems</i> , 32.	636
Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	637
Kristina Toutanova. 2019. Bert: Pre-training of deep	638
bidirectional transformers for language understand-	639
ing.	640
Nanqing Dong and Eric P Xing. 2018. Few-shot seman-	641
tic segmentation with prototype learning. In <i>BMVC</i> ,	642
volume 3, page 4.	643
Richard Fikes and Tom Kehler, 1985. The role of frame-	644

Richard Fikes and Tom Kehler. 1985. The role of framebased representation in reasoning. *Communications of the ACM*, 28(9):904–920.

Srishti Gautam, Ahcene Boubekki, Stine Hansen,

Suaiba Salahuddin, Robert Jenssen, Marina Höhne,

and Michael Kampffmeyer. 2022. Protovae: A

trustworthy self-explainable prototypical variational

model. Advances in Neural Information Processing

Dilan Görür and Carl Edward Rasmussen. 2010. Dirich-

Dat Hong, Tong Wang, and Stephen Baek. 2023.

David W Hosmer Jr, Stanley Lemeshow, and Rodney X

Alon Jacovi, Oren Sar Shalom, and Yoav Gold-

Sarthak Jain and Byron C Wallace. 2019. Attention is

not explanation. arXiv preprint arXiv:1902.10186.

Xisen Jin, Zhongyu Wei, Junyi Du, Xiangyang Xue,

and Xiang Ren. 2019. Towards hierarchical im-

portance attribution: Explaining compositional se-

mantics for neural sequence models. arXiv preprint

Been Kim, Cynthia Rudin, and Julie A Shah. 2014.

The bayesian case model: A generative approach

for case-based reasoning and prototype classification.

Advances in neural information processing systems,

Kamran Kowsari, Kiana Jafari Meimandi, Mojtaba Heidarysafa, Sanjana Mendu, Laura Barnes, and Donald

Junbing Li, Changqing Zhang, Joey Tianyi Zhou,

Huazhu Fu, Shuyin Xia, and Qinghua Hu. 2021.

Deep-lift: Deep label-specific feature learning for

image annotation. IEEE transactions on Cybernetics,

I Loshchilov. 2017. Decoupled weight decay regulariza-

Chiyu Ma, Jon Donnelly, Wenjun Liu, Soroush

Vosoughi, Cynthia Rudin, and Chaofan Chen. 2024.

Interpretable image classification with adaptive

Walaa Medhat, Ahmed Hassan, and Hoda Korashy.

2014. Sentiment analysis algorithms and applica-

tions: A survey. Ain Shams engineering journal,

tion. arXiv preprint arXiv:1711.05101.

prototype-based vision transformers.

Brown. 2019. Text classification algorithms: A sur-

berg. 2018. Understanding convolutional neural

networks for text classification. arXiv preprint

Sturdivant. 2013. Applied logistic regression. John

Protorynet-interpretable text classification via pro-

totype trajectories. Journal of Machine Learning

let process gaussian mixture models: Choice of the

base distribution. Journal of Computer Science and

Systems, 35:17940–17952.

Technology, 25(4):653-664.

Researcyh, 24(264):1–39.

Wiley & Sons.

arXiv:1809.08037.

arXiv:1911.06194.

vey. Information, 10(4):150.

52(8):7732-7741.

5(4):1093-1113.

27.

- 665

- 671
- 672 673

674

- 675 676
- 677 678 679

683

684

- 690

Yao Ming, Panpan Xu, Furui Cheng, Huamin Qu, and Liu Ren. 2019a. Protosteer: Steering deep sequence model with prototypes. IEEE transactions on visualization and computer graphics, 26(1):238–248.

699

700

701

703

705

708

709

710

711

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

749

750

- Yao Ming, Panpan Xu, Huamin Qu, and Liu Ren. 2019b. Interpretable and steerable sequence learning via prototypes. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pages 903-913.
- Saumitra Mishra, Bob L Sturm, and Simon Dixon. 2017. Local interpretable model-agnostic explanations for music content analysis. In ISMIR, volume 53, pages 537-543.
- Christoph Molnar. 2020. Interpretable machine learn*ing*. Lulu. com.
- W James Murdoch, Peter J Liu, and Bin Yu. 2018. Beyond word importance: Contextual decomposition to extract interactions from lstms. arXiv preprint arXiv:1801.05453.
- Mike A Nalls, Jose Bras, Dena G Hernandez, Margaux F Keller, Elisa Majounie, Alan E Renton, Mohamad Saad, Iris Jansen, Rita Guerreiro, Steven Lubbe, et al. 2015. Neurox, a fast and efficient genotyping platform for investigation of neurodegenerative diseases. *Neurobiology of aging*, 36(3):1605–e7.
- Meike Nauta, Jörg Schlötterer, Maurice Van Keulen, and Christin Seifert. 2023. Pip-net: Patch-based intuitive prototypes for interpretable image classification. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 2744-2753.
- Meike Nauta, Ron van Bree, and Christin Seifert. 2021. Neural prototype trees for interpretable fine-grained image recognition.
- Zhongang Qi, Saeed Khorram, and Fuxin Li. 2019. Visualizing deep networks by optimizing with integrated gradients. In CVPR workshops, volume 2, pages 1-4.
- Carl Rasmussen. 1999. The infinite gaussian mixture model. Advances in neural information processing systems, 12.
- Lu Ren, Lan Du, Lawrence Carin, and David B Dunson. 2011. Logistic stick-breaking process. Journal of Machine Learning Research, 12(1).
- Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomáš Kočiský, and Phil Blunsom. 2015. Reasoning about entailment with neural attention. arXiv preprint arXiv:1509.06664.
- Cynthia Rudin. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. Nature machine intelligence, 1(5):206-215.

843

844

845

846

847

848

849

850

851

852

805

806

761

751

752

754

755

771 772 774

775

- 776 777 778 782
- 790
- 794
- 795
- 797
- 798

- 803

- Rory Sayres, Ankur Taly, Ehsan Rahimy, Katy Blumer, David Coz, Naama Hammel, Jonathan Krause, Arunachalam Narayanaswamy, Zahra Rastegar, Derek Wu, et al. 2019. Using a deep learning algorithm and integrated gradients explanation to assist grading for diabetic retinopathy. *Ophthalmology*, 126(4):552–564.
  - Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020a. Mpnet: Masked and permuted pretraining for language understanding.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020b. Mpnet: Masked and permuted pretraining for language understanding. Advances in neural information processing systems, 33:16857-16867
- Pål Sørgaard. 1991. Evaluating expert system prototypes. AI & society, 5:3-17.
- Zhivar Sourati, Darshan Deshpande, Filip Ilievski, Kiril Gashteovski, and Sascha Saralajew. 2023. Robust text classification: Analyzing prototype-based networks. arXiv preprint arXiv:2311.06647.
- Gencer Sumbul, Marcela Charfuelan, Begum Demir, and Volker Markl. 2019. Bigearthnet: A large-scale benchmark archive for remote sensing image understanding. In IGARSS 2019 - 2019 IEEE International Geoscience and Remote Sensing Symposium. IEEE.
- Yee Whye Teh et al. 2010. Dirichlet process. Encyclopedia of machine learning, 1063:280-287.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Michael Tsang, Youbang Sun, Dongxu Ren, and Yan Liu. 2018. Can i trust you more? modelagnostic hierarchical explanations. arXiv preprint arXiv:1812.04801.
- Cinzia Viroli and Geoffrey J McLachlan. 2019. Deep gaussian mixture models. Statistics and Computing, 29:43-51.
- Frank Willard, Luke Moffett, Emmanuel Mokel, Jon Donnelly, Stark Guo, Julia Yang, Giyoung Kim, Alina Jade Barnett, and Cynthia Rudin. 2024. This looks better than that: Better interpretable models with protopnext.
- Mengqi Xue, Qihan Huang, Haofei Zhang, Lechao Cheng, Jie Song, Minghui Wu, and Mingli Song. 2022. Protopformer: Concentrating on prototypical parts in vision transformers for interpretable image recognition. arXiv preprint arXiv:2208.10431.
- Yifei Zhang, Neng Gao, and Cunqing Ma. 2023. Learning to select prototypical parts for interpretable sequential data modeling. In Proceedings of the AAAI

Conference on Artificial Intelligence, volume 37, pages 6612-6620.

Yin Zhang, Rong Jin, and Zhi-Hua Zhou. 2010. Understanding bag-of-words model: a statistical framework. International journal of machine learning and cybernetics, 1:43–52.

#### A DPGMM

To model each similarity distribution as a mixture of Gaussian components, we use a neural network that takes a hidden representation h as input, which is derived from  $\mathbf{m}_k$  via a two-layer MLP: h = $MLP(\mathbf{m}_k)$ . This hidden representation h is then used to generate the parameters of the Gaussian mixture, including the mixture weights  $\pi$ , means  $\mu$ , and standard deviations  $\sigma$ , allowing the model to approximate the similarity distribution effectively. **Means** ( $\mu$ ) and **Standard Deviations** ( $\sigma$ ). The parameters of the Gaussian components are computed as follows:

$$\mu = \text{sigmoid}(\mathbf{W}_{\mu}h + \mathbf{b}_{\mu}) \times T, \qquad (14)$$

$$\sigma = \exp(\mathbf{W}_{\sigma}h + \mathbf{b}_{\sigma}), \tag{15}$$

where  $\mu$  and  $\sigma$  are the mean and standard deviation for each of the M Gaussian components.

**Mixture Weights** ( $\pi$ ). To dynamically determine the mixture weights, we employ the Stick-Breaking Process (Ren et al., 2011), with the Dirichlet Process (DP) (Teh et al., 2010) implicitly implemented through the stick-breaking formulation. The DP provides a nonparametric Bayesian approach that allows the model to determine the appropriate number of components adaptively, which is crucial for handling data with unknown complexity.

We define a maximum number of Gaussian components, G, which represents the potential number of components for approximating the similarity distribution. The mixture weights  $\pi_q$  for each component g are generated as follows:

$$\nu_g = \operatorname{sigmoid}(\mathbf{W}_{\pi}h + \mathbf{b}_{\pi}), \qquad (16)$$

$$\pi_g = \nu_g \prod_{\ell=1}^{g-1} (1 - \nu_\ell), \quad g = 1, \dots, G, \quad (17)$$

Here,  $\nu_g$  is computed by applying a sigmoid function to a linear transformation of the hidden representation h. The Stick-Breaking Process ensures that the mixture weights  $\pi_m$  sum to one and adaptively determine the number of active components, enabling the model to capture complex and potentially multi-modal distributions.

#### **B** Datasets

853

857

864

870

875

876

877

884

894

900

901

902

The IMDB dataset contains 25,000 balanced training and test samples and follows a binary sentiment classification format. The dataset was split into training (90%) and validation (10%) partitions. The Yelp Reviews dataset consists of 580,000 samples, with training and test sets comprising 550,000 and 30,000 samples, respectively. Sentiments were binarized by treating 1-2 stars as negative and 3-4 stars as positive. The Amazon dataset was created by selecting 30,000 random reviews, with 24,000 samples allocated for training and validation and 6,000 for testing. The Hotel dataset includes 20,000 reviews evaluating 1,000 hotels, reduced to a balanced subset of 4,508 reviews (2,254 positive and 2,254 negative). The Steam Reviews dataset consists of 130,000 pre-processed reviews, balanced between positive and negative sentiments. Reviews with fewer than 10 characters or containing less than two sentences were excluded.

> The DBPedia dataset is a multiclass dataset extracted from Wikipedia. For the experiments in this paper, we use only 4 labels: "Person," "Animal," "Building," and "Natural Place." Similarly, the Consumer Complaints dataset is a multiclass dataset. For the experiments, we use only 4 classes: "Checking or Savings Account," "Credit Card or Prepaid Card," "Debt Collection," and "Mortgage."

In all experiments, pre-trained embeddings from the BERT-based language model (Song et al., 2020b) were employed to convert raw text into sentence embeddings, enabling downstream analysis.

#### **C Prototype Interpretation**

To assess the interpretability of the ProtoLens model, we provide prototype-aligned interpretations across multiple datasets. Each figure showcases the top-3 original text sentences from the training set that are most aligned with each prototype. These examples illustrate how ProtoLens associates prototypes with representative samples, making its decision-making process more interpretable and transparent.

For the IMDB dataset, as shown in Figure 9, ProtoLens aligns prototypes with representative training samples that reflect key aspects of movie reviews. Positive prototypes are associated with reviews praising elements such as acting and overall quality, as seen in samples like "He does an excellent job in this movie" and "I deeply enjoyed his performance." Negative prototypes, on the other hand, align with reviews critiquing aspects like plot and execution, exemplified by samples such as "This movie was poorly acted, poorly filmed, poorly written" and "It's talky, illogical, slow, and ultimately boring." These representative samples demonstrate ProtoLens' ability to capture diverse perspectives in sentiment analysis. 903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

In the Yelp dataset, as shown in Figure 10, ProtoLens aligns prototypes with representative samples that capture customer opinions on food, service, and ambiance. Positive prototypes are linked to text such as "The service is impeccable" and "The food is great, good portions and quality," reflecting positive customer experiences. Conversely, negative prototypes correspond to samples highlighting dissatisfaction, such as "The food was horrible" and "The place looked dirty and disorganized." These aligned samples illustrate how ProtoLens effectively represents common patterns in customer feedback.

For the Hotel dataset, as shown in Figure 11, ProtoLens aligns prototypes with representative training samples reflecting both positive and negative experiences. Positive prototypes align with samples such as "Room was clean and good" and "The staff were friendly and helpful," highlighting aspects of comfort and service. Negative prototypes correspond to samples like "The room had no soundproofing" and "The carpet is disgusting and filthy," emphasizing common complaints in hospitality feedback. These representative samples demonstrate ProtoLens' ability to capture recurring themes in hotel reviews.

In the Steam dataset, as shown in Figure 12, ProtoLens identifies prototypes aligned with gaming reviews that reflect both satisfaction and dissatisfaction. Positive prototypes are linked to reviews like "This game is amazing" and "Runs smooth even on low settings," which highlight positive gameplay experiences. Negative prototypes, on the other hand, align with samples such as "The servers are abandoned" and "This game sucks, do not buy it," reflecting technical issues and user frustration. These representative samples demonstrate ProtoLens' ability to adapt to highly specific and technical feedback in gaming.

For the Amazon dataset, as shown in Figure 13, ProtoLens aligns prototypes with representative training samples focusing on product quality, usability, and service. Positive prototypes correspond to samples such as "The decor is beautiful and the ambiance is great" and "I enjoyed this place and will go back," reflecting favorable customer feedback. Negative prototypes align with samples like "The food was uninspired and lacked flavor" and "Horrible management and worse customer service," highlighting dissatisfaction. These examples demonstrate ProtoLens' versatility in capturing meaningful patterns in e-commerce reviews.

955

957

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

978

979

981

983

985

991

995

999

1000

1001

1003

Overall, these results underscore ProtoLens' ability to align prototypes with semantically meaningful training samples, providing interpretable insights into the patterns learned during training. This interpretability is key to understanding the model's reasoning across diverse datasets.

#### **D** Classification Interpretation

ProtoLens explains its classification predictions by aligning input text with prototypes from the training set and computing similarity scores to highlight the most relevant prototypes. Each prototype contributes to the final prediction based on its similarity to the input text and its associated sentiment weight. Below, we discuss how ProtoLens interprets both positive and negative classifications through representative examples.

#### D.1 Positive Sentiment Interpretation

Figure 14 demonstrates a positive sentiment classification. ProtoLens activates three prototypes that correspond to semantically aligned samples from the training set. For instance, \*\*Prototype 10\*\* highlights positive movie reviews with phrases like "In all it is a good movie to see," capturing strong alignment with the input's positive tone. Similarly, \*\*Prototype 14\*\* emphasizes "acting was terrific," contributing further evidence of a positive sentiment. The similarity scores and sentiment weights of these prototypes are combined to determine the final classification as positive. This process underscores how ProtoLens grounds its decisions in interpretable and meaningful text examples.

#### **D.2** Negative Sentiment Interpretation

Figure 15 illustrates a negative sentiment classification. ProtoLens activates prototypes that align with critical text samples from the training set. For example, \*\*Prototype 3\*\* reflects dissatisfaction through statements such as "It's talky, illogical, slow, and ultimately very boring," aligning with the input's description of the movie as "pretty bad." \*\*Prototype 4\*\* further reinforces the negative sentiment by associating with phrases like "poorly acted, poorly filmed, poorly written." These prototypes provide interpretability by grounding the model's negative classification in representative samples that closely match the input text.

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

### **D.3** Interpretability

The examples in Figures 14 and 15 demonstrate ProtoLens' ability to explain its predictions using interpretable prototypes. By aligning input text with training set samples that serve as prototypes, ProtoLens offers a transparent view of how classification decisions are made. The similarity scores and sentiment weights ensure that each activated prototype meaningfully contributes to the overall prediction, enhancing both interpretability and faithfulness of the model.

Overall, these results highlight ProtoLens' capacity to provide human-understandable explanations for sentiment classification tasks, bridging the gap between model interpretability and practical applications.

### E Ablation Study

To demonstrate the effectiveness of the Prototype Alignment and Diversity Constraint, we compare ProtoLens trained with and without these components. Prototype Alignment ensures that prototypes maintain their semantic faithfulness. The Diversity Constraint encourages prototypes to capture distinct, non-redundant features, enhancing generalization and reducing redundancy in representation. The results are shown in Table 2.

**Impact of Diversity Constraints.** The removal of diversity constraints (*w/o Diversity*) leads to a noticeable accuracy decline across all tested datasets, notably on IMDB (from 0.903 to 0.882), Amazon (from 0.937 to 0.926), Yelp (from 0.962 to 0.931) and Hotel (from 0.963 to 0.947). This indicates that the diversity loss plays a crucial role in encouraging distinct and varied prototype representations, which helps the model generalize better across different data points. The drop in accuracy suggests that when prototypes become more redundant, they lose their ability to represent the diversity in the dataset, limiting the model's interpretability and performance.

**Impact of Prototype Alignment.** The ablation results for removing prototype alignment (*w/o Alignment*) show a decline in performance, particularly on the Yelp dataset (from 0.963 to 0.943), highlighting the importance of prototype alignment. Align-

ing prototypes with representative embeddings ensures they remain semantically meaningful, leading to more accurate and interpretable predictions. The slight performance drop across other datasets, such as IMDB and Amazon, further emphasizes that the adaptive update process enabled by prototype alignment promotes more stable and reliable learning, improving the model's interpretability and accuracy.

#### F Hyperparameter

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1066

1067

1068

1070

1071

1072

1073

1074

1075

1076

1078

1079

1080

1081

1082

1083

1084

1086

1087

1088

1089

1090

1091

1092

1094

1095

1096

1097

1098

1099

1100

1101

1102

**Effect of K.** The number of prototypes, denoted by K, plays a crucial role in determining the balance between model interpretability and classification performance. As shown in Figure 7, increasing K generally leads to improved accuracy across most datasets, with the exception of some slight fluctuations. For instance, in the IMDB dataset, increasing K from 10 to 40 boosts the performance from 0.884 to 0.903, while for the Yelp dataset, a similar increase elevates the accuracy from 0.931 to 0.950. The improvements plateau or slightly decrease for higher values of K, suggesting diminishing returns beyond a certain point.

The optimal value of K appears to be datasetdependent. For example, K = 50 yields the best performance on the Amazon and Yelp datasets with 0.937 and 0.962, respectively, while K = 40 provides the best performance on the IMDB dataset (0.903). Meanwhile, for the Hotel dataset, K = 20achieves the highest accuracy at 0.963. This variation indicates that the ideal number of prototypes may depend on the complexity and size of the dataset.

Overall, increasing K allows the model to capture more fine-grained patterns by using a larger set of prototypes, but setting K too high may introduce unnecessary complexity without substantial accuracy gains. Thus, choosing K involves a trade-off between maintaining a manageable number of interpretable prototypes and achieving high predictive performance.

**Effect of n-gram.** An n-gram is a hyperparameter that determines the granularity of text division. As shown in Figure 8, an n-gram size of 5 achieves the best performance across all datasets, with notable improvements on IMDB (0.903), Amazon (0.937), and Hotel (0.963), indicating that n = 5 is the optimal n-gram size, providing the best trade-off between incorporating sufficient context and avoiding unnecessary complexity. For smaller n-gram sizes (e.g., n = 1, 3), performance is slightly lower, 1103 likely due to the model's limited ability to capture 1104 broader contextual information. On the other hand, 1105 a larger n-gram size (n = 7, 9) does not yield im-1106 proved performance and even leads to a decrease 1107 in accuracy on all datasets, as seen with IMDB and 1108 Amazon. This suggests that including too large of 1109 a n-gram introduces noise, which results in slight 1110 performance degradation. 1111

1112

1113

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

#### G Cross-Dataset Prototype Generalization

To assess the generalizability of ProtoLens proto-1114types across datasets, we conducted a cross-dataset1115evaluation. Specifically, we tested the performance1116of ProtoLens on the Hotel dataset using prototypes1117derived from the Yelp and Amazon datasets, which1118also represent customer review domains. Table 31119summarizes the results.1120

Table 3: Cross-dataset evaluation results. ProtoLens performance on the Hotel dataset with prototypes derived from different datasets.

Prototype Source	Accuracy on Hotel Dataset	
Hotel (Original)	0.963	
Yelp	0.954	
Amazon	0.943	

The results demonstrate that ProtoLens maintains strong performance even when using prototypes derived from external datasets. While the accuracy slightly decreases compared to using prototypes generated directly from the target dataset (Hotel), the drop in performance is modest: a 0.9% and 2.0% reduction in accuracy when using Yelp and Amazon prototypes, respectively. This suggests that ProtoLens prototypes capture generalizable patterns that can extend across datasets with similar domains.

These findings underscore the robustness of ProtoLens in leveraging prototypes across related datasets, a desirable property for practical applications where annotated data for prototype derivation may be limited. Furthermore, the ability to generalize across datasets indicates that ProtoLens can identify domain-invariant concepts, making it a promising approach for transfer learning and crossdomain interpretability in prototype-based models.

# IMDB

	Prototype Aligned Interpretation				
pe 0	top-3 representative candidates	Contribution to Positive Class			
Prototy	<ol> <li>He, too does an excellent job in this movie.</li> <li>I have a lot of respect for his acting after viewing his performance in this movie.</li> <li>I was deeply impressed with the character he played.</li> </ol>	0.843			
Prototype 1	<ol> <li>This is supposed to be a horror film, but it's lacking in that area and isn't the least bit scary.</li> <li>I happen to to be a horror movie fan, but this film was just so poor, words fail me.</li> <li>Don't waste your time - even the tried and true horroring trigue classics fail in this movie.</li> </ol>	-0.809			
Prototype 2	<ol> <li>This show gave great laughs in premieres, and it still does during re-runs.</li> <li>When we started watching this series on cable, I had no idea how addictive it would be.</li> <li>it was actually a pretty funny show.</li> </ol>	0.723			
Prototype 3	<ol> <li>It lacks substance and style!</li> <li>It's silly, not thoughtful, and boring.</li> <li>It's talky, illogical, slow, and ultimately very boring.</li> </ol>	-0.956			
Prototype 4	<ol> <li>This movie was poorly acted, poorly filmed, poorly written, and overall horribly executed.</li> <li>The film itself is poorly constructed and acted.</li> <li>The plot is slashed to bits and the acting is horrible.</li> </ol>	-0.854			
Prototype 5	<ol> <li>But sometimes it is interesting to see what goes on through peoples' minds.</li> <li>But, an interesting insight into human nature.</li> <li>There is an underlying theme here.</li> </ol>	0.546			
Prototype 6	<ol> <li>The music and song just fantastic</li> <li>There is great music on the soundtrack.</li> <li>The music is also wonderfully matched and haunting.</li> </ol>	0.494			
Prototype 7	<ol> <li>I'll have to assume that they just didn't have the budget to make a decent film.</li> <li>Seriously, does Hollywood think movies like this are good enough?</li> <li>As a writer I find films this bad making it into production a complete slap in the face.</li> </ol>	-0.398			
Prototype 8	<ol> <li>The overall results is just plain bad.</li> <li>Apart from a couple very entertaining song &amp; dance numbers, this is pretty terrible.</li> <li>It is slow, boring and bordering on pointless.</li> </ol>	-0.421			
Prototype 9	<ol> <li>Stylistically, the film is also beautiful</li> <li>But this film is full of wonderful surprises and performances.</li> <li>Filmed in a theatrical way and excellent acted.</li> </ol>	0.643			

Figure 9: Aligned interpretation of prototypes with corresponding text sentences on the IMDB dataset. Each prototype is associated with specific spans of text and sentiment weights, providing insights into the reasoning behind the model's predictions.

# Yelp

	Prototype Aligned Interpretation			
pe 0	top-3 representative candidates	Contribution to Positive Class		
Prototy	<ol> <li>Yeah, the bar area is nice and colorful.</li> <li>The inside of the location is decent, it mostly has tables a few booths and a new expanded bar.</li> <li>It is always clean, outside patio with mist and a smoking section, TVs throughout the bar area.</li> </ol>	0.797		
Prototype 1	<ol> <li>This place is better than you could imagine based on the concept and is well worth the meal.</li> <li>The food is great, good portions and quality, yummy selections.</li> <li>Not only is the food great, the service is impeccable.</li> </ol>	0.765		
Prototype 2	<ol> <li>The food is horrible the service was bad.</li> <li>The food was terrible and I would not recommend this place to anybody.</li> <li>This is the WORST restaurant I have EVER been to and experienced.</li> </ol>	-0.686		
Prototype 3	<ol> <li>When I came back to ask for a refund they were very rude about it and refused to help.</li> <li>They then got my order wrong and said they wouldn't do anything about it when I told them.</li> <li>I asked for a refund from the company and not a peep out of them.</li> </ol>	-0.388		
Prototype 4	<ol> <li>I have to say I was truly disappointed by the flavors.</li> <li>The flavors are \ok nothing special.</li> <li>The only flavor I didn't care for was the red velvet.</li> </ol>	-0.782		
Prototype 5	<ol> <li>Oh and the noise level was too high.</li> <li>the place looked dirty and disorganized and smelled bad!</li> <li>There was a dingy smell and a security guard wandering the aisles</li> </ol>	-0.433		
Prototype 6	<ol> <li>Service was good, friendly staff.</li> <li>Service was excellent, staff was friendly.</li> <li>The staff was nice and service was prompt.</li> </ol>	0.517		
Prototype 7	<ol> <li>In addition to these items, the bread garlic butter that is served with the meal was also great.</li> <li>We each ordered something different off the menu and everything was just scrumptious!</li> <li>We were teated to a nice spread including vegetarian pasta primavera spicy!</li> </ol>	0.432		
Prototype 8	<ol> <li>I am NEVER going back here and wouldn't recommend anyone to even try the place.</li> <li>I wasn't impressed by this place and I don't think I'll be returning anytime soon.</li> <li>I will never come here again.</li> </ol>	-0.364		
Prototype 9	<ol> <li>Great price at \$5.95 and plenty of food for me.</li> <li>Prices are so reasonable, and with a restaurant.com coupon it was just dirt cheap.</li> <li>The prices are insanely cheap for what you get.</li> </ol>	0.481		

Figure 10: Aligned interpretation of prototypes with corresponding text sentences on the Yelp dataset. The figure highlights the diverse prototypes and their representative candidates, emphasizing interpretability in the sentiment analysis task.

# Hotel

	Prototype Aligned Interpretation			
pe 0	top-3 representative candidates	Contribution to Positive Class		
Prototy	<ol> <li>Its close to restaurants and really any place you want to go</li> <li>The location is outstanding and I suppose you get what you pay for in that aspect.</li> <li>Great clean place to stay.</li> </ol>	0.674		
Prototype 1	<ol> <li>The room had no sound proof.</li> <li>The air conditioning did not work well in either of the rooms in which we stayed.</li> <li>The heater in the room did not work properly.</li> </ol>	-0.730		
Prototype 2	<ol> <li>The carpet is disgusting and filthy.</li> <li>The carpeted floors were very dirty and were not vacuumed.</li> <li>Carpet was dirty smelled and had stains all over.</li> </ol>	-0.862		
Prototype 3	<ol> <li>the bed in my room was also one of the most comfortable hotel beds</li> <li>it was clean beds very comfortable.</li> <li>the beds and pillows were comfortable.</li> </ol>	0.238		
Prototype 4	<ol> <li>Room was clean and good.</li> <li>The room was clean and roomy.</li> <li>The room was clean in very nice condition and everything worked well.</li> </ol>	0.991		
Prototype 5	<ol> <li>My first impression was quite good for the price.</li> <li>All in all a good experience.</li> <li>U get what u pay for</li> </ol>	0.329		
Prototype 6	<ol> <li>I'm thrilled you had a wonderful stay.</li> <li>I am glad that you enjoyed your stay at the hotel.</li> <li>Enjoyed my stay at the hotel.</li> </ol>	0.808		
Prototype 7	<ol> <li>The room we stayed in was smelly dirty and poorly cleaned.</li> <li>The room smelled old and the bathroom was gross.</li> <li>When we entered our room it had a very bad odor.</li> </ol>	-0.926		
Prototype 8	<ol> <li>The staff were delightful and most helpful with special mention of the front desk!</li> <li>Staff were extremely friendly and helpful we felt very welcomed.</li> <li>The staff were friendly and helpful when checking in.</li> </ol>	0.777		
Prototype 9	<ol> <li>The bathtub was peeling and dirty and the mold on the shower curtain was horrible.</li> <li>The pool and hot tub was filthy.</li> <li>Black mold in the shower and lamps that did not work.</li> </ol>	-0.758		

Figure 11: Aligned interpretation of prototypes with corresponding text sentences on the Hotel dataset. The interpretations include both positive and negative sentiment examples, showcasing the model's ability to capture nuanced feedback.

## Steam

	Prototype Aligned Interpretation			
pe 0	top-3 representative candidates	Contribution to Positive Class		
Prototy	<ol> <li>Servers suck devs suck glitches and cheaters run rampant its just not worth the time.</li> <li>The server of game is tooooooooo rubbish.</li> <li>The servers are abandoned always laggy and lots of disconections.</li> </ol>	-0.654		
Prototype 1	<ol> <li>Two more maps posible gamemodes incoming weather and new weapons keeps the game interesting.</li> <li>The constant updates and dlc whick keep the gameplay fresh and original.</li> <li>The addition of new maps vehicles guns and players keeps everything fresh and makes every game a new experience.</li> </ol>	0.826		
Prototype 2	<ol> <li>I got to say that gta 5 is awesome so much you to do on the game itself.</li> <li>Gta v is a great game and its great playing with your friends.</li> <li>There is so much to do in gta v. I recommend this game!</li> </ol>	0.911		
Prototype 3	<ol> <li>The developers have disallowed mods which is simply outrageous.</li> <li>So the games makers have decided to cut off mods.</li> <li>However the developers are ruining their game by removing mod support.</li> </ol>	-0.235		
Prototype 4	<ol> <li>I dont recomend buying this game.</li> <li>If you are a fan of this game type then dont bother buying this game at the moment.</li> <li>This game sucks do not buy it.</li> </ol>	-0.963		
Prototype 5	<ol> <li>Even on my old I702x i can run this game pretty smoothly on normal settings!</li> <li>Runs smooth even on bad pc s on low settings of course.</li> <li>The game runs silky smooth and looks great even on my modest hardware with only 1gb video memory.</li> </ol>	0.577		
Prototype 6	<ol> <li>The game is an amazing game.</li> <li>Personally this is one of my favourite game.</li> <li>It is an amazing game.</li> </ol>	0.935		
Prototype 7	<ol> <li>You basically spend 30 minutes looting just to end up dying by something you cant even do anything about.</li> <li>Sometimes you don't see anyone for 15 minutes and then you die from any sniper anywhere on the map.</li> <li>You literally run around pick up crap shoot at people and usually die very quickly.</li> </ol>	-0.403		
Prototype 8	<ol> <li>It is the worst game every created.</li> <li>It is honestly a terrible game.</li> <li>It's a garbage game plain and simple.</li> </ol>	-0.806		
Prototype 9	<ol> <li>Game is great but full of hackers hackers everywhere!</li> <li>I like this game but atm its full of hackers.</li> <li>So many hackers have appeared in this game.</li> </ol>	-0.359		

Figure 12: Aligned interpretation of prototypes with corresponding text sentences on the Steam dataset. The figure demonstrates how ProtoLens handles diverse feedback in gaming reviews, including issues like performance and user experience.

## Amazon

Prototype Aligned Interpretation			
pe 0	top-3 representative candidates	Contribution to Positive Class	
Prototy	<ol> <li>I won't be going back.</li> <li>I'm sad to say that I won't be going back.</li> <li>I definitely will not be going back.</li> </ol>	-0.579	
Prototype 1	<ol> <li>The decor is nice and there are TV's everywhere, including at every booth.</li> <li>The place is nicely laid out and there are a decent number of tables more than the Tempe location.</li> <li>That aside, the ambience is great and the decor is beautiful.</li> </ol>	0.648	
Prototype 2	<ol> <li>The food did not look appealing.</li> <li>The food was just not good.</li> <li>The food seemed to be uninspired and lacked flavor.</li> </ol>	-0.613	
Prototype 3	<ol> <li>I enjoyed this place and will go back, perhaps today.</li> <li>I've been here numerous times over the years and always had a great time.</li> <li>I went to this place a couple times and took some new friends there today.</li> </ol>	0.582	
Prototype 4	<ol> <li>It's too bad the horrible service outweighed the tasty food.</li> <li>Needless to say, horrible management, even worse customer service and         <ul> <li>I will NOT be returning to this location!</li> <li>Just a really bad experience all around with the slow service and sub par food.</li> </ul> </li> </ol>	-0.865	
Prototype 5	<ol> <li>This is one of the worse restaurants ever!</li> <li>Words can't express how appalled I am about our food experience at this restaurant.</li> <li>The food and service this past week when we dined was awful!</li> </ol>	-0.719	
Prototype 6	<ol> <li>She made sure my order was right and worked closely with me to ensure everything was perfect.</li> <li>I will give kudos to our hostess, who was lovely.</li> <li>They were so positive and they gave us recommendations.</li> </ol>	0.547	
Prototype 7	<ol> <li>The drinks are expensive, but they're made pretty well.</li> <li>but I digress. The drinks are FANTASTIC.</li> <li>The drinks are just what you would expect from a place that is membership only - strong and tasty.</li> </ol>	0.282	
Prototype 8	<ol> <li>So after waiting nearly 15 minutes for anyone to even come and take our order we left.</li> <li>We waiting a good ten minutes before we were even acknowledged, and it was 3pm, restaurant was near empty.</li> <li>We waited 15 minutes, no one came to our table, we watched 3 other servers walk by.</li> </ol>	-0.364	
Prototype 9	<ol> <li>All of the staff were friendly and service was great</li> <li>Service was good and the place was clean.</li> <li>The waiting and serving staff were excellent, they were very helpful.</li> </ol>	0.405	

Figure 13: Aligned interpretation of prototypes with corresponding text sentences on the Amazon dataset. This figure illustrates ProtoLens' interpretability across product reviews, focusing on features such as quality, service, and usability.

#### **Positive Text Instance**

**Good Movie, acting was terrific especially** from Eriq EbouaneyLumumbaand very well directed. It also shows how Lumumba was cornered by the Belgians, U S A and United Nations and how they labelled him a `communist' to scare people as they did to all the Honest True African leaders like Nkrumah, Kenyatta, Nyerere and many others. It shows how western countries preach democracy while they have something else on the back of their minds. It *is a story of injustice, struggle* and brutality. There should have been an explanation why he Lumumba couldn't keep the second largest country in Africa in one piece. And also what was going on with Tshombe and Katanga . Just heads up if you gonna watch the movie Tshombe was controlling the Katanga region which if I am not mistaken is the number one copper producer in the world. In all it is *a good movie to see*. You will learn something new about Africa, it's leaders and it's people and probably will open your eyes why this continent is ridden with wars.

Top-3 Activated Prototype	Prototype 10	Prototype 12	Prototype 14
	1. In all it is a good movie to see.	1. This was the most visually stunning, amaz- ing and incredible story I've ever experienced.	1. Some very good character actors in this fine film.
Aligned Interpretaion	2. Overall, this is a fun film & I highly recommend it.	2. Everything about it was wonderful!	2. The acting by all of these actors is very good.
	3. It's GREAT and a film EVERYONE must see.	3. The story was completely absorbing and entertaining.	3. The actors deliver solid enough performances.
Extracted Span	"good movie to see"	"is a story of injustice, struggle"	"Good Movie, acting was terrific especially"
Similarity	0.732	0.319	0.720
Contribution to Positive Class	0.846	0.247	0.687
Prediction		Positive	

Figure 14: The figure showcases how ProtoLens aligns input text with prototypes to explain a positive sentiment prediction. The extracted spans and similarities for the top-3 activated prototypes are presented, along with sentiment weights contributing to the final prediction.

#### **Negative Text Instance**

I caught this movie on Sci-Fi before heading into work. If you've any interest in seeing Dean Cain dive and avoid being enveloped in flames at least a dozen times, this movie is for you. If that doesn't peak your interest, well, I'm afraid you'll wish that YOU were the one about to be enveloped in flames, because *this movie is pretty bad*. The *acting, to begin with, is awful, awful, awful.* The characters are all completely obnoxious, and the dialogue is worse than your typical Z-grade, Sci-Fi movie. Towards the end, the movie began to remind me of 'Hollow Man' complete with escape via elevator shaft, except with a Dragon, not a naked, invisible man. Unlike other similar flicks, however, this one wasn't even awesomely bad... *it was just plain bad*.

Top-3 Activated Prototype	Prototype 3	Prototype 4	Prototype 8
Aligned Interpretaion	<ol> <li>It lacks substance and style!</li> <li>It's silly, not thoughtful, and boring.</li> <li>It's talky, illogical, slow, and ultimately very boring.</li> </ol>	<ol> <li>This movie was poorly acted, poorly filmed, poorly written, and overall horribly executed.</li> <li>The film itself is poorly constructed and acted.</li> <li>The plot is slashed to bits and the acting is horrible.</li> </ol>	<ol> <li>The overall results is just plain bad.</li> <li>Apart from a couple very entertaining song &amp; dance numbers, this is pretty terrible.</li> <li>It is slow, boring and bordering on pointless.</li> </ol>
Extracted Span	"this movie is pretty bad"	"acting, to begin with, is awful, awful, awful"	"it was just plain bad"
Similarity	0.467	0.558	0.626
Contribution to Positive Class	-0.956	-0.854	-0.421
Prediction		Negative	

Figure 15: The figure shows how ProtoLens aligns input text with prototypes to explain a negative sentiment prediction, supported by similarity scores and sentiment weights.