# Hierarchical Approaches for Domain-Specific Image Captioning: Classification, Distillation, and Optimization

#### Anonymous ACL submission

#### Abstract

This paper presents a novel hierarchical approach to improve image captioning, focusing on enhancing the accuracy, contextual appropriateness, and linguistic diversity of generated captions. We address key limitations of existing multimodal models, including inaccurate object relationships, missed details, and poor domainspecific understanding. Our method incorporates three main components: a classificationguided prompting system that utilizes domainspecific knowledge, a knowledge distillation 011 framework that transfers captioning capabilities from GPT-40 to the LLaVA model, and an iterative Direct Preference Optimization (DPO) approach that refines caption quality. Exten-016 sive experiments demonstrate that our approach outperforms existing methods, achieving near-018 GPT-40 performance while maintaining com-019 putational efficiency. Additionally, we release a high-quality dataset of 9,840 image-caption pairs across 18 categories, providing valuable resources for future research in domain-specific image captioning.

### 1 Introduction

024

037

041

Image captioning, the task of automatically generating natural language descriptions for images, has emerged as a crucial indicator of models' ability to understand visual and language information effectively (Mokady et al., 2021; Li et al., 2023). Despite significant advances in this field, current state-of-the-art models still face substantial challenges in generating accurate, contextually appropriate, and comprehensive descriptions.

Recent multimodal models, while demonstrating impressive capabilities across various tasks (Li et al., 2023; Liu et al., 2024), exhibit notable limitations in image captioning. These limitations manifest in several aspects: inaccurate object relationship descriptions, missed subtle but important details, and inadequate ability to express domainspecific concepts. Moreover, traditional evaluation metrics such as BLEU (Papineni et al., 2002) and METEOR (Banerjee and Lavie, 2005) often fail to capture these nuanced aspects of performance, potentially leading to misleading assessments of model capabilities. 042

043

044

047

048

054

055

058

060

061

062

063

065

066

067

068

069

071

073

074

075

077

To address these limitations, we propose a novel hierarchical approach that combines classificationguided prompting, knowledge distillation, and iterative preference optimization. Our method consists of three key components:

- A classification-first strategy that categorizes images into semantic classes, enabling the use of class-specific prompting templates for more precise and contextually relevant caption generation.
- A knowledge distillation framework that leverages GPT-4o's superior captioning abilities to create high-quality training pairs, which are then used to fine-tune a more efficient LLaVA model<sup>1</sup>.
- An iterative Direct Preference Optimization (DPO) approach that continuously refines the model's output quality by learning from automatically generated preference pairs using a Pixtral (Agrawal et al., 2024) model for evaluation.

Through comprehensive experiments, we demonstrate that this hierarchical approach significantly improves captioning quality across multiple dimensions. The classification-guided prompting ensures domain-specific accuracy, while the knowledge distillation from GPT-40 provides rich linguistic diversity. Furthermore, the iterative DPO refinement helps align the model's outputs with quality standards, leading to more natural and contextually appropriate descriptions.

<sup>&</sup>lt;sup>1</sup>LLaVA serves as our base multimodal model due to its strong vision-language capabilities and computational efficiency.

Our extensive evaluation shows that the proposed approach achieves significant improvements over existing methods. Specifically, our model achieves a Pixtral score of 4.312 after DPO optimization, approaching the performance of GPT-40 (4.473) and GPT-40-mini (4.293), while substantially outperforming the base LLaVA model (3.598). Human evaluation strongly correlates with our automatic metrics (Pearson correlation: 0.779), validating the effectiveness of our evaluation framework. These results demonstrate the success of combining semantic classification, knowledge distillation, and preference optimization in improving image captioning quality.

078

079

084

096

101

102

103

104

106

107

108

109

110

111

112

113

114

115

116

117

118

Our main contributions are summarized as follows:

- We propose a novel classification-guided prompting strategy that leverages domainspecific knowledge for more accurate and contextually appropriate caption generation.
- We present an effective knowledge distillation pipeline that efficiently transfers captioning capabilities from GPT-40 to a more practical LLaVA model.
- We develop an iterative DPO framework that uses automated evaluation for continuous model improvement, demonstrating consistent performance gains over multiple iterations.
- We release a comprehensive dataset of 9,840 high-quality image-caption pairs across 18 diverse categories, where each image is paired with a GPT-40-generated caption. This curated dataset provides valuable training resources for future research in domain-specific image captioning.
  - We introduce a fine-tuned DPO version of the LLaVA model, which demonstrates strong image captioning capabilities, further enhancing the practical application of the model.

### 2 Related Works

### 2.1 Image Captioning Models

119Recent advances in image captioning have been120largely driven by the development of large multi-121modal models. Traditional approaches relied on122encoder-decoder architectures (Xu et al., 2016; An-123derson et al., 2018), while more recent work has

focused on end-to-end multimodal training. Models like BLIP (Li et al., 2023) and LLaVA (Liu et al., 2024) have demonstrated impressive capabilities in understanding and describing visual content. However, these models often struggle with domainspecific details and contextual accuracy, motivating our classification-guided approach. 124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

158

159

160

161

162

163

164

165

166

167

168

169

170

171

#### 2.2 Caption Evaluation Metrics

The evaluation of image captioning systems has evolved through several generations:

**Reference-based Metrics** Traditional metrics such as BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and CIDEr (Vedantam et al., 2015) evaluate captions by comparing them with human-written references. While widely used, these metrics often fail to capture semantic accuracy and contextual appropriateness (Anderson et al., 2016).

**Learning-based Metrics** More recent approaches like BERTScore (Zhang\* et al., 2020) and CLIP-Score (Hessel et al., 2022) leverage pre-trained models to assess caption quality. These metrics show better correlation with human judgments but still lack explainability and domain-specific understanding.

**LLM-based Evaluation** The emergence of large language models has introduced new possibilities for caption evaluation. Recent work has shown that LLMs can provide more nuanced and context-aware evaluations (Liu et al., 2023). Our work builds on this direction by using the Pixtral model for automated evaluation and preference learning.

### 2.3 Knowledge Distillation and Model Alignment

Knowledge distillation has proven effective in transferring capabilities from larger to smaller models (Hinton et al., 2015). In the multimodal domain, recent work has shown success in distilling vision-language capabilities (Li et al., 2023). Our approach extends this to image captioning by distilling knowledge from GPT-40 to a more efficient LLaVA model.

**Direct Preference Optimization** DPO has emerged as a powerful technique for aligning language models with human preferences (Rafailov et al., 2024). Recent work has demonstrated its effectiveness in improving text generation quality (Liu et al., 2023). We adapt this approach to image

219

220

221

223

224

225

226

227

228

230

231

232

233

234

235

236

239

240

241

242

243

244

245

246

247

248

249

250

ScienceQA (Lu et al., 2022), ChartQA (Masry 217

218

$$D_{\text{train}} = (I_i, c_i)_{i=1}^N \tag{1}$$

where  $I_i$  represents an image and  $c_i$  its corresponding category label.

including Flickr, JAAD (Rasouli et al., 2017, 2018),

et al., 2022), etc.:

3) **Model Fine-tuning**: We fine-tune the LLaVA model using our collected dataset:

$$\mathcal{L}_{\text{cls}} = -\sum_{i=1}^{N} \log p_{\theta}(c_i | I_i)$$
(2)

where  $\theta$  represents the model parameters.

**Classification Process** Given a new image *I*, the classification is performed as:

$$c = f_{\text{LLaVA}}(I) \in \mathcal{C} \tag{3}$$

where  $f_{LLaVA}$  is our fine-tuned LLaVA classifier.

**Structured Prompt Design** For each category c, we design a specialized prompting template  $P_c$  to guide the caption generation process. Representative examples of our prompting templates are provided in Appendix A, and the complete set is available in our repository.

**Caption Generation** The final caption generation process involves:

$$C_{\text{final}} = g_{\text{LLaVA}}(I, P_c) \tag{4}$$

where  $g_{\text{LLaVA}}$  is another LLaVA model fine-tuned on GPT-40 distilled data (detailed in Section 3.2).

#### 3.2 Knowledge Distillation from GPT-40

As mentioned in Section 3.1, our caption generation model  $g_{LLaVA}$  is fine-tuned through knowledge distillation from GPT-40. We leverage the same dataset used for classifier training to create high-quality caption pairs.

**High-Quality Caption Generation** For each image I in our training dataset, we use GPT-40 to generate captions:

$$C_{\rm GPT-4o} = g_{\rm GPT-4o}(I, P_c) \tag{5}$$

where  $g_{\text{GPT-40}}$  represents the GPT-40 model's caption generation function. 252

captioning by using automated evaluations to guidethe optimization process.

#### 2.4 Category-Specific Generation

The idea of leveraging category information for improved generation has been explored in various contexts. Previous work has shown that domain-specific prompting can improve text generation quality (Liu et al., 2024). Our work extends this concept to image captioning by introducing a classification-guided prompting strategy that adapts to different visual domains.

### 3 Method

174

175

176

177

178

181

182

184

185

186

187

190

191

193

194

195

196

197

198

200

201

207

210

211

212

213

215

Our approach consists of three main components: (1) a classification-guided prompting system, (2) a knowledge distillation pipeline from GPT-40, and (3) an iterative DPO optimization process. Figure 1 illustrates the complete architecture of our proposed method. In this section, we describe each component in detail.

### 3.1 Classification-Guided Prompting

The effectiveness of image captioning heavily depends on the model's ability to recognize and describe domain-specific content accurately. To address this challenge, we propose a classificationfirst approach that leverages a fine-tuned large multimodal model (LMM) for image classification, followed by category-specific prompting for caption generation.

**Classification Model Training** We fine-tune a LLaVA model to serve as our classifier. The training process consists of several key steps:

1) Category Definition: We define 18 comprehensive categories C covering most common image scenarios:

- Scene types: Animal, Architecture, Art, City, Landscape
- Daily activities: Daily\_Life, Food, Love, Medical, Sports
- Technical: Autonomous\_Driving, Smart\_Cities, Transportation
- Visual content: Cartoon, Chart, Person, Plant
- Others: Other

2)**Training Data Collection**: For each category  $c \in C$ , we collect images from multiple sources,



Figure 1: method

253 Training Data Creation We create training pairs254 by combining:

$$D_{\text{distill}} = \{ (I_i, P_{c_i}, C_{\text{GPT-4o},i}) \}_{i=1}^N$$
 (6)

where N is the total number of images across all categories.

**LLaVA Fine-tuning** We fine-tune the LLaVA model using these high-quality training pairs with both full parameter fine-tuning and LoRA approaches. For both versions, the supervised learning objective is:

$$\mathcal{L}_{\text{SFT}} = -\sum_{i=1}^{N} \log p_{\theta}(C_{\text{GPT-4o},i}|I_i, P_{c_i}) \quad (7)$$

where  $\theta$  represents the parameters of the LLaVA model. The detailed training configurations for both approaches are provided in Appendix B.

#### **3.3** Iterative DPO Optimization

259

260

262

263

265

266

270

To further improve the quality of generated captions, we implement an iterative Direct Preference Optimization (DPO) process using a Pixtral model for preference scoring. We leverage both fullparameter and LoRA fine-tuned models to generate diverse candidate captions. 271

272

273

274

275

276

277

278

279

280

281

282

283

286

287

288

291

**Caption Generation and Scoring** For each image, we use both the full-parameter and LoRA finetuned models to generate five captions each, resulting in ten candidate captions per image. These captions are then evaluated using the Pixtral model, which assigns a score between 1 and 5 to each caption following a structured evaluation prompt (detailed in Appendix C).

**Preference Pair Selection** From the ten candidates, we identify the highest and lowest scoring captions. If their score difference exceeds 2 points (on the 1-5 scale), we use this pair for DPO training. This threshold ensures we only learn from pairs with significant quality differences.

**DPO Training** The DPO objective is defined as:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{(I,C_w,C_l)} \left[ \log \sigma \left( \beta (r_\theta(C_w) - r_\theta(C_l)) \right) \right]$$
(8)

where  $r_{\theta}(C)$  is the reward score for caption  $C, \beta$  is temperature parameter, and  $\sigma$  is sigmoid function.

4

292	Iterative Refinement The DPO process consists
293	of:
294	1. Generate ten captions per image
295	2. Score captions using Pixtral
296	3. Select pairs with significant score differences
297	4. Update model using DPO objective
298	5. Repeat until convergence
299	The detailed training configurations and imple-
300	mentation details are provided in Appendix B.
301	4 Experiments
302	4.1 Experimental Setup
303	Dataset We collect and curate a new dataset
304	from multiple sources including Flickr, JAAD,
305	ScienceQA, ChartQA, etc., covering 18 diverse
306	categories including animal, architecture, art, au-
307	tonomous driving, cartoon, chart, city, daily life,
308	food, landscape, love, medical, other, person, plant,
309	smart cities, sports, and transportation. The dataset
210	consists of 0.840 training images and 706 valids

food, landscape, love, medical, other, person, plant, smart cities, sports, and transportation. The dataset consists of 9,840 training images and 796 validation images, with approximately balanced distribution across categories. All images are manually verified to ensure quality and correct categorization. The validation set maintains a similar category distribution as the training set to ensure fair evaluation.

311

312

313

314

315

316

318

319

320

321

322

323

324

325

326

327

330

332

333

334

337

For each image in the training set, we generate a high-quality caption using GPT-40 following our category-specific prompting strategy. These captions are carefully crafted to capture domainspecific details and contextual information relevant to each category. This results in a rich dataset where each image is paired with a detailed, contextually appropriate caption that can serve as highquality training data for future research.

The dataset is designed with several key features:

- **Category Balance**: Training and validation sets maintain approximately balanced distribution across all categories, ensuring comprehensive coverage across domains.
- Quality Control: All images undergo manual verification for visual quality and category correctness.
- **Domain Diversity**: The 18 categories cover a wide range of scenarios from daily life to specialized domains like medical imaging and autonomous driving.

• **High-Quality Captions**: Each training image is paired with a GPT-40 generated caption that follows category-specific guidelines.

338

339

340

341

342

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

We plan to release this dataset along with our code and models to facilitate future research in domain-specific image captioning.

**Evaluation Metrics** We employ two types of evaluation:

- Automated Evaluation: The Pixtral model serves as our primary metric, assigning scores from 1 to 5 based on caption quality. We evaluate all generated captions on the entire validation set of 796 images.
- Human Evaluation: Three expert annotators evaluate 54 randomly sampled image-caption pairs using the same 1-5 scale. The sample size was chosen to ensure thorough evaluation while maintaining feasible annotation effort.

**Baselines** We compare our method with several strong baselines:

- **GPT-40**: A large-scale language model developed by OpenAI, known for its advanced natural language understanding and generation capabilities. It serves as a benchmark for state-of-the-art performance in various language tasks.
- **GPT-4o-mini**: A more compact version of GPT-4o, designed to offer competitive performance with reduced computational requirements. This model maintains a balance between efficiency and effectiveness, making it suitable for applications with limited resources.
- Claude-3-Haiku: A large language model developed by Anthropic, optimized for speed and affordability. It is designed to process 21,000 tokens per second for prompts under 32,000 tokens, making it suitable for enterprise workloads that require quick analysis of large datasets.
- LLaVA-v1.5-7B: An open-source multimodal model fine-tuned from LLaMA and Vicuna on GPT-generated multimodal instruction-following data. It is an auto-regressive language model based on the transformer architecture, trained to understand and generate 383

our curated dataset of 9,840 image-caption pairs, trained without Direct Preference Optimization (DPO). This model represents the

effectiveness of traditional fine-tuning approaches on our domain-specific dataset.

both text and images. The model was trained in September 2023 and is licensed under the

• LLaVA-FT: A version of LLaVA model that

has undergone full-parameter fine-tuning on

LLAMA 2 Community License.

385

387

388

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427 428

429

430

431

432

- LLaVA-LoRA: A variant of the LLaVA model fine-tuned on our dataset using Low-Rank Adaptation (LoRA) techniques, without DPO. This approach demonstrates the efficacy of parameter-efficient fine-tuning methods on our carefully curated training data while maintaining model performance.
- Qwen-VL-Chat (Bai et al., 2023): A multimodal chatbot model that integrates vision and language understanding, enabling it to process and generate responses based on both textual and visual inputs. This model is designed to handle a wide range of multimodal tasks, including image captioning and visual question answering.

• Qwen2-VL-7B-Instruct (Wang et al., 2024): An advanced multimodal model that builds upon the Qwen-VL-Chat architecture, incorporating instruction-following capabilities to improve its performance on tasks requiring specific guidance. The 7B refers to the model's parameter size, indicating a balance between computational efficiency and performance.

### 4.2 Main Results

The experimental results reveal several comprehensive findings:

**Base Model Comparisons**: Among the models, we observe a clear performance hierarchy. While the closed-source GPT-40 maintains the highest performance (4.473), setting a strong upper bound, other closed-source models like GPT-40-mini (4.293) and Claude-3-Haiku (4.268) also show strong performance. For open-source models, Qwen2-VL-7B-Instruct achieves 4.224 and Qwen-VL-Chat reaches 3.833. LLaVA-v1.5-7B serves as our baseline with a score of 3.598.

**Fine-tuning Effectiveness**: Our fine-tuning approaches show significant improvements over the

Model	<b>Pixtral Score</b>	$\Delta$	#Params
GPT-40	4.473	-	-
GPT-4o-mini	4.293	-	-
Claude-3-Haiku	4.268	-	-
LLaVA-v1.5-7B (Base)	3.598	-	7B
Qwen-VL-Chat	3.833	+0.235	7B
Qwen2-VL-7B-Instruct	4.224	+0.626	7B
LLaVA-FT (Full)	4.078	+0.480	7B
LLaVA-LoRA	4.034	+0.436	7B
DPO (1 iteration)	4.113	+0.515	7B
DPO (2 iterations)	4.168	+0.570	7B
DPO (3 iterations)	4.290	+0.692	7B
DPO (4 iterations)	4.312	+0.714	7B

Table 1: Performance comparison on the validation set (796 images). Scores range from 1 to 5.  $\Delta$  shows absolute improvement over the base model. #Params shows the total number of parameters.

base model. Full parameter fine-tuning achieves a score of 4.078 (+0.480 over base), demonstrating the effectiveness of comprehensive model adaptation. LoRA fine-tuning reaches 4.034 (+0.436 over base), nearly matching full fine-tuning while being parameter-efficient. Both approaches outperform Qwen-VL-Chat but still trail behind more advanced models like Claude-3-Haiku and Qwen2-VL-7B-Instruct, suggesting room for further improvement.

433

434

435

436

437

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

465

466

**DPO Improvement**: The iterative DPO process shows consistent and meaningful improvements. First iteration achieves +0.515 over base model, surpassing both fine-tuning approaches. Subsequent iterations show continuous improvement, with the second iteration adding +0.055 (reaching 4.168), the third iteration showing the largest incremental gain of +0.122 (achieving 4.290), and the fourth iteration providing a modest +0.022 improvement (reaching 4.312). The diminishing returns after the third iteration suggest convergence of the optimization process. Notably, our final DPO model outperforms both Qwen2-VL-7B-Instruct by +0.088 points and Claude-3-Haiku by +0.044 points, demonstrating the effectiveness of our approach.

**Gap Analysis with Leading Models**: Our final model (DPO 4 iterations) achieves a score of 4.312, which is particularly noteworthy as it surpasses both closed-source models (GPT-4omini and Claude-3-Haiku) and open-source models (Qwen2-VL-7B-Instruct and Qwen-VL-Chat). While there remains a small gap of 0.161 points from GPT-4o, our model outperforms GPT-4omini by +0.019 points, Claude-3-Haiku by +0.044

points, and the best open-source model Qwen2-VL-467 7B-Instruct by +0.088 points. This achievement is 468 especially significant considering we maintain the 469 computational efficiency of a 7B parameter model, 470 demonstrating that our approach effectively nar-471 rows the performance gap between open-source 472 and closed-source models in the field of image cap-473 tioning. 474

**Computational Efficiency**: All our models maintain the original 7B parameter count, offering practical deployment capabilities compared to larger models like GPT-40 and Claude-3-Haiku. This enables efficient inference time while achieving competitive performance, providing a scalable solution for real-world applications.

### 4.3 Human Validation Study

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

To validate the reliability of Pixtral as an automated evaluation metric, we conducted a comprehensive human evaluation study. The study included 54 image-caption pairs randomly sampled from the validation set, evaluated by three experts with ML/CV backgrounds and experience in image captioning. Each annotator independently scored all pairs following detailed evaluation criteria, using a 1-5 scale with 1 point increments, which aligns with the Pixtral scoring scale.

**Correlation Analysis** The comparison between human and Pixtral scores shows strong alignment with a Pearson correlation of 0.762 (p-value: 0.000), Kendall's Tau of 0.545, and Spearman correlation of 0.595. These strong correlations across different statistical measures validate Pixtral's effectiveness as an automated evaluation metric.

500 Error Analysis The difference between Pixtral 501 and human scores demonstrates reasonable agree-502 ment, with a Mean Absolute Error (MAE) of 0.414 503 and Root Mean Square Error (RMSE) of 0.572. 504 These values indicate that Pixtral's scores typically 505 deviate from human consensus by less than half a 506 point on the 5-point scale, suggesting strong practi-507 cal reliability.

**Inter-annotator** Agreement The Cohen's Kappa coefficients between annotators are 0.507, 0.617, and 0.466 for pairs 1-2, 1-3, and 2-3, 510 511 respectively. These values reflect moderate to good inter-annotator agreement, highlighting that 512 while there are some subjective differences in 513 caption evaluation, the overall consistency is still 514 acceptable. These results further underscore the 515

advantage of using a consistent automated metric like Pixtral, and the importance of averaging multiple human judgments to mitigate individual biases.

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

### 4.4 Qualitative Analysis

We demonstrate our model's capability through several test cases on images not present in our training set. As shown in Fig. 1, our classification-guided prompting effectively handles various domainspecific characteristics. For animal images, our model accurately captures physical attributes (e.g., "fluffy gray and white kitten") and behavioral aspects (e.g., "lounges comfortably"). For person descriptions, it successfully combines human attributes with environmental context (e.g., "standing against a vibrant blue wall"). In plant-related cases, it shows capability in describing fine-grained details (e.g., "green soybean seed"). These examples demonstrate our model's robust performance on novel images. However, challenges remain in several aspects. The model sometimes struggles with complex scenes involving multiple objects or actions, may have difficulty with rare or unusual visual concepts, and occasionally misses subtle contextual details. These limitations suggest directions for future improvements in our approach.

# 5 Discussion and Conclusion

### 5.1 Key Findings

Our experimental results support several key conclusions:

1) The effectiveness of our approach in improving caption quality, evidenced by both automated and human evaluation. Through classificationguided prompting, knowledge distillation, and DPO refinement, our method demonstrates consistent improvements over the base model.

2) The validity of Pixtral as an automated evaluation metric, supported by strong correlation with human judgments. This validation enables efficient automatic evaluation of large-scale image captioning systems.

3) The efficiency of our pipeline, with our final model (DPO 4 iterations) achieving strong performance (4.312) while maintaining practical deployment requirements. This demonstrates the effectiveness of our knowledge distillation approach in transferring capabilities from larger models.

4) The complementary benefits of different components: classification-guided prompting for lan-

646

647

648

649

650

651

652

653

654

655

656

657

658

659

612

613

614

615

guage quality improvement, knowledge distillation for model capability transfer, and DPO for continuous optimization. Each component contributes to the overall performance gain, leading to consistent improvements over the base model.

> 5) We introduce a fine-tuned DPO version of the LLaVA model, which demonstrates strong image captioning capabilities, further enhancing the practical application of the model.

### 5.2 Dataset Contribution

565

566

567

570

571

573

574

576

577

578

582

583

584

585

587

590

591

593

598

599

601

606

607

608

611

A significant outcome of this work is our curated dataset, which includes 9,840 training images and 796 validation images, spanning 18 diverse categories. This dataset offers several unique advantages:

• High-quality GPT-4o-generated captions: Each image is paired with a detailed caption generated by the GPT-4o model, which accurately captures domain-specific details. We have carefully designed category-specific prompting strategies to ensure that the descriptions reflect the unique context and information of each category.

 Balanced category distribution: We ensured that both the training and validation sets have a balanced distribution of categories, ensuring comprehensive coverage across all categories. This balance helps the model perform well across different types of images, avoiding bias towards any specific category and improving generalizability.

• Manual verification ensuring data quality and category accuracy: All images have undergone manual verification to ensure their quality and the accuracy of their category labels. Expert reviewers manually checked each image to ensure it is correctly categorized and meets the required standards, eliminating the potential impact of incorrect labels.

• Diverse domain coverage: The dataset spans a wide range of domains, from common scenarios to specialized fields, including animal, architecture, art, daily life, food, medical, smart cities, and more. These categories include everyday scenes as well as specialized fields such as smart transportation, healthcare, and autonomous driving. • Rich Contextual Information and Details: The captions provided with each image extend beyond mere visual descriptions, incorporating comprehensive contextual information. These captions transcend simple visual representation, integrating broader background and nuanced details to provide profound insights into the scene, environment, and wider context.

We believe that this dataset will serve as a valuable resource for future research in domain-specific image captioning and multimodal understanding. It provides researchers with a high-quality, manually verified data foundation that can significantly advance research in image understanding, caption generation, and cross-modal learning. Additionally, the dataset's diversity and high-quality captions offer a rich source of experimental data for training and evaluating domain-specific generation models, contributing to the further development of more refined models for image captioning.

### 5.3 Conclusion

We have presented a hierarchical approach to image captioning that combines classification-guided prompting, knowledge distillation, and iterative preference optimization. Our method demonstrates significant improvements over baseline approaches, achieving performance comparable to much larger models while maintaining practical deployment requirements. We also introduce a fine-tuned DPO version of the LLaVA model, which further enhances the practical application of the model for image captioning. Together with our released dataset, this work contributes both methodological advances and valuable resources to the field of image captioning. We believe these contributions will facilitate future research in domain-specific visual understanding and caption generation.

### 5.4 Limitations and Future Work

Although our approach performs well, it has some limitations. The current classification system with 18 categories may not capture all image scenarios, and the model faces challenges with complex scenes containing multiple objects. Additionally, the DPO method shows slow convergence in later iterations. Future work should focus on improving classification systems and enhancing the model's ability to handle complex scenes.

764

765

766

767

768

769

770

#### References

660

667

676

679

687

701

705

706

710 711

712

713

714

- Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, Soham Ghosh, Amélie Héliou, Paul Jacob, Albert Q. Jiang, Kartik Khandelwal, Timothée Lacroix, Guillaume Lample, Diego Las Casas, Thibaut Lavril, Teven Le Scao, Andy Lo, William Marshall, Louis Martin, Arthur Mensch, Pavankumar Muddireddy, Valera Nemychnikova, Marie Pellat, Patrick Von Platen, Nikhil Raghuraman, Baptiste Rozière, Alexandre Sablayrolles, Lucile Saulnier, Romain Sauvestre, Wendy Shang, Roman Soletskyi, Lawrence Stewart, Pierre Stock, Joachim Studnia, Sandeep Subramanian, Sagar Vaze, Thomas Wang, and Sophia Yang. 2024. Pixtral 12b. Preprint, arXiv:2410.07073.
  - Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2016. Spice: Semantic propositional image caption evaluation. *Preprint*, arXiv:1607.08822.
  - Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. *Preprint*, arXiv:1707.07998.
  - Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv* preprint arXiv:2308.12966.
  - Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
  - Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2022. Clipscore: A referencefree evaluation metric for image captioning. *Preprint*, arXiv:2104.08718.
  - Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *Preprint*, arXiv:1503.02531.
  - Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: Bootstrapping Language-Image Pretraining With Frozen Image Encoders and Large Language Models. *arXiv preprint*.
  - Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning. *Preprint*, arXiv:2310.03744.
  - Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval:

NLG evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.

- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Ahmed Masry, Do Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263– 2279, Dublin, Ireland. Association for Computational Linguistics.
- Ron Mokady, Amir Hertz, and Amit H Bermano. 2021. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn.
  2024. Direct preference optimization: Your language model is secretly a reward model. *Preprint*, arXiv:2305.18290.
- Amir Rasouli, Iuliia Kotseruba, and John K Tsotsos. 2017. Are they going to cross? a benchmark dataset and baseline for pedestrian crosswalk behavior. In *ICCVW*, pages 206–213.
- Amir Rasouli, Iuliia Kotseruba, and John K Tsotsos. 2018. It's not all about size: On the role of data properties in pedestrian detection. In *ECCVW*.
- Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. *Preprint*, arXiv:1411.5726.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, and Yoshua Bengio. 2016. Show, attend and tell: Neural image caption generation with visual attention. *Preprint*, arXiv:1502.03044.

Tianyi Zhang\*, Varsha Kishore\*, Felix Wu\*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Eval-uating text generation with bert. In *International Conference on Learning Representations*. 

783

784

787

789

790

791

793

794

796

797

799

800

801

808

809

810

811

813

814

815

# A Category-Specific Prompts

This appendix provides the complete set of
category-specific prompts used in our research.
Each prompt follows the structured design described in Section 3.1. Below we present the
prompts for all 18 categories, illustrating our comprehensive approach to domain-specific image captioning.

A.1 Animal Category: Biological Entity Focus

System Prompt: Generate a detailed and accurate caption for the provided animal image. The caption should include the following elements: 1. Accuracy: - Ensure the description is precise and free of errors. 2. Animal Species: - Mention the species of the animal if identifiable 3. Physical Characteristics: - Describe the physical features of the animal (e.g., color, size, distinguishing marks) 4. Behavior: - Describe the behavior or activity of the animal in the image 5. Habitat: - Mention the environment or setting where the animal is located (e.g., forest, desert, ocean) 6. Contextual Details: - Provide any additional relevant information (e.g., weather conditions, time of

> day) User Text: "You are given an image of an animal. Please generate a detailed and accurate caption for the image without any additional information."

### A.2 Architecture Category: Built Environment Focus

### 816 System Prompt:

817Generate a detailed and accurate caption818for the provided architecture image. The

	caption should include the following ele- ments:	819 820
	1. Accuracy:	821
	- Ensure the description is precise and free of errors.	822 823
	2. Building Type:	824
	- Mention the type of building (e.g., residential, commercial, historical)	825 826
	3. Architectural Style:	827
	- Describe the architectural style (e.g., modern, Gothic, Baroque)	828 829
	4. Physical Characteristics:	830
	- Describe the physical features of the building (e.g., height, materials, distinc- tive elements)	831 832 833
	5. Surroundings:	834
	- Mention the building's surroundings (e.g., urban area, countryside)	835 836
	6. Contextual Details:	837
	- Provide any additional relevant infor- mation (e.g., weather conditions, time of day)	838 839 840
	<b>User Text</b> : "You are given an image of an architectural structure. Please gener- ate a detailed and accurate caption for the image without any additional infor- mation."	841 842 843 844 845
A.3	Art Category: Creative Work Focus	846
	System Prompt:	847
	Generate a detailed and accurate caption for the provided art image. The caption should include the following elements:	848 849 850
	1. Accuracy:	851
	- Ensure the description is precise and free of errors.	852 853
	2. Art Form:	854
	- Mention the form of art (e.g., painting, sculpture, digital art)	855 856
	3. Style:	857
	- Describe the style of the artwork (e.g., abstract, realism, impressionism)	858 859
	4. Elements and Composition:	860

861 862	- Describe the elements and composition of the artwork (e.g., colors, shapes, sub-		- Mention speed and relative position to other vehicles	905 906
863	jects)		6. Potential Hazards:	907
864	5. Context:		- Point out any potential dangers or situa-	908
865	- Mention any relevant context (e.g.,		tions requiring special attention	909
866	artist, period, cultural significance)		7. Autonomous Driving Relevant Ele-	910
867	6. Emotional and Interpretative Ele-		ments:	911
868	ments:		- Describe elements particularly impor-	912
869	- Provide interpretation or emotional ele-		tant for the autonomous driving system	913
870	ments conveyed by the artwork		- Mention situations that may require spe-	914
871	User Text: "You are given an image of		cial handling	915
872	an artwork. Please generate a detailed		User Text: "You are given an image re-	916
873	out any additional information "		lated to autonomous driving. Please gen-	917
07.1			erate a detailed and accurate caption for	918
875	A.4 Autonomous Driving Category: Vehicle		the image without any additional infor-	919
876	Environment Focus		mation.	920
877	System Prompt:	A.5	Cartoon Category: Animated Content	921
878	Generate a detailed and accurate caption		Focus	922
879	for the provided autonomous driving im-		System Prompt:	923
880	age. The caption should include the fol-		Generate a detailed and accurate caption	924
001			for the provided cartoon image. The	925
882	1. Accuracy:		caption should include the following ele-	926
883	- Ensure the description is precise and free of errors		ments:	927
004			1. Accuracy:	928
885	2. Scene Overview:		- Ensure the description is precise and	929
886	- Briefly describe the overall scene (e.g.,		free of errors.	930
887	urban street, nighway, suburban area)		2. Character Description:	931
888	- Mention weather conditions and time of		- Mention the main characters and their	932
009			characteristics	933
890	3. Road Environment:		3. Content Description:	934
891	- Describe road type (e.g., two-way four-		- Describe the content of the image	935
892	lane, one-way sheet)		4 Scene Description	036
893	- Mention road surface conditions (e.g.,		Describe the same or setting (a.g. loss	000
094	Describered and increase tractions		tion time of day)	937
895	- Describe road markings and traffic signs		5 Actions and Interactions	550
896	4. Traffic Conditions:		5. Actions and Interactions:	939
897	- Describe positions and behaviors of sur-		- Describe the actions and interactions of	940
898	rounding vehicles		the characters	941
899	- Mention pedestrians, cyclists, and other		6. Visual Style:	942
900	road users		- Mention the visual style of the cartoon	943
901	- Describe traffic light status (if visible)		(e.g., color palette, drawing style)	944
902	5. Ego Vehicle Status:		7. Contextual Details:	945
903	- Describe the autonomous vehicle's po-		- Provide any additional relevant informa-	946
904	sition and driving status		tion (e.g., mood, tone)	947

948	User Text: "You are given an image of a	- Mention notable landmarks or buildings	992
949	cartoon. Please generate a detailed and	4. Activity and Movement:	993
950 951	accurate caption for the image without any additional information."	- Describe any visible activity or move- ment (e.g., traffic, pedestrians)	994 995
952	A.6 Chart Category: Data Visualization	5. Weather and Time of Day:	996
953	Focus	- Mention weather conditions and time of	997
954	System Prompt:	day	998
955	Generate a detailed and accurate cap-	6. Contextual Details:	999
956	tion for the provided chart image. The	- Provide any additional relevant infor-	1000
957	caption should include the following ele-	mation (e.g., cultural or historical signifi-	1001
900		cance)	1002
959	1. Accuracy:	User Text: "You are given an image of	1003
960	- Ensure the description is precise and	a city. Please generate a detailed and	1004
961	free of errors.	accurate caption for the image without	1005
962	2. Details:	any additional information."	1006
963	- Include relevant specifics such as the	A.8 Daily Life Category: Everyday Scene	1007
964	type of chart, data points, trends, and key	Focus	1008
900		System Prompt:	1009
966	3. Consistency:	Generate a detailed and accurate caption	1010
967	- Maintain a consistent style and format.	for the provided daily life image. The	1011
968	4. Readability:	caption should include the following ele-	1012
969	- Write in clear, concise language, avoid-	ments:	1013
970	ing unnecessary complexity.	1. Accuracy:	1014
971	5. Insights or Recommendations:	- Ensure the description is precise and	1015
972	- If appropriate, provide simple insights	tree of errors.	1016
973	or recommendations based on the chart	2. Scene Overview:	1017
974	data.	- Briefly describe the overall scene (e.g.,	1018
975	User Text: "You are given an image of a	home, office, street)	1019
976 977	and accurate caption for the image with-	3. Activities:	1020
978	out any additional information."	- Describe the activities of people in the	1021
		scene	1022
979	A.7 City Category: Urban Environment Focus	4. Objects and Environment:	1023
981	System Prompt:	- Mention notable objects and environ- mental details	1024 1025
982	Generate a detailed and accurate caption	5. Emotional and Social Context:	1026
983	for the provided city image. The caption	- Describe the emotional and social con-	1027
984	should include the following elements:	text of the scene	1027
985	1. Accuracy:	6 Contextual Details	1029
986	- Ensure the description is precise and	Drouide any additional relevant informa	1020
987	free of errors.	tion (e.g., time of day, weather)	1030
988	2. Cityscape Overview:	User Text: "You are given an image de-	1022
989	- Briefly describe the overall cityscape	picting daily life. Please generate a de-	1032
990	(e.g., skyline, street view)	tailed and accurate caption for the image	1034
991	3. Landmarks and Buildings:	without any additional information."	1035

1037System Prompt:rivers, trees, cliffs)1038Generate a detailed and accurate caption for the provided food image. The caption should include the following elements:4. Weather and Time of Day: - Mention weather conditions and time of day1040should include the following elements:- Mention weather conditions and time of day10411. Accuracy:5. Flora and Fauna:1042- Ensure the description is precise and free of errors Describe any visible flora and fauna10436. Contextual Details:10442. Dish Name: able- Provide any additional relevant informa- tion (e.g., location, season)1045- Mention the name of the dish if identifi- ableUser Text: "You are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."	
1038Generate a detailed and accurate caption for the provided food image. The caption should include the following elements:4. Weather and Time of Day:1039for the provided food image. The caption should include the following elements:- Mention weather conditions and time of day1040should include the following elements:- Mention weather conditions and time of day10411. Accuracy:5. Flora and Fauna:1042- Ensure the description is precise and free of errors Describe any visible flora and fauna1043free of errors.6. Contextual Details:10442. Dish Name: able- Provide any additional relevant informa- tion (e.g., location, season)1045- Mention the name of the dish if identifi- ableUser Text: "You are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."	
1039for the provided food image. The caption should include the following elements:- Mention weather conditions and time of day1040should include the following elements:- Mention weather conditions and time of day10411. Accuracy:5. Flora and Fauna:1042- Ensure the description is precise and free of errors Describe any visible flora and fauna1043free of errors.6. Contextual Details:10442. Dish Name: - Mention the name of the dish if identifi- able- Provide any additional relevant informa- tion (e.g., location, season)1045- Mention the name of the dish if identifi- ableUser Text: "You are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."	
1040should include the following elements:day10411. Accuracy:5. Flora and Fauna:1042- Ensure the description is precise and free of errors Describe any visible flora and fauna1043free of errors.6. Contextual Details:10442. Dish Name: - Mention the name of the dish if identifi- able- Provide any additional relevant informa- tion (e.g., location, season)1045- Mention the name of the dish if identifi- ableUser Text: "You are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."	
10411. Accuracy:5. Flora and Fauna:1042- Ensure the description is precise and free of errors Describe any visible flora and fauna1043free of errors.6. Contextual Details:10442. Dish Name: - Mention the name of the dish if identifi- able- Provide any additional relevant informa- tion (e.g., location, season)1045- Mention the name of the dish if identifi- able- User Text: "You are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."	
1042- Ensure the description is precise and free of errors Describe any visible flora and fauna1043free of errors.6. Contextual Details:10442. Dish Name: - Mention the name of the dish if identifi- able- Provide any additional relevant informa- tion (e.g., location, season)1045- Mention the name of the dish if identifi- able- User Text: "You are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."	
1043free of errors.6. Contextual Details:10442. Dish Name:- Provide any additional relevant informa- tion (e.g., location, season)1045- Mention the name of the dish if identifi- able- Provide any additional relevant informa- tion (e.g., location, season)1046ableUser Text: "You are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."1048- Describe the main ingredients used	
10442. Dish Name:- Provide any additional relevant informa- tion (e.g., location, season)1045- Mention the name of the dish if identifi- abletion (e.g., location, season)1046ableUser Text: "You are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."1048- Describe the main ingredients used	
1045- Mention the name of the dish if identifi- ableHon (e.g., focation, season)1046ableUser Text: "You are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."1047- Describe the main ingredients usedany additional information."	
1046ableOser Text.Total are given an image of a landscape. Please generate a detailed and accurate caption for the image without any additional information."10473. Ingredients:accurate caption for the image without any additional information."	
10473. Ingredients:accurate caption for the image without1048- Describe the main ingredients usedany additional information."	
1048- Describe the main ingredients usedany additional information."	
10494. Presentation:A.11Love Category: Emotional Focus	
1050- Describe the presentation and plating ofSystem Prompt:	
1051 the dish Generate a detailed and accurate caption	
10525. Appearance:for the provided love-themed image. The	
1053- Describe the appearance of the dishcaption should include the following ele-	
1054 6. Context:	
1055 - Mention any relevant context (e.g., cui-	
1056 sine type, occasion) - Ensure the description is precise and free of errors	
1057     7. Sensory Details:     2. Scene Overview:	
1058 - Provide sensory details (e.g., aroma, Briefly describe the overall scene	
1059 taste, texture) 3 Emotional Elements:	
1060 User Text: "You are given an image of Describe the emotions depicted in the	
1061 food. Please generate a detailed and ac-	
1062curate caption for the image without any21063additional information."4. Activities:	
- Mention any specific activities or inter-	
1064     A.10     Landscape Category: Natural Scene     actions       1065     Focus     actions	
5. Visual Details:	
- Provide visual details that enhance the	
1067 Generate a detailed and accurate capiton 1068 for the provided landscape image. The	
1069 caption should include the following ele- User Text: "You are given an image with	
1070 ments: a love theme. Please generate a detailed	
1071 1. Accuracy: out any additional information."	
- Ensure the description is precise and	
1073     free of errors.     A.12     Medical Category: Healthcare Focus	
1074   2. Scene Overview:   System Prompt:	
1075 - Briefly describe the overall landscape Generate a detailed and accurate caption for the provided medical image. The	
1076 (e.g., mountains, beach, forest) caption should include the following ele-	
10773. Natural Features:ments:	

1122	1. Accuracy:	User Text: "You are given an image of
1123	- Ensure the description is precise and	a person. Please generate a detailed and
1124	free of errors.	accurate caption for the image without any additional information "
1125	2. Scene Overview:	ung udditional miormation.
1126	- Briefly describe the overall scene	A.14 Plant Category: Botanical Focus
1127	3. Medical Procedures:	System Prompt:
1128	- Describe any visible medical proce-	Generate a detailed and accurate cap-
1129	dures or activities	tion for the provided plant image. The
1130	4. Medical Equipment:	ments.
1131	- Mention notable medical equipment and	
1132	instruments	1. Accuracy.
1133	5. Contextual Details:	- Ensure the description is precise and free of errors.
1134	- Provide any additional relevant informa-	2 Plant Species
1135	tional atmosphere)	2. That Speeks.
1137	User Text: "You are given an image de-	- Mention the species of the plant if iden-
1138	picting a medical scene. Please generate	3 Physical Characteristics
1139	a detailed and accurate caption for the	Describe the abusical features of the
1140	image without any additional informa-	- Describe the physical features of the plant (e.g. color size distinctive ele-
1141	tion.	ments)
1142	A.13 Person Category: Human Subject Focus	4. Environment:
1143	System Prompt:	- Mention the environment or setting
1144	Generate a detailed and accurate cap-	where the plant is located (e.g., garden,
1145	tion for the provided person image. The	forest)
1140	ments:	5. Contextual Details:
1148	1. Accuracy:	- Provide any additional relevant informa-
1149	- Ensure the description is precise and	tion (e.g., season, time of day)
1150	free of errors.	User Text: "You are given an image of
1151	2. Physical Appearance:	a plant. Please generate a detailed and
1152	- Describe the person's physical appear-	accurate caption for the image without
1153	ance (e.g., age, gender, clothing)	any additional information.
1154	3. Actions and Activities:	A.15 Smart Cities Category: Urban
1155	- Mention any actions or activities the	Technology Focus
1156	person is engaged in	<b>System Prompt</b> : Generate a detailed and
1157	4. Emotional State:	city image. The caption should include
1158	- Describe the person's emotional state or	the following elements: 1. Accuracy:
1159	expression	- Ensure the description is precise and
1160	5. Context and Setting:	free of errors. 2. <b>Cityscape Overview</b> :
1161	- Provide context about the setting or	(e.g. skyline street view) 3 Smart In-
1162	background	frastructure: - Mention notable smart
1163	6. Contextual Details:	infrastructure elements (e.g., smart build-
1164	- Provide any additional relevant informa-	ings, IoT devices) 4. Technological Ele-
1165	tion (e.g., weather, time of day)	<b>ments:</b> - Describe visible technological

211	elements (e.g., sensors, automated sys-
212	tems) 5. Traffic and Transportation:
213	- Mention traffic and transportation sys-
214	tems (e.g., autonomous vehicles, smart
215	traffic lights) 6. Contextual Details: -
216	Provide any additional relevant informa-
217	tion (e.g., time of day, weather)

**User Text**: "You are given an image of a smart city. Please generate a detailed and accurate caption for the image without any additional information."

1218

1219

1220

1221

1222

1223

1224

1226

1227

1228

1229

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253 1254

1255

1256

1257

1258

#### A.16 Sports Category: Athletic Activity Focus

System Prompt: Generate a detailed and accurate caption for the provided sports image. The caption should include the following elements: 1. Accuracy: - Ensure the description is precise and free of errors. 2. Sport Type: - Mention the type of sport being played 3. Players and Actions: - Describe the players and their actions 4. Equipment: - Mention any visible sports equipment 5. Setting: -Describe the setting (e.g., stadium, field) 6. Contextual Details: - Provide any additional relevant information (e.g., score, spectators, weather)

**User Text**: "You are given an image of a sports scene. Please generate a detailed and accurate caption for the image without any additional information."

### A.17 Transportation Category: Mobility Focus

System Prompt: Generate a detailed and accurate caption for the provided transportation image. The caption should include the following elements: 1. Accuracy: - Ensure the description is precise and free of errors. 2. Vehicle Type: -Mention the type of vehicle(s) shown 3. Scene Overview: - Briefly describe the overall scene (e.g., road, railway, airport) 4. Activity and Movement: - Describe any visible activity or movement of the vehicles 5. Environment: - Mention the environment or setting (e.g., urban, rural) 6. Contextual Details: - Provide any additional relevant information (e.g., weather conditions, time of day)

User Text: "You are given an image of<br/>a transportation scene. Please generate a1259detailed and accurate caption for the im-<br/>age without any additional information."1261

1263

1291

1292

1293

#### A.18 Other Category: General Content Focus

System Prompt: Generate a detailed and 1264 accurate caption for the provided image. 1265 The caption should include the following 1266 elements: 1. Accuracy: - Ensure the de-1267 scription is precise and free of errors. 2. 1268 **Content Description**: - Provide a clear and relevant description of the content 1270 in the image. 3. In-depth Understand-1271 ing: - Go beyond surface-level details 1272 by interpreting the underlying meaning 1273 or significance of the scene, including 1274 any implied relationships, emotions, or 1275 actions. 4. Contextual Details: - In-1276 clude any necessary contextual details 1277 (e.g., cultural, historical, social signifi-1278 cance) to enhance the understanding of 1279 the image. 5. Visual and Environmen-1280 tal Elements: - Describe notable visual and environmental elements present in 1282 the image, paying attention to subtle fea-1283 tures that may influence the overall inter-1284 pretation. 1285

User Text:"You are given an image.1286Please generate a detailed and accurate1287caption for the image without any addi-1288tional information, focusing on both sur-1289face details and deeper interpretations."1290

We plan to make these prompts and associated implementations available to facilitate future research in domain-specific image captioning.

B Training Details	1294
B.1 Model Configurations	1295
Classification Model	1296
• Base model: LLaVA-v1.5-7B	1297
• Training strategy: Full parameter fine-tuning	1298
• Learning rate: 2e-5	1299
• Batch size: 16 per device	1300
• Number of epochs: 3	1301

302	<b>Caption Model (Full Parameter)</b>	Based on the criteria provided, does the	1338
303	• Base model: LLaVA-v1.5-7B	Plasse rate the ception on a scale of 1 to	1040
304	• Training data: GPT-40 generated captions	5, where:	1340
305	• Learning rate: 2e-5	• 1 - Poor: Does not meet the criteria	1342
306	• Batch size: 16 per device	• 2 - Fair: Meets a few criteria but has	1343
307	• Training epochs: 3	significant issues.	1345
308	Caption Model (LoRA)	• 3 - Good: Meets most criteria with minor issues	1346
309	• LoRA rank: 128	• 4 - Very Good: Meets all criteria	1348
310	• LoRA alpha: 256	<ul> <li>with very few issues.</li> <li>5 - Excellent: Perfectly meets all</li> </ul>	1349
311	• LoRA dropout: 0	criteria.	1351
312	• Learning rate: 2e-4	Be as objective as possible. After provid-	1352
313	• Target modules: all	response on a scale of 1 to 5 by strictly	1353 1354
314	DPO Training	following this format: "Rating: [[rat- ing]]" for example: "Rating: [[5]]"	1355 1356
315	• LoRA rank: 8	Here is the image and the corresponding	1357
316	• LoRA alpha: 16	caption.	1358
317	• Learning rate: 2e-6	Ear detailed evaluation implementations and ari	1005
318	• Batch size: 2 per device	teria, we will make these publicly available in our	1361
319	• Gradient accumulation steps: 8	upcoming repository to facilitate future research in this area.	1362 1363
320	B.2 Computing Infrastructure		
321	• Hardware: 8 NVIDIA A100 GPUs (80GB)		
322	• Framework: LlamaFactory with DeepSpeed		
323	• Mixed precision: BF16		
324	For detailed training scripts and configurations,		
325	we will make these publicly available in our up-		
326 327	this area.		
328	C Pixtral Evaluation Implementation		
329	The Pixtral evaluation system uses a standardized		
330	prompt structure for evaluating captions:		
331	You are a helpful and precise assistant		
332	for evaluating the quality of captions. A		
333	caption will be provided for a particular		
335	tion:		
000	(Critaria Dania) (astronomic i fa ari		
336 337	<criteria begin=""> {category-specific cri- teria} <criteria end=""></criteria></criteria>		