

Seeking Universal Shot Language Understanding Solutions

Anonymous CVPR submission

Paper ID 4

Abstract

001 *Shot language understanding (SLU) is crucial for inter-*
002 *preting narrative, emotion, and aesthetic style in filmmak-*
003 *ing. Although vision-language models (VLMs) demonstrate*
004 *strong general capabilities, they struggle with the complex,*
005 *multi-dimensional nature of cinematography, primarily due*
006 *to a severe lack of high-quality training data. To bridge this*
007 *gap, we introduce SLU-SUITE, a comprehensive dataset*
008 *comprising 490K human-annotated image and video QA*
009 *pairs that cover 33 distinct tasks across six cinematic di-*
010 *mensions. Leveraging SLU-SUITE, we originally identify*
011 *that the core bottleneck of VLMs for SLU is semantic align-*
012 *ment with expert boundaries and terminology, rather than*
013 *raw visual perception. Guided by this insight, we propose a*
014 *balanced data scheduling and parameter-efficient strategy*
015 *to train a universal SLU model, UniShot-8B. Comprehen-*
016 *sive evaluations, including in-domain and out-of-domain*
017 *settings, demonstrate the superiority of UniShot-8B.*

018 1. Introduction

019 High-quality filmmaking relies not only on the script, but
020 also on *shot language* [2, 8, 12] to convey narrative mean-
021 ing, emotion, and aesthetic style. Shot language under-
022 standing (SLU) is highly complex: (1) it spans multiple as-
023 pects of cinematography, such as composition, shot scale,
024 camera motion, lighting, colors, and inter-shot connections;
025 and (2) even within a single aspect, different experts of-
026 ten adopt different taxonomies (e.g., different composition
027 rules or different levels of granularity for camera motion).

028 Vision-language models (VLMs) show strong perfor-
029 mance on general visual understanding. However, recent
030 studies [4, 7, 8, 10, 11] reveal significant discrepancies be-
031 tween VLM judgments and film experts on SLU tasks. Con-
032 sequently, there is a pressing need to develop universal mod-
033 els for SLU that can align well with film experts, handle
034 diverse SLU tasks, and generalize to unseen scenarios.

035 However, a major obstacle is the lack of large-scale,
036 high-quality training data. As shown in Table 1, existing
037 human-annotated SLU datasets [7–10] usually cover only a
038 small subset of tasks and are limited in scale. Meanwhile,

current VLM-labeled datasets[6, 12], such as using GPT-
4V, exhibits underlying misalignment issues. Thus, we in-
troduce **SLU-SUITE**, a comprehensive training and eval-
uation suite for general SLU. SLU-SUITE organizes data
into six film-grounded dimensions [2–4] (*composition, cov-
erage, viewpoint, motion, lighting, and cuts*) and covers **33**
detailed tasks with **490K human-labeled QA pairs**.

Building upon SLU-SUITE, we further study **VLM fine-
tuning strategies for building universal SLU models**. We
propose a simple data scheduling strategy to balance learn-
ing across different dimensions, together with a parameter-
efficient adaptation strategy, i.e., finetuning the language
model only. Specifically, we empirically demonstrate that
the SLU performance of VLMs is mainly bottlenecked by
semantic alignment rather than visual perception.

Based on these insights, we efficiently fine-tune Qwen3-
VL-8B to obtain our universal SLU model, **UniShot-8B**.
Through **evaluations on 23 SLU tasks**, UniShot-8B shows
significant superiority: (1) in **in-domain (ID)** evaluation,
one single UniShot-8B surpasses 12 task-specific fine-tuned
VLMs in 9/12 cases; and (2) in **out-of-domain (OOD)** eval-
uation, UniShot-8B outperforms Gemini-3.0-Flash and Pro
by an average 5% accuracy improvement.

In short, our contributions are four-fold: the construc-
tion of the comprehensive dataset suite, the study of precise
VLM fine-tuning strategies, the training of a universal SLU
model, and comprehensive ID and OOD evaluations.

2. SLU-SUITE: Training and Evaluation Suite for General Shot Language Understanding

2.1. The Datasets

To mitigate the scarcity of data in shot language understand-
ing (SLU), we build **SLU-SUITE**, a unified suite designed
for *general* SLU training and evaluation. Compared with
existing human-labeled cinematography datasets (Table 1),
SLU-SUITE covers a broader set of tasks and explicitly
supports both in-domain (ID) and out-of-domain (OOD)
evaluation. As detailed in Table 2, SLU-SUITE spans six
high-level dimensions, includes **33 detailed task**, and con-
tains **~490K human-labeled QA pairs (including 53,607**

Human-Labeled Dataset	Scale	Training Supp.	OOD Eval.
ShotBench [8]	8 tasks ~70k QA 5 tasks	✓	×
CameraBench [7]	~150k labels 7 tasks	✓	×
CineTechBench [10]	~0.7k QA 1 task	×	N/A
MovieCuts [9]	~110k clips	✓	×
SLU-SUITE (Ours)	33 tasks ~490k QA	✓	✓

Table 1. **Comparison with representative human-labeled cinematography datasets.** **Scale** reports the number of tasks and the approximate data size. **Training Supp.** indicates whether official splits support fine-tuning/training. **OOD Eval.** indicates whether the benchmark supports controlled out-of-domain evaluation by design (e.g., unseen sources/taxonomies/tasks). CineTechBench is evaluation-oriented without training partitions, thus OOD is marked as N/A. **Our SLU-SUITE** covers a broader set of tasks, provides with large-scale human-labeled data, and explicitly supports both in-domain (ID) and out-of-domain (OOD) evaluation to support general shot language understanding research, especially model finetuning/training.

078 image QA pairs and 433,718 video QA pairs.

079 Organizing 33 Tasks into 6 Film-Grounded Dimensions.

080 Shot language understanding data are highly heteroge-
081 neous: different sources use different label spaces (e.g.,
082 coarse vs. fine-grained shot scale), some tasks are multi-
083 label (option combinations), and some tasks are binary or
084 free-form (caption). Training and evaluating on a flat list of
085 tasks makes it hard to control coverage, prevent taxonomy
086 confusion, and analyze transfer.

087 We therefore organize SLU-SUITE into **six high-level**
088 **dimensions** that follow standard cinematography controls
089 used in film analysis and on-set communication (e.g.,
090 framing/composition, scale and lensing, viewpoint, camera
091 movement, lighting/color, and editing/cuts) and are consis-
092 tent with structured filmmaking taxonomies [2–4].

093 As summarized in Table 2, each dimension contains mul-
094 tiple *task variants*:

- 095 • *Composition* (**3** tasks) captures how subjects and vi-
096 sual mass are arranged in the frame, including layout
097 templates and weight/placement cues (*LayoutTemplate_**,
098 *WeightPlacement*).
- 099 • *Coverage* (**9** tasks) models what a shot covers and how
100 coverage is constructed, including shot scale taxonomies
101 (*Scale_**), staging-based coverage patterns (*Staging*),
102 mixed staging–scale labels (*StagingScaleMix*), and lens-
103 related cues (*FocalLength*).
- 104 • *Viewpoint* (**6** tasks) describes camera placement relative
105 to the subject and scene, including angle taxonomies (*Angle_**)
106 and camera height (*Height*, *Height.Cartoon*).
- 107 • *Motion* (**10** tasks) covers both movement taxonomies
108 and movement attributes: coarse vs. compound move-
109 ments (*Move_**), motion description (*Move.Caption_**),

and attribute-focused VQA such as speed, shaking, and
complexity (*VQA_**).

- *Lighting* (**4** tasks) captures illumination and color style,
including source/condition, lighting style/attributes, and
palette cues (*SourceCondition*, *Style*, *Attribute*, *Color-
Palette*).
- *Cuts* (**1** task) focuses on editing transitions within a shot-
level context (*InerShotCutType*).

Note that, we consider semantically similar tasks as dif-
ferent detailed tasks when their option sets differs ((e.g.,
Scale_Basic vs. *Scale_Finegrained*).

Constructing via 11 Sources and 4 Stages.

We build SLU-SUITE through a multi-source pipeline. We first iden-
tify 11 dataset sources that satisfy two requirements: (i) an-
notations are produced by humans (or verified by humans),
rather than generated by LLMs; (ii) the sources jointly cover
the six dimensions above, and each dimension is supported
by at least two sub-tasks or has sufficient data volume (e.g.,
≥50K). For each source, we collect the corresponding me-
dia and labels, and remove samples with missing files, bro-
ken links, corrupted content, or decoding failures. We then
filter sensitive content using *Gemini-3.0-Flash*. Finally,
we convert all sources into a unified QA format, includ-
ing both classification-style QA and captioning-style QA.
For classification-style QA, we provide explicit candidate
options in the question and allow multi-label answers.

2.2. Evaluation Protocols

Mitigating media-level leak. Due to curating from mul-
tiple sources, we perform de-duplication across images and
videos to ensure that evaluation media do not appear in the
training set. After de-duplication, we obtain **about 53k**
unique images and 154k unique video clips.

ID and OOD evaluation. We evaluate generalization in
two settings: *ID*, where test tasks are included in training,
and *OOD*, where test tasks are excluded from training.

For OOD evaluation, we hold out **11** classification-
style QA tasks from two existing multi-task SLU datasets,
CineTechBench [10] and *CameraBench* [7], and use them
only for testing.

For the remaining **22** tasks, we select **12** tasks with more
than 3K samples and split each of them into **80%** training
and **20%** ID evaluation. We add all samples from the other
10 tasks to the training set. To further avoid cross-source
leakage, we remove any OOD test sample whose media ap-
pear in the training set.

To keep it simple, we use captioning-style QA only for
training, while both ID and OOD tests use classification-
style QA with explicit answer options. Finally, we use
around **410K** QA pairs for training and around **50K** QA
pairs for evaluation.

High-Level Dimension	Detailed Tasks	Sources Num.	Samples Num.
Composition	3 tasks: <i>CompositionRule; CompositionPattern; VisualWeightPlacement</i>	3	4,841
Coverage	9 tasks: <i>Scale_Basic; Scale_Classic; Scale_Extended; Scale_Finegrained; Scale_Cartoon; Scale_Historical; Staging; StagingScaleMix; FocalLength</i>	7	51,693
Viewpoint	6 tasks: <i>Angle_Basic; Angle_Extended; Angle_Finegrained; Angle_Cartoon; Height; Height_Cartoon</i>	4	37,742
Motion	10 tasks: <i>Move_Coarse; Move_Compound_A; Move_Compound_B; Move_Captioning; VQA_Complexity; VQA_Movement; VQA_Shaking; VQA_Speed; VQA_Presence; VQA_MixedType</i>	5	282,297
Lighting	4 tasks: <i>SourceCondition; Style; Attribute; ColorPalette</i>	2	925
Cuts	1 task: <i>InterShotCutType</i>	1	109,827
Total	33 tasks	11 sources	487,325 QA pairs

Table 2. **SLU-SUITE summary by high-level dimension.** Each row groups task variants under one dimension. **Detailed Tasks** lists task variants in that dimension (count in bold). We define a task variant by a specific option set (label space/taxonomy); thus, semantically similar tasks from different sources are treated as different variants when their option sets differ. **Sources Num.** denotes the number of contributing datasets, and **Samples Num.** denotes the number of labeled QA pairs. Overall, SLU-SUITE includes 33 task variants from 11 datasets, totaling 487,325 human-labeled QA pairs and covering 53,569 unique images and 154,451 unique videos after de-duplication.

160

3. Fine-Tuning Strategies for Universal SLU

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

Universal shot language understanding (SLU) requires learning many task variants with different taxonomies and formats, while keeping the underlying vision-language model (VLM) generalizable. We focus on two practical challenges: (i) **data strategy**—SLU-SUITE is highly imbalanced across dimensions and task variants (Table 2), so naive sampling can over-train a few large groups and hurt balanced performance; and (ii) **parameter strategy**—although modern VLMs are strong on generic vision-language tasks, recent benchmarks suggest they remain weak on SLU, where success often depends on expert terminology and taxonomy boundaries. Blind full-parameter fine-tuning is expensive and may also degrade the model’s general capabilities, which can further reduce generalization. We therefore design a simple data schedule to balance training, and study where parameter-efficient adaptation should be applied to improve SLU most effectively.

178

179

180

181

182

183

184

185

186

187

188

189

Data strategy: balancing performance across dimensions. Sampling examples in proportion to dataset size is dominated by large dimensions (e.g., motion and cuts), which weakens balanced SLU competence. Leveraging the hierarchical organization of SLU-SUITE, we use a **dimension-balanced sampler**. For each mini-batch, we (1) sample one of the six high-level dimensions *uniformly*, and then (2) sample a QA instance from that dimension by first choosing a task variant at random and then choosing an instance from that variant. This schedule prevents large task groups from overwhelming training and encourages the model to revisit all dimensions throughout training.

190

191

192

193

194

195

196

197

Parameter strategy: precisely identifying the bottleneck of VLM. Given the same data strategy above, we study which component limits SLU performance in a typical VLM pipeline: the **vision encoder** (visual perception), the **multimodal projector** (vision–language connector), or the **language model** (semantic decision making). We conduct a controlled **adapter placement** study on Qwen3-VL-8B [1] using *Low-Rank Adaptation (LoRA)* [5]. Specifi-

Table 3. **SLU accuracy with different LoRA adapter placements in Qwen3-VL-8B.** All settings use the same data strategy and hyperparameters; only the adapter target differs.

	None	Vision Encoder	Projector	Language Model	All
Overall Acc. (↑)	0.467	0.567	0.571	<u>0.699</u>	0.702

cally, we compare five settings by enabling LoRA in: (i) none (no adaptation), (ii) vision encoder only, (iii) projector only, and (iv) language model only. (v) all layers. We control variables by using the *same training data, the same LoRA rank, and the same batch size and optimizer*, and only changing the trainable module.

As shown in Table 3, updating the **language model only** yields the largest gain, almost the same performance with updating all layers, but with significantly less parameters. This result suggests that VLMs’ SLU performance is largely bottlenecked by *semantic alignment*—i.e., aligning with film experts’ decision boundaries and taxonomy-specific terms—rather than capturing visual signals. Based on this finding, we adopt **language-model-only LoRA** as our default parameter-efficient fine-tuning strategy.

Implementation details of our SFT solution. Based on these, we fine-tune Qwen3-VL-8B with supervised fine-tuning (next-token prediction loss) for 3 epochs. We train on $2 \times A100$ 80GB GPUs, with batch size 4 per GPU and LoRA rank 32, within 2 days. We denote the resulting universal SLU model as **UniShot-8B**.

4. Experiments

We study two questions: (1) **ID performance:** can our universal multi-task model match or exceed other general or even task-specific solutions? (2) **OOD generalization:** does the model generalize well to unseen tasks?

We compare our **UniShot-8B** with: **Five Universal (one-for-all) Solutions.** We evaluate: (i) Qwen3-VL-8B [1] as an open-source VLM baseline; (ii) *Gemini-2.5-Flash*, *Gemini-3.0-Flash*, *Gemini-3.0-Pro* as powerful commercial VLMs; (iii) *ShotVL-7B* [8] as a SOTA multi-task SLU model via finetuning Qwen2.5-VL-7B; and **Task-specific**

Table 4. **In-domain (ID) results on 12 tasks.** We compare *Universal* (one model for all tasks) against *Specific* (one model per task). **Count #1** reports how many task variants each method ranks first.

12 In-Domain Tasks	Universal (one-for-all)					Specific (one-for-one)	
	Qwen3-VL-8B	Gemini-2.5-Flash	Gemini-3.0-Flash	Gemini-3.0-Pro	ShotVL-7B	UniShot-8B (Ours)	12 task-specific-SFT Qwen3-VL-8B
CompositioRule	0.020	0.005	0.040	0.040	0.775	0.791	0.784
ScaleExtended	0.145	0.165	0.125	0.170	0.115	0.290	0.320
ScaleHistorical	0.670	0.770	0.850	0.825	0.740	0.900	0.855
ScaleCartoon	0.715	0.675	0.650	0.660	0.460	0.940	0.875
ScaleBasic	0.650	0.465	0.595	0.555	0.425	0.875	0.865
Staging	0.430	0.165	0.280	0.295	0.290	0.785	0.685
AngleBasic	0.330	0.420	0.440	0.485	0.245	0.940	0.960
AngleCartoon	0.350	0.515	0.450	0.545	0.260	0.690	0.710
Height	0.630	0.690	0.705	0.670	0.650	0.905	0.875
HeightCartoon	0.540	0.385	0.555	0.580	0.600	0.740	0.720
MoveCoarse	0.745	0.650	0.710	0.585	0.600	0.850	0.840
InterShotCutType	0.230	0.215	0.295	0.270	0.210	0.400	0.390
Avg. Acc. (↑)	0.455	0.427	0.475	0.473	0.448	0.759	0.740
Count #1 (↑)	0/12	0/12	0/12	0/12	0/12	9/12	3/12

Table 5. **Out-of-domain (OOD) results on 11 unseen task variants** (Universal solutions only). OOD tasks use unseen task forms and/or new taxonomies from different sources.

11 Out-of-Domain Tasks	Qwen3-VL-8B	Gemini-2.5-Flash	Gemini-3.0-Flash	Gemini-3.0-Pro	ShotVL-7B	UniShot-8B (Ours)
CompositionPattern	0.083	0.208	0.158	0.525	0.133	0.558
ScaleClassic	0.371	0.486	0.600	0.436	0.507	0.414
FocalLength	0.217	0.333	0.467	0.433	0.300	0.367
AngleExtended	0.650	0.575	0.825	0.792	0.608	0.608
MoveCompound	0.815	0.750	0.810	0.585	0.805	0.880
VQAMovement	0.366	0.452	0.570	0.441	0.301	0.634
VQAComplexity	0.535	0.555	0.570	0.355	0.490	0.800
VQAShaking	0.640	0.615	0.565	0.520	0.615	0.720
VQASpeed	0.510	0.680	0.605	0.645	0.730	0.675
LightingStyle	0.582	0.591	0.600	0.691	0.582	0.700
ColorPalette	0.483	0.600	0.617	0.600	0.550	0.633
Avg. Acc. (↑)	0.478	0.531	0.581	0.547	0.511	0.636
Count #1 (↑)	0/11	0/11	3/11	0/11	1/11	7/11

230 **(one-for-one) Solutions.** We also consider 12 *task-specific-*
 231 *SFT Qwen3-VL-8B models* in ID evaluations, each fine-
 232 tuned on a single training task corresponding to the same
 233 test task.

234 All evaluations use classification-style QA. Accuracy is
 235 reported over three runs with decoding temperature as 0.
 236 For multi-label questions, we use strict exact-match.

237 **ID Results: one UniShot-8B beats 12 task-specific**
 238 **fine-tuning models.** Table 4 summarizes results on 12 ID
 239 task variants. UniShot-8B achieves the best average ac-
 240 curacy (**0.759**) and ranks first on **9/12** tasks, substantially
 241 improving over the base Qwen3-VL-8B and the general-
 242 purpose VLMs. Notably, UniShot-8B also slightly ex-
 243 ceeds the task-specific upper bound on average (0.759 vs.
 244 0.740), suggesting that *multi-task training can provide pos-*
 245 *itive transfer across SLU tasks, where learning one task can*
 246 *strengthen others compared to training them independently.*
 247 Overall, the ID results support our goal of a universal SLU
 248 model that covers many tasks without sacrificing accuracy.

249 **OOD Results: UniShot-8B generalizes well to unseen**

250 **SLU tasks.** Table 5 reports performance on 11 OOD task
 251 variants, which include unseen task forms and/or new tax-
 252 onomies from different sources. UniShot-8B achieves the
 253 best average accuracy (**0.636**) and ranks first on **7/11** tasks,
 254 outperforming the open-source baselines and the stronger
 255 commercial model Gemini-3.0-Flash on average. *The OOD*
 256 *results demonstrate that large-scale multi-task training en-*
 257 *ables UniShot-8B with strong generalization ability, thereby*
 258 *broadly supporting unseen tasks.*

259 **Other Insights.** (I) We also observe that **composi-**
 260 **tion** remains a clear weakness for general-purpose VLMs:
 261 most general VLM baselines fail on the two composition
 262 tasks, with accuracies staying below 20%. (II) In addition,
 263 ShotVL-7B is limited by its narrower SLL task coverage:
 264 although its overall SLU performance is comparable to that
 265 of Gemini-2.5-Flash, it shows notable generalization fail-
 266 ures on several uncovered tasks and even underperforms the
 267 similarly sized Qwen3-VL-8B. These results further moti-
 268 vate the need for a comprehensive and large-scale suite like
 269 our SLU-SUITE to support universal SLU models building.

270

References

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

- [1] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, et al. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025. 3
- [2] David Bordwell, Kristin Thompson, and Jeff Smith. *Film art: An introduction*. McGraw-Hill New York, 2008. 1, 2
- [3] Blain Brown. *Cinematography: theory and practice: image making for cinematographers and directors*. Routledge, 2016.
- [4] Agneet Chatterjee, Rahim Entezari, Maksym Zhuravinskyi, Maksim Lapin, Reshinh Adithyan, Amit Raj, Chitta Baral, Yezhou Yang, and Varun Jampani. Stable cinematics: Structured taxonomy and evaluation for professional video generation. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*. 1, 2
- [5] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Liang Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *Iclr*, 1(2):3, 2022. 3
- [6] Xuan Ju, Yiming Gao, Zhaoyang Zhang, Ziyang Yuan, Xintao Wang, Ailing Zeng, Yu Xiong, Qiang Xu, and Ying Shan. Miradata: A large-scale video dataset with long durations and structured captions. *Advances in Neural Information Processing Systems*, 37:48955–48970, 2024. 1
- [7] Zhiqiu Lin, Siyuan Cen, Daniel Jiang, Jay Karhade, Hewei Wang, Chancharik Mitra, Yu Tong Tiffany Ling, Yuhan Huang, Rushikesh Zawat, Xue Bai, et al. Towards understanding camera motions in any video. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*. 1, 2
- [8] Hongbo Liu, Jingwen He, Yi Jinn, Dian Zheng, Yuhao Dong, Fan Zhang, Ziqi Huang, Yinan He, Weichao Chen, Yu Qiao, et al. Shotbench: Expert-level cinematic understanding in vision-language models. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*. 1, 2, 3
- [9] Alejandro Pardo, Fabian Caba Heilbron, Juan León Alcázar, Ali Thabet, and Bernard Ghanem. Moviecuts: A new dataset and benchmark for cut type recognition. In *European Conference on Computer Vision*, pages 668–685. Springer, 2022. 2
- [10] Xinran Wang, Songyu Xu, Shan Xiangxuan, Yuxuan Zhang, Muxi Diao, Xueyan Duan, Kongming Liang, Zhanyu Ma, et al. Cinetechbench: A benchmark for cinematographic technique understanding and generation. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*. 1, 2
- [11] Hang Wu, Yujun Cai, Haonan Ge, Hongkai Chen, Ming-Hsuan Yang, and Yiwei Wang. Refineshot: Rethinking cinematography understanding with foundational skill evaluation. *arXiv preprint arXiv:2510.02423*, 2025. 1
- [12] Weijia Wu, Mingyu Liu, Zeyu Zhu, Xi Xia, Haoen Feng, Wen Wang, Kevin Qinghong Lin, Chunhua Shen, and Mike Zheng Shou. Moviebench: A hierarchical movie level dataset for long video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 28984–28994, 2025. 1