A Structured Framework for Evaluating and Enhancing Interpretive Capabilities of Multimodal LLMs in Culturally Situated Tasks

Anonymous ACL submission

Abstract

This study aims to test and evaluate the capabilities and characteristics of current mainstream Visual Language Models (VLMs) in generating critiques for traditional Chinese painting. To achieve this, we first developed a quantitative framework for Chinese painting critique. This framework was constructed by multi-dimensional extracting evaluative features-including evaluative stance, core focal points, and argumentative qualityfrom human expert critiques using a zero-shot classification model. Based on these features, several representative critic personas were defined and quantified. This framework was then employed to evaluate selected VLMs (e.g., Gemini 2.5 Pro). The experimental design involved persona-guided prompting to assess the VLM's ability to generate critiques from diverse perspectives. Our findings reveal the current performance levels, strengths, and areas for improvement of VLMs in the domain of art critique, offering insights into their potential and limitations in complex semantic understanding and content generation tasks. The code used for our experiments can be publicly accessed at: https://github.com/anon user/ anon repo¹.

1 Introduction

006

800

011

012

015

017

019

027

032

041

language models (LLMs) have Large demonstrated remarkable performance on general NLP benchmarks, yet their applicability in culturally embedded, humanistic domains remains limited. In high-context interpretive tasks such as art criticism, clinical narrative analysis, or historical commentary, model performance depends not only on linguistic fluency or factual accuracy, but also on deeper forms of cognitive alignment-epistemic sensitivity, rhetorical coherence, and cultural adaptability.

A representative and particularly demanding testbed for such capabilities is Chinese art This genre, especially when commentary. analyzing works like traditional landscape or court paintings, involves symbolic interpretation, aesthetic judgment, and deeply situated cultural discourse. Existing multimodal LLMs are rarely evaluated in this space. Standard benchmarks such as MME (Fu et al., 2024) and MMBench (Liu et al., 2024) focus on object recognition or task-oriented vision-language reasoning, while frameworks like ArtGPT (Chen et al., 2024) emphasize captioning and factual These methods largely overlook grounding. interpretive nuance and disciplinary diversity.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Meanwhile, humanistic commentary often exhibits non-linear logic, specialized lexicons, and varied stylistic conventions, particularly in Chinese art contexts where rhetorical strategies such as *yijing* (意境, artistic conception) or *qiyun shengdong* (气韵生动, spiritual resonance) are essential but difficult to quantify. Without appropriate grounding, LLMs risk producing synthetic outputs that mimic surface patterns but fail to demonstrate epistemic alignment (Guo et al., 2023; Mishra et al., 2024). This growing mismatch calls for new paradigms in evaluation and adaptation.

To address these challenges, we introduce Vision-Understanding VULCA-the and Language-based Cultural Adaptability Framework. VULCA is a structured evaluation and enhancement framework designed to assess how well MLLMs align with domain-specific interpretive practices in culturally situated tasks. Our work centers on Chinese art commentary, but the methodology generalizes to other multimodal and epistemically rich domains such as religion, medicine, or history. VULCA combines three core components: (1) a multi-dimensional human expert benchmark (MHEB) constructed from 163

¹Repository will be linked upon paper acceptance.

art commentaries annotated across five cultural capability dimensions; (2) a persona-guided 084 recontextualization mechanism using eight interpretive personas and a domain-specific knowledge base; and (3) a joint evaluation pipeline integrating vector-space semantic alignment with rubric-based capability scoring. Commentaries are generated from annotated traditional Chinese paintings, and their alignment with expert patterns is evaluated with and without interventions. As a result, we produce five contributions: (i) the definition of VULCA, a new structured framework for assessing and enhancing MLLMs in culturally grounded, multimodal reasoning tasks; (ii) we construct D1, 097 а high-quality human benchmark of Chinese art commentary annotated across five capability dimensions; (iii) we develop and evaluate 100 persona-guided recontextualization interventions 101 using eight expert personas and a domain-specific 102 knowledge base; (iv) we demonstrate over 20% 103 improvement in symbolic reasoning and over 30% improvement in argumentative coherence on 105 Gemini 2.5 Pro using our proposed method; and 106 (v) we establish the generalizability of our 107 evaluation methodology to other epistemically rich domains such as religion, history, and 109 education. 110

Together, our work highlights the need for new evaluation paradigms that go beyond benchmark metrics and toward measuring how well LLMs can adapt to the interpretive demands of real-world, interdisciplinary contexts.

2 Related Work

111

112

113

114

115

116

Missing Evaluation Dimensions for Cultural 117 Reasoning. Existing benchmarks for large or 118 multimodal language models, such as (Fu et al., 119 2024; Li et al., 2023), emphasize factual accuracy 120 or instruction following, seldom addressing 121 symbolic interpretation or epistemic alignment. 122 ArtGPT (Chen et al., 2024), for instance, 123 evaluates stylistic generation but lacks formal 124 metrics for interpretive depth. While prior work 125 explores aesthetic reasoning (Qi, 2024; Wang, 126 2024), these studies rarely offer structured, multi-capability evaluation. Our work introduces 128 cultural adaptability, operationalized through a 129 multi-dimensional human expert benchmark 130 (MHEB) with capability rubrics, enabling 131 quantitative comparison in high-context domains 132

like Chinese art.

Limitations of Persona Conditioning Without Grounding. Persona use in LLM evaluation shows promise for style control (Wang et al., 2023a, 2024), yet most methods lack structured knowledge grounding, especially in epistemically rich domains. Our method combines persona simulation with curated domain-specific knowledge to guide generation towards symbolic reasoning and cultural interpretation, not just stylistic alignment, offering a controlled intervention mechanism. 133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

157

158

159

160

161

162

163

164

165

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

Gap in Multimodal Input–Interpretation Evaluation. Current multimodal frameworks like MMBench or LLaVA-Bench (Zhang et al., 2023) primarily focus on classification, question answering, or instruction following, rarely requiring grounded interpretation. Our pipeline links annotated symbolic elements with structured prompts for art commentary, evaluating MLLM outputs for semantic alignment with MHEB using vector-space and rubric-based metrics, addressing a gap in assessing image-conditioned cultural reasoning.

Lack of Comparative Cultural Interventions Across Models. Surveys (Mishra et al., 2024; Guo et al., 2023) discuss LLM limitations in nuanced discourse, but few studies compare model responsiveness to structured cultural interventions. Our empirical evaluation shows persona and knowledge base intervention improves symbolic reasoning and argumentative coherence by over 20-30%, highlighting epistemic alignment's role beyond fluency. This cross-model, capability-specific analysis distinguishes our work.

3 Methodology

This research aims to comprehensively evaluate Visual Language Models (VLMs) capabilities in generating critiques for traditional Chinese painting, assessing their understanding of image content, commentary quality, and adaptability to guided perspectives. The workflow involves: Framework Construction, developing a quantitative analytical framework from human commentaries, including expert defining evaluative dimensions and critic personas; VLM Evaluation Experiment Design, creating structured protocols for VLM critique generation



Figure 1: Overview of the VULCA framework, illustrating its components and their interactions for structured evaluation and intervention in art criticism.

under conditions like persona-based and baseline prompting; and Experimentation and Result Analysis, implementing experiments, collecting VLM critiques, and analyzing them with the developed framework to assess capabilities and intervention impacts. Figure 1 provides an overview of this framework and its components.

183

184

185

186

187

188

A cornerstone is the quantitative framework 189 benchmark for VLM critiques, built upon human expert commentaries on Chinese art. To ensure 191 objective, reproducible, and fine-grained evaluation, an automated capability assessment 193 framework was developed. This involves feature 194 extraction, multi-dimensional capability scoring, profile assignment, and visualization, using a 196 zero-shot classification model for fine-grained 197 evaluative labels. The scoring covers painting 198 element recognition, Chinese painting 199 understanding, and language usage, each with a dedicated rubric . This structured, rule-based 201 approach enhances objectivity and facilitates large-scale benchmarking (Jiang and Chen, 2025; Hayashi et al., 2024). 204



Figure 2: T-SNE visual representation of human expert art commentaries.

3.1 Feature Engineering from Human Expert Critiques

205

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

226

227

228

229

230

231

232

234

235

236

237

238

239

Framework foundation relies on human expert commentaries, significantly from Giuseppe Castiglione's (Lang Shining) "Twelve Months" (十二月令图) series—Qing imperial court paintings fusing Chinese and Western traditions. To enhance model training and evaluation, a sliding window cropping strategy (640×640 pixel sub-images) was applied to these high-resolution images, augmenting data diversity and granularity for improved VLM detail recognition and evaluation accuracy, a common practice in computer vision (e.g., (Lin et al., 2014; Krishna et al., 2017)).

We employed a zero-shot classification model evaluative to systematically extracted characteristics. This model, proficient in multilingual text classification without task-specific training, objectively identified and scored texts against predefined labels across three dimensions: Evaluative Stance (e.g., "Historical Research"), Core Focal Points (e.g., "Use of Color"), and Argumentative Quality (e.g., "Profound Insight"). This process created a multi-dimensional structured. feature representation for each expert commentary. Appendix D.5 lists these labels. Figure 2 visualizes the MHEB semantic distribution from these features.

The zero-shot classification model serves as an analytical tool for deconstructing expert texts and building our evaluation framework, distinct from the VLMs (e.g., Gemini 2.5 Pro, Qwen-VL) evaluated later.

3.2 Evaluation Dimensions and Label System

240

241

243

244

245

246

247

253

255

256

261

263

265

267

269

270

271

272

275

276

279

To systematically analyze commentary content, we define a structured annotation scheme based on three major dimensions: **Stance**, **Feature Focus**, and **Commentary Quality**. These dimensions were derived from domain-specific literature and refined through pilot studies with expert annotations.

Stance characterizes the rhetorical or evaluative position taken by the commentator (e.g., historical interpretation, praise, or critique). Feature Focus identifies the specific visual or contextual aspects discussed in the commentary line quality, symbolism, (e.g., spatial composition). Commentary Quality captures the analytical depth and logical structure of the commentary, ranging from clear, well-argued insights to superficial or biased remarks.

Each dimension comprises a set of fine-grained subcategories with bilingual English–Chinese mappings. Full definitions and label lists are provided in Appendix D.5.

To illustrate the system, we briefly explain one representative label from each dimension:

Stance -Aesthetic Appreciation (美学鉴赏 型): Commentary focuses on the beauty and expressive power of the painting, often using evocative or poetic language to highlight visual elegance or emotional resonance. Feature Focus -Brushwork Technique (笔法技巧): The analysis emphasizes the artist's brushstroke styles, control, or variation, such as fine lines, dry brush texture, or fluid ink application. Commentary Quality -Profound Insight (见解深刻独到): The argument demonstrates deep understanding, originality, and relevance, going beyond surface observations to offer meaningful interpretations.

These representative sub-dimensions help bridge formal annotation with art historical reasoning. The full taxonomy serves as the foundation for profile construction, persona classification, and performance evaluation.

3.3 Construction and Definition of Critic Personas

To capture holistic critique style and depth beyond granular features, we constructed "critic personas" representing archetypal critical perspectives. Their development was data-driven, analyzing features from human expert commentaries, complemented by art history domain expertise. Five core personas were defined: Comprehensive Analyst (博学通论型), Historically Focused Critic (历史考据型), Technique & Style Focused Critic (技艺风格型), Theory & Comparison Focused Critic (理论比较 型), and General Descriptive Profile (泛化描述 型). 290

291

292

293

295

296

297

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

327

328

329

330

331

332

333

334

335

336

337

338

339

Each persona is quantitatively defined by rules and thresholds based on zero-shot classification This rule-based matching feature scores . objectively assigns commentaries (human or Persona definition and VLM) to personas. matching rely on explicit features and rule-based logic, not primarily direct semantic embedding of raw text. Dimensionality reduction (t-SNE/UMAP) visualizes commentary and persona distribution in the feature space, not for initial persona vector generation.

3.4 Value and Application of the Framework

The resulting quantitative framework, which integrates fine-grained feature analysis with the abstracted critic personas, offers a multi-dimensional, quantifiable, and empirically grounded benchmark. Rooted in the discernible characteristics of human expert critiques, this framework provides a structured and robust foundation for the subsequent systematic evaluation and comparative analysis of Chinese painting commentaries generated by Visual Language Models.

3.5 Experimental Design for VLM Evaluation

This quantitative framework guided experiments evaluating selected VLMs (e.g., Gemini 2.5 Pro, Qwen-VL). The core task required VLMs to generate commentary on provided traditional Chinese painting images. Experiments typically involved structured, multi-round interactions for each VLM per image, including persona-based and baseline Q&A rounds (Zhou et al., 2024; Wang et al., 2023b).

Inputs were multifaceted: high-definition "Monthly Images" (sometimes segmented); predefined "Persona Cards" (Wang et al., 2023b) guiding analysis—Mama Zola (佐拉妈妈), Professor Elena Petrova (埃琳娜·佩特洛娃教授), Okakura Kakuzō (冈仓天心), Brother Thomas (托马斯修士), John Ruskin (约翰·罗斯金), Su Shi (苏轼), Guo Xi (郭熙), Dr. Aris Thorne (阿 里斯·索恩博士); standardized prompt templates (Zhou et al., 2024); and an optional JSON knowledge base (Zhang et al., 2024b; Bin et al., 2024). Persona guidance aimed to assess VLM capability to simulate diverse perspectives and analytical styles (Zhang et al., 2024a).

340

341

342

347

351

353

354

358

359

361

363

367

371

373

374

375

377

379

389

All VLM-generated texts were recorded and systematically organized. These outputs were then analyzed using the quantitative framework (Section 3.2), applying zero-shot classification to extract feature scores and matching critiques against predefined "Critic Personas" to assess alignment, especially under specific persona guidance.

Key VLM evaluation dimensions include: Painting Element Recognition (5-point scale); Chinese Painting Understanding (7-point scale); and Chinese Language Usage (5-point scale). particularly for structured Prompt design, commentary, targeted these dimensions .

Vector Space Representation and 3.6 Visualization

To compare human and VLM critiques, we converted feature scores (evaluative stance, focal points, argumentative quality) from both into numerical vectors. These vectors were projected into a 2D space using t-SNE for visualisation (Reimers and Gurevych, 2019), enabling assessment of semantic similarity and distributional differences. Figure 2 illustrates such a comparative visualisation, showing the semantic distribution of human expert commentaries versus baseline MLLM-generated commentaries, highlighting their initial semantic gap.

These visualizations help analyze how MLLM outputs align with human expert benchmarks, identify specific MLLM strengths/weaknesses, and assess persona/knowledge interventions' impact on aligning MLLM critiques with desired expert profiles.

3.7 Automated Workflow

This research implemented a modular, automated experimental pipeline for profile scoring, dimensionality reduction, and dataset preparation for visualizations.

Experimental benchmarking involved MLLM commentary generation using a curated artwork dataset with varied prompts (baseline, persona-specific, knowledge-enhanced). MLLM outputs were logged, versioned, and organized by

model, persona, and prompt. Subsequent 390 automated analysis involved feature extraction, and comparative metrics persona scoring, 392 generation. This systematic approach facilitated 393 large-scale, reproducible evaluation of MLLM 394 performance in Chinese art critique. 395

391

396

397

399

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

3.8 Multi-Model Comparative Evaluation

To comprehensively assess the capabilities of state-of-the-art large language and vision-language models, we conducted a systematic comparative evaluation across four representative models: Google Gemini 2.5 Pro, Meta Llama-3.1-8B-Instruct, Meta Llama-4-Scout-17B-16E-Instruct, and Qwen-2.5-VL-7B. All models were evaluated using the same experimental protocol, dataset splits, and evaluation metrics to ensure fair and reproducible comparison.

All models were accessed via their official APIs or open-source checkpoints, with inference settings kept consistent. For multimodal tasks, only models supporting both text and image inputs were included in the corresponding benchmarks. The evaluation covers a range of tasks, including argumentative quality, core focal points, stance analysis, and semantic space visualization, as detailed in Section D.

Quantitative Modeling and Formalisms 3.9

This section details the key mathematical formulations used in our analytical framework, covering semantic representation, comparative metrics, and the profile matching algorithm.

Semantic Embedding. Conceptually:

 $\mathbf{v}_d = \text{SentenceTransformer}(\text{document}_d)$ (1)

Where $(\mathbf{v}_d \in \mathbb{R}^N)$ (e.g., (N = 1024) for BAAI/bge-large-zh-v1.5).

Average Quality Score for Radar Chart ($\bar{q}_{i,G}$). For a quality dimension j and a group of documents G (e.g., Human Experts, MLLM Baseline):

$$\bar{q}_{j,G} = \frac{1}{|N_G|} \sum_{d \in N_G} s_{j,d} \tag{2}$$

Where $s_{j,d}$ is the score of document d on quality 431 dimension j, and $|N_G|$ is the number of documents 432 in group G. 433

473

474

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

500

501

502

503

505

457

458

459

Centroid Calculation in Dimensionality Reduced Space (c $_p$). For a profile/condition p, its centroid in a 2D space (e.g., t-SNE):

$$\mathbf{c}_p = (\bar{x}_p, \bar{y}_p)$$
$$= \left(\frac{1}{|D_p|} \sum_{d \in D_p} x_d, \frac{1}{|D_p|} \sum_{d \in D_p} y_d\right) \quad (3)$$

Where (x_d, y_d) are the 2D coordinates of document *d* belonging to profile/condition *p*, and $|D_p|$ is the number of documents in profile/condition *p*.

Cohen's d (Effect Size). To measure the standardized difference between two group means (\bar{X}_1, \bar{X}_2) :

$$d = \frac{\bar{X}_1 - \bar{X}_2}{s_p} \tag{4}$$

Where s_p is the pooled standard deviation:

s

$$_{p} = \sqrt{\frac{(n_{1} - 1)s_{1}^{2} + (n_{2} - 1)s_{2}^{2}}{n_{1} + n_{2} - 2}}$$
 (5)

And here n_1 , n_2 are the sample sizes of group 1 and group 2, while s_1^2 , s_2^2 are the variances of group 1 and group 2.

Stance Contribution Formula (S_C) . We compute the stance contribution S_C using the following conditions:

 S_C

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

$$= \begin{cases} \frac{s_{actual} - s_{min_rule}}{s_{max_rule} - s_{min_rule}}, & \text{if } L_{actual} = L_{rule}, \\ s_{actual} \ge s_{min_rule}, \\ s_{max_rule} \ne s_{min_rule}, \\ 1, & \text{if } L_{actual} = L_{rule}, \\ s_{actual} \ge s_{min_rule}, \\ s_{max_rule} = s_{min_rule}, \\ 0, & \text{otherwise} \end{cases}$$

Where S_C is the stance contribution score, L_{actual} is the actual stance label of the text, L_{rule} is the required stance label in the profile rule, s_{actual} is the actual stance score, and s_{min_rule} , s_{max_rule} represent the required range.

4 Results

We present results from semantic alignment, capability profiling, and the effects of persona-guided interventions on MLLMs. All evaluations are made with respect to the multi-dimensional human expert benchmark (MHEB), using both vector-space analysis and rubric-based scoring.

4.1 Semantic Divergence from Expert Commentary

Baseline MLLM outputs exhibit significant divergence from human expert commentaries. As shown in Figure 3 (left), expert texts cluster tightly in semantic space, while MLLM outputs are more dispersed and form distinct clusters. Profile-based visualizations (Figure 4 (right)) further confirm this divergence: baseline models frequently align with generic or technique-oriented profiles, rarely matching complex expert personas.

4.2 Capability Profile Differences

Human expert commentaries, as quantified by our ZSL analysis (see Table 4 in Appendix E.3 for full data which Figure 4 (left) visualizes), emphasize symbolic and historical interpretation (e.g., average scores of 0.676 in Historical Context and 0.661 in Symbolism) but notably less on technical aspects like Brushwork Technique (0.199). They also exhibit high subjectivity and non-linear reasoning (e.g., 0.674 in Subjective View, 0.093 in Clear Logic, as detailed in Table 7).

In contrast, baseline MLLMs show varied performance. For instance, Llama-4-Scout-17B-16E-Instruct achieves high scores in Historical Context (0.710) and Symbolism (0.758), comparable to or exceeding human experts. Qwen-2.5-VL-7B also performs well in these areas (0.650 and 0.773 respectively) and particularly excels in Artistic Conception (0.891) and Brushwork Technique (0.937) – the latter being dramatically higher than the human expert average of 0.199 for this feature (see 'Table 4'). Gemini-2.5pro shows strength in Structure (0.874),Lavout and while Meta-Llama-3.1-8B-Instruct generally presents lower scores across several nuanced dimensions like Historical Context (0.366) and Symbolism (0.529).These differences are summarized in Figure 4 (left) and supported by the radar plots in Figure 3 (right).

4.3 Effectiveness of Persona-Guided Interventions

Persona-guided prompting, especially when supported by domain knowledge, substantially



Figure 3: Impact of Persona and Knowledge Base Interventions on MLLM Critiques: A comprehensive analysis comparing intervened MLLM outputs with a human expert benchmark. Left: t-SNE and KDE plots visualize the semantic distribution of critiques from different sources (human experts, baseline MLLMs, intervened MLLMs). Right: A radar chart compares average capability scores across dimensions like Profound Insight and Logical Clarity.



Figure 4: Profiling Summary: A comparative visualization of Human Experts vs. MLLMs across key textual features (left), mean profile alignment scores (center), and t-SNE projection of profile vectors (right).

improves MLLM outputs. Figure 3 (right) illustrates that Qwen-2.5-VL improves scores across key dimensions-e.g., Profound Insight (from 0.31 to 0.61), Strong Argumentation (0.33 to 0.66), and Detailed Analysis (0.33 to 0.70), with full details available in 'Table 7'. These results indicate stronger alignment with expert-style reasoning. Alignment improvements are also visible in profile scores (Figure 4 (center)), with intervened outputs matching sophisticated expert types like "Comprehensive Analyst" (e.g., Qwen-2.5-VL-7B achieving an alignment score of 0.778 for this profile, as detailed in 'Table 5') more closely than baseline.

507

508

510

511

512

513

514

515

517

519

523

525

4.4 Cross-Model Comparison and Configurations

Qwen-2.5-VL and LLaMA-4-Scout-17B demonstrate strong performance under intervention. In Figure 4 (left), which visualizes data from Table 4, both models demonstrate high scores in areas like Artistic Conception (Qwen: 0.891, Llama-4: 0.851), Brushwork Technique (Qwen: 0.937, Llama-4: 0.903), and Layout and Structure (Qwen: 0.895, Llama-4: 0.916). Their profile alignment in Figure 4 (center) confirms their ability to emulate multiple expert types. The overall performance rankings, detailed in 'Table 1', reveal that the Qwen-2.5-VL-7B model, when guided by the Mama Zola persona and an external knowledge base, achieved the top composite score (9.2/10) and expert alignment (100%).

526

527

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

These results show that interpretive capability in MLLMs can be substantially improved by structured prompting and domain-specific conditioning. Culturally aligned personas are particularly effective, highlighting the potential of the VULCA framework to guide MLLMs toward expert-level reasoning in specialized domains. The distribution of MLLM outputs in semantic

Table 1: Overall Rankings: Top performing model and persona combinations across capability dimensions.

Rank	Configuration	Composite Score	Expert Alignment
1	Qwen-2.5-VL-7B + Mama Zola (佐拉妈妈) + KB	9.2/10	100%
2	meta-llama_Llama-4-Scout-17B-16E-Instruct + John Ruskin (约翰·罗斯金) + KB	8.9/10	97%
3	meta-llama_Llama-4-Scout-17B-16E-Instruct + Mama Zola (佐拉妈妈) + KB	8.7/10	95%
4	meta-llama_Llama-4-Scout-17B-16E-Instruct + Brother Thomas (托马斯修士) + KB	8.5/10	92%
5	meta-llama_Llama-4-Scout-17B-16E-Instruct + Su Shi (苏轼) + KB	8.5/10	92%
-	Human Expert Benchmark (avg)	9.2/10	100%

space, based on their profile scores (centroids detailed in Appendix Table 3), also shifts with interventions, indicating changes in their overall analytical posture.

5 Discussion

546

547

550

551

552

553

554

556

559

560

564

565

566

567

571

572

573

575

577

579

585

This study demonstrates that while baseline VLMs exhibit a notable semantic and capability gap compared to human experts in Chinese art critique, targeted interventions using personas and knowledge bases can significantly improve alignment. The VULCA framework provides a robust methodology for quantifying these changes. Our findings highlight VLMs' potential in specialized domains but also underscore the need for culturally aware prompting and knowledge integration for nuanced understanding. The observed 20-30% capability enhancement in some models via our interventions is a promising step.

contributions include the **VULCA** Kev multi-dimensional framework itself as а evaluation tool and the empirical demonstration of intervention effectiveness. This offers pathways for developing more culturally attuned The critic personas, and expert-like VLMs. derived from human expert data, provide a practical mechanism for guiding VLMs towards desired analytical styles.

Limitations include the specific set of VLMs and artworks; future work could broaden this scope. The definition of "expert critique" is also culturally situated and can be further explored. Investigating more sophisticated knowledge integration techniques and dynamic persona adaptation are promising future directions. Further research could also explore cross-cultural VLM critique capabilities.

6 Conclusion

This research introduced VULCA, a quantitative framework for evaluating VLM-generated

critiques of traditional Chinese painting, and demonstrated its utility in assessing baseline VLM capabilities and the impact of persona and knowledge-based interventions. We found that such interventions significantly enhance VLM performance, moving their outputs closer to human expert standards in terms of semantic alignment and critical depth. The study the importance underscores of culturally grounded approaches for developing VLMs capable of nuanced engagement with specialized domains like art criticism. Future work will continue to refine these methods and explore their applicability across diverse cultural contexts and aiming to foster more artistic traditions, sophisticated AI-assisted cultural analysis.

587

588

589

590

591

593

594

596

597

599

600

601

602

603

604

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

Acknowledgments

A Dataset Details

A.1 Lang Shining's "Twelve Months" Dataset

Our study centers on Giuseppe Castiglione's "Twelve Months" series (十二月令图), 12 paintings showing seasonal activities in the Qing imperial court. These paintings fuse Chinese and Western artistic traditions, ideal for cross-cultural interpretation study. We compiled digital images (6 million pixels) from the National Palace Museum (Taiwan) digital archives under CC BY 4.0 license. The dataset includes historical texts and scholarly analyses in both Chinese and English, from Qing Dynasty sources and modern scholarship.

B Experimental Setup Details

B.1 Automated Data Processing and Analysis Workflow

The quantitative analysis in this research is supported by a series of automated scripts. The workflow is divided into three main phases:

- 625

633

635

637

641

647

650

654

664

666

667

670

B.1.1 Phase 1: Feature Extraction from **Human Expert Texts**

This phase is handled by an automated script.

• Purpose: То automatically extract predefined textual features from a collection of human expert commentaries on Chinese art. These features include evaluative stances, core analytical focal points (such as "Use of Color", "Artistic Conception"), and argumentative quality aspects (such as "Profound Insight", "Clear Logic").

• Input:

- A base directory containing .txt files of human expert critiques.
- Predefined lists of English candidate labels for stance, features, and quality (as listed in Appendix D.5), along with their Chinese translations used internally for reporting if needed.

• Key Processing Steps:

- 1. Recursively scans the input directory for all .txt files.
- 2. Loads a zero-shot classification model.
- 3. For each critique text:
 - Predicts the primary evaluative stance (single-label classification from the stance label set).
 - Predicts multiple core focal points (multi-label classification from the feature label set).
 - Predicts multiple argumentative features quality (multi-label classification from the quality label set).
- 4. Stores the determined English label and its confidence score for the primary stance.
- 5. Stores the English labels and confidence scores for all identified features and quality aspects. These are typically stored as dictionary-like structures mapping the label to its score.

6. Outputs:

- A consolidated master CSV file where each row represents a critique: a unique file id, a text preview, the predicted stance and its score, all predicted focal points and their scores, and all predicted quality features and their scores.

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

- Individual CSV files for each (derived scholar/work from sub-directory names), containing the same information for critiques within that specific scope.

Phase 2: Exploratory Data Analysis B.1.2 and Feature Visualization of Human and MLLM Data

This phase is handled by an automated script.

- Purpose: To perform exploratory data analysis (EDA) on the extracted features from both human experts (output from Phase 1) and a baseline MLLM, and to visualize the combined feature space using dimensionality reduction.
- Input:
 - The consolidated human expert features CSV from Phase 1.
 - A consolidated MLLM features CSV, assumed to have a compatible structure, particularly for 'features' and 'quality' columns.

• Key Processing Steps:

- 1. Loads and combines the human and MLLM feature data, adding а source type column to differentiate origins.
- 2. Parses stringified 'features' and 'quality' columns (if stored as strings in CSV) back into dictionary objects.
- 3. Performs EDA, including:
 - Calculating distributions for predicted stances and source types.
 - For each identified feature and quality item: calculating overall mention counts, mention frequency within comments that have such data, and descriptive statistics (mean, median, std dev, min, max) of their scores across all relevant texts.

4. Constructs unified feature vectors for 716 each commentary by concatenating all 717 individual feature and quality scores. 718 719 720 721 722 feature vectors to 2D for visualization. 724 • Output: 725 726 727 730 731 with file id and source type. 732 734 source type. 735 736 737 739 740 741 742 743 745 746 747 data. 748 • Input: 749 750 751 752 753 Micro 755 757

758

761

(such Profiles" as "Comprehensive Analyst," "Historically Focused") and "General

- on Human and MLLM Data This phase is handled by an automated script. • Purpose: To apply a more sophisticated analytical layer by scoring texts against predefined expert profiles, performing dimensionality reduction on these profile scores, calculating profile proportions, and preparing a rich dataset for composite

Neighbor Embedding) and UMAP (Uniform Manifold Approximation and

5. Applies t-SNE (t-distributed Stochastic

Projection, if the library is available) to reduce the dimensionality of these

- An EDA summary CSV file detailing the
 - statistical findings from the EDA. coordinates (x, y) from t-SNE and
- A CSV file containing the 2D UMAP for each commentary, along
- PNG image files for the t-SNE plot and UMAP plot of the combined feature space, typically colored by stance or

B.1.3 Phase 3: Profile Scoring, Candidate Selection, and Advanced Visualization

- visualizations. This phase also focuses on comparing human expert data with MLLM
 - The consolidated human expert features CSV from Phase 1.
 - The consolidated MLLM features CSV.
- Predefined criteria for "Specialized Descriptive Profiles." These profiles are defined by rules that consider stance labels, specific feature scores (e.g., "Historical Context" score > 0.5), and quality scores.

Key Processing Steps:

1. Loads and combines human and 763 MLLM feature data, performing 764 necessary preprocessing like parsing 765 feature/quality dictionaries. 766 2. For each text, calculates a match score 767 (typically 0-1) against each predefined 768 expert profile. This involves: 769 - Checking if a primary stance 770 requirement is met. 771 - Evaluating if a minimum number 772 of flexible rules (based on 773 feature/quality score thresholds) 774 are satisfied. 775 - Combining these into an overall 776 profile match score, potentially 777 with weighting for stance and 778 feature contributions. Specialized 779 logic is used for "Comprehensive 780 Analyst" and "General Descriptive 781 Profile." 782 3. These profile scores are added as new 783 columns to the dataset. 784 4. Constructs new feature vectors based 785 on these profile scores (and potentially 786 other score * columns). 787 5. Applies t-SNE and UMAP for 2D 788 visualization of these 789 profile-score-based vectors. 790 6. Calculates the proportional contribution 791 of each specialized profile score to the 792 sum of all specialized profile scores for a 793 given text. Determines a primary profile 794

762

796

797

798

799

800

801

802

803

804

805

806

• Output:

- A primary CSV file designed for external visualization tools. This includes file id, text previews, original profile status, source type, t-SNE/UMAP coordinates derived from profile scores, original filenames, all profile score columns, and the calculated profile proportion columns.

based on the highest proportion.

 A secondary, more comprehensive CSV file containing the entire processed

849

807 dataframe with all original and derived columns, including the profile-based dimensionality reduction coordinates.
810 C Persona Definitions
811 The following eight persona cards were utilized in this study, each detailed in a separate subsection:

813 C.1 Mama Zola (佐拉妈妈)

814

815

816

817

818

819

821

822

823

824

827

828

830

832

834

838

841

842

847

- **Basic Information:** Elderly West African oral historian and textile artist (female, born 1955, Senegalese village). Guardian of tribal wisdom.
- Key Influences/Background: Grew up in a culture without written records, learning history and wisdom through oral traditions, songs, dances, and rituals. Textile skills passed down through generations; her works are themselves carriers of narrative and history. Critical of Western museums' plunder and misinterpretation of African art.
- Analytical Style and Characteristics: Interprets art from the perspective of community function, ritual significance, and ancestral connection. Emphasizes the practicality, locality, and collective creativity of art. Values the symbolic meaning of materials and the spiritual infusion during the crafting process. Believes art is part of life, not an isolated "artwork."
- Numeric Attributes (Scale: 1-10):
 - Community Culture Perspective: 10
 - Oral Tradition Connection: 9
 - Decolonization Awareness: 8
 - Sensitivity to Craft and Materials: 9
 - Spirituality and Rituality: 7
 - Acceptance of Western Art Theory: 2
 - Language and Expression Style: Language is simple, vivid, full of storytelling and life wisdom. Often uses proverbs and metaphors. Critiques as if telling an ancient story, emphasizing emotional connection and collective memory. Tone is gentle but firm.
 - Sample Phrases:

- "Every pattern on this cloth tells the story of our ancestors, more truly than any book."
- "What you call 'artworks,' we use to celebrate harvests and connect the living with the dead. It is alive, breathing with us."
- "Those masks in museums, separated from their dances and songs, are like fish out of water, soulless."
- "To dye this indigo thread requires the moon's blessing and the earth's gift; this color holds the memory of our people."
- "True beauty is what makes the whole village feel warmth and strength, not something hung on a wall for individual admiration."

C.2 Okakura Kakuzō (冈仓天心)

- Basic Information: Prominent Japanese Meiji era art activist, thinker, and educator (male, 1863-1913, Yokohama). A founder of the Tokyo School of Fine Arts (now Tokyo University of the Arts) and Head of the Chinese and Japanese Art Department at the Museum of Fine Arts, Boston.
- Key Influences/Background: Dedicated to reviving and promoting Japanese and Eastern traditional arts, resisting the blind Westernization of the early Meiji Restoration. Deeply influenced by Eastern philosophy (especially Zen and Daoism). Authored English works such as "The Ideals of the East" and "The Book of Tea," introducing Eastern culture and aesthetics to the West.
- Analytical Style and Characteristics: Emphasized the cultural concept of "Asia is one." Valued the spirituality and symbolic meaning of art, believing the core of Eastern art lies in the "rhythm of life." Advocated for an aesthetic of simplicity, subtlety, and harmony with nature. Possessed a deep understanding of Western art and conducted comparative studies.

• Numeric Attributes (Scale: 1-10):

- Emphasis on Eastern Spirituality: 10
- Cross-Cultural Comparative 894
 Perspective: 9
 895

- Focus on Materials and Craft: 6
• Language and Expression Style: Language is poetic and philosophical, reflecting both Eastern and Western cultural literacy. Elegant prose, adept at interpreting art from a macro-cultural perspective. When introducing to Western readers, often used vivid metaphors and insightful discussions.
• Sample r nrases:
 "Asia is one. The Himalayas divide, only to accentuate, two mighty civilisations, the Chinese with its communism of Confucius, and the Indian with its individualism of the Vedas."
 "Teaism is a cult founded on the adoration of the beautiful among the sordid facts of everyday existence."
 "The Art of life lies in a constant readjustment to our surroundings."
- "In the trembling grey of a breaking dawn, when the birds were whispering in mysterious cadence among the trees, have you not felt that they were talking to their mates about the untold mystery of waking life?"
 "True beauty could be discovered only by one who mentally completed the incomplete."
C.3 Professor Elena Petrova (埃琳娜·佩特洛 娃教授)
• Basic Information: Rigorous Russian Formalist art critic (female, born 1965, St. Petersburg). Professor in the Department of Comparative Literature and Art Theory at a university.
• Key Influences/Background: Deeply influenced by Russian Formalist literary theory (e.g., Shklovsky, Eikhenbaum). Believes the essence of art lies in its formal techniques and "defamiliarization" effect,

- Awareness of Traditional Revival: 8

- Understanding of Western Art: 7

- Interpretation of Symbolic Meaning: 7

rather than social content or the artist's biography.

• Analytical Style and Characteristics: Focuses on the "literariness" of artworks (or "artisticness" itself for visual arts). Analyzes the structure, devices (privom), and media-specific properties of works, and how these elements interact to produce aesthetic effects. Rejects viewing art as a simple reflection of social, historical. or psychological phenomena.

• Numeric Attributes (Scale: 1-10):

- Depth of Formal Analysis: 10
- Focus on Defamiliarization Effect: 9 953
- Sensitivity to Media Properties: 8
- Rejection of Historical/Social Context: 955
 7 956
- Disregard for Authorial Intent: 8
- Restraint in Emotional Interpretation: 6
- Language and Expression Style: Precise, objective language, like scientific analysis. Extensive use of Formalist terminology. Arguments are logically rigorous, with layered dissection. Tone is calm and devoid of personal emotion.

• Sample Phrases:

- "The device is the content of art. We are concerned not with *what* the artist says, but *how* it is said, i.e., its 'device' (priyom)."
- "This painting, through its distortion of conventional perspective, successfully creates a 'defamiliarization' (ostranenie) effect, compelling the viewer to re-examine familiar objects."
- "We must treat the work as a self-sufficient system of signs, analyzing the tensions and harmonies among its internal elements, rather than resorting to external biographical or psychological factors."
- "So-called 'themes' or 'ideas' are merely motivations for stringing
 981

- 983 984
- 985
- 987
- 98
- 99
- 9
- 9
- 995
- 0
- 997 998
- 999
- 1000 1001
- 1002 1003
- 1003
- 1004
- 1005
- 1007
- 1008
- 1009
- 1011
- 1012
- 101
- 1014

1015

1016 – Adherence to Traditional Techniques: 8

10

- Focus on Image Archetypes: 7
 - Evaluation of Secular Art: 3
 - Receptiveness to Innovation: 2
- Language and Expression Style: Language
 is devout, tranquil, and full of religious
 metaphors. Often quotes Scripture and
 Patristic texts. Commentary focuses on
 revealing the divine reality and spiritual
 guidance behind images. Tone is peaceful,
 humble, with mystical overtones.

• Sample Phrases:

together various artistic devices; they

are not the core of artistic analysis

- "The artistic merit of this piece lies in

its clever orchestration of fundamental

'devices' (ustanovka) such as color, line,

and composition, not in the narrative

• Basic Information: Contemplative hermit monk and iconographer (male, born 1970, a

monastery on Mount Athos). Dedicated to

preserving ancient Byzantine icon painting

spiritual and artistic training within the

Eastern Orthodox monastic tradition. Deeply

the

Neoplatonism, and icon theology (e.g., St.

John of Damascus). Believes art is a window

• Analytical Style and Characteristics:

Interprets art from theological and spiritual

meaning of artworks, archetypes, and their

function in liturgy and prayer. Emphasizes

fasting, prayer, and spiritual concentration

during the creative process. Believes true

- Theological Symbolism Interpretation:

- Emphasis on Spiritual Function: 9

Desert

Focuses on the symbolic

Received

Fathers,

itself."

scene it depicts."

C.4 Brother Thomas (托马斯修士)

techniques and theology.

influenced

to the divine.

perspectives.

• Key Influences/Background:

by

beauty points to divine beauty.

• Numeric Attributes (Scale: 1-10):

"This icon is not merely a 'depiction'; it
is itself a 'revelation' of the divine
presence, a window to the unseen
world."

1027

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1043

1044

1045

1046

1047

1048

1049

1050

1051

1054

1055

1056

1057

1058

1059

1060

1062

1064

1065

1066

1067

1068

1070

- "One should view an icon with a prayerful heart. The direction of lines, the use of color, all follow ancient patristic norms, guiding the soul upwards."
- "When creating, the iconographer must fast and pray, becoming a pure conduit for the divine light to flow through the brush."
- "The gold background symbolizes eternal light; the figures' 'inverse perspective' is not 'unrealistic' but transcends worldly vision to present the heavenly order."
- "Every detail, from the folds of a robe to the gesture of a finger, carries profound theological meaning, a silent sermon."

C.5 John Ruskin (约翰·罗斯金)

- **Basic Information:** Leading English art critic of the Victorian era, social reformer, writer, and poet (male, 1819-1900, London). Slade Professor of Fine Art at the University of Oxford.
- Key Influences/Background: Influenced by Romantic views of nature and Christian ethical thought. Championed the Pre-Raphaelite Brotherhood, emphasizing the moral and didactic function of art and fidelity to nature. Had a deep understanding of Gothic architecture.
- Analytical Style and Characteristics: Emphasized "truth to nature." Believed that beauty was intrinsically linked with truth and goodness. Focused on the detailed depiction in artworks, craftsmanship, and the social and moral meanings they conveyed. Held a critical stance towards the social problems and artistic alienation brought by industrialization.

• Numeric Attributes (Scale: 1-10):

- Emphasis on Fidelity to Nature: 10

1073	- Moral/Didactic Function: 9	• A
1074	- Acuity of Detail Observation: 8	۷۵ "۲
1075	- Evaluation of Craftsmanship: 7	ez de
1076	- Social Critical Awareness: 8	cł cr
1077	 Acceptance of Formalism: 3 	na
1078	• Language and Expression Style: Eloquent	51
1079	and powerful language, full of passion and	• N
1080	moral appeal. Ornate writing style, rich in	
1081	literary description. Often used complex long	
1082	sentences and abundant rhetoric. Sharp in	
1083	criticism, fervent in praise.	
1084	Sample Phrases:	
1085	- "Go to Nature in all singleness of heart,	
1086	and walk with her laboriously and	
1087	trustingly, having no other thought but	
1088	how best to penetrate her meaning, and	
1089	remember her instruction."	
1090	- "All great art is praise. And the greatest	• L
1091	art is that which praises the highest	pı
1092	things."	th
1000	"The munder and most thoughtful minde	cr
1093	- The purest and most thoughtful minds	re
1094	are mose which love colour the most.	m hı
1095	- "Fine art is that in which the hand, the	
1096	head, and the heart of man go together."	• Sa
1097	- "To see clearly is poetry, prophecy, and	
1098	religion, —all in one."	
1099	C.6 Su Shi (苏轼)	
1100	• Basic Information: Chinese Northern Song	
1101	Dynasty writer, calligrapher, painter, and art	
1102	theorist (male, 1037-1101, Meishan,	
1103	Meizhou). Courtesy name Zizhan,	
1104	pseudonym Dongpo Jushi. A key founder of	
1105	literati painting theory.	
1106	 Key Influences/Background: Deeply 	
1107	influenced by Confucianism, Daoism, and	
1108	Chan (Zen) Buddhism. Advocated for	
1109	"scholar-official painting" $(\pm \land \blacksquare)$,	
1110	emphasizing the integration of poetry,	
1111	calligraphy, and painting, and the expression	
1112	of inner spirit. His artistic ideas had a	
1113	protound impact on the development of later	
1114	meran panning.	

An electrical State and Channet entriet	
Analytical Style and Characteristics:	1115
Values the "spiritual resonance" (神韵) and	1116
"artistic interest" (意趣) of artworks over	1117
external formal likeness. Emphasizes the	1118
decisive role of the artist's personal	1119
character, knowledge, and cultivation in	1120
creation. Esteems an aesthetic realm of	1121
natural innocence, plainness, and distanced	1122
simplicity.	1123
Numeric Attributes (Scale: 1-10):	1124
– Literary Integration: 10	1125

- Emphasis on Brushwork Interest: 9 1126
- Subjective Spiritual Expression: 9
- Requirement for Formal Accuracy: 3 1128
- Importance of Historical Tradition: 8 1129

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

- Theoretical Innovation: 7
- Language and Expression Style: Elegant prose, rich in philosophical and poetic thought. Often uses poetry as analogy; critiques are profound yet accessible, with refined and insightful language. Tone is moderate, balanced, and imbued with humanistic concern.
- Sample Phrases:
 - "The way to view a painting is to first observe its spiritual resonance, not to seek formal likeness; formal likeness is the business of artisans."
 - "To judge painting by formal likeness is to see with the eyes of a child. To insist a poem must be *this* poem, means one certainly doesn't know poets."
 - "Savoring Mojie's (Wang Wei) poetry, there is painting within the poetry; viewing Mojie's painting, there is poetry within the painting."
 - "One must have the bamboo fully formed in one's chest before applying it to the brush and paper; this is beyond those who do not have the bamboo formed in their chests."

1156- "This painting deeply captures the
meaning of creation; the brushwork is
simple yet the meaning is complete.1159This is what is meant by 'the height of
brilliance returns to plainness.""

C.7 Guo Xi (**郭熙**)

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1189

1191

- Basic Information: Outstanding Chinese Northern Song Dynasty landscape painter and painting theorist (male, c. 1023-c. 1085, Wen County, Heyang). Served as an Erudite (艺 学) in the imperial painting academy during Emperor Shenzong's reign.
- Key Influences/Background: Inherited and developed the traditions of the Northern school of landscape painting, emphasizing observation and experience of nature. His theoretical work "The Lofty Message of Forests and Streams" (林泉高致) is a seminal text in Chinese landscape painting theory.
- Analytical Style and Characteristics: Emphasized that landscape paintings should be "walkable, viewable, wanderable, and habitable" (可行、可望、可游、可居). Proposed methods for observing and depicting landscapes such as the "Three Distances" (三远: high distance, deep distance, level distance). Valued the influence of seasons and climate on scenery, striving for majestic and varied artistic conceptions (意境).

• Numeric Attributes (Scale: 1-10):

- Depth of Nature Observation: 9
 - Spatial Representation Skill: 10
- Creation of Landscape Atmosphere: 9
 - Theoretical System Construction: 8
- 1192 Diversity of Brushwork Techniques: 7
- Connection to Humanistic Spirit: 6
- Language and Expression Style: Language is simple, concrete, and rich with summaries of practical experience. Adept at using vivid metaphors to describe landscape forms and the artist's insights. Discourse is systematic and clear, possessing both theoretical depth and practical guidance.

• Sample Phrases:

 "Landscapes can be those one can walk through, those one can gaze upon, those one can wander in, and those one can dwell in. When a painting achieves this, it is a masterpiece." 1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

- "Mountains have three distances: looking up at the peak from the foot of a mountain is called high distance; peering into the back from the front of a mountain is called deep distance; looking from a near mountain towards a distant mountain is called level distance."
- "In real landscapes of rivers and valleys, observe them from afar to capture their 势 (shi overall configuration/momentum), and observe them up close to capture their 质 (zhi substance/texture)."
- "Spring mountains are delicately charming as if smiling; summer mountains are lush green as if dripping; autumn mountains are clear and bright as if adorned; winter mountains are bleak and somber as if sleeping."
- "Mountains take water as their blood vessels, vegetation as their hair, and mist and clouds as their spirit and radiance."

C.8 Dr. Aris Thorne (阿里斯·索恩博士)

- **Basic Information:** Futurist digital art historian and ethicist (non-binary, born 2042, Neo-Kyoto). Specializes in AI-generated art, bio-art, and the philosophical implications of post-human creativity.
- Key Influences/Background: Raised in a highly technological society but trained in classical art history. Deeply influenced by cybernetics, post-humanism, and existentialist philosophy. Dedicated to building bridges between rapidly developing techno-art and core human values.
- Analytical Style and Characteristics: 1243
 Examines emerging techno-art forms with a critical eye. Focuses on ethical issues such as algorithmic bias, authorship, and the authenticity and originality of art. When 1247

1248	analyzing works, explores both their	D.1 Capability Assessment Framework		
1249	technological innovation and their reflection	Our three-dimensional capability assessment		
1250	on and questioning of the human condition.	framework is designed to evaluate MLLM		
1251	• Numeric Attributes (Scale: 1-10):	performance in Chinese art commentary through both vector space analysis and specific capability		
1252	- Focus on Tech Ethics: 10	metrics:		
1253	 Insight into Future Trends: 9 	• Painting Element Recognition (5-point		
1254	- Critical Thinking: 8	visual elements, compositional features, and		
1255	– Interdisciplinary Integration: 9	technical aspects.		
1256	– Traditional Art Literacy: 6	 Level 1: Minimal recognition of basic elements significant errors or omissions 		
1257	- Emotional Resonance: 5	- Level 2: Basic recognition of major		
1258	• Language and Expression Style: Precise,	elements, but with notable inaccuracies		
1259	Often uses emerging scientific and	Lavel 3: Accurate identification of		
1260	technological terms and philosophical	major compositional elements and		
1201	concents Arguments are rigorous tending	techniques		
1202	to pose open-ended questions rather than	teeninques		
1264	providing definitive answers	- Level 4: Detailed recognition of both		
120-	Course Discourse	major and minor elements with few		
1265	• Sample Phrases:	errors		
1266	- "When algorithms become paintbrushes,	Level 5. Comprehensive and nuanced		
1267	how do we define the creator? When	recognition of subtle visual elements and		
1268	code generates beauty, where does the	technical features		
1269	boundary of originality lie?"	technical reatures		
1070	"This AI concreted image is its 'style'	• Chinese Painting Understanding (7-point		
1270	- This Al-generated image, is its style	scale): Evaluates depth of understanding		
1271	data or an emerging 'machine	cultural meanings, historical contexts, and		
1272	intuition'?"	symbolic references specific to Chinese painting traditions.		
1274	- "Bio-art challenges the traditional			
1275	dichotomy of life and non-life, forcing	– Level 1: Minimal recognition of		
1276	us to rethink what is 'natural' and what	obvious symbols, significant cultural		
1277	is 'artificial.'"	misinterpretations		
1278	- "Under the post-human gaze, does this	- Level 2: Basic recognition of common		
1279	work enhance our humanity, or does it	symbols but limited understanding of		
1280	herald its dissolution?"	their significance		
1281	- "In evaluating such works, we must not	– Level 3: Moderate understanding of		
1282	only ask 'what is it.' but more	major symbols with some contextual		
1283	importantly, 'what does it make us	awareness		
1284	think,' and 'where will it lead us?'"			
		– Level 4: Accurate interpretation of		
1285	D Evaluation Framework and Prompts	major cultural symbols with appropriate historical context		
1286	This section details the evaluation framework,			
1287	including the multi-dimensional capability	- Level 5: Detailed understanding of both		
1288	assessment rubric and the standardized prompts	common and specialized symbolic		
1289	used for eliciting commentaries from MLLMs.	elements		

1333	- Level 6: Sophisticated analysis of
1334	symbolic relationships with strong
1335	historical contextualization
1336	– Level 7: Expert-level analysis of
1337	symbolic networks with nuanced
1338	cultural and historical insights
1339	• Chinese Language Usage (5-point scale):
1340	including terminology accuracy stylictic
1341	appropriateness and fluency in Chinese art
1343	discourse.
1344	- Level 1: Significant terminology errors,
1345	inappropriate style for art commentary
1346	- Level 2: Basic fluency but frequent
1347	terminology errors and stylistic
1348	inconsistencies
1349	 Level 3: Generally appropriate language
1350	with occasional specialized terminology
1351	errors
1352	- Level 4: Accurate terminology usage
1353	with appropriate stylistic features for art
1354	commentary
1355	- Level 5: Expert-level language usage
1356	with precise terminology and
1357	stylistically sophisticated expression
1358	D.2 Structured Commentary Evaluation
1309	
1360	Our evaluation of structured commentaries
1361	tollows a detailed rubric designed specifically for
1362	the two-part format (paragraph-form analysis and
1363	structured assessment). This rubric maps specific
1304	three core canability dimensions:
1366	Mapping to Core Capabilities:
1367	- Painting Element Recognition is
1368	evaluated primarily through:
1369	* Accuracy in identifying visual
1370	elements from predefined lists in
1371	the structured template
1372	* Correct classification of
1373	compositional techniques from
1374	multiple-choice options
1375	* Precision in describing spatial
1376	relationships using standardized
4077	terminology

* Recognition of brushwork	1378
techniques from a predefined	1379
taxonomy	1380
Chinasa Dainting Understanding is	1001
- Chinese Fainting Understanding is	1381
evaluated primarily through:	1382
* Correct matching of symbols with	1383
their cultural meanings from	1384
provided options	1385
* Appropriate selection of historical	1386
context categories from a	1387
predefined list	1388
* Accurate identification of	1389
philosophical concepts relevant to	1390
the painting	1391
* Proper classification of the work	1392
within Chinese painting traditions	1393
- Chinese Language Usage is evaluated	1394
primarily through:	1395
* Correct use of specialized Chinese	1396
art terminology from a provided	1397
glossary	1398
* Appropriate stylistic features for	1399
Chinese art commentary	1400
* Proper application of Chinese	1401
aesthetic concepts in context	1402
* Fluency and naturalness in Chinese	1403
language expression	1404
Structured Template Scoring:	1405
_ Primary Visual Flomonts (Painting	1/06
- IT many Visual Elements (Laming Flement Recognition):	1400
Element Recognition).	1407
* 0 points: Fails to identify any	1408
correct elements from the	1409
predefined list	1410
* 1 point: Identifies 1-2 basic	1411
elements correctly	1412
* 2 points: Identifies 3-4 elements	1413
correctly with minor errors	1414
* 3 points: Identifies 5+ elements	1415
correctly with proper	1416
categorization	1417
* 4 points: Identifies all major and	1418
several minor elements with	1419
precise descriptions	1420
* 5 points: Comprehensive	1421
identification with nuanced	1422
understanding of relationships	1423

1424	– Symbolic Content (Chinese Painting
1425	Understanding):
1426	* 0 points: Fails to match any
1427	symbols with their cultural
1428	meanings
1429	* 1-2 points: Matches basic symbols
1430	with simplified meanings
1431	* 3-4 points: Matches multiple
1432	symbols with appropriate
1433	meanings and basic context
1434	* 5-6 points: Matches complex
1435	symbols with detailed cultural
1436	explanations
1437	* 7 points: Sophisticated matching
1438	with interconnected symbolic
1439	networks and philosophical depth
1440	– Key Terminology (Chinese Language
1441	Usage):
1442	* 0 points: Uses incorrect or
1443	inappropriate terminology
1444	throughout
1445	* 1 point: Uses basic terminology
1446	with frequent errors
1447	* 2-3 points: Uses standard
1448	terminology with occasional errors
1449	<pre>* 4 points: Uses specialized</pre>
1450	terminology accurately and
1451	appropriately
1452	* 5 points: Demonstrates mastery of
1453	specialized terminology with
1454	nuanced application
4455	The structured template includes marific
1455	sections with predefined options, multiple choice
1450	selections and classification tasks that allow for
1458	objective scoring. For example:
	$\mathbf{T} = \mathbf{T} = $
1459	• The "Primary visual Elements" section
1460	20+ elements
1401	
1462	• The "lechnical Approach" section uses
1463	multiple-choice classification of techniques
1464	• The "Symbolic Content" section requires
1465	matching symbols to meanings from
1466	provided options
1467	• The "Historical Context" section uses
1468	categorical classification from predefined
1469	traditions

•	The	"Key	Terminology"	section	requires	1470
	selec	tion fro	om a specialized	l glossary	y	1471

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

This structured approach enables direct comparison with annotated ground truth and 1473 provides a standardized framework for evaluating 1474 all three core capabilities across different models 1475 and personas.

Structured Commentary Prompt Design D.3

We developed a standardized structured elicit prompting approach to consistent commentaries across all models. The core prompt given to the MLLMs is detailed below. For persona-enhanced prompts, the respective persona card information (see Section C) was prepended to this core prompt, with an additional instruction to adopt the persona's perspective, knowledge base, and communication style.

Please assume the role of a Hello! professional art critic. Next, you will receive an image of a Chinese painting and any associated textual annotations (if available). Please provide a detailed, insightful, and well-structured critique of this

Your output should consist of two parts:

- 1. The complete commentary text. 1496
- 2. A JSON object summarizing 1497 your core evaluation points. 1498

Part One: Commentary Text

artwork and information.

Please write one or more coherent paragraphs to thoroughly analyze multiple aspects of the artwork. It is recommended that you consider and cover at least the following points (but you are not limited to them):

- Composition and Layout: Evaluate the overall structure of the painting, the organization of elements, the creation of space, visual guidance, etc.
- Brushwork and **Technique:** 1511 Analyze the use of lines (such as 1512 thickness, speed, turns, strength), 1513 the variations in ink tones (dense, 1514 light, wet, dry), texture strokes (皴 1515

法), moss dots (点苔), coloring, and other specific painting techniques and their effects.

1516

1517

1518

1519

1520

1522

1523

1524

1525

1526

1527

1528

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1559

1560

1561

- Use of Color (if applicable): Discuss the paintings color palette, the coordination and contrast between colors, and the emotions or symbolic meanings conveyed by the colors.
- Theme and Content: Interpret the subject matter depicted in the artwork (such as landscapes, figures, flowers and birds, etc.), specific objects, potential storylines or narrative elements, and any underlying symbolic meanings or cultural connotations.
- Artistic Conception and Emotion (意境): Elaborate on the overall atmosphere, aesthetic taste, and artistic style conveyed by the painting, as well as the emotional resonance or philosophical reflections it might evoke in the viewer.
 - Style and Heritage: Analyze the artistic style characteristics of the artwork, its connections to major historical painting schools, traditional techniques, or specific artists, and its potential innovations based on inherited traditions.

Please strive for meticulous analysis, clear viewpoints, and support your statements with specific visual elements from the artwork and any provided textual information.

Part Two: Structured Evaluation in JSON Format

After your commentary text, please start a new line and provide a JSON object strictly adhering to the following structure and key names. Fill in your evaluation results into the corresponding values.

1562Please ensure the JSON format is1563correct, and all string values use double

quotes. Do not add any extra markers1564or explanations before or after the1565JSON object. Your commentary text1566and this JSON object will be your1567complete response to this artwork.1568

1569

1570

1571

1572

1573

1574

1575

1576

1577

1578

1579

1580

1581

1582

1583

1584

1585

1586

1587

1588

1589

1590

1591

1592

1593

1594

1595

1596

1597

1598

1599

1600

1601

1602

1603

1604

1605

1606

1608

1609

1610

D.4 Vector Space Analysis Methods

Our vector space analysis employed several complementary methods:

- Embedding We Model: used the BAAI/bge-large-zh-v1.5 model. а specialized multilingual sentence transformer. This model generates 1024-dimensional vectors that capture between semantic relationships commentaries.
- Similarity Metrics: We primarily used cosine similarity to measure semantic closeness between vectors, supplemented by Earth Mover's Distance (EMD) to capture distribution differences between vector spaces.
- Dimensionality Reduction: For visualization purposes, we employed UMAP (Uniform Manifold Approximation and Projection) and t-SNE (t-distributed Stochastic Neighbor Embedding) to reduce the high-dimensional vectors to two or three dimensions while preserving semantic relationships. The resulting coordinates were also saved for detailed analysis (Table 6).
- Clustering Analysis: We applied hierarchical clustering to identify patterns in the vector spaces, particularly to analyze grouping by persona, painting subject, or capability level.

All vector space analyses were conducted using consistent parameters across comparisons to ensure valid results.

D.5 Zero-Shot Classification Labels for Feature Extraction

The initial feature extraction from textual commentaries (both human expert and MLLM-generated) employed а zero-shot classification model with the following predefined candidate label sets, derived from the extraction scripts.

D.5.1 Evaluative Stance Labels

• Historical Research (历史考证型) 1611

1612	• Aesthetic Appreciation (美学鉴赏型)	general descriptive profiles used in this study.
1613	• Socio-cultural Interpretation (社会文化解读	Scores for features and qualities are generally on a
1614	型)	0-1 scale, derived from the zero-shot classification
1615	• Comparative Analysis (比较分析型)	model.
1616	• Theoretical Construction (理论建构型)	
1617	• Critical Inquiry (质疑与思辨型)	D.6.1 Specialized Profile Criteria
1618	• High Praise (高度赞扬与推崇)	(Micro-Level)
1619	• Objective Description (客观中性描述)	These profiles aim to capture more specific
1620	• Mild Criticism (温和批评与保留)	anlytical tendencies.
1621	• Strong Negation (强烈否定与驳斥)	• 博学通论型 (Comprehensive Analyst):
1622	D.5.2 Core Focal Point Labels	- Description: Characterized by a broad
1623	• Use of Color (色彩运用)	engagement with numerous facets of the
1624	• Brushwork Technique (笔法技巧)	artwork. This profile does not rely on a
1625	• Texture Strokes (Chunfa) (皴法特点)	single dominant stance but requires high
1626	• Line Quality (线条质量)	scores (e.g., ≥ 0.6) across a significant
1627	• Ink Application (墨法变化)	number (e.g., at least 10) of diverse feature
1628	• Layout and Structure (布局与结构)	labels (e.g., "Use of Color", "Brushwork
1629	• Spatial Representation (空间营造)	Technique", "Historical Context",
1630	• Artistic Conception (意境表达)	"Symbolism", etc.).
1631	• Emotional Expression (情感传递)	– Example Rule Logic:
1632	• Subject Matter (主题内容)	min flexible rules to pass:
1633	• Genre (题材选择)	10, where each rule is feature score
1634	• Symbolism (象征意义)	>= 0.6 for a wide range of features listed
1635	• Historical Context (历史背景)	in
1636	• Artist Biography (画家生平)	ALL_POSSIBLE_FEATURE_LABELS.
1637	• Style/School (风格流派)	
1638	• Technique Inheritance & Innovation (技法传	• 历史考据空 (Historically Focused):
1639	承与创新)	- Description: Emphasizes the historical and
1640	• Cross-cultural Influence (跨文化影响)	biographical aspects of the artwork and
1641	D.5.3 Argumentative Quality Labels	artist.
16/0	• Drofound Insight (田韶深刻) • Drofound Insight (田韶深刻)	- Example Rule Logic: Requires at least 2
1642	• Strong Argumentation (论证本公右力)	flexible rules to pass, such as:
1643	• Clear Logic (逻辑洁晰严密)	★ Feature "Historical Context": score >
16/5	• Detailed Analysis (细节分析目体)	$\%$ Teature Thistoriear context : score \geq
16/6	• Classical Citations (引田经曲佐证)	* Feature "Artist Biography": score
1647	• Objective Viewpoint (观占安观公分)	* reature Artist Diography : score \geq 0.40
1648	• Superficial Treatment (论状流于表面)	0.40
1649	• Overly General Content (内容较为宽浮)	* Teature Style/School . score ≥ 0.40
1650	• Lacks Examples (缺乏具体例证)	* Quality Classical Citations : score \geq
1651	• Logical Gaps (逻辑存在跳跃)	0.25
1652	• Subjective/Biased View (观点主观片面)	• 技艺风格型 (Technique & Style Focused):
1652	D.6 Export Profile Definitions for	- Description: Focuses on the aesthetic
1654	Commentary Analysis	appreciation of technical skills, artistic
1034	Commentary Analysis	style, and expressive qualities.
1655	To further categorize and understand the nuanced	- Example Rule Logic: Main stance is
1656	styles of art commentaries, a rule-based profiling	"Aesthetic Appreciation" (score > 0.40)
1657	system was developed. This system assigns texts	AND at least 2 flexible rules nass such as:
1658	to predefined profiles based on their stance, focal	
1659	points (features), and argumentative quality scores.	* Feature "Technique Inheritance &
1660	Below are the definitions for key specialized and	Innovation": score ≥ 0.30

1707 1708	* Feature "Artistic Conception": score ≥ 0.20
1709	•理论比较型 (Theory & Comparison Focused):
1710	- Description: Characterized by comparative
1711	analysis, theoretical framing, and critique.
1712	often examining structural and symbolic
1713	elements.
1714	- Example Rule Logic: Requires at least 3
1715	flexible rules to pass, such as:
1716	* Feature "Stylistic Analysis": score \geq
1717	0.30
1718	* Feature "Cross-cultural Comparison":
1719	$score \ge 0.40$
1720	* Feature "Theoretical Construction":
1721	score ≥ 0.30
1722	* Feature "Layout and Structure": score
1723	≥ 0.50
1724	* Feature "Symbolism": score ≥ 0.50
1725	D.6.2 General Descriptive Profile Criteria
1726	This profile captures texts that provide broader
1727	descriptions without a highly specialized focus.
1728	• 泛化描述型 (General Descriptive Profile):
1729	- Description: Applies when a commentary
1730	discusses several common aspects of an
1700	1
1731	artwork with moderate scores and holds a
1731 1732	artwork with moderate scores and holds a generally common stance (e.g., Objective
1731 1732 1733	artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation)
1731 1732 1733 1734	artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria
1731 1732 1733 1734 1735	artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles.
1731 1732 1733 1734 1735 1736	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is
1731 1732 1733 1734 1735 1736 1737	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description",
1731 1732 1733 1734 1735 1736 1737 1738	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic
1731 1732 1733 1734 1735 1736 1737 1738 1739	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with acces > 0.15 AND st bast 2 features from the standard standa
1731 1732 1733 1734 1735 1736 1737 1738 1739 1740	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15, AND at least 3 features from a predefined profile.
1731 1732 1733 1734 1735 1736 1736 1737 1738 1739 1740 1741	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15, AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Lea of Color")
1731 1732 1733 1734 1735 1736 1737 1738 1739 1740 1741 1742	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15, AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an every stance start
1731 1732 1733 1734 1735 1736 1737 1738 1739 1740 1741 1742 1743	artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. – <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15 , AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20
1731 1732 1733 1734 1735 1736 1736 1737 1738 1739 1740 1741 1742 1743 1744	artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. – <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15 , AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20 .
1731 1732 1733 1734 1735 1736 1737 1738 1739 1740 1741 1742 1743 1744	artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. – <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15 , AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20 . E Detailed Results
1731 1732 1733 1734 1735 1736 1736 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15, AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20. E Detailed Results E.1 Detailed Persona Capability Scores
1731 1732 1733 1734 1735 1736 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746 1747	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15, AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20. E Detailed Results E.1 Detailed Persona Capability Scores Table 2 shows distinct capability score patterns
1731 1732 1733 1734 1735 1736 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746 1747 1748	artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. – <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15 , AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20 . E Detailed Results E.1 Detailed Persona Capability Scores Table 2 shows distinct capability score patterns across personas:
1731 1732 1733 1734 1735 1736 1736 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746 1747 1748 1749	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15, AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20. E Detailed Results E.1 Detailed Persona Capability Scores Table 2 shows distinct capability score patterns across personas: Personas with Chinese cultural backgrounds
1731 1732 1733 1734 1735 1736 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746 1745 1746 1747 1748 1749 1750	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15, AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20. E Detailed Results E.1 Detailed Persona Capability Scores Table 2 shows distinct capability score patterns across personas: Personas with Chinese cultural backgrounds (e.g., Mama Zola, Okakura Kakuzō)
1731 1732 1733 1734 1735 1736 1736 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746 1747 1748 1749 1750 1751	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15, AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20. E Detailed Results E.1 Detailed Persona Capability Scores Table 2 shows distinct capability score patterns across personas: Personas with Chinese cultural backgrounds (e.g., Mama Zola, Okakura Kakuzō) generally scored higher in Chinese Painting
1731 1732 1733 1734 1735 1736 1736 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746 1747 1748 1749 1750 1751 1752	 artwork with moderate scores and holds a generally common stance (e.g., Objective Description, Socio-cultural Interpretation) but does not meet the more stringent criteria of specialized profiles. <i>Example Rule Logic:</i> Primary stance is one of ("Objective Description", "Socio-cultural Interpretation", "Aesthetic Appreciation", "Historical Research") with score ≥ 0.15, AND at least 3 features from a predefined pool (e.g., "Historical Context", "Symbolism", "Use of Color") are mentioned with an average score ≥ 0.20. E Detailed Results E.1 Detailed Persona Capability Scores Table 2 shows distinct capability score patterns across personas: Personas with Chinese cultural backgrounds (e.g., Mama Zola, Okakura Kakuzō) generally scored higher in Chinese Painting Understanding and Chinese Language

Personas with Western art backgrounds (e.g.,	1754
Professor Elena Petrova, Brother Thomas)	1755
performed well in Painting Element	1756
Recognition but were weaker in Chinese	1757
Painting Understanding and Language	1758
Usage.	1759
The cross-cultural expert persona (John	1760

1762

1763

1764

1765

1766

1767

1768

1769

1770

1771

1772 1773 1774

1775

1785

1786

1787

1788

1789

- Ruskin) demonstrated balanced capabilities, excelling in Chinese Painting Understanding, suggesting knowledge base support can bridge cultural gaps.
- The technology-oriented persona (Dr. Aris Thorne) achieved the highest in Painting Element Recognition but was less proficient in cultural understanding and language.
- The contemporary Chinese persona (Guo Xi) showed strong Painting Element Recognition and good Chinese Painting Understanding.

E.2 Prompt Sensitivity Analysis

Semantic s	similarity	scores	between	responses	to
different fo	ormulation	ns:			

Positive/Negative Formulations:

7
8
9
0
1

- Okakura Kakuzō: 0.86 1782
- Professor Elena Petrova: 0.67 1783
- Shen Mingtang: 0.89 1784
- Data Provenance and Licensing: The Twelve Months Series paintings were accessed through the National Palace Museum (Taiwan) digital archives under CC BY 4.0 license.
- Computational Resources: Our vector 1790 space analysis approach requires significant 1791 computational resources, which may limit 1792 accessibility for some researchers or 1793 institutions. 1794

Table 2: Mean Capability Scores Across Different Personas (5-point scale for Painting Element Recognition and Chinese Language Usage, 7-point scale for Chinese Painting Understanding)

Model	Persona	Painting Elements	Cultural Understanding	Argumentation	Profile Match
google gemini-2.5pro	Brother Thomas (托马斯修士)	-0.2	0.5	0.1	+6
google gemini-2.5pro	Unknown Persona	-0.2	-0.1	0.0	+-1
google gemini-2.5pro	Guo Xi (郭熙)	-0.1	-0.1	0.2	+-7
google gemini-2.5pro	John Ruskin (约翰·罗斯金)	-0.2	0.5	0.2	+1
google_gemini-2.5pro	Mama Zola (佐拉妈妈)	-0.3	-0.0	0.1	+-2
google_gemini-2.5pro	Su Shi (苏轼)	0.4	0.5	0.4	+6
google_gemini-2.5pro	Okakura Kakuzō (冈仓天心)	0.1	0.3	0.1	+6
meta-llama_Llama-4-Scout-17B-16E-Instruct	Brother Thomas (托马斯修士)	-0.1	0.1	-0.2	+6
meta-llama_Llama-4-Scout-17B-16E-Instruct	Unknown Persona	-0.5	-0.4	-0.6	+-6
meta-llama_Llama-4-Scout-17B-16E-Instruct	Guo Xi (郭熙)	-0.3	-0.0	-0.4	+-3
meta-llama_Llama-4-Scout-17B-16E-Instruct	John Ruskin (约翰·罗斯金)	0.1	0.3	0.4	+0
meta-llama_Llama-4-Scout-17B-16E-Instruct	Mama Zola (佐拉妈妈)	-0.1	0.4	0.1	+2
meta-llama_Llama-4-Scout-17B-16E-Instruct	Su Shi (苏轼)	-0.2	0.2	0.2	+-2
meta-llama_Llama-3.1-8B-Instruct	Brother Thomas (托马斯修士)	-0.2	-0.2	-0.0	+0
meta-llama_Llama-3.1-8B-Instruct	Unknown Persona	0.2	0.2	0.0	+2
meta-llama_Llama-3.1-8B-Instruct	Guo Xi (郭熙)	0.0	-0.9	-0.3	+-11
meta-llama_Llama-3.1-8B-Instruct	John Ruskin (约翰·罗斯金)	-0.3	0.1	0.2	+-6
meta-llama_Llama-3.1-8B-Instruct	Mama Zola (佐拉妈妈)	-0.5	-0.4	-0.1	+-15
meta-llama_Llama-3.1-8B-Instruct	Su Shi (苏轼)	0.4	0.7	0.7	+10
Qwen-2.5-VL-7B	Brother Thomas (托马斯修士)	0.6	1.6	1.4	+19
Qwen-2.5-VL-7B	Unknown Persona	0.6	1.3	0.9	+18
Qwen-2.5-VL-7B	Guo Xi (郭熙)	0.5	1.2	1.0	+12
Qwen-2.5-VL-7B	John Ruskin (约翰·罗斯金)	0.7	1.7	1.3	+24
Qwen-2.5-VL-7B	Mama Zola (佐拉妈妈)	0.9	2.4	2.1	+22
Qwen-2.5-VL-7B	Su Shi (苏轼)	0.8	1.5	1.5	+16

Table 3: Mean Centroid Coordinates in Reduced Dimensions (t-SNE/UMAP) for Evaluated MLLM Sources

Source	t-SNE X (Mean)	t-SNE Y (Mean)	UMAP X (Mean)	UMAP Y (Mean)
Qwen-2.5-VL-7B	-2.1547577	-0.667885	2.5803347	1.209615
gemini-2.5pro	-1.7324703	-1.3018972	1.8234636	1.2407658
meta-llama_Llama-3.1-8B-Instruct	-2.4183042	-1.4762617	2.4776638	1.8536302
meta-llama_Llama-4-Scout-17B-16E-Instruct	0.0048952624	-0.812603	0.3323455	-1.037882

• Expert Knowledge Access: The development of effective persona cards requires access to specialized knowledge, which may create barriers to implementing similar approaches in other cultural domains.

E.3 Supplementary Quantitative Data Tables

1795

1796

1798

1799

1800

1801

1802

1803

1804

1805

1806

1808

1810

1811

1812

1813

1814

This section provides supplementary tables detailing the quantitative data underlying some of the figures and analyses presented in the main paper. The mean centroid coordinates for evaluated MLLM sources in the reduced dimensional space are detailed in Table 3. For a detailed breakdown of the key feature scores that underpin the visualizations in Figure 4A, please refer to Table 4. Similarly, the mean profile alignment scores visualized in Figure 4B are presented in detail in Table 5. The specific capability scores used to generate the radar chart in Figure 3B can be found in Table 7.

F Representative Output Samples

1815The examples in Table 8 demonstrate the1816differences in content generated by MLLMs

under basic prompts versus different persona 1817 prompts. Through comparison, we can observe: 1. 1818 **Basic Prompt Outputs**: Without persona 1819 guidance, models tend to generate more 1820 generalized, descriptive content, primarily 1821 focusing on visible elements in the image, and 1822 often exhibit quality issues such as Logical Gaps and Subjective/Biased assessments. 2. **Chinese 1824 Artist Persona Outputs**: Under the guidance of 1825 personas like Mama Zola and Okakura Kakuzō, 1826 the output content demonstrates stronger historical research tendencies and aesthetic 1828 appreciation capabilities, with significantly higher 1829 Classical Citations scores and better performance 1830 on features such as Historical Context. 3 1831 **Language and Style Differences**: Outputs 1832 guided by Chinese personas often begin in Chinese, use more professional terminology, and 1834 reference classical literature more frequently, which is closely related to the relevant knowledge points contained in the persona knowledge base. 1837

Further analysis reveals a significant shift in1838semantic space under different persona guidance,1839validating the substantial impact of persona1840

Table 4: Key Feature Scores for Human Experts and MLLMs. These scores correspond to data visualized in Figure 4A.

Source	Hist. Context	Art. Conception	Symbolism	Brush. Tech.	Layout Struct.	Use of Color	Line Quality	Subject Matter
human_expert	0.676	0.599	0.661	0.199	0.549	0.395	0.496	0.691
gemini-2.5pro	0.4261660233	0.6015897764	0.6935903973	0.6399750158	0.8743446511	0.6952415214	0.7324248211	0.5401486428
meta-llama Llama-3.1-8B-Instruct	0.3659920343	0.5850531087	0.5293492947	0.5909547665	0.7457691074	0.6573745586	0.4430214438	0.4339093090
meta-llama Llama-4-Scout-17B-16E-Instruct	0.7100048551	0.8508161700	0.7583027472	0.9033655355	0.9164849845	0.9357454672	0.8192868597	0.7891201358
Qwen-2.5-VL-7B	0.6504738033	0.8907955483	0.7733450871	0.9369910086	0.8949400724	0.9436663414	0.7946821108	0.6997969688

Table 5: Mean Profile Alignment Scores for Human Experts and MLLMs. These scores correspond to data visualized in Figure 4B.

Source	Comprehensive Analyst	Historically Focused	Technique Style Focused	Theory Comparison Focused	General Descriptive Profile
human expert	0.709	0.623	0.518	0.431	0.665
gemini-2.5pro	0.6066217268	0.4645543554	0.5805458927	0.7892081424	0.6725181508
meta-llama Llama-3.1-8B-Instruct	0.4859600855	0.3351432514	0.4807204770	0.7763639851	0.5595579955
meta-llama Llama-4-Scout-17B-16E-Instruct	0.7796032621	0.6908934862	0.8188009710	0.8516423824	0.8236625996
Qwen-2.5-VL-7B	0.7783469856	0.6530052284	0.8566955672	0.8481851482	0.7842983472

Table 6: Sam	ple Data from	t-SNE and KDE	Analysis	(underlying	Figure 3A).
--------------	---------------	---------------	----------	-------------	-------------

Model Name	Source Type	Intervention	t-SNE X	t-SNE Y	File ID
gemini-2.5pro gemini-2.5pro gemini-2.5pro	model model model	baseline baseline baseline	-8.245 -0.607 -2.392	-7.489 -15.201 -1.717	august_八月(basic).txt august_八月(with_Dong_Qichang).txt august_八月(with_Dr_Evelyn_Reed).t
gemini-2.5pro gemini-2.5pro	model model	baseline baseline	-12.369 -7.852	-5.803 -6.419	august_八月(with_Li_Ruoyun).txt august_八月(with_Marcus_Fabius).tx t
human_expert	human	ground_truth	3.451	-0.876	列文森 (JosephLevenson)中国早期绘画中的政治和个人因素.txt

Table 7: Capability Scores for Radar Chart Dimensions (underlying Figure 3B).

Model Name	Intervention	Profound Insight	Strong Arg.	Detailed Analysis	Clear Logic	Objective Viewpoint	Class. Citations	Logical Gaps	Subjective/ Biased View
HumanAvg	Human Expert	0.396	0.448	0.540	0.093	0.327	0.419	0.465	0.674
Gemini-2.5-Pro	Baseline	0.458	0.486	0.527	0.318	0.461	0.334	0.409	0.483
Gemini-2.5-Pro	Intervened	0.569	0.643	0.689	0.227	0.601	0.492	0.388	0.536
meta-llama Llama-3.1-8B-Instruct	Baseline	0.342	0.371	0.388	0.451	0.305	0.253	0.521	0.399
meta-llama Llama-3.1-8B-Instruct	Intervened	0.495	0.573	0.612	0.274	0.549	0.427	0.417	0.580
meta-llama_Llama-4-Scout-17B-16E-Instruct	Baseline	0.511	0.539	0.583	0.367	0.524	0.399	0.367	0.445
meta-llama_Llama-4-Scout-17B-16E-Instruct	Intervened	0.647	0.701	0.735	0.201	0.676	0.581	0.312	0.502
Qwen-2.5-VL-7B	Baseline	0.311	0.338	0.329	0.515	0.262	0.219	0.599	0.341
Qwen-2.5-VL-7B	Intervened	0.608	0.660	0.695	0.301	0.629	0.518	0.591	0.666

Table 8: Representative MLLM Output Samples with Feature Scores. Fields: Commentary Preview, Stance Label, Feature Scores (excerpt), Quality Assessment (excerpt). Shows differences between basic prompts and personaguided prompts (Mama Zola, Okakura Kakuzō).

File ID	Commentary Preview	Stance Label	Feature Scores (excerpt)	Quality Assessment (excerpt)
April				
四月(basic).txt	This artwork, evidently a section from the	Socio-cultural	Brushwork: 0.99; Layout: 0.98;	Logical Gaps: 0.58; Strong
	"Fourth Month"	Interpretation	Line: 0.93	Argumentation: 0.50
August				
八月(basic).txt	This analysis delves into a magnificent	Comparative Analysis	Line: 0.90; Layout: 0.90;	Detailed Analysis: 0.59; Lacks
	example of Qing Dynasty court painting		Spatial: 0.78	Examples: 0.44
December		-	•	
十二月(basic).txt	This magnificent scroll, a segment from the	Socio-cultural	Brushwork: 0.99; Cross-	Subjective/Biased: 0.89;
	"Twelve Months Paintings"	Interpretation	cultural: 0.98; Layout: 0.98	Logical Gaps: 0.74
April		-	·	
四月(with_Mama_Zola	从这幅《四月令图》中,我们可见郎世宁	Historical Research	Historical Context: 0.96;	Classical Citations: 0.78;
).txt	融合中西方画法的独特成就 (如果这里		Brushwork: 0.92; Cross-	Profound Insight: 0.62
	提及了 Li Ruoyun,则改为 Mama Zola)		cultural: 0.87	
May				
五月(with_Okakura_Ka	此《五月令图》乃郎世宁为乾隆皇帝所作,	Aesthetic Appreciation	Brushwork: 0.95; Style/School:	Classical Citations: 0.82; Strong
kuzo).txt	笔法精妙,构图宏大… (如果这里提及了		0.88; Historical Context: 0.85	Argumentation: 0.55
	Dong Qichang, 则改为 Okakura Kakuzō)			

- 1842 1843
- 1845 1846
- 1847

1849

1852 1853

1854

1855

1858

1859

1860

1863

- 1850
- 1851
- 1856 1857

1861

- 1864 1865
- 1866 1867
- 1868 1869
- 1870
- 1871

1875 1876

- 1879
- 1880 1881

1882 1883 1884

1885 1886

1888

intervention on model output characteristics.

G **Limitations and Ethics Statement**

G.1 Limitations

This study has several limitations that should be considered when interpreting our results:

Bevond the specific points enumerated below. this study confronts broader limitations inherent in current AI capabilities and evaluation methodologies. Models, despite interventions, may still reflect biases from their foundational training data or struggle with true generalization to vastly different cultural artifacts or artistic forms beyond the Chinese paintings studied. The very tools of our framework, such as the zero-shot classifier for feature extraction and the predefined granularity of persona cards and knowledge bases, introduce their own constraints and potential blind spots, possibly failing to capture the full spectrum of expert nuance or the entirety of relevant domain knowledge. Furthermore, the sensitivity of LLMs to prompt engineering and the finite scope of our dataset could influence the observed outcomes. At a more fundamental level. a significant challenge remains in distinguishing between genuine understanding or deep cultural adaptability and sophisticated pattern matching or role-play by the models. The rich, often tacit, knowledge that informs human expert critiquesubtleties of intuition, embodied experience, and deeply internalized cultural schemas-largely eludes current computational approaches and quantitative metrics, posing an ongoing frontier for research in culturally-situated AI.

Vector Space Model Limitations:

- Embedding Model Specificity: Our vector space analysis relies on a specific embedding model, and results might vary with different models. While we selected a model fine-tuned for Chinese art commentary, it may still have limitations in capturing certain cultural nuances.
- Dimensionality **Reduction:** Visualizations using dimensionality reduction techniques inevitably lose some information from the original high-dimensional space, potentially obscuring subtle relationships.

- Semantic Similarity Metrics: Cosine 1889 similarity and other metrics provide useful quantitative comparisons but 1891 may not perfectly align with human 1892 judgments of semantic similarity in 1893 specialized domains. 1894

1895

1901

1902

1903

1904

1905

1911

Structured Commentary Limitations:

- Format Constraints: The structured 1896 format may artificially constrain both 1897 human and MLLM expression patterns, 1898 potentially reducing stylistic diversity 1899 and creative interpretation.
- Scaffolding **Effects:** The template-based section may artificially performance boost MLLM by providing explicit categories and prompts that guide responses.
- Human Expert Adaptation: 1906 Converting existing human expert 1907 commentaries to our structured format 1908 required interpretation and adaptation, 1909 potentially introducing biases. 1910

• Model and Evaluation Limitations:

- Model Selection: While we selected 1912 diverse and state-of-the-art models, 1913 results might differ with other MLLMs. Our API-based approach precludes 1915 deep analysis of models' internal 1916 mechanisms. 1917
- Prompt Sensitivity: Despite our 1918 standardized approach, MLLMs may 1919 exhibit sensitivity to minor variations in 1920 prompt phrasing or structure, affecting 1921 the consistency of results. 1922
- Temporal Limitations: Our study 1923 represents a snapshot of current MLLM 1924 capabilities, which are rapidly evolving. 1925
- Evaluation Subjectivity: Despite our structured rubrics, the evaluation of 1927 cultural understanding and language 1928 quality necessarily involves some subjective judgment. 1930
- Dataset and Domain Limitations: 1931

- 1932 Dataset Specificity: Our analysis
 1933 focuses on specific collections of
 1934 Chinese paintings, and findings may
 1935 not generalize to other cultural artifacts
 1936 or artistic traditions.
 - Annotation Influence: The annotations on input images may influence MLLM outputs in ways that differ from how they would process unannotated images.
 - Cross-Lingual Considerations: Cultural nuances may be lost in translation between Chinese and English, particularly for specialized art terminology.
 - Generalizability: The methodology's effectiveness may vary across different cultural domains and content types.
 - G.2 Ethics Statement

1939

1941

1942

1943

1944

1945

1946

1947

1948

1949

1951

1952

1953

1954

1955

1956

1957

1958

1959

1961

1962

1963

1964

1965

1967

1968

1969

1970

1971

1972

1973

1974

1976

This research raises several ethical considerations:

Cultural Representation and Respect:

- Cultural Authority: We acknowledge the ethical complexities of computational systems interpreting culturally significant artifacts. particularly when these systems are developed primarily in Western contexts.
- Interpretive Plurality: We recognize that there is no single "correct" interpretation of cultural symbols, and that diverse perspectives have validity within their cultural contexts.
- Persona Construction: Our persona designs inevitably reflect our own understanding and conceptualization of different expert roles, which may contain unintentional biases or oversimplifications.

AI Application Considerations:

 Potential Misuse: The structured prompting and persona intervention techniques could potentially be misused to generate misleading or biased interpretations if not implemented responsibly. Algorithmic Bias: The underlying
 MLLMs may contain biases that affect
 their interpretations, particularly across
 cultural contexts, which our
 interventions might not fully address.

1983

1984

1985

1986

1988

1989

1990

1991

1992

1993

1994

1995

1997

1998

1999

2004

2005

2006

2010

2011

2013

2014

2015

2016

2018

2021

- **Transparency:** We emphasize the importance of transparency when deploying AI systems for cultural interpretation, including clear disclosure of the use of persona interventions.
- Human Oversight: While our methods can enhance MLLM capabilities, we advocate for maintaining human expert oversight in sensitive cultural heritage applications.

• Data and Resource Considerations:

- Data Provenance and Licensing: The Twelve Months Series paintings were accessed through the National Palace Museum (Taiwan) digital archives under CC BY 4.0 license.
- Computational Resources: Our vector space analysis approach requires significant computational resources, which may limit accessibility for some researchers or institutions.
- Expert Knowledge Access: The development of effective persona cards requires access to specialized knowledge, which may create barriers to implementing similar approaches in other cultural domains.

We have designed our research to contribute to more culturally sensitive AI applications while acknowledging the limitations of computational approaches to cultural interpretation. Our quantitative framework and structured evaluation methods aim to provide transparent and reproducible results while respecting the complexity and diversity of cultural interpretation.

G.3 Standardized Prompt Design 2019

H Model Details for Multi-Model Comparative Evaluation

This section provides detailed information on the
models included in the multi-model comparative2022
2023

evaluation, including architecture, parameter modality context 2025 count, support, length, knowledge cutoff, licensing, and access method. Detailed specifications for each evaluated model are provided in Table 9. 2028

Notes:

2024

2027

2029

2030

2031

2032 2033

2035

2036

2037

2038

2039

2040

2042

2044

2045

2046

2048

2049

2050

2052

2053

2054

2056

2057

2058

2059

2061

2067

2070

- Gemini 2.5 Pro: Google proprietary model, supports multimodal input (text, image, code), commercial API only, knowledge cutoff 2024.
 - Llama-3.1-8B-Instruct: Open-source, textonly, 8B parameters, 128K context, supports 8+ languages, weights available on Meta and Hugging Face.

• Llama-4-Scout-17B-16E-Instruct:

Open-source, natively multimodal (text, image, code), 17B activated parameters (MoE, 16 experts, 109B total), 10M context, 12 languages, weights on Meta/HF, knowledge cutoff Aug 2024.

• Qwen-2.5-VL-7B: Open-source, supports text, image, video, 7B parameters, 32K+ context, 12+ languages, Apache 2.0 license, weights on Alibaba/HF, knowledge cutoff 2025.

For further details on model usage, inference settings, and prompt templates, see the main text and project documentation.

Knowledge Base Content Ι

This section contains the full content of the knowledge base.json file used to provide structured domain knowledge to the MLLMs during certain experimental conditions.

• Chinese Landscape Painting Concepts (中 国山水画概念):

- Core Concept (核心理念): The core of Chinese landscape painting is "spirit resonance" (qi yun sheng dong), the foremost principle of Xie He's "Six Canons", referring to the vitality, spirit, and verve presented in a work, emphasizing the unity of inner spirit and outer expression. Another core concept is "artistic conception" (vi jing), which is the emotion, atmosphere, and profound meaning conveyed by the painting beyond the objects themselves, pursuing an artistic effect of fused 情景 2071 (emotion/scene) 境 and 2072 (milieu/boundary), inspiring 2073 contemplation. Landscape painting also 2074 embodies the idea of "harmony between man and nature" (tian ren he vi), entrusting philosophical thoughts 2077 and emotions through the depiction of nature.

- Main Features (主要特点): The main features of Chinese landscape painting Subject Matter: Primarily 1. are: natural mountains and rivers, forests, clouds, and water, often imbued with literati sentiments such as reclusion and spiritual refreshment. 2. Brush and Ink (Bi Mo): Utilizes a brush, ink, and Xuan paper, emphasizing the "bone method in brushwork" (gu fa yong bi), shaping the texture of objects and expressing emotions through variations in the strength of lines and the density, wetness, and dryness of ink (e.g., outlining, texturing, rubbing, dotting, dyeing). 3. Composition (Zhang Fa): Focuses on the interplay of void and solid, appropriate density, echoing openings and closings, and leaving blank spaces to create profound artistic conception and pictorial momentum, often using perspective methods like "level distance" (ping yuan), "high distance" (gao yuan), "deep and distance" (shen yuan). 4. Pursuit of Artistic Conception: Seeks not complete formal resemblance but rather spiritual likeness, emphasizing the 2107 integration of poetry, calligraphy, 2108 painting, and seals, and pursuing 2109 meaning beyond the painted image. 2110
- Brief History (简史): Chinese 2111 landscape painting originated in the 2112 Jin, Southern and Northern Wei, 2113 Dynasties, and became an independent 2114 genre in the Sui and Tang Dynasties. 2115 The Five Dynasties to the Northern 2116 Song (907-1127) was its "great era", 2117 with numerous famous artists (e.g., Jing 2118 Hao, Guan Tong, Dong Yuan, Ju Ran, Li Cheng, Fan Kuan, Guo Xi), forming 2120

Table 9: Detailed Specifications of Evaluated Models (Corresponds to Table 1 in 'list.md' outline)

Model	Parameters	Architecture	Modality	Context Length	Knowledge Cutoff	License	Access
Gemini 2.5 Pro	Proprietary	Proprietary (Google)	Text, Image, Code	1M+	2024	Commercial	API (Google Cloud)
Llama-3.1-8B-Instruct	8B	Transformer	Text	128K	Dec 2023	Llama 3.1 Community	Open Weights (Meta, HF)
Llama-4-Scout-17B-16E-Instruct	17B (MoE, 16 experts)	MoE Transformer	Text, Image, Code	10M	Aug 2024	Llama 4 Community	Open Weights (Meta, HF)
Owen-2.5-VL-7B	7B	Transformer	Text, Image, Video	32K+	2025	Apache 2.0	Open Weights (Alibaba, HF)

distinct northern and southern styles: northern landscapes were majestic, while southern water towns were gentle. The Southern Song period placed more emphasis on poetic meaning and personal emotional expression (e.g., Ma Yuan, Xia Gui). Literati painting rose in the Yuan Dynasty, emphasizing the interest of brush and ink and subjective expression (e.g., Zhao Mengfu, the Four Masters of Yuan). The Ming and Qing Dynasties saw further development and a divergence of schools based on inherited traditions, with court painting and literati painting coexisting.

2121

2122

2123

2124

2125

2126

2127

2128

2129

2130

2131

2133

2134

2135

2137

2138

2139

2140

2141

2142

2143

2144

2145

2146

2147

2148

2149

2150

2151

2152

2153

2154

2155

2156

2157

2158

2159

2160

2161

2162

2163

2164

・ Qing Court Painting (清代宫廷绘画):

- Overview (概述): Qing Dynasty court painting was managed by the Imperial Household Department. During the Qianlong era, specialized institutions such as the Ruyi Guan (Palace Ateliers) and the Painting Academy Office were established. Painters were strictly managed, with systems for examination, ranking, rewards and punishments, and work review. It primarily served the imperial family, with functions including recording the appearance and life of emperors and empresses, documenting major state events and ceremonies (e.g., Southern Inspection Tours, battle scenes), decorating palaces and gardens, religious propaganda, and historical reference. Its development is divided into three periods: Shunzhi-Kangxi (initial phase), Yongzheng-Qianlong (peak, with a complete system and numerous famous artists), and post-Jiaqing (decline), synchronized with the rise and fall of national strength.

 Characteristics (特点): Qing Dynasty court painting covered a wide range of subjects, including portraits of emperors, empresses, and meritorious officials, 'scenes of pleasure' (xingletu), major historical events (Southern Inspection Tours, wars, ceremonies), religious paintings, decorative landscapes and flower-and-bird paintings, and documentary-style depictions of tribute animals and plants. The overall style was meticulous, detailed, richly colored, and regal. The most prominent characteristic was the fusion of Chinese and Western styles: influenced by European missionary it emphasized light and painters, shadow, three-dimensionality, employed linear perspective ("xianfa hua"), and introduced oil painting and copperplate engraving. Simultaneously, landscape traditional ("the Four Wangs" school) and flower-and-bird (Yun Shouping's school) painting styles also continued.

2165

2166

2167

2168

2169

2170

2171

2172

2173

2174

2175

2176

2177

2178

2179

2180

2181

2182

2183

2184

2185

2186

2187

2188

2189

2190

2191

2192

2193

2194

2195

2196

2197

2198

2199

- Representative Figures (代表人物): Representative painters include: early figures such as Jiao Bingzhen, Leng Mei, Tang Dai; peak period Chinese painters like Chen Mei, Ding Guanpeng, Jin Tingbiao, Xu Yang, Yao Wenhan, Zhang Zongcang; European painters (excluding Lang Shining) such as Jean Denis Attiret (Wang Zhicheng), Ignatius Sickeltart (Ai Qimeng), etc. Additionally, there were court official painters like Dong Bangda, Jiang Tingxi, etc.

• Giuseppe Castiglione (郎世宁):

- Biography Summary (生平简介): Giuseppe Castiglione (Lang Shining, 2202 1688-1766), an Italian from Milan, was 2203 a Jesuit. He came to China in the 54th 2204 year of Kangxi (1715) and entered the 2205 court around the Kangxi-Yongzheng transition. serving the Kangxi, Yongzheng, and Qianlong emperors. 2208

2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
0005
2235
2230
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2231

2210

2211

2212

His main activities included creating paintings, participating in the design of the Western-style buildings in the Old Summer Palace (Yuanmingyuan), teaching Western painting techniques, and assisting Nian Xiyao in writing 'Shi Xue' (The Study of Vision). He was favored during the Qianlong era and was posthumously granted the title of Vice Minister.

- Artistic Style Overview (艺术风格概 述): In his early period, Lang Shining's style was typically Western. Later, to adapt to the aesthetic tastes of the Chinese imperial family, he integrated Chinese painting techniques, forming a style that blended Chinese and Western elements. His paintings emphasized realism, focusing on light and shadow, perspective, and anatomical structure, but also adopted Chinese painting methods such as even lighting and a focus on line work. Although his style was praised by the court, it was not recognized by the literati painting school.

- Major Contributions (主要贡献): He introduced systematically Western painting techniques such as oil painting and linear perspective (xianfa hua) to the Qing court and taught them, promoting the fusion of Chinese and Western art and forming a new look for Qing court painting. He assisted in the completion of 'Shi Xue' (The Study of advancing the spread of Vision), perspective studies. His documentary-style paintings are important historical materials.
- Representative Works Mention (代表 作列举): Besides the 'Twelve Months Paintings', his representative works include 'One Hundred Horses'. 'Assembled Auspicious Objects', 'Pine, Rock, and Auspicious Fungus', 'Ayusi Attacking Bandits with a Spear', 'Emperor Qianlong's Spring Message of Peace', etc. He also participated in large-scale documentary creating paintings such as 'Banquet in the

Garden of Ten Thousand Trees' and 'Equestrian Skills'.

2259

2260

2261

2263

2265

2266

2267

2268

2271

2272

2274

2275

2276

2277

2278

2281

2282

2284

2287

2288

2289

2290

2291

2292

2294

2295

2296

2297

2299

2301

2302

2303

2304

2305

2306

2307

• Twelve Months Paintings (十二月令图):

- Theme Content (主题内容): The 'Twelve Months Paintings' is a series of 12 works on silk with colors, created by Lang Shining, depicting representative seasonal activities and life scenes in the Qing Dynasty court for each month of the year, such as viewing lanterns in the first month, dragon boat racing in the fifth month, and moon gazing in the eighth month, meticulously showcasing figures, costumes, architecture, and natural scenery.
- Artistic Significance (艺术意义): This series is a mature representative work of Lang Shining's style blending Chinese and Western elements, integrating Western perspective and light/shadow with traditional Chinese composition and aesthetics. It is not only a precious pictorial historical material for studying Qing Dynasty court life and culture but also an important testament to Sino-Western artistic exchange in the 18th century.
- Dataset Source Annotation (数据集 来源与标注): The images for this research dataset are primarily sourced from the National Palace Museum (Taiwan) digital archives (600dpi, CC BY 4.0). Each painting has been annotated in three layers: visual elements, cultural symbols, and artistic techniques, to support AI evaluation and cultural-aesthetic analysis.

References

- Yi Bin, Wenhao Shi, Yujuan Ding, Zhiqiang Hu, Zheng Wang, Yang Yang, See-Kiong Ng, and Heng Tao Shen. 2024. Gallerygpt: Analyzing paintings with large multimodal models. *arXiv preprint arXiv:2408.00491*.
- Zhengqing Chen, Jiayi Jiang, Yifan Jiang, Yiyang Yin, and Nanyun Peng. 2024. Artgpt-4: Artistic visionlanguage understanding with adapter-enhanced mllm. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition (CVPR), pages 1–10.

2413

2414

2415

Chaoyou Fu, Peixian Chen, Yunhao Shen, Yunjie Lin, Shuhuai Zhao, Fangyun Zhang, Baobao Zhao, Weizhu Xie, and Yu Qiao. 2024. Mme: A comprehensive evaluation benchmark for multimodal large language models. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 1–10.

2308

2309

2311

2313

2314

2315

2316

2317

2318

2319

2320

2324

2326

2327

2328

2338

2339

2340

2341

2345

2354

2356

2357

2358

2359

2360

2361

2364

- Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, and Deyi Xiong. 2023. Evaluating large language models: A comprehensive survey. *Preprint*, arXiv:2310.19736.
- Kazuki Hayashi, Kazuma Onishi, Toma Suzuki, Yusuke Ide, Seiji Gobara, Shigeki Saito, Yusuke Sakai, Hidetaka Kamigaito, Katsuhiko Hayashi, and Taro Watanabe. 2024. Irr: Image review ranking framework for evaluating vision-language models. *arXiv preprint arXiv:2402.12121*.
- Ruixiang Jiang and Changwen Chen. 2025. Multimodal llms can reason about aesthetics in zero-shot. *arXiv preprint arXiv:2501.09012*.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *Preprint*, arXiv:1602.07332.
- Haoxuan Li, Shounan An, Zhe Geng, Haotian Liu, Qing Lian, He Kuan, Wentao Wu, Xizhou Zhu, Yong Jae Lee, and Chunyuan Li. 2023. Ferret: Refer and ground anything anywhere at any granularity. *Preprint*, arXiv:2310.07702.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. 2014. Microsoft coco: Common objects in context. *ArXiv*, abs/1405.0312.
- Yifan Liu, Yijia Shao, Yisi Jay Lin, Hwaran Lee, Arjun Kingdon, Sourodeep Jash, Zhaohary HUTHERN Andrianome, Tong Sun, Ning Wang, Kunal Karia, Parth Patwa, Haoyu Wang, Vignesh Pagadala, Bryan Hsueh, Wenxi Chen, Yuan Fang, Peng Gao, Liqun Shao, Rogerio Feris, and 5 others. 2024. Towards multimodal llm benchmarking: A practical guide based on real-world use cases. *Preprint*, arXiv:2402.18060.
- Tanisha Mishra, Edward Sutanto, Rini Rossanti, Nayana Pant, Anum Ashraf, Akshay Raut, Germaine Uwabareze, Ajayi Oluwatomiwa, and Bushra Zeeshan. 2024. Use of large language models as artificial intelligence tools in academic research and publishing among global clinical researchers. *Scientific Reports*, 14(1):31672.
- Zhixiang Qi. 2024. The Spirit of Traditional Chinese Aesthetics. *Asian Studies*, 12(1):281–300. Publisher: Springer Nature Singapore.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical*

Methods in Natural Language Processing, pages 3982–3992.

- Junfeng Wang, Qiushi Jiang, Tianyi Wang, Xiaoya Liu, Xipeng Qiu, Lei Li, and Xuanjing Huang. 2023a. Personallm: Investigating the ability of Ilms to understand and represent different personas. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 14092–14105.
- Junfeng Wang, Qiushi Jiang, Tianyi Wang, Xiaoya Liu, Xipeng Qiu, Lei Li, and Xuanjing Huang. 2023b. Personallm: Investigating the ability of llms to understand and represent different personas. *arXiv preprint arXiv:2307.10188*.
- Yuhan Wang. 2024. The Changes of "Shen" and "Yi" in Chinese Painting Aesthetics: From Gu Kaizhi to Ni Zan. *Open Access Library Journal*, 11(1-6).
- Zekun Wang, Zhiwei Zhang, Tianyi Chen, Xin Wang, Zijian Zhao, Haoran Huang, and William Yang Wang. 2024. RoleLLM: Benchmarking, Eliciting, and Enhancing Role-Playing Abilities of Large Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 814–830.
- Mingi Zhang, Yicheng Li, Xufang Tian, Yuda Zhang, Zihan Huang, Yuxuan Lu, Banghua Chen, Xueyuan Cao, and Dong Liu. 2023. Can llm-augmented autonomous agents cooperate?, an evaluation of their cooperative capabilities through melting pot. *Preprint*, arXiv:2311.17933.
- Tuo Zhang, Tiantian Feng, Yibin Ni, Mengqin Cao, Ruying Liu, Katharine Butler, Yanjun Weng, Mi Zhang, Shrikanth S. Narayanan, and Salman Avestimehr. 2024a. Creating a lens of chinese culture: A multimodal dataset for chinese pun rebus art understanding. *arXiv preprint arXiv:2406.10318*.
- Wei Zhang, Wong Kam-Kwai, Biying Xu, Yiwen Ren, Yuhuai Li, Minfeng Zhu, Yingchaojie Feng, and Wei Chen. 2024b. Cultiverse: Towards cross-cultural understanding for paintings with large language model. *arXiv preprint arXiv:2405.00435*.
- Zijian Zhou, Ling Jiang, Xi Yin, Xu Ren, Radu Soricut, Christoph Feichtenhofer, and Florian Strub. 2024. Culturally diverse prompting for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–10.