

# DreamFrame: Enhancing Video Understanding via Automatically Generated QA and Style-Consistent Keyframes

This supplementary document is organized as follows:

- Section A contains more details about the model architecture of DreamFrame-7B.
- Section B contains more details about training parameters.
- Section C contains prompt examples for the stage one of our approach.
- Section D contains type distribution and more examples of our generated data.
- Section E contains comparison results on movie understanding.

## A MORE DETAILS ON MODEL ARCHITECTURE

Here, we introduce the fundamental modules used in LLaMA-VID. LLaMA-VID basically consists of a visual encoder, a text decoder, a projector and a LLM. For visual encoder, LLaMA-VID use EVA-G [2] as its ViT-based backbone and the patch size is set to 14. For text decoder, Qformer-7b [1] is used in our experiment. For projector, one-layer MLP is used to transform the embedding into context token. For LLM, we use vicuna-7b.

## B MORE TRAINING DETAILS ON DREAMFRAME-7B

As mentioned, we conduct experiments after the second stage of LLaMA-VID. During the model training phase, we employ the original LLaMA-VID configuration as the foundation for our training process. We utilize 2 NVIDIA A100 GPUs. To conserve GPU memory, we employ deepspeed with zero3 during model training, disabling tf32 and opting for fp16. The remaining parameters are shown in Table A

## C PROMPT EXAMPLE

We illustrate the movie plot generation prompt in Figure A. By providing specific elements such as themes, overview, and styles, we guide GPT-4 to produce movie-level key frame descriptions tailored to the latter generation process.

## D MORE DETAILS AND RESULTS OF GENERATED DATA

We show more results of our generated data in Figure B and Figure C.

## E COMPARISON RESULTS ON MOVIE UNDERSTANDING

We demonstrate several comparison results for the comprehension of movie clips to show how we improve our baseline model LLaMA-VID in Figure D. Words in blue color shows the alignment with ground truth while words in red color mean wrong prediction part. Model trained on our dataset demonstrates more reasonable answers.

Settings	Parameter	
Batch size	4	64
Learning rate	2e-5	65
Learning schedule	Cosine decay	66
Warmup ratio	0.03	67
Weight decay	0	68
Optimizer	AdamW	69
Vision encoder	Freeze	70
Text decoder	Freeze	71
Max token	65536	72

Table A: Training settings.

---

### Algorithm 1: The Paradigm of DreamFrame

---

```

Require: In-context learned LLM  $G$ , Diffusion model  $\epsilon_\theta$ .
Input: movie description  $D$ , prompts  $P$ .
Output: Style descriptions  $S$ , Frame descriptions  $FD$ , Init style image set  $I$ , Overview  $O$ , Style embedding  $v$ , Visual frames  $VF$ , QA pairs  $Q$ .
1  $S, O \leftarrow G(P, D);$ 
2  $FD, Q \leftarrow G(P, O)$  using story expanding;
3 for  $i \leftarrow 1$  to 10 do
4    $x_t \sim \mathcal{N}(0, I);$ 
5   for  $t \leftarrow T$  to 1 do
6      $| x_{t-1} \leftarrow \epsilon_\theta(x_t, c_\theta(S), t);$ 
7   end
8   Add  $x_0$  to  $I$ ;
9 end
10 for  $j \leftarrow 1$  to  $\text{len}(I)$  do
11   for  $t \leftarrow T$  to 1 do
12      $| v_t \leftarrow \text{Optimized by LDM Loss with gt noise } \epsilon(I_j);$ 
13   end
14 end
15 for  $k \leftarrow 1$  to  $\text{len}(FD)$  do
16    $| FD_k \leftarrow \text{Prefix with a style word + } FD_k;$ 
17    $| x_t \sim \mathcal{N}(0, I);$ 
18   for  $t \leftarrow T$  to 1 do
19      $| x_{t-1} \leftarrow \epsilon_\theta(x_t, c_\theta\{v\}(FD_k), t);$ 
20   end
21   Add  $x_0$  to  $VF$ ;
22 end

```

---

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174



## CHATGPT PROMPT FOR MOVIE PLOT GENERATION

You are a specialist in creating stories. Generate a story which theme is "[theme]". Requirements: First, generate the title, overview and "story tone and the style" of the story. I will use the story to generate images, so "Story tone and style" should describe the artistic style of the pictures, and should use only adjectives and short phrases, not sentences, with the most important phrases included in "("")", and a single one-word style-keyword included in "{}"(Note that you must provide the style-keyword like the example). Second, list all the locations that the story will take place in, at least 12 locations. Only common places are allowed and use at most 2 words for the location. Do not mention character names in the locations. Do not use location like "Mary's" or "Mary's house" to indicate someone's home. Third, generate all the characters' name, and use the name of a Western celebrity to describe them, with clothes color, in the format of "<xxx: a man/woman looks like xxx in xxx>". For clothes color, a simple one-word color is enough. If a character appears in the story for at least twice, put it in "Character" part. A main character appears most frequently, and a supporting character appears in part of the story. Except for the characters listed in "Character" part, all the other characters and objects that will appear in the story should only appear once. For the new story, I need 1 or 2 main characters, and 0 or 1 supporting character. Fourth, generate 11 substories based on the overview. "Scene" part should pick 1,2 or 3 locations in the "Locations" part, and the whole substory will only take place in these locations. "QA" part is a question-answer pair based on each substory. You can only ask one of the four types of questions"what", "why", "how" and "where" after each substory. The answer should be one or two sentences, no need to be too long. Below is an example:

-----  
Story Title: The Symphony of Love

Overview: 'The Symphony of Love' is a heartwarming tale set in the bustling city of New York, painted in a vibrant, contemporary style. The story revolves around Amelia, a talented violinist with a passion for music, and Ethan, a successful businessman with a hidden love for the arts. Their lives intersect as they navigate the complexities of love, ambition, and the pursuit of dreams.

Story tone and style: {Romantic}, ((Urban chic)), ((modern elegance, bohemian charm)), pastel palette, minimalist aesthetic, soft lighting, candid moments, architectural beauty, cityscape, dynamic, emotional, intimate, dreamy, nostalgic

Style-keyword: Romantic

Locations: ['Concert Hall', 'Office', 'Cafe', 'Art Gallery', 'Park', 'Restaurant', 'Apartment', 'Subway', 'Rooftop', 'Street', 'Lake', 'Classroom', 'Cinema']

Character:

Main characters:

<Amelia: a woman looks like Anne Hathaway in red>  
<Ethan: a man looks like Chris Hemsworth in blue>

Supporting character:

<Oliver: a man looks like Robert Downey Jr. in brown>

Substories:

Substory1: The Melody of Chance

Scene: Subway, Street

Content: Amelia and Ethan have a chance encounter in the subway, and walk on the street.

QA: Question: Where did Amelia and Ethan meet each other? Answer: In the subway

Substory2: The Harmony of Friendship

Scene: Cafe

Content: Amelia and Ethan start to build a friendship over shared coffees and conversations at their favorite cafe.

QA: Question: How did Amelia and Ethan know each other? Answer: They meet each other in cafe, drinking coffee.

I didn't list all the substories, but you should generate the complete version.

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

**Figure A: Prompt for Movie Plot Generation.** Our prompt for movie plot generation comprises basic elements, including overview, characters, style and so on.

233		291
234		292
235		293
236		294
237		295
238		296
239		297
240		298
241		299
242	generate an image in Comedic style: In the newspaper office, a man looks like Hugh Laurie in black holds a portfolio labeled 'Wedding of the Century', while standing near a large glass window, looking thoughtful.	300
243		301
244		302
245		303
246		304
247		305
248		306
249		307
250		308
251		309
252		310
253		311
254		312
255		313
256		314
257	generate an image in tense style: In the coffee shop, a man looks like Hugh Jackman in black is still sitting across from a woman looks like Jessica Chastain in navy, looking at her, waiting for a response.	315
258		316
259		317
260		318
261		319
262		320
263		321
264		322
265		323
266		324
267		325
268		326
269		327
270		328
271		329
272		330
273	generate an image in Dynamic style: In the modern office building, a man looks like Robert Downey Jr. in blue is sitting across from the investor, looking at her with hope, showing her his groundbreaking concept.	331
274		332
275		333
276		334
277		335
278		336
279		337
280	<b>Figure B: Examples of generated long video instruction data We use GPT-4 and guided text-to-image generation models to generate consistent key frames of move-level video with reasonable lines and corresponding question-answer pairs. These data are used to train multimodal large language models on video understanding.</b>	338
281		339
282		340
283		341
284		342
285		343
286		344
287		345
288		346
289		347
290		348



358 generate an image in Gritty style:  
359 In the desert, a man looks like Clint  
360 Eastwood in black is sitting on the  
361 desert sand, his gaze fixed on the  
362 hidden faces in the outlaw camp,  
363 with a grave face.

358 generate an image in Gritty style:  
359 In the desert, Monroe is walking  
360 cautiously, approaching the outlaw  
361 camp while holding a small stick,  
362 his face serious.

358 generate an image in Gritty style:  
359 In the desert, a man looks like Clint  
360 Eastwood in black is squatting,  
361 using the stick to draw a map in the  
362 sand, with a determined expression.

358 generate an image in Gritty style:  
359 In the desert, Monroe is standing,  
360 scanning the camp once more, a  
361 mixture of anticipation and concern  
362 etched on his face.

358 generate an image in Gritty style:  
359 In the desert, Monroe is walking  
360 toward the camp, a saddlebag  
361 across his shoulder and eyes on the  
362 horizon, his face hardened.



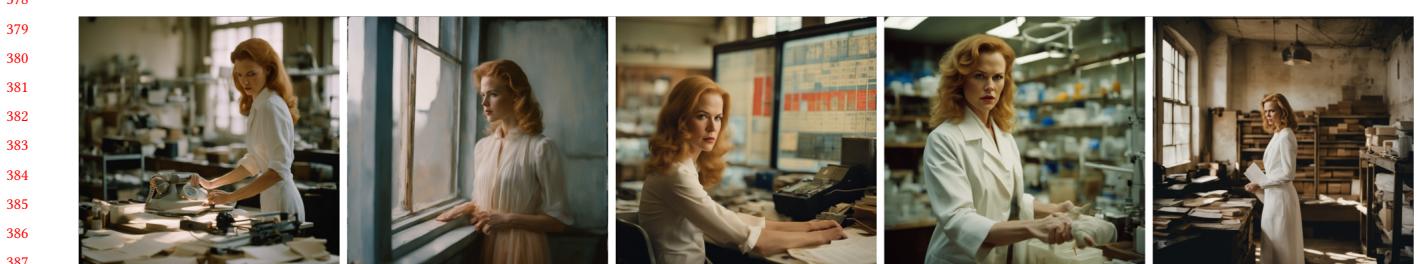
373 generate an image in sad style: At  
374 home, a boy looks like Timothee  
375 Chalamet in black is standing in the  
376 living room, looking out of the  
377 window, with an anguished face.

373 generate an image in sad style: At  
374 home, a boy looks like Timothee  
375 Chalamet in black is sitting in his  
376 room, holding a map of the  
377 countryside, with a thoughtful face.

373 generate an image in sad style: In  
374 the countryside, a boy looks like  
375 Timothee Chalamet in black is  
376 walking across an open field, with  
377 a relieved face.

373 generate an image in sad style: In  
374 the countryside, a boy looks like  
375 Timothee Chalamet in black is  
376 sitting against a haystack, looking  
377 at the setting sun, with a peaceful  
378 face.

373 generate an image in sad style: In  
374 the countryside, a boy looks like  
375 Timothee Chalamet in black is  
376 standing near a stream, watching  
377 the water flow, with a calm face.



388 generate an image in Thrilling style:  
389 In the lab, a woman looks like  
390 Nicole Kidman in white is walking  
391 towards a telephone on her desk, a  
392 look of resolve returning to her  
393 face.

388 generate an image in Thrilling style:  
389 In the lab, a woman looks like  
390 Nicole Kidman in white is standing  
391 by the window, looking out, a  
392 worried look clouding her face.

388 generate an image in Thrilling style:  
389 In the lab, a woman looks like  
390 Nicole Kidman in white is sitting  
391 down, looking at the countdown  
392 clock on her desk, her eyes filled  
393 with determination.

388 generate an image in Thrilling style:  
389 In the lab, a woman looks like  
390 Nicole Kidman in white is standing,  
391 holding her lab coat, looking ready  
392 to fight, a determined expression  
393 hardening her face.

388 generate an image in Thrilling style:  
389 In the lab, a woman looks like  
390 Nicole Kidman in white is walking  
391 out the door, holding her file of  
392 evidence, her face full of  
393 determination for the task ahead.

394 **Figure C: Examples of generated long video instruction data We use GPT-4 and guided text-to-image generation models to**  
395 **generate consistent key frames of move-level video with reasonable lines and corresponding question-answer pairs. These data**  
396 **are used to train multimodal large language models on video understanding.**

397  
398  
399  
400  
401  
402  
403  
404  
405  
406

452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463

465

466

467

468

469

470

471

472

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

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

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

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

## 581 REFERENCES

- 582 [1] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao,  
583 Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. InstructBLIP:  
584 Towards General-purpose Vision-Language Models with Instruction Tuning.  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638
- 639 *arXiv:2305.06500* (2023).
- 640 [2] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun  
641 Huang, Xinlong Wang, and Yue Cao. 2023. Eva: Exploring the limits of masked  
642 visual representation learning at scale. In *CVPR*.
- 643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
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