

VIDEOAGENT: ALL-IN-ONE AGENTIC FRAMEWORK FOR VIDEO UNDERSTANDING AND EDITING

Anonymous authors

Paper under double-blind review

ABSTRACT

Video editing has become essential in digital media creation, yet existing automated systems are restricted to short segment processing and domain-specific tasks. They face two critical limitations: i) inability to handle diverse video comprehension and editing operations, and ii) lack of long-video understanding for coherent narrative creation. We propose **VideoAgent**, an all-in-one agentic framework addressing these challenges through two key innovations. First, we develop automated video shot creation with shot planning agents for coherent narratives and cross-modal retrieval for aligned visual content. Second, we design a multi-agent orchestration framework integrating over thirty specialized editing agents. Intent parsing filters relevant tools while self-reflective graph orchestration assembles complex editing pipelines. Extensive experiments on our newly-proposed **VideoEdit** benchmark and public datasets demonstrate VideoAgent’s superiority over existing multimodal LLMs and agentic systems. VideoAgent achieves 87-98% orchestration success rates while reducing API costs by 60%. Human evaluation across six video categories shows VideoAgent produces professional-quality content approaching human-level performance, with ratings only 4% below human-created videos.

1 INTRODUCTION

Video editing has become indispensable in modern digital media creation (Esser et al., 2023), as evidenced by the explosive growth of platforms like YouTube and TikTok where millions of creators produce content daily (Chai et al., 2023). However, traditional video editing requires complex professional tools, specialized skills, and intricate multi-step workflows that create significant barriers for ordinary users. This gap between widespread content creation demand and technical accessibility has driven the need for automated video editing systems (Wang et al., 2025) that can interpret natural language instructions and execute sophisticated editing tasks without requiring professional expertise.

Recent advances in AI technology have made it increasingly feasible to leverage Large Language Models (LLMs) (Minaee et al., 2024; Jiang et al., 2024; Guo et al., 2025) and Vision Language Models (VLMs) (Zhang et al., 2024b; Li et al., 2022) for multimodal data understanding and editing. Deep Video Discovery (DVD) employs LLM-coordinated tool ensembles to enable autonomous video exploration through systematic "observe-reason-act" cycles (Zhang et al., 2025). VideoRAG introduces a sophisticated dual-channel architecture combining knowledge graphs with multimodal encoding for complex reasoning over long videos (Ren et al., 2025). mPLUG-2 demonstrates modular multimodal foundations using dual visual encoders for unified text, image, and video understanding (Xu et al., 2023b). Though these works have made pioneering explorations on multimodal content understanding and editing, none of them have adequately solved the complex requirements of automated video editing. These existing approaches suffer from two critical limitations: first, they typically focus on specific understanding or editing tasks and cannot handle the diverse range of video comprehension and editing operations required in real-world scenarios. Second, they often lack the long-video understanding and planning capabilities essential for human-level video editing, preventing them from creating videos with sufficient length and coherent narrative structure.

To address these limitations and achieve human-level automated video creation capabilities, we identify two fundamental challenges that must be overcome:

- **Coherent Long-Form Video Planning.** Professional video creation demands coherent narratives that span multiple shots and extended durations, requiring sophisticated top-down shot planning

with multi-step decomposition and strong text-based reasoning capabilities. For instance, creating a movie trailer requires understanding character arcs, identifying climactic moments, and sequencing emotional beats to build narrative tension, a process that involves complex hierarchical planning from story structure down to individual shot specifications. Critically, such shot planning must be conducted with comprehensive awareness of available visual materials, demanding global understanding and precise retrieval across extensive content libraries. Consider the task of creating a highlight reel from a 2-hour film: the system must possess deep comprehension of both granular visual details and overarching narrative elements to select and sequence appropriate clips. Current approaches predominantly focus on isolated segment understanding and editing, lacking the holistic story creation capabilities essential for professional-quality content. Addressing this challenge is crucial as it represents the fundamental difference between automated clip assembly and genuine creative video production that rivals human expertise.

- **Multi-Agent Workflow Orchestration.** Video editing requires orchestrating diverse tools spanning foundational utilities (audio extraction, rhythm detection) to sophisticated creative agents (dialogue creation, narrative structuring) across visual, auditory, and textual modalities with complex interdependencies. Effective production demands systems that can parse natural language instructions, identify creative patterns, and dynamically compose non-linear workflow networks coordinating multi-modal processing. For instance, creating a music video synchronized to beat drops requires seamlessly integrating audio analysis, rhythm detection, visual retrieval, and temporal alignment in orchestrated sequences where each component’s output feeds subsequent operations. Existing methods are constrained to specific domains or predefined workflows, lacking comprehensive tool libraries and adaptive orchestration intelligence for general-purpose creation. This challenge is paramount as it determines whether AI systems can achieve general-purpose video creation or remain limited to domain-specific automated assistance.

To tackle the above challenges, this paper proposes an all-in-one agentic framework VideoAgent for video understanding and editing. Our approach first develops an automated video shot creation system that employs shot planning agents to generate coherent narrative structures, coupled with cross-modal retrieval mechanisms to identify and extract semantically aligned visual content from extensive material libraries, effectively addressing coherent long-form video planning. To achieve general-purpose video creation across diverse genres, we further design a comprehensive multi-agent orchestration framework that integrates over thirty specialized editing agents through dynamic workflow composition, where intent parsing mechanisms efficiently filter relevant tools and self-reflective graph orchestration adaptively assembles complex editing pipelines, enabling seamless execution of sophisticated creative tasks that span multiple modalities and production workflows.

To summarize, this paper makes the following contributions:

- **All-in-one Agentic Framework.** We propose VideoAgent, a comprehensive agentic framework that unifies video understanding and editing capabilities within a single system, enabling automated video creation across diverse genres and production workflows.
- **Novel Technical Solutions.** We introduce global-aware video shot creation and self-reflective agent graph orchestration to address coherent long-form video planning and multi-agent workflow orchestration challenges that have hindered fully automated video editing and remaking tasks.
- **Comprehensive Evaluation.** We construct the VideoEdit benchmark with high-quality instructions and human-aligned evaluation. Experiments demonstrate VideoAgent’s superiority over existing systems, achieving 87-98% success rates while reducing API costs by 60%. Human evaluation shows quality ratings only 4% below professional human-created videos. To facilitate reproducibility, we release our code at <https://anonymous.4open.science/r/VideoAgent-DC33>.

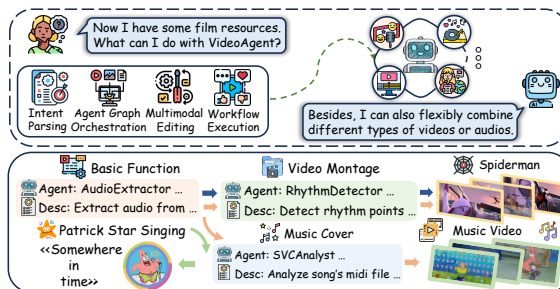


Figure 1: Automated video editing with VideoAgent.

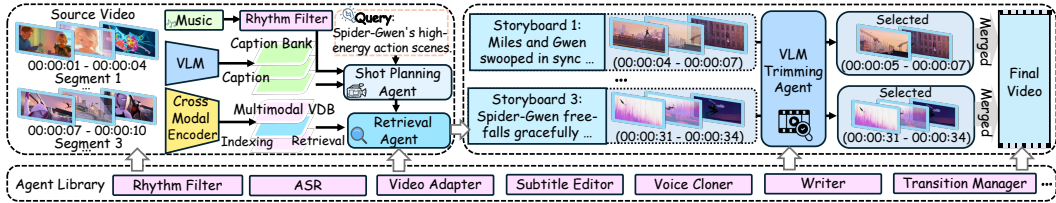


Figure 2: Automated video shot creation with shot planning, video retrieval and trimming.

2 METHODOLOGY

2.1 VIDEO EDITING TASK FORMALIZATION

We formally define the video editing task as follows. **Given i)** a natural language instruction I describing the desired video editing objectives, and **ii)** a collection of multimodal content resources $\mathcal{M} = \{m_1, m_2, \dots, m_{|\mathcal{M}|}\}$ (e.g., raw video clips, audio tracks, and textual documents for narration), the system **produces** a complete edited video V that satisfies the specified instruction I . The output video V comprises a temporally ordered sequence of video shots $V = (v_1, v_2, \dots, v_{|V|})$, where each shot v_i represents a coherent video segment containing synchronized visual and audio components. In essence, the video editing task is to build an agentic system f as follows:

$$f : (I, \mathcal{M}) \rightarrow V = (v_1, v_2, \dots, v_{|V|}) \quad (1)$$

2.2 AUTOMATED VIDEO SHOT CREATION

2.2.1 SHOT PLANNING AGENT

To generate the video shot sequence $(v_1, v_2, \dots, v_{|V|})$, we develop a shot planning agent that produces shot-level storyboards in natural language. Each storyboard provides a textual description specifying the visual content for the corresponding shot, ensuring all shots collectively form a coherent narrative fulfilling the video creation instruction. This process requires global awareness: the shot planning agent must understand both the overall story structure specified by the user instruction I and the available visual resources \mathcal{M} . Given these requirements, the shot planning agent operates as follows:

$$\begin{aligned} S &= (s_1, s_2, \dots, s_{|V|}) = \text{ShotPlanning}(I, \mathcal{M}, |V|) = \text{LLM}(I, C, |V|; p_s) \\ C &= \text{LLM}(c_1, c_2, \dots, c_{|\mathcal{M}|}; p_c), \quad c_i = \text{Caption}(\text{KeyFrames}(m_i)) \end{aligned} \quad (2)$$

where S is the generated shot sequence, C is the compressed visual context summarizing all available materials, c_i is the caption for the i -th visual material, p_s and p_c are prompts for shot planning and visual context generation respectively, and $|V|$ is the desired number of shots.

The shot planning agent operates through two sequential LLM processing steps. First, it processes all keyframe captions to generate a compressed visual summary C that captures the available video materials. Second, it combines this visual summary with the user instruction I to generate the required number of shot descriptions that form the overall narrative structure.

2.2.2 SHOT VISUAL CONTENT RETRIEVAL

Given the shot-level storyboards S , VideoAgent employs paired video indexing and retrieval modules to identify and fetch appropriate video clips for each shot in an efficient manner. The indexing module first extracts keyframes from each video segment using the KeyFrames method mentioned above, then employs cross-modal representation methods such as ImageBind and CLIP to generate embeddings for these keyframes within a unified text-visual representation space. For retrieval, VideoAgent generates embeddings for each shot’s storyboard s_i using the same cross-modal representation approach, then selects the visual content with the highest cosine similarity to the storyboard embedding for use in the current shot. Formally, this indexing and retrieval process works as follows:

$$r_i = \arg \max_j \cos(e_j, \text{Embed}(s_i)), \quad e_i = \text{Embed}(\text{KeyFrames}(m_i)), \quad i = 1, 2, \dots, |V| \quad (3)$$

where r_i represents the retrieved material index for the i -th shot, e_i denotes the embedding for the i -th material’s keyframes, and $\cos(\cdot, \cdot)$ denotes cosine similarity. With this method, our visual content retrieval module fetches semantically aligned visual content to form a coherent narrative following the foregoing text-based shot planning.

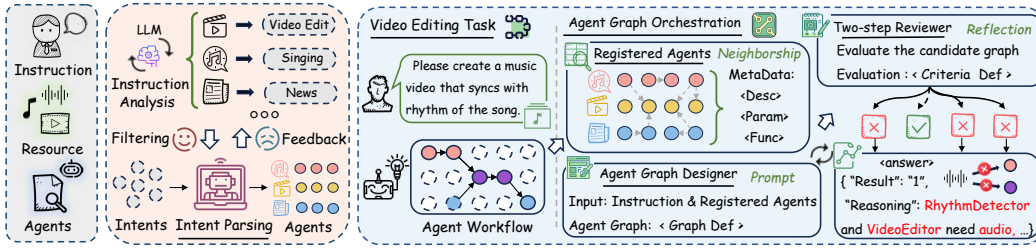


Figure 3: Video editing with agent graph orchestration and execution.

2.2.3 FINE-GRAINED VIDEO TRIMMING

Though the previous modules fetch semantically aligned visual materials, these raw materials are trimmed with uniform length. In real creative work, shot lengths are variable to match the rhythm of background music and dialogue duration (see Listing 9 for shot length decision agents). To address this issue, VideoAgent further designs a fine-grained video trimming agent that employs Vision Language Models (VLMs) with visual understanding capabilities to select the most suitable sub-segment from the retrieved video clip v_{r_i} based on the current shot’s storyboard text, serving as the final visual material for that shot, formally as follows:

$$v_i = m_{r_i}[t_{\text{start}_i} : t_{\text{end}_i}], \quad [t_{\text{start}_i}, t_{\text{end}_i}] = \text{VLM}(m_{r_i}, s_i, t_{i+1} - t_i; p_i) \quad (4)$$

where v_i is the final trimmed video segment for the i -th shot, m_{r_i} is the retrieved raw material, t_{start_i} and t_{end_i} are the start and end timestamps for trimming, $t_{i+1} - t_i$ represents the target duration for the shot, and p_i is the prompt for video trimming guidance. Through this adaptive trimming design, the system achieves precise temporal alignment between visual content and narrative requirements.

2.3 VIDEO EDITING WITH AGENT GRAPH ORCHESTRATION

2.3.1 MULTIMODAL EDITING AGENT LIBRARY

To enhance VideoAgent with advanced video editing capabilities beyond simple concatenation, we develop a multi-agent system for sophisticated video editing operations. This system enables complex creative transformations requiring deep understanding of content semantics, artistic styles, and technical production workflows. The system encompasses two primary categories of tool agents:

- **Foundational multimodal processing tools** that provide essential technical operations such as music beat extraction, voice synthesis, face swapping, lip-sync generation, and text translation.
- **Creation prompts with human expertise** that offer guidance for content generation including style-specific scriptwriting, lyrical composition, and meme-based humor creation. For example, the system can transform English stand-up routines into Chinese crosstalk, generate derivative meme content, or synchronize video cuts with background music rhythms.

These diverse tools are essential for enabling flexible and creative video production workflows, and to support this broad spectrum of editing requirements, we develop a comprehensive Agent Library covering a large array of creation genres. Notably, the functionalities used in aforementioned video shot creation (Section 2.2) are also wrapped as tool agents, ensuring all capabilities are **unified within this multi-agent system framework**. Details about all our tool agents are given in Appendix A.5.

2.3.2 EDITING AGENT GRAPH

To capture the complexity of editing workflows required by advanced video editing, our VideoAgent proposes to compose the tool agents into an agent graph. In this graph, nodes represent individual agents while edges capture information flow, representing the sequential dependencies of agents processing data. Formally, the tool agent library is denoted as $A = \{A_0, A_1, \dots, A_n\}$, and the graph-structured agent workflow is represented as $\mathcal{G} = \{A, E, \delta\}$, where $E \subseteq A \times A$ denotes directed edges encoding inter-agent dependencies (e.g., $A_i \rightarrow A_j$), and $\delta : E \rightarrow \Theta$ defines the state transition function for parameter passing between agents, where Θ denotes the parameter space.

2.3.3 EFFICIENT AGENT SELECTION VIA INTENT PARSING

To achieve efficient agent selection for each creation instruction I , VideoAgent employs an intent parsing mechanism to filter and prioritize relevant agents. Given instruction I , the system reads the agent registry R containing parameter specifications and functional descriptions, and uses LLMs to decompose both explicit and latent sub-intents within I based on a predefined intent-to-agent mapping table M . This process can be defined as:

$$\mathcal{S}_0 = \text{LLM}(I, R, M; p_{iparse}), \quad R = \Psi(A_i | A) \quad (5)$$

Here \mathcal{S}_0 denotes intents extracted from the user instruction, which correspond to a set of agents usually used to solve these intents. p_{iparse} is the prompt for intent parsing, and Ψ is the registry function that generates JSON schema formatted registration tables for LLM comprehension.

2.3.4 SELF-REFLECTIVE AGENT GRAPH ORCHESTRATION AND EXECUTION

To more accurately compose suitable agent graphs for solving user instruction I , VideoAgent adopts a self-reflective workflow for agent graph Orchestration. The process begins with the **intent parsing stage**, where it performs iterative self-reflection and revises its intent parsing outputs. This self-reflection process considers whether the extracted intents can form an agent graph sufficient to solve the user instruction. Formally, it works iteratively as follows:

$$\mathcal{S}_{i+1}, F_{i+1} = \text{LLM}(I, \mathcal{S}_i, F_i; p_{iref}) \quad (6)$$

where \mathcal{S}_i represents the intent set at iteration i , F_i denotes the feedback from the previous iteration, and p_{iref} is the prompt for iterative reflection. Through this iterative self-reflection, more accurate intent parsing results are obtained for efficient agent pre-selection. The process has a maximum iteration limit, with the LLM controlling whether continued reflection is necessary.

Subsequently, for the **workflow planning** stage, we design a two-step self-evaluation mechanism: the first step assesses whether the current graph requires further modification, while the second step reflects and generates a new iteration of \mathcal{G} . This approach enables precise identification of graph deficiencies and addresses potential failure modes, which is defined as follows:

$$r_i = \text{LLM}(F'_i, I, R, \mathcal{G}_i; p_{gref}), \quad \mathcal{G}_{i+1}, F'_{i+1} = \text{LLM}(I, \mathcal{S}, F'_i, \mathcal{G}_i; p'_{gref}) \quad (7)$$

where r_i represents the assessment result at iteration i , F'_i denotes the feedback for graph reflection, R is the agent registry, and p_{gref}, p'_{gref} are the prompts for graph assessment and revision respectively.

Agent Graph Execution. After constructing the agent graph \mathcal{G} , the execution follows the graph-defined workflow through sequential agent invocation. The execution process begins by identifying entry nodes (agents with no incoming edges) and proceeds through topological ordering of the graph structure. Formally, the execution process can be described as:

$$\mathcal{O}_j = A_j(\mathcal{I}_j, \theta_j), \quad \mathcal{I}_j = \bigcup_{(A_i, A_j) \in E} \delta((A_i, A_j))(\mathcal{O}_i) \quad (8)$$

where A_i, A_j represent agents in the graph, $\mathcal{I}_*, \mathcal{O}_*$ refer to input and output data, respectively. θ_j represents the agent-specific configuration parameters. The state transition function $\delta((A_i, A_j))$ transforms the output from agent A_i into the appropriate input format for agent A_j according to the edge specification. The execution continues until all agents in the graph have been processed and the final video editing outputs are generated from the terminal nodes (agents with no outgoing edges).

3 EVALUATION

We conduct extensive experiments to validate VideoAgent’s effectiveness across five research questions: **RQ1** evaluates VideoAgent’s performance compared to various baselines, **RQ2** examines the effect of key modules, **RQ3** analyzes hyperparameter sensitivity, **RQ4** assesses consistency between self-evaluation and human assessment, and **RQ5** demonstrates superiority in real-world scenarios.

3.1 EXPERIMENTAL SETTINGS

Datasets and Evaluation Protocols. We evaluate VideoAgent on video understanding and workflow orchestration using two datasets: **Shot2Story** (Han et al., 2023) for video retrieval evaluation, and our

Table 1: Workflow orchestration performance comparison in terms of *Success Rate*.

Backbone	Claude-Sonnet-4		Claude-Sonnet-3.7		GPT-4o		Deepseek-v3	
Data	Audio	Video	Audio	Video	Audio	Video	Audio	Video
CoT	0.78	0.77	0.80	0.72	0.68	0.65	0.69	0.65
Debate	0.66	0.52	0.67	0.66	0.76	0.69	0.61	0.50
Step Back	0.79	0.85	0.68	0.77	0.63	0.38	0.72	0.70
ExpertPrompting	0.71	0.76	0.78	0.73	0.69	0.68	0.75	0.70
Intelligent Go-Explore	0.69	0.63	0.70	0.75	0.88	0.86	0.58	0.53
Flow	0.62	0.64	0.84	0.83	0.68	0.60	0.66	0.61
VideoAgent	0.95	0.87	0.98	0.95	0.92	0.88	0.94	0.89

new **VideoEdit** benchmark for workflow orchestration with high-quality video creation instructions and human-aligned judge evaluation. See Appendix A.2 for details.

Baseline Methods. VideoAgent is compared with a comprehensive list of baseline methods, including **i) Pure-Language Agentic Systems:** Chain-of-Thought (Wei et al., 2022), LLM Debate (Du et al., 2023), Step Back (Zheng et al., 2023), ExpertPromptingXu et al. (2023a), Intelligent Go-ExploreLu et al. (2024), FlowNiu et al. (2025) **ii) Multimodal Understanding LLMs and Agents:** Qwen2.5-VL (Bai et al., 2025), VideoRAG (Ren et al., 2025), VideoMind (Liu et al., 2025) **iii) Video Generation Agents:** NoteBookLM (Google Labs, 2023), Director (VideoDB, 2024), FunClip (ModelScope, 2024), NarratoAI (linyqh, 2025). Details about baselines are provided in Appendix A.3.

Implementation Details of our VideoAgent and baseline methods are provided in Appendix A.4.

3.2 OVERALL PERFORMANCE COMPARISON (RQ1)

We first compare VideoAgent with baselines on orchestrating effective video editing workflows and video retrieval. The results are shown in Tables 1 and 2, respectively. We observe the following:

- **Superiority in Video Editing Orchestration.** VideoAgent consistently outperforms baselines across all backbone LLMs, achieving success rates of 0.87-0.98 compared to baselines including CoT-based prompting methods and general-purpose agent frameworks. Our method surpasses the best baseline by 2-25% across different modalities and models. This superiority stems from video editing’s inherent complexity involving multiple modalities and diverse functionalities, which existing agentic methods struggle to identify and orchestrate. Our flexible agent graph design and self-reflective orchestration method effectively handles these complex multimodal requirements.
- **Improves Backbones in Video Understanding.** VideoAgent significantly enhances video retrieval performance across foundational models. For example, GPT-4o’s Recall improves from 31.22% to 48.85%, and Gemini-2.5-flash from 30.04% to 44.93%. These improvements result from our shot planning agent’s superior global information perception, generating coherent shot-level storyboards for accurate retrieval. Compared to existing RAG and agent frameworks (VideoRAG, VideoMind-7B), VideoAgent achieves stronger visual-language alignment and overall performance.
- **Cost-Efficiency of VideoAgent.** VideoAgent reduces MLLM API costs by 60% across backbone models while maintaining superior performance. Unlike multimodal LLM baselines that place all candidate shots in context for direct judgment, our strategic shot planning and targeted retrieval minimize redundant processing and focus computational resources on relevant content segments, achieving both better results and significantly lower costs.

Table 2: Video understanding and retrieval performance comparison in terms of Recall@1 (%), Embedding Matching score (EM) (%), Intersection over Union (IoU) (%), and API calling Costs.

Method	Recall	EM	IoU	Cost	Method	Recall	EM	IoU	Cost
Claude-Sonnet-3.7	46.03	27.95	23.91	0.374	Ours-Claude-Sonnet-3.7	44.27	28.18	24.81	0.147
Claude-Sonnet-3.5	27.28	27.35	12.72	0.375	Ours-Claude-Sonnet-3.5	38.70	28.45	21.32	0.147
Gemini-2.5-pro	45.98	27.78	25.91	0.349	Ours-Gemini-2.5-pro	47.24	28.21	25.74	0.136
Gemini-2.5-flash	30.04	28.04	16.69	0.070	Ours-Gemini-2.5-flash	44.93	28.25	25.07	0.028
GPT-4o	31.22	27.64	18.53	0.253	Ours-GPT-4o	48.85	28.26	26.99	0.099
VideoRAG	31.03	15.84	14.35	0.100	Qwen-2.5-VL-72B-Instruct	18.89	27.99	10.51	-
VideoMind-7B	38.26	27.75	19.67	-					

Table 3: Performance of different ablated VideoAgent, in terms of Success Rate (SR).

Backbone	Claude-Sonnet-3.7				Deepseek-v3				GPT-4o			
Data	Audio		Video		Audio		Video		Audio		Video	
Metric	SR	Cost	SR	Cost	SR	Cost	SR	Cost	SR	Cost	SR	Cost
-I	0.96	0.15	0.95	0.19	0.90	0.16	0.87	0.10	0.91	0.15	0.85	0.17
-G	0.70	0.08	0.69	0.12	0.56	0.06	0.61	0.06	0.50	0.11	0.68	0.09
-IG	0.58	0.22	0.66	0.11	0.52	0.15	0.66	0.10	0.46	0.21	0.66	0.13
Origin	0.98	0.09	0.95	0.05	0.94	0.05	0.89	0.09	0.92	0.11	0.88	0.10

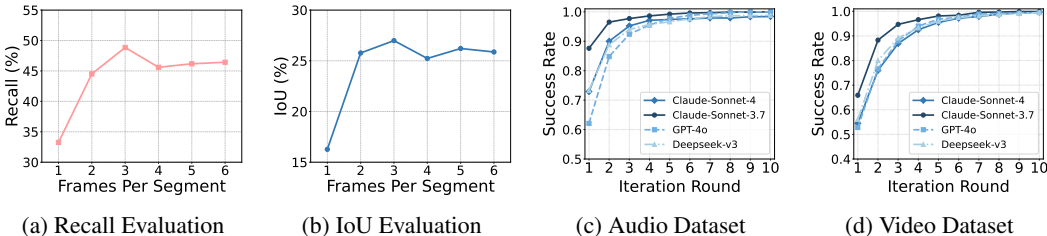


Figure 4: Hyperparameter study for the proposed VideoAgent.

3.2.1 HUMAN EVALUATION ON CREATED VIDEO QUALITY

To assess the quality of videos produced by VideoAgent, we conduct a comprehensive human evaluation study. We generate creation instructions covering six diverse video types including tutorials, entertainment, and promotional content, such as "Please create a commentary video about the provided novel. The video should have a voiceover and appropriate visuals that match commentary content." A total of 26 participants evaluate the videos using a 5-point scoring system based on criteria including visual quality, content coherence, and overall effectiveness. We compare VideoAgent against four baseline methods and human-created videos. The results are shown in Figure 5.

VideoAgent successfully creates videos across all six categories, while existing methods often fail on certain video types, producing no evaluable output and demonstrating VideoAgent’s superior **generalization capability**. Across all categories, VideoAgent achieves **significantly higher scores than baseline methods**, reflecting its ability to effectively orchestrate diverse editing tools with accurate visual content understanding. Notably, VideoAgent’s performance **approaches or even exceeds human-created videos** in several categories, as VideoAgent consistently produces high-quality results while human creators exhibit varying skill levels.

3.3 ABLATION STUDY (RQ2)

We validate the effectiveness of key design components in VideoAgent through systematic ablation studies. We first assess our video shot creation component using the Shot2Story dataset (Figure 6).

- When the shot planning agent is removed and original Shot2Story queries are used directly for retrieval, performance drops significantly across all metrics, with Recall declining from 48.85 to 45.01. This demonstrates the **importance of global shot planning** even when users provide shot-level descriptions.
- Additionally, we replace our cross-modal representation (CM Rep.) method with a caption-then-encoding approach using a pure-language embedding model (all-MiniLM-L6-v2). Results show that our **cross-modal representation paradigm better preserves visual fidelity**, achieving superior retrieval performance at the shot level.

Furthermore, we evaluate the impact of Intent Parsing (-I) and Agent Graph (-G) on workflow orchestration, with results shown in Table 3.

- Removing the graph-based agent orchestration substantially reduces the orchestration success rate from over 90% to below 55%, demonstrating the complexity of video editing workflows and the **critical importance of graph-based orchestration**.
- While removing intent parsing does not significantly affect performance, it leads to a notable increase in computational cost, highlighting the **cost-efficiency of our intent parsing approach**.

3.4 HYPERPARAMETER STUDY (RQ3)

This section studies the impact of two key hyperparameters in our VideoAgent. The evaluation results are shown in Figure 4. From the results we make the following analysis:

Frames per Segment in shot planning affects the number of key frames extracted for each raw video resource m_i . Using only 1 frame results in a substantial decrease in recall performance, indicating insufficient visual information extraction capability. When the frame count is excessively increased, performance slightly declines, potentially due to over-sampling introducing noisy frames.

Iteration Round in workflow orchestration shows that increasing iteration rounds consistently improves the success rate of workflow orchestration, demonstrating typical performance scaling during inference. The results also indicate that, with our graph-based orchestration method, sufficient iterations can significantly reduce performance gaps between different underlying LLM models.

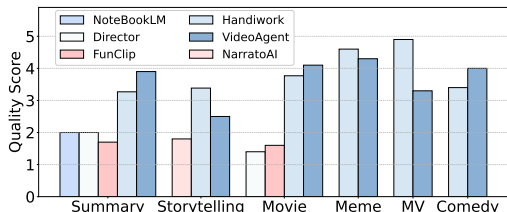


Figure 5: Human-rated video quality assessment.

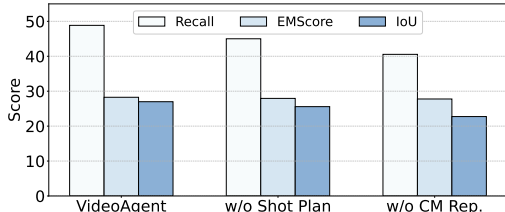


Figure 6: VideoAgent ablation study results.

3.5 LLM-HUMAN JUDGMENT CONSISTENCY (RQ4)

To validate the reliability of LLM-based evaluation on orchestration success rate, we compare the consistency between LLM judge results and human annotations. Specifically, we randomly generate 30 agent graphs using VideoAgent for different creation instructions with iteration rounds set to 1, manually annotate their success status, and calculate accuracy, precision, recall, and F1 scores for LLM judgments compared to human ground truth. Results are shown in Figure 7.

Our analysis reveals that Claude-Sonnet-3.7 demonstrates the most reliable evaluation performance, achieving scores between 0.85-1.0 across all metrics, indicating high consistency with human judgment. We observe that recall consistently exceeds precision across model-modality combinations, suggesting that LLM evaluators tend to adopt more permissive evaluation stances with higher rates of false positives than false negatives. This pattern indicates that our evaluation framework provides conservative assessments that rarely miss actual successful orchestrations.

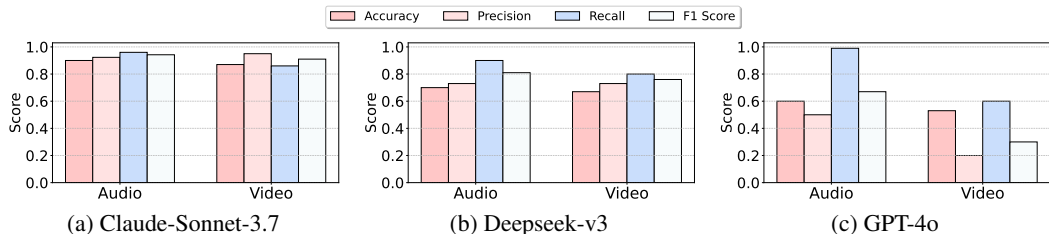


Figure 7: Consistency study on LLM self-evaluation

3.6 CASE STUDY ON CREATED VIDEOS (RQ5)

To demonstrate VideoAgent’s practical capabilities and workflow flexibility, we conduct a comprehensive real-world case study. As illustrated in Figure 8, when a user provides the instruction about creating a rhythm-synced Spiderman movie montage video, the system autonomously parses the instruction’s intent to efficiently orchestrate the appropriate agent graph, subsequently invoking agents such as RhythmDetector that identifies the high-energy, dynamic nature of the content to handle fast-paced visual editing requirements. This showcases VideoAgent’s ability to handle creative scenarios through intelligent agent orchestration. Most notably, when addressing complex video montage creation, VideoAgent demonstrates advanced multimodal understanding, for example, synchronizing visual cuts with musical rhythms and retrieving target scenes or characters based on shot text descriptions. The case study particularly highlights VideoAgent’s strength in multimodal integration and flexible agent combination. Rather than following rigid predefined workflows, the system dynamically adapts its agent graph structure based on user intents and available resources. This flexibility enables seamless transitions between different content types—from basic audio extraction to sophisticated rhythm-synchronized montages—all within a unified framework. The results

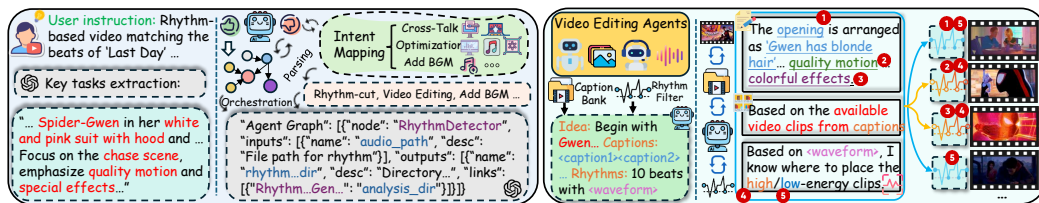


Figure 8: Case study: Creating a rhythm-synced Spiderman movie montage with VideoAgent.

demonstrate VideoAgent’s practical utility in real-world creative scenarios, providing users with intuitive yet powerful tools for diverse video editing requirements.

4 RELATED WORK

Automated Multimodal Content Editing integrates textual, visual, and temporal modalities to enable intelligent content manipulation. Recent advances in diffusion-based models have facilitated text-guided image editing through cross-modal attention mechanisms (Shuai et al., 2024; Fang et al., 2025), while video editing research has addressed temporal consistency via spatiotemporal frameworks incorporating multiple conditional inputs (He et al., 2025; Yang et al., 2025). Multimodal large language models have further enabled systems that parse complex editing instructions and execute corresponding visual operations (Zhang et al., 2024a; Cheng et al., 2025). However, contemporary research primarily emphasizes editing granularity and temporal coherence while overlooking workflow orchestration challenges. Our work addresses comprehensive workflow orchestration by integrating over thirty specialized agents that adaptively assemble editing workflows based on user instructions, enabling seamless execution of complex multi-step tasks without manual coordination.

Multi-Agent Systems. Recent advances in multi-agent systems have shifted from monolithic AI models to collaborative frameworks where multiple autonomous agents interact to solve complex tasks (Zhuge et al., 2024a; Tang et al., 2025a). Key developments include the use of LLMs to enable natural language reasoning and tool invocation. This paradigm demonstrates remarkable flexibility through dynamic role specialization, where agents adaptively assume functions such as evaluation and refinement (Zhuge et al., 2024b), program synthesis and execution (Tang et al., 2025b), and seamless interfacing with external APIs, databases, and computational resources (Shi et al., 2024). However, existing multi-agent systems primarily target text-based applications like coding and mathematics. Our work extends multi-agent frameworks to multimodal video editing workflows.

Video Understanding and Retrieval are intrinsically connected pillars of visual intelligence, with advances in one domain driving progress in the other. Modern research in video understanding focuses on developing specialized models for challenges like visual question answering (VQA) (Qian et al., 2024; Sima et al., 2024; Li et al., 2025), leveraging complex fusion mechanisms to extract spatial and temporal semantics. Representative works such as VideoMind (Liu et al., 2025) and Qwen-2.5-VL (Bai et al., 2025) demonstrate strong capabilities in long-video scenarios, enabling precise temporal localization and spatio-temporal reasoning. Complementarily, retrieval-oriented methodologies leverage cross-modal video-text models that learn aligned representations in shared embedding spaces. Techniques like VideoRAG (Ren et al., 2025) construct comprehensive knowledge graphs integrating diverse long-video sources. Our work differs by designing a framework with strong video understanding and retrieval capabilities specifically for video editing workflows.

5 CONCLUSION

We presented VideoAgent, a comprehensive agentic framework that unifies video understanding and editing capabilities to enable automated video creation across diverse genres and production workflows. Our approach addresses two fundamental challenges in automated video editing: coherent long-form video planning through global-aware shot creation, and multi-agent workflow orchestration via self-reflective graph composition. VideoAgent integrates over thirty specialized editing agents and demonstrates significant superiority over existing methods, achieving 87-98% orchestration success rates while reducing MLLM API costs by 60%. Human evaluation confirms that VideoAgent produces high-quality videos comparable to human creators across multiple categories. These results establish VideoAgent as an effective solution for bridging the gap between widespread content creation demand and technical accessibility, paving the way for more sophisticated automated creative systems.

ETHIC STATEMENT

VideoAgent democratizes video creation by lowering technical barriers, enabling broader creative expression and content production across diverse user communities. However, we acknowledge that automated video editing capabilities could potentially be misused to create misleading or manipulative content that spreads misinformation. To mitigate these risks, we recommend implementing content moderation frameworks and encourage responsible usage practices among users and platforms. LLMs were used sparingly for language refinement in the preparation of this manuscript.

REPRODUCIBILITY STATEMENT

To facilitate reproducibility and open scientific research, we have open-sourced VideoAgent at an anonymous link: <https://anonymous.4open.science/r/VideoAgent-DC33>. The repository includes the complete implementation code, evaluation datasets with corresponding evaluation scripts, and all 30+ specialized tool agents comprising our multimodal editing framework. Beyond this submission, our project has been open-sourced on GitHub and has established a thriving community of over 200 users for deployment and collaboration. The codebase described in this paper has been successfully deployed and utilized by a substantial number of users in real-world scenarios, demonstrating its practical viability and robustness across diverse video editing workflows.

USAGE OF LLMs

In this work, large language models (LLMs) serve as key components across multiple editing agents to enhance understanding, planning, and execution for adaptive video generation aligned with user goals. Specifically, 1) LLMs compress visual keyframe captions and integrate them with user instructions to generate coherent shot-level storyboards, ensuring global narrative consistency; 2) Vision-language models adaptively select precise video segments guided by storyboard text and target durations, achieving fine-grained temporal alignment with narrative rhythm. Additionally, 3) LLMs parse user instructions to extract intents that facilitate efficient agent selection and prioritization from a registry, while 4) dual-stage self-reflective agent graph orchestration iteratively refines intent understanding and evaluates workflows to improve agent coordination and mitigate failures. Furthermore, LLMs are employed in a limited capacity for language refinement during manuscript preparation.

REFERENCES

- Anthropic. Claude sonnet 3.7 technical report, 2025a. URL <https://www.anthropic.com/news/claude-3-7-sonnet>.
- Anthropic. Claude sonnet 4, 2025b. URL <https://claude.ai/new>.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Siboz Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Wenhao Chai, Xun Guo, Gaoang Wang, and Yan Lu. Stablevideo: Text-driven consistency-aware diffusion video editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 23040–23050, 2023.
- Dabing Cheng, Haosen Zhan, Xingchen Zhao, Guisheng Liu, Zemin Li, Jinghui Xie, Zhao Song, Weiguo Feng, and Bingyue Peng. Text-to-edit: Controllable end-to-end video ad creation via multimodal llms. *arXiv preprint arXiv:2501.05884*, 2025.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first International Conference on Machine Learning*, 2023.
- Zhihao Du, Yuxuan Wang, Qian Chen, Xian Shi, Xiang Lv, Tianyu Zhao, Zhifu Gao, Yexin Yang, Changfeng Gao, Hui Wang, Fan Yu, Huadai Liu, Zhengyan Sheng, Yue Gu, Chong Deng, Wen Wang, Shiliang Zhang, Zhijie Yan, and Jing-Ru Zhou. Cosyvoice 2: Scalable

- 540 streaming speech synthesis with large language models. *ArXiv*, abs/2412.10117, 2024. URL
541 <https://api.semanticscholar.org/CorpusID:274762932>.
- 542
- 543 Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis.
544 Structure and content-guided video synthesis with diffusion models. In *Proceedings of the*
545 *IEEE/CVF international conference on computer vision*, pp. 7346–7356, 2023.
- 546 Rongyao Fang, Chengqi Duan, Kun Wang, Linjiang Huang, Hao Li, Shilin Yan, Hao Tian, Xingyu
547 Zeng, Rui Zhao, Jifeng Dai, et al. Got: Unleashing reasoning capability of multimodal large
548 language model for visual generation and editing. *arXiv preprint arXiv:2503.10639*, 2025.
- 549
- 550 Gemini Team. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long
551 context, and next generation agentic capabilities. *ArXiv*, abs/2507.06261, 2025. URL <https://api.semanticscholar.org/CorpusID:280151524>.
- 552
- 553 Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand
554 Joulin, and Ishan Misra. Imagebind one embedding space to bind them all. *2023 IEEE/CVF*
555 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15180–15190, 2023. URL
556 <https://api.semanticscholar.org/CorpusID:258564264>.
- 557
- 558 Google Labs. NotebookLM (AI-powered research assistant). <https://notebooklm.google>, 2023.
559 [Experimental Preview; accessed 2025-08-11].
- 560 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
561 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
562 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 563 Mingfei Han, Linjie Yang, Xiaojun Chang, and Heng Wang. Shot2story20k: A new benchmark for
564 comprehensive understanding of multi-shot videos. *arXiv preprint arXiv:2312.10300*, 2023.
- 565
- 566 Yangfan He, Sida Li, Jianhui Wang, Kun Li, Xinyuan Song, Xinhang Yuan, Keqin Li, Kuan Lu,
567 Menghao Huo, Jingqun Tang, et al. Enhancing low-cost video editing with lightweight adaptors
568 and temporal-aware inversion. *arXiv preprint arXiv:2501.04606*, 2025.
- 569
- 570 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
571 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.
572 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 573
- 574 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-
575 training for unified vision-language understanding and generation. In *International conference on*
576 *machine learning*, pp. 12888–12900. PMLR, 2022.
- 577
- 578 Zhangbin Li, Jinxing Zhou, Jing Zhang, Shengeng Tang, Kun Li, and Dan Guo. Patch-level sounding
579 object tracking for audio-visual question answering. In *Proceedings of the AAAI Conference on*
580 *Artificial Intelligence*, volume 39, pp. 5075–5083, 2025.
- 581
- 582 Shijia Liao, Yuxuan Wang, Tianyue Li, Yifan Cheng, Ruoyi Zhang, Rongzhi Zhou, and Yijin
583 Xing. Fish-speech: Leveraging large language models for advanced multilingual text-to-speech
584 synthesis. *ArXiv*, abs/2411.01156, 2024. URL <https://api.semanticscholar.org/CorpusID:273811499>.
- 585
- 586 linyqh. NarratoAI. <https://github.com/linyqh/NarratoAI>, 2025. [Online; accessed 2025-08-
587 11].
- 588
- 589 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
590 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*
591 *arXiv:2412.19437*, 2024.
- 592
- 593 Jinglin Liu, Chengxi Li, Yi Ren, Feiyang Chen, and Zhou Zhao. Diffsinger: Singing voice synthesis
594 via shallow diffusion mechanism. In *AAAI Conference on Artificial Intelligence*, 2021. URL
595 <https://api.semanticscholar.org/CorpusID:235262772>.
- 596
- 597 Songting Liu. Zero-shot voice conversion with diffusion transformers. *ArXiv*, abs/2411.09943, 2024.
598 URL <https://api.semanticscholar.org/CorpusID:274117146>.

- 594 Ye Liu, Kevin Qinghong Lin, Chang Wen Chen, and Mike Zheng Shou. Videomind: A chain-of-lora
595 agent for long video reasoning. *arXiv preprint arXiv:2503.13444*, 2025.
- 596
597 Cong Lu, Shengran Hu, and Jeff Clune. Intelligent go-explore: Standing on the shoulders of giant
598 foundation models. *arXiv preprint arXiv:2405.15143*, 2024.
- 599
600 Brian McFee, Colin Raffel, Dawen Liang, Daniel P. W. Ellis, Matt McVicar, Eric Battenberg, and
601 Oriol Nieto. librosa: Audio and music signal analysis in python. In *Proceedings of the 14th python
602 in science conference*, volume 8, 2015.
- 603
604 Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier
605 Amatriain, and Jianfeng Gao. Large language models: A survey. *arXiv preprint arXiv:2402.06196*,
2024.
- 606
607 ModelScope. FunClip. <https://github.com/modelscope/FunClip>, 2024. [Online; accessed
2025-08-11].
- 608
609 Boye Niu, Yiliao Song, Kai Lian, Yifan Shen, Yu Yao, Kun Zhang, and Tongliang Liu. Flow:
610 Modularized agentic workflow automation. *arXiv preprint arXiv:2501.07834*, 2025.
- 611
612 OpenAI. Gpt-4o system card. *ArXiv*, abs/2410.21276, 2024. URL <https://api.semanticscholar.org/CorpusID:273662196>.
- 613
614 Tianwen Qian, Jingjing Chen, Linhai Zhuo, Yang Jiao, and Yu-Gang Jiang. Nuscenes-qa: A multi-
615 modal visual question answering benchmark for autonomous driving scenario. In *Proceedings of
616 the AAAI Conference on Artificial Intelligence*, volume 38, pp. 4542–4550, 2024.
- 617
618 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever.
619 Robust speech recognition via large-scale weak supervision. In *International Conference on
Machine Learning*, 2022. URL <https://api.semanticscholar.org/CorpusID:252923993>.
- 620
621 Xubin Ren, Lingrui Xu, Long Xia, Shuaiqiang Wang, Dawei Yin, and Chao Huang. Vide-
622 orag: Retrieval-augmented generation with extreme long-context videos. *arXiv preprint
arXiv:2502.01549*, 2025.
- 623
624 Yaya Shi, Xu Yang, Haiyang Xu, Chunfen Yuan, Bing Li, Weiming Hu, and Zhengjun Zha. Emscore:
625 Evaluating video captioning via coarse-grained and fine-grained embedding matching. *2022
626 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 17908–17917,
2021. URL <https://api.semanticscholar.org/CorpusID:244270365>.
- 627
628 Zhengliang Shi, Shen Gao, Xiuyi Chen, Yue Feng, Lingyong Yan, Haibo Shi, Dawei Yin, Pengjie
629 Ren, Suzan Verberne, and Zhaochun Ren. Learning to use tools via cooperative and interactive
630 agents. *arXiv preprint arXiv:2403.03031*, 2024.
- 631
632 Xincheng Shuai, Henghui Ding, Xingjun Ma, Rongcheng Tu, Yu-Gang Jiang, and Dacheng Tao. A
633 survey of multimodal-guided image editing with text-to-image diffusion models. *arXiv preprint
arXiv:2406.14555*, 2024.
- 634
635 Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Hanxue Zhang, Chengen Xie, Jens
636 Beißwenger, Ping Luo, Andreas Geiger, and Hongyang Li. Drivelm: Driving with graph vi-
637 sual question answering. In *European conference on computer vision*, pp. 256–274. Springer,
638 2024.
- 639
640 Jiabin Tang, Tianyu Fan, and Chao Huang. Autoagent: A fully-automated and zero-code framework
641 for llm agents. *arXiv preprint arXiv:2502.05957*, 2025a.
- 642
643 Jiabin Tang, Lianghao Xia, Zhonghang Li, and Chao Huang. Ai-researcher: Autonomous scientific
644 innovation. *arXiv preprint arXiv:2505.18705*, 2025b.
- 645
646 Vidi Team, Celong Liu, Chia-Wen Kuo, Dawei Du, Fan Chen, Guang Chen, Jiamin Yuan, Lingxi
647 Zhang, Lu Guo, Lusha Li, Longyin Wen, Qingyu Chen, Rachel Deng, Sijie Zhu, Stuart Siew,
Tong Jin, Wei Lu, Wen Zhong, Xiaohui Shen, Xin Gu, Xing Mei, and Xueqiong Qu. Vidi: Large
multimodal models for video understanding and editing. *ArXiv*, abs/2504.15681, 2025. URL
<https://api.semanticscholar.org/CorpusID:277994371>.

- 648 VideoDB. Director. <https://github.com/video-db/Director>, 2024. [Online; accessed 2025-
649 08-11].
- 650
- 651 Yukun Wang, Longguang Wang, Zhiyuan Ma, Qibin Hu, Kai Xu, and Yulan Guo. Videodirector:
652 Precise video editing via text-to-video models. In *Proceedings of the Computer Vision and Pattern
653 Recognition Conference*, pp. 2589–2598, 2025.
- 654 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
655 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
656 neural information processing systems*, 35:24824–24837, 2022.
- 657
- 658 Benfeng Xu, An Yang, Junyang Lin, Quan Wang, Chang Zhou, Yongdong Zhang, and Zhendong
659 Mao. Expertprompting: Instructing large language models to be distinguished experts. *arXiv
660 preprint arXiv:2305.14688*, 2023a.
- 661 Haiyang Xu, Qinghao Ye, Ming Yan, Yaya Shi, Jiabo Ye, Yuanhong Xu, Chenliang Li, Bin Bi,
662 Qi Qian, Wei Wang, et al. mplug-2: A modularized multi-modal foundation model across text,
663 image and video. In *International Conference on Machine Learning*, pp. 38728–38748. PMLR,
664 2023b.
- 665
- 666 Dingyi Yang, Chunru Zhan, Ziheng Wang, Biao Wang, Tiezheng Ge, Bo Zheng, and Qin Jin. Synchron-
667 ized video storytelling: Generating video narrations with structured storyline. In *Annual Meeting
668 of the Association for Computational Linguistics*, 2024. URL [https://api.semanticscholar.
669 org/CorpusID:269983517](https://api.semanticscholar.org/CorpusID:269983517).
- 670 Xiangpeng Yang, Linchao Zhu, Hehe Fan, and Yi Yang. Videograin: Modulating space-time
671 attention for multi-grained video editing. In *The Thirteenth International Conference on Learning
672 Representations*, 2025.
- 673 Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li,
674 Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint
675 arXiv:2408.01800*, 2024.
- 676
- 677 David Junhao Zhang, Dongxu Li, Hung Le, Mike Zheng Shou, Caiming Xiong, and Doyen Sahoo.
678 Moonshot: Towards controllable video generation and editing with multimodal conditions. *arXiv
679 preprint arXiv:2401.01827*, 2024a.
- 680 Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A
681 survey. *IEEE transactions on pattern analysis and machine intelligence*, 46(8):5625–5644, 2024b.
- 682
- 683 Xiaoyi Zhang, Zhaoyang Jia, Zongyu Guo, Jiahao Li, Bin Li, Houqiang Li, and Yan Lu. Deep
684 video discovery: Agentic search with tool use for long-form video understanding. *arXiv preprint
685 arXiv:2505.18079*, 2025.
- 686 Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen, Heng-Tze Cheng, Ed H Chi, Quoc V Le, and
687 Denny Zhou. Take a step back: Evoking reasoning via abstraction in large language models. *arXiv
688 preprint arXiv:2310.06117*, 2023.
- 689 Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen
690 Schmidhuber. Gptswarm: Language agents as optimizable graphs. In *Forty-first International
691 Conference on Machine Learning*, 2024a.
- 692
- 693 Mingchen Zhuge, Changsheng Zhao, Dylan Ashley, Wenyi Wang, Dmitrii Khizbullin, Yongyang Xiong,
694 Zechun Liu, Ernie Chang, Raghuraman Krishnamoorthi, Yuandong Tian, et al. Agent-as-a-judge:
695 Evaluate agents with agents. *arXiv preprint arXiv:2410.10934*, 2024b.
- 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

A APPENDIX

For a deeper understanding of the VideoAgent’s operational mechanism, we provide a comprehensive technical description of the multi-modal agentic framework in the appendix.

- A1. Video Showcase 14
- A2. Evaluation Datasets and Protocols 21
- A3. Baseline Method 26
- A4. Implementation Details 27
- A5. Tooluse Agents 28
- A6. Visual Processing Agents 29
- A7. Audio Processing Agents 34
- A8. Knowledge Analysis Agents 61
- A9. Prompt of Intent Analysis 83
- A10. Graph Construction Rules 84
- A11. Two-step Reviewer 85
- A12. Cases of Multi-modal Agent Workflow 87
- A13. MLLM/VLM for Video Editing Evaluation 90

A.1 VIDEO SHOWCASE

Below are video showcases we created with VideoAgent, which can be grouped into six main categories: **Beat-synced Edits**, **Video Overviews**, **Storytelling**, **Meme Video Edits**, **Song Remixes**, and **Cross-lingual Adaptations** (including stand-up comedy and crosstalk adaptations).

Beat-synced Edits: The user provides background music and film footage. VideoAgent automatically detects beat and tempo changes in the music and aligns high-energy film visual shots with strong beat moments. **Video Overview:** The user provides video footage of a news item or event. VideoAgent transcribes the speech, writes a summary of the event, matches the meaning of each summary sentence to corresponding visual segments from the original video, and generates voiceover for each sentence to complete the edit. **Storytelling:** The user provides plain text of a novel or story along with video footage for visual matching. VideoAgent scripts the story in a specific narration style, matches each line of the script to suitable video segments based on semantic meaning, and generates voiceover for each script segment to produce the final video. **Meme Video Edits:** The user supplies a source video and a custom script or narrative. VideoAgent extracts and transcribes the original audio, generates new speech, precisely syncs the replacement audio to video frames, and outputs a professionally dubbed, meme-style edit. **Song Remixes:** The user provides a MIDI file, lyrics, background music, and a target voice sample. VideoAgent generates a cover, clones the target voice, aligns timing with the arrangement, and integrates the result into the video pipeline for a polished, synced remix. **Cross-lingual Adaptations:** The user provides source audio and target voice samples. VideoAgent adapts the script into the desired cultural format (e.g., Chinese crosstalk or English talk show), synthesizes target voices, applies appropriate effects, and integrates the result into the video editing pipeline for a polished, synced adaptation.

Listing 1: VideoAgent Showcase 1 with Prompts

```

Categories:
1. Beat-synced Edits
  Description:
    - User provides background music and film footage.
    - Detects beat/tempo changes and aligns high-energy shots to strong beats.
↔
Prompt:
```

756 Begin with Gwen, who has blond hair, sitting at a dining table in front of
 757 ↪ a window, then transition to her playing drums with pop textures and
 758 ↪ musical notes in the background. Include action sequences featuring Miguel
 759 ↪ O'Hara in a dark blue suit with red accents, sharp red claws, and black-and
 760 ↪ -red eye lenses; Spider-Gwen in a white-and-pink suit with a hood and
 761 ↪ ballet shoes; Miles Morales with curly hair and a red spider logo on his
 762 ↪ chest; and the Spot in a black suit covered in white spots using portal
 763 ↪ powers. Focus on the chase scene against a blue sky with trains, and
 764 ↪ emphasize high-quality motion throughout--web-swinging, combat, and vibrant
 765 ↪ special effects.

766 2. Video Overviews

767 Description:

768 - User provides a news/event video.
 769 - Transcribe speech, summarize, align summary sentences to visuals,
 770 ↪ generate voiceover.

771 Prompt:

772 - Short tech news, colloquial expression within 250 words, check the
 773 ↪ accuracy of key terms, e.g. the GPT model name should be 40 instead of 4.0

774 3. Storytelling

775 Description:

776 - User provides story/novel text and footage.
 777 - Script in chosen narration style, match lines to visuals, generate
 778 ↪ voiceover.

779 Prompt:

780 - A verbal interpretation copy of no less than 1,000 words

781 4. Meme Video Edits

782 Description:

783 - User supplies a source video and custom script.
 784 - Extract/transcribe audio, generate new speech, sync precisely to frames
 785 ↪ , output dubbed edit.

786 Prompt:

787 - Create a humorous narrative about two PhD students seeking advice from
 788 ↪ Master Ma. For the two PhD students, one of them is known for high citation
 789 ↪ counts and the other for numerous publications. Transform martial arts
 790 ↪ terms into AI research terminology while keeping phrase lengths similar (
 791 ↪ length difference should be less than two Chinese characters). The story
 792 ↪ highlights their academic rivalry and ends with Master Ma advising against
 793 ↪ "internal competition". Keep signature phrases like "wasn't cautious enough
 794 ↪ " and "achieving great results with minimal effort" while avoiding mentions
 795 ↪ of real institutions. The word combinations should be logical and
 796 ↪ appropriate for an academic context.

797 5. Song Remixes

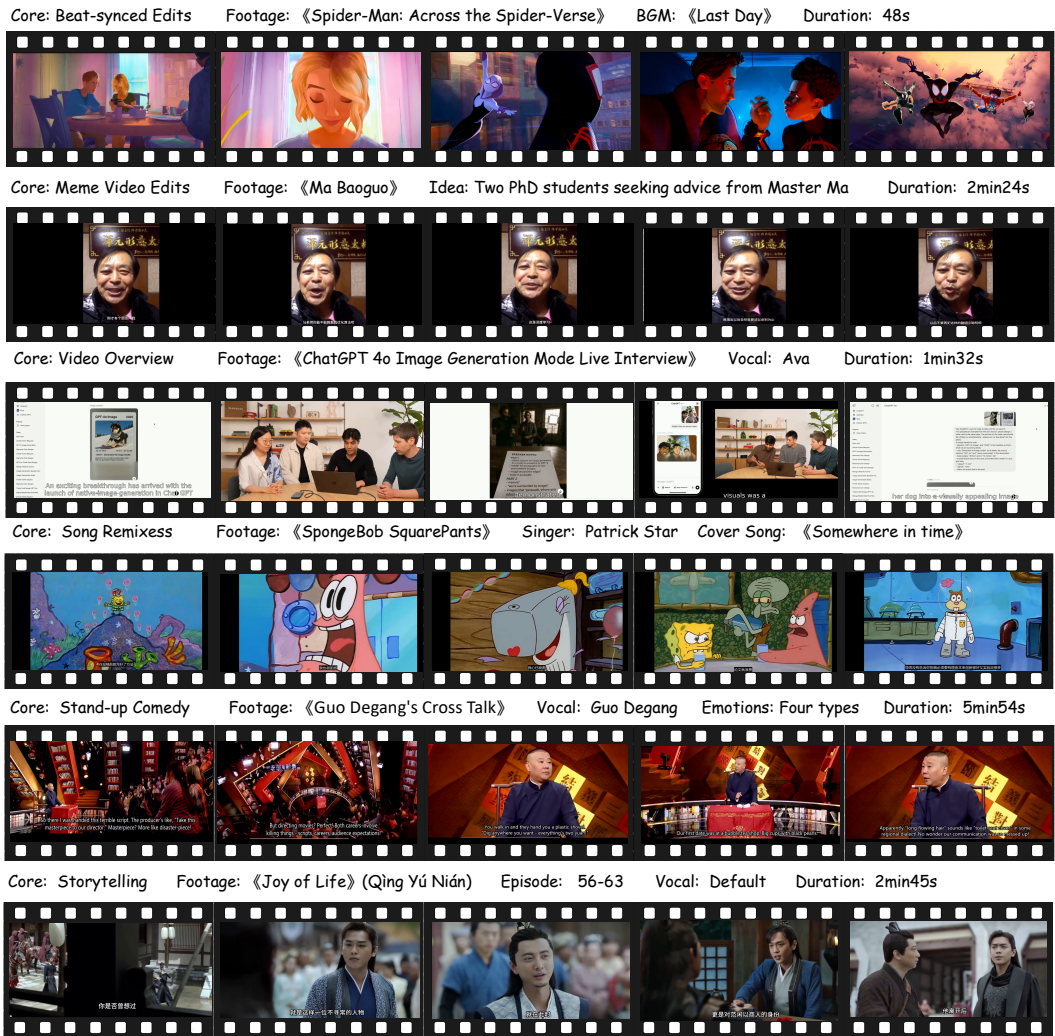
798 Description:

799 - User provides MIDI, lyrics, background music, and a target voice sample.
 800 ↪
 801 - Generate cover, clone voice, align timing, integrate into video
 802 ↪ pipeline.

803 Prompt:

804 - The song is performed by Patrick Star, focusing on the theme of "the
 805 ↪ struggles of manuscript submission and dealing with overly critical
 806 ↪ reviewers", following the original lyrics' sentence structure while
 807 ↪ replacing specific content. It incorporates elements of reviewer nitpicking
 808 ↪ (e.g., questioning innovation, demanding redundant experiments) and
 809 ↪ expresses frustration with lines like "If only I could swap reviewers, this

810 ↪ academic fate is too cruel” to highlight the emotional toll of peer review.
 811 ↪
 812
 813 6. Cross-lingual Adaptations
 814 Description:
 815 - User provides source audio and target voice samples.
 816 - Adapt script to cultural format (e.g., crosstalk/talk show), synthesize
 817 ↪ voices, apply effects, integrate into pipeline.
 818 Prompt:
 819 - Adapting to the cultural context of talk shows and localizing humor



855 Figure 9: Video cases of VideoAgent in real-world scenarios - Case 1
 856
 857
 858
 859
 860
 861
 862
 863

Listing 2: VideoAgent Showcase 2 with Prompts

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

```

Categories:
1. Beat-synced Edits
  Description:
    - User provides background music and film footage.
    - Detects beat/tempo changes and aligns high-energy shots to strong beats.
  ↪
  Prompt:
    - Conflicts between Nezha, Dragon Prince Ao Bing (blue-robed and blue
  ↪ hair) and Shen Gongbao (black-robed).

2. Video Overviews
  Description:
    - User provides a news/event video.
    - Transcribe speech, summarize, align summary sentences to visuals,
  ↪ generate voiceover.
  Prompt:
    - Short movie podcast, colloquial expression within 300 words; identify
  ↪ which actor or host is speaking; do not mention movie-ticket availability.

3. Storytelling
  Description:
    - User provides story/novel text and footage.
    - Script in chosen narration style, match lines to visuals, generate
  ↪ voiceover.
  Prompt:
    - Write a fluent commentary script of 1,500 words.

4.1 Meme Video Edits (Fan Zhiyi)
  Description:
    - User supplies a source video and custom script.
    - Extract/transcribe audio, generate new speech, sync precisely to frames
  ↪ , output dubbed edit.
  Prompt:
    - Emphasize the positive impact of IShowSpeed's China tour, which is not
  ↪ only an online carnival but also a milestone event in cultural
  ↪ communication in the digital age.

4.2 Meme Video Edits (Xiao Ming Live)
  Description:
    - User supplies a source video and custom script.
    - Extract/transcribe audio, generate new speech, sync precisely to frames
  ↪ , output dubbed edit.
  Prompt:
    - Based on the following scenario, create an angry rebuttal from Zhuge
  ↪ Liang (me):
    - Speaker: Zhuge Liang (me)
    - Start with "Why don't you look at your own problems for the failure,"
  ↪ followed by "...look at your own problems" pattern sentences that all
  ↪ reference anime events
    - Anime examples must mention specific characters
    - Only the last "...look at your own problems" should return to the
  ↪ failure scenario
    - Use colloquial language and diverse anime references

5. Cross-lingual Adaptations
  Description:
    - User provides source audio and target voice samples.

```

918 - Adapt script to cultural format (e.g., crosstalk/talk show), synthesize
 919 ↪ voices, apply effects, integrate into pipeline.
 920 Prompt:
 921 - To generate a crosstalk-style piece, the story must conform to
 922 ↪ objective reality and be about forty to fifty sentences.

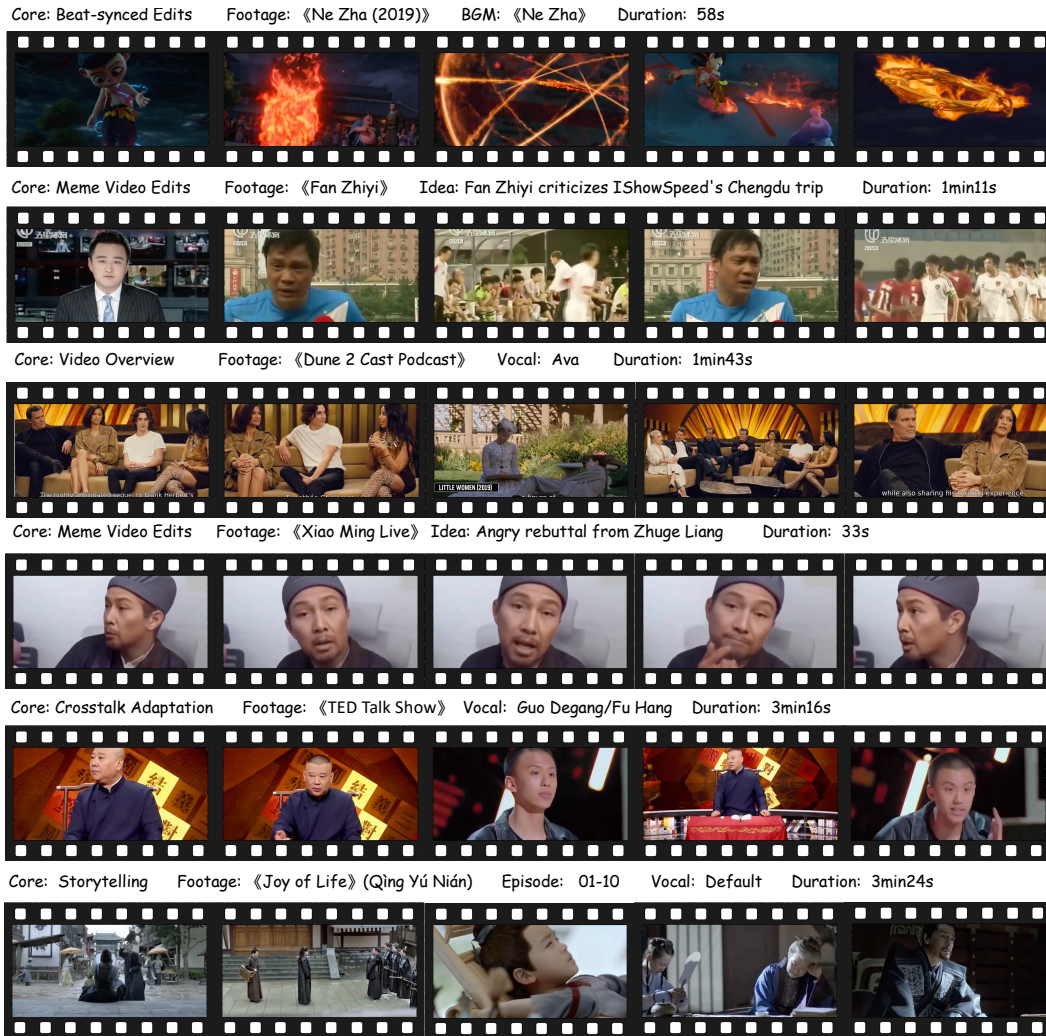


Figure 10: Video cases of VideoAgent in real-world scenarios - Case 2

Listing 3: VideoAgent Showcase 3 with Prompts

962 Categories:
 963 1.1 Beat-synched Edits (Titanic)
 964 Description:
 965 - User provides background music and film footage.
 966 - Detects beat/tempo changes and aligns high-energy shots to strong beats.
 967 ↪ Prompt:
 968 - A romantic and sweet love story about Jack and Rose meeting on the
 969 ↪ Titanic. It cannot include the part where the ship is in distress, nor the
 970 ↪ night scene. In the first section, Rose, wearing a purple hat and a white
 971 ↪ shirt, walks out of a white car with a purple umbrella, looking
 ↪ thoughtfully.

972
 973 1.2 Beat-synced Edits (Interstellar)
 974 Description:
 975 - User provides background music and film footage.
 976 - Detects beat/tempo changes and aligns high-energy shots to strong beats.
 977 ↩
 978 Prompt:
 979 - Celebrate humanity's courage in space exploration. Include scenes
 980 ↩ featuring spaceships, wormholes, black holes, space station docking
 981 ↩ maneuvers, ocean planets, and glacial worlds. Show astronauts in their
 982 ↩ distinctive white spacesuits as they venture into the unknown, highlighting
 983 ↩ humanity's drive to explore the cosmos.

984 1.3 Beat-synced Edits (Interstellar)
 985 Description:
 986 - User provides background music and film footage.
 987 - Detects beat/tempo changes and aligns high-energy shots to strong beats.
 988 ↩
 989 Prompt:
 990 - Love can transcend time and space.

991 1.4 Beat-synced Edits (Ne Zha)
 992 Description:
 993 - User provides background music and film footage.
 994 - Detects beat/tempo changes and aligns high-energy shots to strong beats.
 995 ↩
 996 Prompt:
 997 - Capture more scenes of conflicts and battles between Nezha and Shen
 998 ↩ Gongbao (black-robed) and Dragon Prince Ao Bing (blue-robed).

1000 2.1 Meme Video Edits (Xiao Ming Live)
 1001 Description:
 1002 - User supplies a source video and custom script.
 1003 - Extract/transcribe audio, generate new speech, sync precisely to frames
 1004 ↩ , output dubbed edit.
 1005 Prompt:
 1006 - Background: Mixue Ice Cream is a national chain brand focusing on ice
 1007 ↩ cream and tea beverages. On March 15 (Consumer Rights Day), they were
 1008 ↩ reported to be using overnight lemons. However, compared to other exposures
 1009 ↩ , using overnight lemons is not considered a particularly serious violation
 1010 ↩ and is somewhat understandable.
 1011 - Speaker: Snow King (Mixue's representative)
 1012 ↩ - Purpose: Emphasize that the "overnight lemon" situation is not too
 1013 ↩ serious, highlighting Mixue's good reputation
 1014 - Must preserve the phrases "Look in my eyes tell me why" and "tell me"
 1015 - Must end with the word "say it"
 1016 ↩ - If the original text contains awkward phrasing, such as redundant words
 1017 ↩ or confused semantics, do not imitate that style or sentence structure
 1018 - Ensure natural and fluent sentences

1018 2.2 Meme Video Edits (Xiao Ming Live)
 1019 Description:
 1020 - User supplies a source video and custom script.
 1021 - Extract/transcribe audio, generate new speech, sync precisely to frames
 1022 ↩ , output dubbed edit.
 1023 Prompt:
 1024 - Based on the following scenario, create an angry rebuttal from Zhuge
 1025 ↩ Liang (me):
 - Speaker: Zhuge Liang (me)

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

- Zhuge Liang (me) is challenged about why a certain Three Kingdoms
↪ character has a higher rating than him and launches a fierce rebuttal
- Later rating comparisons should show stark differences (can be
↪ exaggerated)
- Use colloquial language, align with historical facts, and only
↪ replace specific content

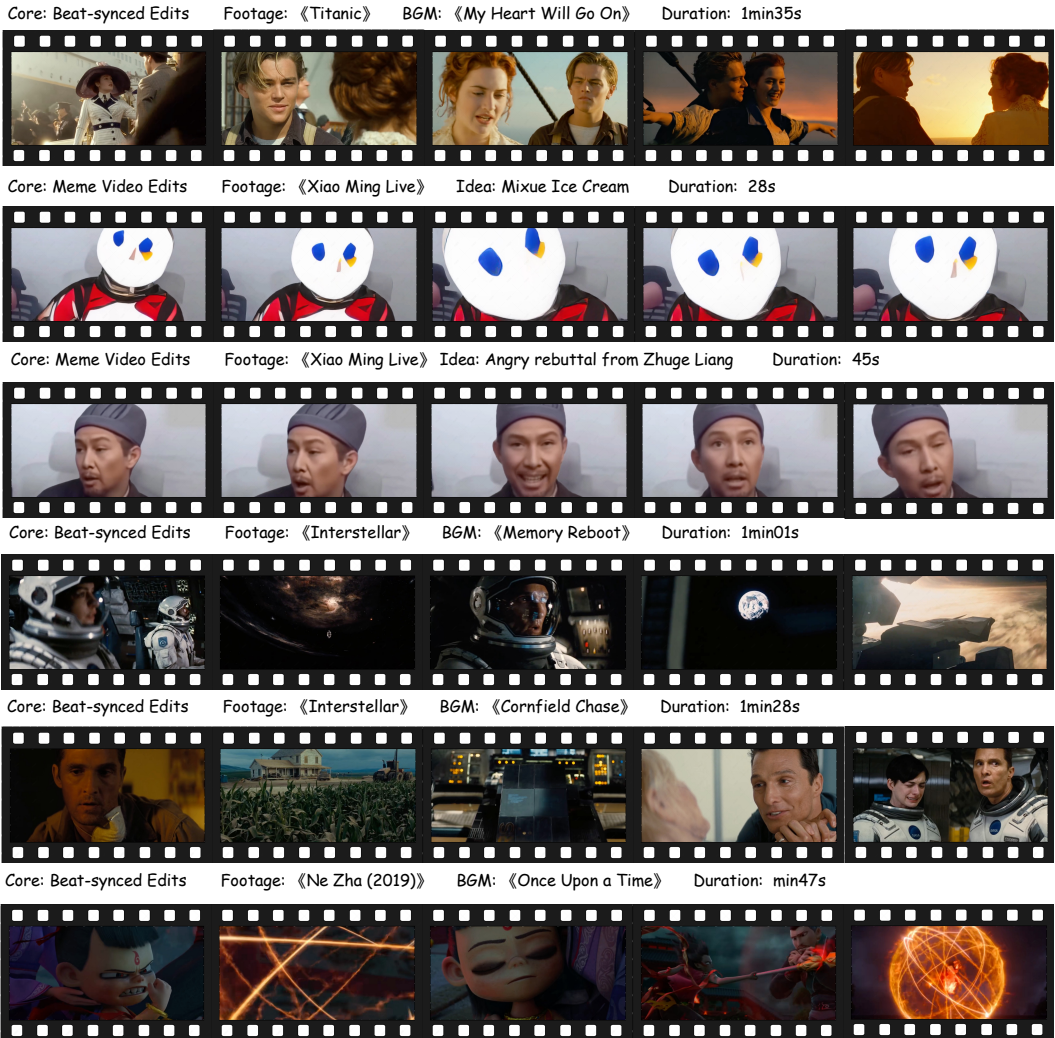


Figure 11: Video cases of VideoAgent in real-world scenarios - Case 3

A.2 EVALUATION DATASETS AND PROTOCOLS

To evaluate VideoAgent in terms of both video understanding and workflow orchestration, we conduct experiments from the following two perspectives.

- **Video Understanding and Retrieval Evaluation.** We utilize the **Shot2Story** dataset to evaluate whether VideoAgent can effectively retrieve appropriate video clips according to user inputs. This dataset contains video clips paired with corresponding queries that describe each clip’s content. We randomly select 100 videos with more than 5 queries each, using the midpoint frames as keyframes for retrieval evaluation. In total, there are 599 individual queries. Note that achieving such global understanding and precise retrieval across multiple video segments is computationally intensive. Processing all queries for a multimodal LLM or agent typically requires around 3-9 hours. The distribution of evaluation video durations in our sampled dataset is shown in Figure 12. Videos are grouped into fixed intervals 10–15s, 15–20s, 20–25s, 25–30s, and 30s+ based on their total segment lengths. The caption queries per video distribution in our dataset is illustrated in Figure 13. Each video in the dataset contains between 5 and 8 captions, providing multiple query perspectives for evaluation. For video categories distribution in our sampled dataset is shown in Figure 14. The dataset exhibits a diverse range of content categories, with Entertainment, Film & Animation, Autos & Vehicles, Sports, News & Politics, Howto & Style, Science & Technology, Comedy, Education, Travel & Events, People & Blogs, Gaming and Pets & Animals of the sampled videos.

We use three metrics to compare retrieved video clips and groundtruth shots: **i) Recall:** This is the percentage of groundtruth videos being retrieved with the first rank for each query. **ii) Embedding Matching Score:** Following Shi et al. (2021); Yang et al. (2024); Cheng et al. (2025), we compute the coarse-grained embedding matching score between retrieved video content and LLM-generated captions (see List 6) using ImageBind Girdhar et al. (2023) cross-modal encoder. This evaluates whether retrieved videos align with user intents at a semantic level. **iii) Intersection over Union:** Following Team et al. (2025), we compute temporal IoU between predicted and ground-truth video segments. IoU measures the overlap between retrieved and target clips, ranging from 0 to 1.

- **Workflow Orchestration Evaluation.** To assess the model’s performance in orchestrating workflows for video editing instructions, we construct the **VideoEdit** dataset containing 2000 creation instructions encompassing diverse visual and audio editing requirements. Starting with authentic video editing demands, such as creating prompts for ChatGPT-4o image generation live stream overviews like "Short movie podcast, colloquial expression within 300 words, notice to identify which actor or host is talking, don’t mention movie tickets available issue", details of prompts can see Appendix A.1. We then use Claude-Sonnet-3.7 for data augmentation to generate additional creative requirements, the audio editing and video editing system prompts are presented in List 4 and List 5. For audio editing, examples include requirements such as "I want to create a funny stand-up comedy audio with audience reactions. The content should be adapted from a script I provide, and I’d like the comedian’s voice to match a specific target vocal style." For video editing, examples include "I have a video that I need to edit with a different script while maintaining the original speaker’s voice. I want to keep the visuals exactly the same, but change what the person is saying to better match my needs. The new content should sound natural and match the original speaker’s voice and speaking style." Through this approach, we yield 2000 distinct instructions covering the full spectrum of video editing tasks.

We employ Claude-Sonnet-3.7 to build a **Judge Agent** that evaluates whether VideoAgent and baseline models successfully orchestrate agent workflows for these complex requirements. This judge agent verifies instruction fulfillment and validates data flow compatibility across agents, ensuring outputs from preceding agents are properly formatted as inputs for subsequent agents. The judge agent demonstrates high alignment with human evaluators in comparative testing (see Section 3.5). This evaluation uses **Success Rate** as the metric, indicating the percentage of samples where the orchestrated workflow fulfills the user instruction.

Listing 4: Audio Editing Instruction Prompt

```
You are an Autonomous Agent System Designer. I will provide Registered Agents
↔ (Name, Description and Parameters information):

Your task is to:
1. Generate Feasible User Requirement:
```

```

1134   - Only consider generating feasible User Requirement for **audio-related**
1135   ⇨ aspects (Finally generate audio instead of video) based on metadata of
1136   ⇨ Registered Agents.
1137   - Avoid mentioning specific agent names in User Requirement, but clearly
1138   ⇨ describe the needs.
1139   - Avoid mentioning the concrete steps and details of implementation.
1140   - Reflect real-world needs, and be phrased conversationally.
1141 2. Design Executable Agent Graph:
1142   - Format: List
1143   - Agent Graph shall contain metadata for each Agent Node including:
1144     * name: (string)
1145     * inputs: (list of input parameter objects with):
1146       - parameter: input parameter name
1147       - description: brief parameter description
1148     * outputs: (list of output parameter objects with):
1149       - parameter: output parameter name
1150       - description: brief parameter description
1151       - links: (list of dictionaries) where each dictionary specifies:
1152         - key of dictionaries: target agent name
1153         - value of dictionaries: target agent's input parameter name that
1154   ⇨ this output connects to
1155   - Agent Graph case:
1156     [
1157       {
1158         "node": Agent Name,
1159         "inputs": [
1160           {
1161             "name": input parameter name1,
1162             "description": ...
1163           },
1164           {
1165             "name": input parameter name2,
1166             "description": ...
1167           },
1168         ],
1169         "outputs": [
1170           {
1171             "name": output parameter name1,
1172             "description": ...,
1173             "links": [
1174               {"next_agent1": next_agent1's input parameter name"},
1175               {"next_agent2": next_agent2's input parameter name"}
1176             ],
1177           },
1178           {
1179             "name": output parameter name2,
1180             "description": ...,
1181             "links": []
1182           },
1183           ...
1184         ],
1185       },
1186       ...
1187     ]
1188 3. Generate Agent Chain:
1189   - Format: List
1190   - Based on the description of the Agent and the sequential information
1191   ⇨ contained in the designed Agent Graph, generate the Agent Chain
1192   - Agent Chain case:

```

```

1188     ["agent1", "agent2", ...]
1189 4. Generate User Input Graph
1190   - Format: List
1191   - Parameter nodes with no in-degree (no incoming edges) are uniformly
1192   ↪ considered to require user input.
1193   - Parameter nodes with no in-degree may have different names but share the
1194   ↪ same user input, meaning a single user input parameter can point to
1195   ↪ multiple such nodes.
1196   - Parameter nodes with no in-degree that are linked to user input should be
1197   ↪ represented in the format **AgentName.input_parameter**
1198   - Generate the User Input Graph based on the Agent descriptions and
1199   ↪ parameter passing information in the designed Agent Graph.
1200   - User Input Graph case:
1201     [
1202       {
1203         "node": User input parameter name,
1204         "description": brief description of the parameter
1205         "links": [
1206           {"agent1": agent1.input_parameter1},
1207           {"agent2": agent2.input_parameter2}
1208         ]
1209       },
1210       ...
1211     ]
1212 5. Output Reasoning:
1213   - Provide concise reasoning (<200 words) explaining the entire workflow
1214   ↪ logic
1215
1216 In addition to the above formatting requirements, please also note the
1217 ↪ following:
1218 1. For each element of **outputs** in each Agent Node:
1219   - Ensure that the **links** in the **outputs** point to an input parameter
1220   ↪ that actually exists in the next Agent Node.
1221   - Also, make sure the description of the output parameter matches the
1222   ↪ description of the input parameter in the next Agent Node.
1223 2. Final JSON Output Format Specification:
1224 {
1225   "User Requirement": ...,
1226   "Agent Graph": ...,
1227   "Agent Chain": ...,
1228   "User Input Graph": ...,
1229   "Reasoning": ...
1230 }
1231 Strictly follow JSON output format!
1232
1233 Metadata of registered agents:
1234 {registry_agents}
1235
1236 Real-life Requirement Examples:
1237 {real-life_examples}

```

Listing 5: Video Editing Instruction Prompt

```

1238 You are an Autonomous Agent System Designer. I will provide Registered Agents
1239 ↪ (Name, Description and Parameters information):
1240
1241 Your task is to:
1242 1. Generate Feasible User Requirement:

```

1242 - Only consider generating feasible User Requirement for ****video-related****
1243 ↪ aspects based on metadata of Registered Agents.
1244 - Avoid mentioning specific agent names in User Requirement, but clearly
1245 ↪ describe the needs.
1246 - Avoid mentioning the concrete steps and details of implementation.
1247 - Reflect real-world needs, and be phrased conversationally.

1248 2. Design Executable Agent Graph:
1249 - Format: List
1250 - Agent Graph shall contain metadata for each Agent Node including:
1251 * name: (string)
1252 * inputs: (list of input parameter objects with):
1253 - parameter: input parameter name
1254 - description: brief parameter description
1255 * outputs: (list of output parameter objects with):
1256 - parameter: output parameter name
1257 - description: brief parameter description
1258 - links: (list of dictionaries) where each dictionary specifies:
1259 - key of dictionaries: target agent name
1259 - value of dictionaries: target agent's input parameter name that
1260 ↪ this output connects to
1261 - Agent Graph case:
1262 [

```

1263     {
1264       "node": Agent Name,
1265       "inputs": [
1266         {
1267           "name": input parameter name1,
1268           "description": ...
1269         },
1270         {
1271           "name": input parameter name2,
1272           "description": ...
1273         }
1274       ],
1275       "outputs": [
1276         {
1277           "name": output parameter name1,
1278           "description": ...,
1279           "links": [
1280             {"next_agent1": next_agent1's input parameter name"},
1281             {"next_agent2": next_agent2's input parameter name"}
1282           ]
1283         },
1284         {
1285           "name": output parameter name2,
1286           "description": ...,
1287           "links": []
1288         },
1289         ...
1290       ]
1291     }

```

1291 3. Generate Agent Chain:
1292 - Format: List
1293 - Based on the description of the Agent and the sequential information
1294 ↪ contained in the designed Agent Graph, generate the Agent Chain
1295 - Agent Chain case:
["agent1", "agent2", ...]

```

1296 4. Generate User Input Graph
1297   - Format: List
1298   - Parameter nodes with no in-degree (no incoming edges) are uniformly
1299   ↪ considered to require user input.
1300   - Parameter nodes with no in-degree may have different names but share the
1301   ↪ same user input, meaning a single user input parameter can point to
1302   ↪ multiple such nodes.
1303   - Parameter nodes with no in-degree that are linked to user input should be
1304   ↪ represented in the format **AgentName.input_parameter**
1305   - Generate the User Input Graph based on the Agent descriptions and
1306   ↪ parameter passing information in the designed Agent Graph.
1307   - User Input Graph case:
1308     [
1309       {
1310         "node": User input parameter name,
1311         "description": brief description of the parameter
1312         "links": [
1313           {"agent1": agent1.input_parameter1},
1314           {"agent2": agent2.input_parameter2}
1315         ]
1316       },
1317       ...
1318     ]
1319 5. Output Reasoning:
1320   - Provide concise reasoning (<200 words) explaining the entire workflow
1321   ↪ logic
1322
1323 In addition to the above formatting requirements, please also note the
1324 ↪ following:
1325 1. For each element of **outputs** in each Agent Node:
1326   - Ensure that the **links** in the **outputs** point to an input parameter
1327   ↪ that actually exists in the next Agent Node.
1328   - Also, make sure the description of the output parameter matches the
1329   ↪ description of the input parameter in the next Agent Node.
1330 2. Final JSON Output Format Specification:
1331 {
1332   "User Requirement": ...,
1333   "Agent Graph": ...,
1334   "Agent Chain": ...,
1335   "User Input Graph": ...,
1336   "Reasoning": ...
1337 }
1338 Strictly follow JSON output format!
1339
1340 Metadata of registered agents:
1341 {registry_agents}
1342
1343 Real-life Requirement Examples:
1344 {real-life_examples}

```

Listing 6: Video Caption to High-Level Summary Transformation Prompt

```

1344 prompt = f"""
1345 Transform the following detailed video caption into a concise, high-level
1346 ↪ summary that captures the main themes, setting, and overall atmosphere in
1347 ↪ no more than 50 words, no point form, use pure sentence.
1348
1349 Focus on:
1350 - Setting identification: Where does the action take place?

```

- Key themes: What are the main activities or purposes being shown?
- Atmosphere and tone: What mood or feeling does the scene convey?
- Overall narrative: What story or message emerges?

Instructions:

1. Condense specific details into broader concepts
2. Identify recurring elements that suggest themes
3. Synthesize individual actions into a cohesive narrative about the environment
4. Use descriptive language that captures essence rather than listing events
5. Maintain the core message while abstracting from granular details
6. Preserve emotional tone and workplace/environmental culture
7. Keep response under 50 words

Input format: whole_caption"[detailed description]"

Output format: high_level_caption"[concise thematic summary]"

Rules:

- Extract thematic elements from specific actions
- Focus on overarching narrative rather than individual events
- Capture the purpose
- Ensure output maintains fidelity to original content intent

Transform the provided detailed caption following these guidelines." ""

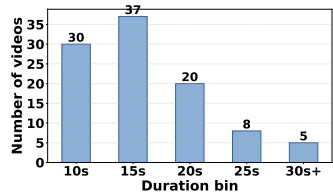


Figure 12: Evaluation video durations distribution.

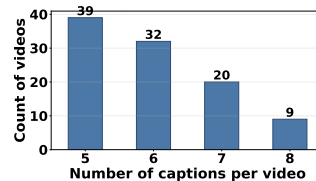


Figure 13: Evaluation video captions/queries distribution.

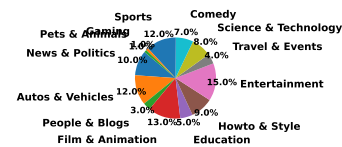


Figure 14: Evaluation video categories distribution.

A.3 BASELINE METHOD

To ensure the comprehensiveness of the benchmark, we conducted a multi-dimensional comparison with n baselines encompassing different mechanisms.

Agentic Systems on Different Domains.

- **Chain-of-Thought** Wei et al. (2022): This method proposes a chain of thought to improve LLMs' ability to perform complex reasoning by integrating explicit, structured intermediate steps.
- **LLM Debate** Du et al. (2023): It thoughtfully and systematically integrates debate results from multiple independent language model instances across several iterative rounds of refinement.
- **Step Back** Zheng et al. (2023): This method enables LLMs to abstract from instances containing specific details to generate high-level concepts and foundational first principles efficiently.
- **ExpertPrompting** Xu et al. (2023a): It leverages the potential of LLMs by delegating complex tasks to domain-specific expert agents and automatically synthesizing their respective opinions.
- **Intelligent Go-Explore** Lu et al. (2024): This method is designed to replace handcrafted heuristics with intelligence, iteratively returning to observed states and searching for optimal solutions.
- **Flow** Niu et al. (2025): This framework introduces modular, parallel workflows to boost efficiency and adaptability, enabling real-time updates that preserve global coherence during challenges.

MLLM and Agents.

- **GPT-4o** OpenAI (2024): GPT-4o is a unified multimodal model for text, audio, image, and video, with near-human latency, Turbo-level text/code, better non-English, faster, and stronger vision.
- **Deepseek-v3** Liu et al. (2024): This large language model excels in code generation, mathematical reasoning, and multilingual understanding, combining high performance with efficient inference for a wide range of technical and general-purpose applications.
- **Gemini-2.5** Gemini Team (2025): Gemini 2.X spans Pro and Flash models with strong multimodal reasoning; 2.5 Pro offers SoTA coding/reasoning and 3-hour video context, while 2.5/2.0 Flash and Flash-Lite deliver high performance at much lower latency/cost, enabling agentic workflows.
- **Claude-Sonnet-3.7** Anthropic (2025a): This efficient multimodal language model features hybrid reasoning and extended thinking capabilities, optimized for coding, customer service, and general-purpose tasks with strong performance at low latency.
- **Claude-Sonnet-4** Anthropic (2025b): This high-performance language model supports up to 1 million tokens of context, excelling in complex reasoning, code generation, and long-context analysis, making it suitable for large-scale document understanding and planning tasks.
- **Qwen2.5-VL** Bai et al. (2025): This vision-language model is designed for long-video comprehension spanning hours, with precise second-level event localization and multimodal reasoning.
- **VideoRAG** Ren et al. (2025): It introduces a RAG framework for long-context video understanding with graph-grounded and multimodal encoding, enabling unlimited retrieval and beating SOTA.
- **VideoMind** Liu et al. (2025): This method offers a role-based workflow for temporally grounded video understanding, using Chain-of-LoRA for efficient role-switching via lightweight adapters.

Video Generation Agents.

- **NoteBookLM** Google Labs (2023): NotebookLM is a personalized AI knowledge assistant that leverages advanced LLMs to analyze and interact with user-provided documents collaboratively.
- **Director** VideoDB (2024): This framework is capable of reasoning through video-related tasks while simultaneously streaming the processed results in real-time to facilitate video synthesis.
- **FunClip** ModelScope (2024): FunClip is an open-source tool for automated video clipping. Users can quickly extract video segments by selecting recognized text or speaker segments.
- **NarratoAI** linyqh (2025): NarratoAI is a video narration platform, combining script writing, editing, voice-over generation, and subtitles into a unified workflow for content production.

A.4 IMPLEMENTATION DETAILS

This section will provide detailed descriptions of the implementation details and parameter configurations for each baseline in different experiments.

All baseline methods are configured using their default or recommended parameter settings and implemented following the original papers. All experiments were conducted on NVIDIA RTX 3090 GPUs. **Workflow Orchestration Evaluation:** In the workflow orchestration performance comparison and workflow orchestration ablation study, for fairness, we use Claude-Sonnet-3.7 as the Judge Agent to evaluate Agent Graphs and Agent Chains from multiple perspectives, including execution sequence validation, parameter routing correctness, agent functional redundancy assessment, and requirement fulfillment analysis. We also set the number of reflection rounds to 2 for all self-reflective methods. In the LLM-human judgment consistency experiment, to reduce the error rate of human judgment, we first execute agent graph checks for execution sequence validation and parameter routing correctness issues, then proceed with agent functional redundancy assessment and requirement fulfillment analysis. **Video Editing Evaluation:** All video inputs are segmented into three-second segments, under fps1, one frame is sampled per second. For recall calculations, the midpoint timestamp of each retrieved video segments serves as the reference point. Subsequently, two-second clips are extracted, one second before and one second after the midpoint timestamp, yielding a total duration of two seconds per retrieved clip. These extracted clips are then concatenated sequentially to generate the final output video for computing the Embedding Matching score and Intersection over Union metrics, all ours VLM agent/caption use MiniCPM-V 2.6 int4 Yao et al. (2024). In the video editing ablation study and hyperparameter experiments, all our base models utilize GPT-4o, and continue to keep the video segments divided into three seconds for input.

A.5 TOOLUSE AGENTS

To accommodate diverse user needs, we have designed a variety of multimodal tool-use agents capable of handling tasks such as audio preprocessing, video overview, storytelling, beat-synced edits, meme video remaking, song remixes, cross-lingual adaptations as well as video edits. The agent names and detailed descriptions of these tools are presented in Table 4.

Table 4: The detailed information of tooluse agents.

Agent Name	Type	Description
AudioExtractor	Audio Preprocessing	To extract audio from a single video or all videos in a directory
LoudnessNormalizer	Audio Preprocessing	Audio loudness normalization tool
Merge	Audio Preprocessing	To merge video and audio tracks
Mixer	Audio Preprocessing	To mix audio with the BGM
Resampler	Audio Preprocessing	Audio Resampling tool
Separator	Audio Preprocessing	Audio source separation tool
VoiceGenerator	Speech Generation	To generate speech based on scene content.
CrossTalkAdapter	Cross-lingual Adaptations	Adapt a reference script into segmented cross talk.
CrossTalkSynth	Cross-lingual Adaptations	Segment-by-segment cross talk audio synthesis with final merge.
StandUpAdapter	Cross-lingual Adaptations	Adapt a reference script into segmented stand-up comedy.
StandUpSynth	Cross-lingual Adaptations	Segment-by-segment stand-up comedy audio synthesis with final merge.
SVCAdapter	Song Remixess	Adapt the original lyrics.
SVCAnalyst	Song Remixess	Analyze the original song’s MIDI file to extract information such as notes, note durations, etc.
SVCAnalyst	Song Remixess	Analyze the original song’s MIDI file to extract information such as notes, note durations, etc.
SVCCoverist	Song Remixess	Source-to-target voice timbre cloning for singing voice synthesis.
SVCSingle	Song Remixess	Split the adapted lyrics into segments, then synthesize each segment with the default vocal singing voice.
TTSInfer	Meme Video	Take the sliced audio clips as target voice references, then conduct Text-To-Speech synthesis using their rewritten text segments.
TTSReplace	Meme Video	Replace audio of the original video with derivative audio segments from the sliced clips.
TTSSlicer	Meme Video	Slice the audio into segments. The audio segments often need to be transcribed afterward.
TTSWriter	Meme Video	Rewrite the transcript of the sliced audio segments based on user requirements.
CommentaryContent Generator	Storytelling	To generates commentary content based on user ideas and text source materials with specialized formatting for video presentations.
NewsContentGenerator	Video Overview	To generate news summary based on user ideas and reference materials.

Agent Name	Type	Description
RhythmDetector	Beat-synced Edits	To create cut points for video editing based on music rhythms.
RhythmDetector	Beat-synced Edits	To create cut points for video editing based on music rhythms.
RhythmContent Generator	Beat-synced Edits	To extract video segment content, create scene-focused narrative summaries, and generate rhythm-aware storyboards.
VideoConversion	Video Edits	converts audio content with JSON timestamps into visual scene descriptions for video generation.
VideoEditor	Video Edits	Ultimately merging the clips and adding audio.
VideoPreloader	Video Edits	To initialize environment and preprocess video files.
VideoSearcher	Video Edits	To retrieve matching video clips based on timestamp file.
VideoQA	Video Interaction	Transcribes all videos in a directory and provides interactive Q&A session.
VideoSummarization Generator	Video Interaction	Agent that generates summarization content based on user ideas and reference materials.
FaceSwapping	Video Processing	Replace the faces appearing in the video with those of the target person.
LipSynchronization	Video Processing	Synchronize the speaker's lip movements in the video with the audio track.

A.6 VISUAL PROCESSING AGENTS

Visual processing agents are specialized components that handle video content analysis, scene caption, and visual element manipulation for multimedia production workflows. These agents bridge the gap between textual content and visual representation by converting visual scenes into actionable narrative descriptions, managing video file operations, and coordinating the integration of visual and audio elements.

A.6.1 VIDEOPRELOADER

The VideoPreloader serves as the foundational initialization agent for video processing pipelines, establishing the necessary directory structure and preprocessing operations required for video editing workflows. It integrates with VideoRAG systems to enable efficient video content indexing and retrieval.

Listing 7: VideoPreloader Structure

```

Input: video_directory
Output: preprocessing_status

Algorithm:
1. System initialization:
  a. Configure multiprocessing with spawn method
  b. Set up logging and suppress warnings
  c. Initialize project root and tools directory paths
  d. Add necessary modules to system path

2. Directory structure setup:

```

```

1566     a. Create dataset base directory structure
1567     b. Initialize video_edit working directory
1568     c. Create required subdirectories:
1569         - audio_analysis: for audio processing outputs
1570         - scene_output: for scene analysis results
1571         - videosource-workdir: for VideoRAG operations
1572         - writing_data: for text and script storage
1573         - video_output: for final video products
1574     d. Set working directory for VideoRAG operations
1575
1576 3. Video source validation:
1577     a. Verify video source directory existence
1578     b. Check for MP4 video files in source directory
1579     c. Generate list of valid video file paths
1580     d. Report video discovery statistics
1581
1582 4. VideoRAG module loading:
1583     a. Dynamically import VideoRAG and QueryParam classes
1584     b. Handle import errors with fallback mechanisms
1585     c. Initialize VideoRAG content processing system
1586
1587 5. Video preprocessing pipeline:
1588     a. Initialize VideoRAG with working directory
1589     b. Execute video insertion process:
1590         - Extract video metadata and features
1591         - Generate content embeddings for retrieval
1592         - Create searchable video index database
1593         - Store preprocessed data in working directory
1594     c. Validate preprocessing completion
1595
1596 6. Status reporting and cleanup:
1597     a. Generate processing statistics
1598     b. Report successful video preprocessing
1599     c. Return execution status for pipeline coordination
1600     d. Prepare for subsequent VideoSearcher operations
1601
1602 Error Handling:
1603     - Directory creation failures with permission checks
1604     - VideoRAG import failures with module path resolution
1605     - Video file access errors with format validation
1606     - Processing interruption with graceful cleanup

```

1605 A.6.2 VIDEOSearcher

1607 The VideoSearcher retrieves matching video clips from preprocessed video databases based on scene
1608 semantic descriptions. It leverages VideoRAG content indexing to perform intelligent video segment
1609 matching for automated video assembly workflows.

1610 Listing 8: VideoSearcher Structure

```

1612 Input: video_scene_path
1613 Output: search_status
1614
1615 Algorithm:
1616 1. Path and directory setup:
1617     a. Initialize dataset and video_edit directory paths
1618     b. Set scene_output and working directory locations
1619     c. Create necessary directories if not existing
1620     d. Add tools directory to system path

```

```

1620
1621 2. VideoRAG module initialization:
1622   a. Dynamically import VideoRAG and QueryParam classes
1623   b. Handle import failures with error logging
1624   c. Configure VideoRAG content processing system
1625
1626 3. Scene file processing:
1627   a. Load JSON scene file with UTF-8 encoding
1628   b. Extract segment_scene content for queries
1629   c. Validate scene content availability
1630   d. Prepare query string for VideoRAG processing
1631
1632 4. Query parameter configuration:
1633   a. Initialize QueryParam with videoragcontent mode
1634   b. Configure reference inclusion settings:
1635     - wo_reference = False: include video clip references
1636     - wo_reference = True: exclude references from response
1637   c. Set multiprocessing spawn method for compatibility
1638
1639 5. VideoRAG content search:
1640   a. Initialize VideoRAG with working directory
1641   b. Execute semantic query against video database:
1642     - Process scene descriptions into embeddings
1643     - Perform similarity matching against indexed videos
1644     - Retrieve relevant video segments with timestamps
1645   c. Generate search response with matching clips
1646
1647 6. Result processing and validation:
1648   a. Validate VideoRAG query completion
1649   b. Log search statistics and performance metrics
1650   c. Handle search failures with appropriate error messages
1651   d. Return structured search status
1652
1653 Error Handling:
1654 - JSON file not found or invalid format errors
1655 - VideoRAG import and initialization failures
1656 - Empty or malformed scene content handling
1657 - Query processing exceptions with detailed logging

```

1657 A.6.3 VIDEOEDITOR

1658 The VideoEditor performs final video assembly by integrating matched video segments with synchro-
1659 nized audio tracks. It utilizes multimodal analysis to optimize clip selection and ensures temporal
1660 alignment between visual content and audio/storyboards narratives.
1661

1662 Listing 9: VideoEditor Structure

```

1663 Input: video_directory, audio_path, timestamp_path
1664 Output: final_video_path
1665
1666 Algorithm:
1667 1. Initialization and data loading:
1668   a. Configure directory structure and model paths
1669   b. Load MiniCPM-V-2_6-int4 model for multimodal analysis
1670   c. Parse visual_retrieved_segments.json for video segments
1671   d. Load kv_store_video_segments.json for timing metadata
1672   e. Initialize beat synchronization and storyboard processing
1673
1674 2. Time period generation:

```

```

1674     a. Parse rhythm_points.json for beat timestamps
1675     b. Create time periods from consecutive beat intervals
1676     c. Load video_scene.json for storyboard descriptions
1677     d. Align time periods with narrative segments
1678
1679 3. Frame extraction and preprocessing:
1680   For each video segment:
1681     a. Extract frames at 1fps starting from segment start time
1682     b. Include exact start timestamp plus subsequent whole seconds
1683     c. Convert frames to RGB PIL Images (224x224 resolution)
1684     d. Handle RGBA to RGB conversion for consistent processing
1685
1686 4. Multimodal scene matching with VLM prompt and frame insertion:
1687   System context: Video segment analysis for optimal clip selection
1688
1689   Message structure construction:
1690   msgs = [{'role': 'user', 'content': [prompt_text]}]
1691   For each extracted frame:
1692     msgs[0]['content'].append(frame_image)
1693
1694   User prompt template embedded with frames:
1695   "You are analyzing a video segment from {start_time:.3f}s to {end_time:.3f}
1696   ↪ s.
1697   You see {frame_count} consecutive video frames, each representing 1 second.
1698   The required clip duration is {duration:.3f} seconds.
1699   Find the best sequence matching this description: '{description}'
1700
1701   Requirements:
1702   1. Analyze ALL {frame_count} frames for best consecutive sequence
1703   2. Choose starting frame (0-{max_frame}) allowing {duration:.3f}s clip
1704   3. Maximum starting frame: {max_start} to fit duration constraints
1705   4. Return ONLY single number - starting frame number
1706   5. Select frames with optimal visual quality and scene consistency
1707   6. Prioritize alignment with scene description content"
1708
1709   Frame insertion process:
1710   - Text prompt inserted as first content element
1711   - Sequential frame images appended to content array
1712   - Model processes text prompt with visual frame sequence
1713   - VLM analyzes prompt context alongside embedded frame data
1714
1715 5. Clip extraction and temporal alignment:
1716   a. Parse VLM response to extract optimal starting frame
1717   b. Calculate precise clip timing: start_time + frame_offset
1718   c. Validate clip boundaries within segment constraints
1719   d. Extract video clips with exact duration matching beat intervals
1720   e. Maintain audio track based on keep_original_audio parameter
1721
1722 6. Audio integration pipeline:
1723   a. Load background audio and resize to match video duration
1724   b. Configure audio mixing strategy:
1725     - keep_original_audio=True: composite with background at mix_ratio
1726     - keep_original_audio=False: replace with background music only
1727   c. Apply volume adjustments and create composite audio tracks
1728   d. Synchronize audio timeline with video concatenation
1729
1730 7. Final video assembly and export:
1731   a. Concatenate video clips using compose method for seamless transitions
1732   b. Apply composite audio track with proper codec configuration

```

- 1728 c. Export with optimized encoding parameters:
- 1729 - fps=24, codec=libx264, audio_codec=aac
- 1730 - preset=medium, threads=4 for performance optimization
- 1731 d. Clean up temporary resources and close video handles

- 1732
- 1733 Error handling and fallback mechanisms:
- 1734 - Frame extraction failures with segment boundary validation
- 1735 - VLM response parsing with regex-based number extraction
- 1736 - Video file access errors with path verification
- 1737 - Audio synchronization failures with default timeline alignment
- 1738 - Memory management with proper resource cleanup

1739
1740
1741

1742 A.6.4 FACESWAPPING

1743
1744 The FaceSwapping agent is capable of replacing a specified face in a video with that of a target
1745 person. It uses the Viggie AI tool to achieve precise face replacement, enhancing the overall viewing
1746 experience of the video.

1747
1748

Listing 10: FaceSwapping Agent Structure

1749 Input: source_video_path, target_face_image
1750 Output: swapped_video_path
1751

1752 Algorithm:

- 1753 1. Initialization:
 - 1754 a. Set up file paths and working directories
 - 1755 b. Initialize access to Viggie AI face swapping tool
- 1756 2. Source video processing:
 - 1757 a. Load source video and extract frames
 - 1758 b. Detect and track the specified face across frames
 - 1759 c. Prepare frames containing the target face for swapping
- 1760 3. Face swapping operation:

1761 For each frame with detected target face:

 - 1762 a. Send frame and target_face_image to Viggie AI tool
 - 1763 b. Receive the frame with the specified face replaced by the target face
 - 1764 c. Replace the original frame with the swapped output
- 1765 4. Video reconstruction:
 - 1766 a. Reassemble processed frames preserving original video properties
 - 1767 b. Synchronize original audio track if available
 - 1768 c. Output the final face-swapped video file
- 1769 5. Error handling and cleanup:
 - 1770 a. Handle failures in calling Viggie AI or frame processing gracefully
 - 1771 b. Log issues and fall back to original frames if needed
 - 1772 c. Release resources and temporary files

1773 Notes:

- 1774 - Input video and target face image are required inputs
- 1775 - Output is a video with the specified face replaced by the target
- 1776 - Reliant on Viggie AI for core face swapping capabilities
- 1777 - Designed for seamless integration within video pipelines

1778
1779
1780
1781

1782 A.6.5 LIPSYNCHRONIZATION

1783

1784

1785

1786

1787

1788

1789

1790

1791

1792

1793

1794

1795

1796

1797

1798

1799

1800

1801

1802

1803

1804

1805

1806

1807

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

1819

1820

1821

1822

1823

1824

1825

1826

1827

1828

1829

1830

1831

1832

1833

1834

1835

The LipSynchronization agent is capable of synchronizing the speaker’s lip movements in a video with a target audio track. It uses the Kling AI tool to achieve precise lip-sync alignment, enhancing the realism and coherence of the video’s audio-visual experience.

Listing 11: LipSynchronization Agent Structure

```

Input: source_video_path, target_audio_path
Output: synced_video_path

Algorithm:
1. Initialization:
  a. Configure input/output file paths and working directories
  b. Initialize connection and access to Kling AI lip synchronization tool

2. Data preparation:
  a. Load source video and extract frames and original audio if present
  b. Load target audio track for synchronization

3. Lip synchronization process:
  a. Send source video frames and target audio to Kling AI tool
  b. Receive processed video frames with speaker's lip movements synchronized
  ↪ to target audio
  c. Replace original frames with synchronized frames

4. Video reconstruction:
  a. Combine synchronized frames preserving original video resolution and
  ↪ frame rate
  b. Integrate the target audio track as the output video's audio
  c. Output final lip-synced video file

5. Error handling and cleanup:
  a. Handle communication errors with Kling AI tool gracefully
  b. Provide fallback to original video if synchronization fails
  c. Log all processing steps and encountered issues
  d. Release temporary resources and close file streams

Notes:
- Supports video and separate audio inputs
- Output video shows visually consistent lip motions aligned with target audio
- Relies fully on Kling AI external service for synchronization
- Designed for smooth integration with video editing workflows

```

1827 A.7 AUDIO PROCESSING AGENTS

1828

1829

1830

1831

1832

1833

1834

1835

Audio processing agents constitute a fundamental component of the multimedia production pipeline, handling diverse audio manipulation tasks. These agents leverage advanced signal processing techniques and machine learning models Du et al. (2024); Liao et al. (2024); Liu (2024); Radford et al. (2022); McFee et al. (2015); Liu et al. (2021) to ensure high-quality audio output while maintaining compatibility across different formats and sampling rates. The audio processing workflow encompasses both basic preprocessing operations and sophisticated content analysis capabilities, with agents designed to handle everything from simple format conversions to complex multi-track audio synthesis.

1836 A.7.1 AUDIOEXTRACTOR
1837

1838 The AudioExtractor agent provides comprehensive audio extraction capabilities from video sources,
1839 supporting both single file and batch directory processing. It utilizes FFmpeg for high-quality audio
1840 conversion while maintaining consistent output format specifications across all processed files.

1841 Listing 12: AudioExtractor Structure
1842

```

1843 Input: video_path (file or directory)
1844 Output: audio_paths, data_directory
1845
1846 Algorithm:
1847 1. Input validation and path analysis:
1848     a. Validate video_path existence and accessibility
1849     b. Determine input type (single file vs directory)
1850     c. Initialize supported video format extensions:
1851         { .mp4, .avi, .mov, .mkv, .wmv, .flv, .webm, .m4v }
1852
1853 2. Video file discovery:
1854     For directory input:
1855     a. Scan directory for video files with supported extensions
1856     b. Filter files by extension validation (case-insensitive)
1857     c. Generate sorted list of video file paths
1858     d. Validate file accessibility and integrity
1859
1860 3. Audio extraction process:
1861     For each video file:
1862     a. Generate output audio path: video_name.wav
1863     b. Configure FFmpeg command with parameters:
1864         - Input: source video file
1865         - Video processing: disabled (-vn flag)
1866         - Audio codec: PCM 16-bit LE (pcm_s16le)
1867         - Sample rate: 44.1 kHz (-ar 44100)
1868         - Channels: stereo (-ac 2)
1869         - Error handling: error-level logging only
1870
1871 4. FFmpeg execution pipeline:
1872     Command structure:
1873     ffmpeg -y -i input_video -vn -acodec pcm_s16le
1874     -ar 44100 -ac 2 -loglevel error output_audio.wav
1875
1876     a. Execute subprocess with error handling
1877     b. Monitor conversion progress and status
1878     c. Handle CalledProcessError for conversion failures
1879     d. Detect FFmpeg availability with FileNotFoundError handling
1880
1881 5. Batch processing workflow:
1882     a. Process each video file sequentially
1883     b. Track successful vs failed extractions
1884     c. Generate progress reporting for each file
1885     d. Collect all successful audio file paths
1886     e. Report final statistics: success_count/total_count
1887
1888 6. Output generation and validation:
1889     a. Return audio_paths as list (batch) or string (single)
1890     b. Provide data_dir as containing directory path
1891     c. Handle empty results for failed extractions
1892     d. Ensure output path consistency and accessibility
1893
1894 Error Handling:

```

- 1890 - FFmpeg availability detection and user guidance
- 1891 - Individual file conversion error isolation
- 1892 - Directory access permission validation
- 1893 - Graceful degradation for partial batch failures
- 1894 - Comprehensive error reporting with specific failure reasons

1895
1896

1897 A.7.2 RHYTHMDETECTOR

1898 The RhythmDetector agent provides comprehensive music rhythm analysis for rhythm-cut video
1899 creation workflows. It employs advanced signal processing techniques using librosa and scipy to
1900 detect beat patterns, analyze temporal distributions, and generate precise cut points for video editing
1901 synchronization with musical rhythms.

1902
1903

Listing 13: RhythmDetector Structure

1904
1905
1906

Input: audio_file_path
Output: rhythm_analysis_directory

1907

Algorithm:

1908

1. Audio loading and preprocessing:

1909

- a. Load audio file using librosa.load() with original sample rate
- b. Store audio data, sample rate, and extract base filename
- c. Configure analysis parameters:
 - frame_length=2048, hop_length=512 for STFT analysis
 - Initialize analysis state variables for processing

1910

1911

1912

1913

1914

1915

2. RMS energy calculation and normalization:

1916

- a. Calculate RMS energy using librosa.feature.rms():
 - Apply frame_length and hop_length for temporal resolution
 - Extract energy envelope from audio signal
- b. Normalize RMS values: rms_normalized = rms / max(rms)
- c. Apply smoothing with convolution kernel:
 - kernel = ones(smoothing_window) / smoothing_window
 - Smooth energy curve to reduce noise artifacts

1917

1918

1919

1920

1921

1922

1923

3. Peak detection with temporal constraints:

1924

Function _detect_rhythm_points():

1925

- a. Apply scipy.signal.find_peaks() with parameters:
 - height=energy_threshold (default 0.4) for minimum peak amplitude
 - distance=min_samples_interval for minimum beat separation
- b. Convert peak indices to timestamps using librosa.frames_to_time()
- c. Apply masking for unwanted detection regions (e.g., intro/outro)

1926

1927

1928

1929

1930

4. Mask-based rhythm filtering:

1931

1932

1933

1934

1935

- a. Define mask_ranges for temporal exclusion zones
- b. Filter detected timestamps within masked periods:
 - Check each timestamp against mask start/end boundaries
 - Separate filtered and masked timestamps for reporting
- c. Generate rhythm points with sequential ID and precise timing

1936

1937

5. Comprehensive visualization generation:

1938

Function _plot_rhythm_detection():

1939

- a. Create three-panel visualization:
 - Waveform with rhythm point markers and masked regions
 - RMS energy curve with threshold line and detection points
 - Spectrogram with frequency-domain rhythm visualization
- b. Apply color coding: red for rhythm points, gray for masked areas
- c. Export high-resolution plots (300 DPI) for documentation

1940

1941

1942

1943

```

1944 6. Rhythm interval analysis and statistics:
1945   Function _analyze_rhythm_distribution():
1946   a. Calculate inter-beat intervals using np.diff(timestamps)
1947   b. Generate statistical measures:
1948       - mean, median, min, max interval durations
1949       - standard deviation for rhythm consistency analysis
1950   c. Create distribution visualizations:
1951       - Histogram of interval frequencies
1952       - Temporal plot showing interval variations over time
1953
1954 7. Output generation and data export:
1955   a. Structure rhythm data in JSON format:
1956       beat_data: {
1957           count: total_rhythm_points,
1958           beats: [{id: index, timestamp: precise_time}]
1959       }
1960   b. Include mask information if applied for analysis transparency
1961   c. Export to dataset/video_edit/audio_analysis/ directory:
1962       - cut_points.json: rhythm timing data
1963       - rhythm_detection.png: comprehensive analysis visualization
1964       - rhythm_distribution.png: statistical analysis plots
1965
1966 Technical Specifications:
1967   - Signal processing: librosa with STFT-based energy analysis
1968   - Peak detection: scipy.signal.find_peaks with configurable thresholds
1969   - Temporal resolution: Frame-based analysis with hop_length=512 samples
1970   - Visualization: matplotlib with multi-panel comprehensive displays
1971   - Precision: 3 decimal places for timestamp accuracy
1972
1973 Rhythm Analysis Features:
1974   - Adaptive energy thresholding for various music styles
1975   - Temporal masking for intro/outro exclusion
1976   - Smoothing filters for noise reduction in energy detection
1977   - Statistical analysis for rhythm consistency evaluation
1978   - Multi-domain visualization (time, frequency, energy)
1979
1980 Error Handling:
1981   - Audio file format validation with librosa compatibility
1982   - Path existence verification for input and output directories
1983   - Graceful handling of insufficient rhythm points for analysis
1984   - Exception management for file I/O operations
1985   - Comprehensive error reporting for debugging workflows

```

1984 A.7.3 LOUDNESSNORMALIZER

1986 The LoudnessNormalizer agent provides comprehensive audio loudness normalization capabilities
 1987 for maintaining consistent audio levels across multiple files. It serves as a critical preprocessing
 1988 component in audio production workflows, utilizing the FAP (Fast Audio Processing) tool for batch
 1989 loudness standardization.

1990 Listing 14: LoudnessNormalizer Structure

```

1992 Input: data_directory
1993 Output: processing_status
1994
1995 Algorithm:
1996 1. Input validation and path processing:
1997   a. Validate data_dir parameter using InputSchema
   b. Convert input path to resolved Path object

```

```
1998     c. Verify directory existence and accessibility
1999     d. Ensure input is directory (not single file)
2000     e. Generate absolute path for consistent processing
2001
2002 2. Path preprocessing workflow:
2003     a. Apply Path.resolve() for absolute path conversion
2004     b. Validate directory permissions and read access
2005     c. Handle path existence verification with detailed error reporting
2006     d. Ensure target directory contains processable audio files
2007
2008 3. FAP command construction:
2009     a. Build command array: ["fap", "loudness-norm", input_dir,
2010        output_dir, "--overwrite", "--recursive"]
2011     b. Configure overwrite mode for existing file replacement
2012     c. Enable recursive processing for subdirectory traversal
2013     d. Filter empty arguments from command array
2014
2015 4. Subprocess execution with real-time monitoring:
2016     a. Initialize subprocess.Popen with pipe configuration:
2017        - stdout=PIPE, stderr=PIPE for output capture
2018        - bufsize=1 for line-buffered output
2019     b. Create threading.Thread for stdout and stderr monitoring
2020     c. Apply real-time output reading with UTF-8 decoding
2021     d. Handle decoding errors with 'replace' error handling
2022
2023 5. Real-time output processing:
2024     Function _read_output(pipe):
2025     a. Iterate through pipe.readline() until empty
2026     b. Decode bytes to UTF-8 with error replacement
2027     c. Strip whitespace and format with [FAP] prefix
2028     d. Print real-time progress updates to console
2029     e. Ensure proper pipe closure after processing
2030
2031 6. Process completion and result handling:
2032     a. Wait for subprocess completion with process.wait()
2033     b. Join stdout and stderr threads for cleanup
2034     c. Evaluate return_code for success/failure determination
2035     d. Generate success status for return_code == 0
2036     e. Raise RuntimeError for non-zero return codes
2037
2038 Error Handling:
2039 - Path existence and type validation with specific error messages
2040 - Directory permission verification and access control
2041 - FAP tool availability detection with installation guidance
2042 - Real-time error capture and user feedback
2043 - Graceful subprocess termination and resource cleanup
2044 - Comprehensive exception handling with status reporting
```

2043 A.7.4 MERGE

2045 The Merge agent provides direct video and audio track combination capabilities without intermediate
2046 video clip processing. It distinguishes itself from other multimedia agents by performing complete
2047 file merging rather than segmented video editing, utilizing FFmpeg for high-quality encoding and
2048 synchronization.

2049 Listing 15: Merge Structure

```
2050 Input: video_path, audio_path
2051 Output: merged_video_path
```

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

Algorithm:

1. Input validation and parameter processing:
 - a. Validate video_path and audio_path using InputSchema
 - b. Verify file existence and accessibility for both inputs
 - c. Set output_path to overwrite original video file
 - d. Ensure input files are in compatible formats
2. FFmpeg command construction:
 - a. Configure dual input streams:
 - Primary input: source video file (-i video_path)
 - Secondary input: replacement audio file (-i audio_path)
 - b. Set video encoding parameters:
 - Codec: libx264 for H.264 compression
 - Preset: fast for speed/quality balance
 - CRF: 23 for optimal quality (18-28 range)
 - c. Set audio encoding parameters:
 - Codec: aac for broad compatibility
 - Bitrate: 192k for high-quality audio
3. Stream mapping configuration:
 - a. Video stream selection: -map 0:v:0 (first video from input 0)
 - b. Audio stream selection: -map 1:a:0 (first audio from input 1)
 - c. Duration control: -shortest flag for sync with shorter input
 - d. Optimization: +faststart movflags for streaming compatibility
4. Encoding execution pipeline:

Command structure:

```
ffmpeg -i video_path -i audio_path -c:v libx264 -preset fast  
-crf 23 -c:a aac -b:a 192k -map 0:v:0 -map 1:a:0  
-shortest -movflags +faststart -y output_path
```

 - a. Execute subprocess.run with check=True for error detection
 - b. Monitor process completion and return code validation
 - c. Handle real-time progress feedback through console output
5. Quality assurance and output validation:
 - a. Verify successful merge completion through return code
 - b. Validate output file creation and size consistency
 - c. Ensure audio-video synchronization maintenance
 - d. Generate success confirmation with output path
6. Error handling and recovery:
 - a. Capture CalledProcessError for FFmpeg execution failures
 - b. Handle file permission and access errors gracefully
 - c. Provide detailed error reporting for troubleshooting
 - d. Ensure partial file cleanup on processing failure

Technical Specifications:

- Video encoding: H.264 with CRF 23 quality level
- Audio encoding: AAC at 192 kbps bitrate
- Synchronization: Shortest input duration matching
- Output optimization: Fast-start flag for progressive download
- File handling: Automatic overwrite with -y flag

2106 A.7.5 MIXER

2107
2108 The Mixer agent provides sophisticated audio mixing capabilities for combining foreground audio
2109 with background music tracks. It utilizes PyDub for comprehensive audio processing, supporting
2110 dynamic volume adjustment, length synchronization, and high-quality audio overlay operations.

2111 Listing 16: Mixer Structure

```

2112
2113 Input: bgm_path, audio_path
2114 Output: mixed_audio_path
2115
2116 Algorithm:
2117 1. Input validation and path processing:
2118     a. Validate bgm_path and audio_path using InputSchema
2119     b. Verify file existence and format compatibility
2120     c. Generate output filename: "mixed_" + original_audio_name
2121     d. Set output directory to match input audio directory
2122     e. Determine output format from file extension or default to WAV
2123
2124 2. Audio loading and preprocessing:
2125     a. Load BGM file using AudioSegment.from_file(bgm_path)
2126     b. Load vocal audio using AudioSegment.from_file(audio_path)
2127     c. Support multiple audio formats (WAV, MP3, FLAC, etc.)
2128     d. Handle audio format conversion internally through PyDub
2129
2130 3. Volume adjustment and normalization:
2131     a. Apply BGM volume reduction: bgm_audio + bgm_volume
2132     b. Default BGM volume adjustment: -1 dB for balanced mixing
2133     c. Maintain vocal audio at original volume level
2134     d. Report volume adjustment settings to user
2135
2136 4. Length synchronization and BGM extension:
2137     a. Compare audio lengths: len(bgm_audio) vs len(vocal_audio)
2138     b. For shorter BGM:
2139         - Initialize repeated_bgm = bgm_audio
2140         - Loop BGM extension: repeated_bgm += bgm_audio
2141         - Continue until repeated_bgm >= vocal_audio length
2142         - Trim to exact vocal length: repeated_bgm[:len(vocal_audio)]
2143     c. For longer BGM: use original BGM without modification
2144
2145 5. Audio overlay and mixing process:
2146     a. For extended BGM case:
2147         - Apply overlay: vocal_audio.overlay(repeated_bgm)
2148         - Vocal audio maintains prominence in mix
2149     b. For original BGM case:
2150         - Apply overlay: bgm_audio.overlay(vocal_audio)
2151         - BGM provides extended background beyond vocal length
2152     c. Overlay operation preserves audio quality and dynamic range
2153
2154 6. Export and output generation:
2155     a. Determine output format from file extension analysis
2156     b. Export mixed audio: mixed_audio.export(output_path, format)
2157     c. Support format options: WAV, MP3, FLAC, AAC, OGG
2158     d. Maintain original audio quality settings during export
2159     e. Return output path for downstream processing
2160
2161 Technical Specifications:
2162 - Volume adjustment: -1 dB default BGM reduction
2163 - Length handling: Automatic BGM looping for shorter backgrounds
2164 - Overlay method: Vocal-priority mixing with background preservation

```

- 2160 - Format support: Multi-format input/output through PyDub
- 2161 - Quality preservation: No quality degradation during mixing process
- 2162

2163 Error Handling:

- 2164 - File format compatibility validation
- 2165 - Audio loading error detection and reporting
- 2166 - Export format verification with fallback to WAV
- 2167 - Exception handling with detailed error messaging
- 2168 - Graceful degradation for unsupported audio formats
- 2169

2170 A.7.6 RESAMPLER

2172 The Resampler agent provides comprehensive audio resampling capabilities for standardizing sample
 2173 rates across multiple audio files within a directory. It leverages the FAP (Fast Audio Processing)
 2174 framework to perform batch resampling operations with real-time progress monitoring and threaded
 2175 output handling.

2176 Listing 17: Resampler Structure

```

2178 Input: data_directory
2179 Output: processing_status
2180
2181 Algorithm:
2182 1. Input validation and path processing:
2183   a. Validate data_dir parameter using InputSchema
2184   b. Convert input path to resolved Path object using Path.resolve()
2185   c. Verify directory existence with detailed error reporting
2186   d. Ensure input is directory type (not single file)
2187   e. Generate absolute path for consistent processing
2188
2189 2. Path preprocessing workflow:
2190   Function _process_path(input_path):
2191   a. Initialize Path object from input string
2192   b. Validate path.exists() with FileNotFoundError handling
2193   c. Verify path.is_dir() with type validation
2194   d. Return resolved absolute path for processing
2195   e. Handle permission and access control validation
2196
2197 3. FAP command construction and configuration:
2198   a. Build command array: ["fap", "resample", input_dir,
2199   output_dir, "--overwrite"]
2200   b. Configure input and output directories (same directory)
2201   c. Enable overwrite mode for in-place file replacement
2202   d. Filter empty arguments from command array
2203   e. Generate command string for execution logging
2204
2205 4. Subprocess execution with threading:
2206   a. Initialize subprocess.Popen with pipe configuration:
2207     - stdout=PIPE, stderr=PIPE for output capture
2208     - bufsize=1 for line-buffered real-time output
2209   b. Create separate threads for stdout and stderr monitoring
2210   c. Start threads for concurrent output processing
2211   d. Execute process.wait() for completion synchronization
2212
2213 5. Real-time output monitoring system:
2214   Function _read_output(pipe):
2215   a. Implement iter(pipe.readline, b'') for line iteration
2216   b. Decode bytes to UTF-8 with 'replace' error handling
2217   c. Strip whitespace and apply [FAP] prefix formatting
  
```

2214 d. Print real-time progress updates to console
 2215 e. Ensure proper pipe closure in finally block
 2216
 2217 6. Process completion and result validation:
 2218 a. Wait for subprocess completion with return_code capture
 2219 b. Join stdout_thread and stderr_thread for cleanup
 2220 c. Evaluate return_code for success/failure determination
 2221 d. Return success status for return_code == 0
 2222 e. Raise RuntimeError with detailed code for failures
 2223
 2224 Technical Specifications:
 2225 - Resampling engine: FAP framework with optimized algorithms
 2226 - Processing mode: In-place file replacement with --overwrite
 2227 - Output handling: Real-time threaded monitoring system
 2228 - Error detection: Return code validation with detailed reporting
 2229 - Concurrency: Separate threads for stdout/stderr processing
 2230
 2231 Error Handling:
 2232 - Path existence and type validation with specific messages
 2233 - Directory permission verification and access control
 2234 - FAP tool availability detection with installation guidance
 2235 - Process execution error capture with return code analysis
 2236 - Thread synchronization and resource cleanup
 2237 - Comprehensive exception handling with status reporting

2238 A.7.7 SEPARATOR

2240 The Separator agent provides advanced audio source separation capabilities for isolating vocal and
 2241 instrumental components from mixed audio tracks. It utilizes machine learning-based separation
 2242 algorithms through the FAP framework to perform high-quality vocal extraction and background
 2243 music isolation across multiple audio files.

2244 Listing 18: Separator Structure

2245 Input: data_directory
 2246 Output: processing_status
 2247
 2248 Algorithm:
 2249 1. Input validation and directory processing:
 2250 a. Validate data_dir parameter using InputSchema
 2251 b. Convert input path to resolved Path object
 2252 c. Verify directory existence and accessibility
 2253 d. Ensure input is directory type for batch processing
 2254 e. Generate absolute path for consistent file handling
 2255
 2256 2. Path preprocessing and validation:
 2257 Function _process_path(input_path):
 2258 a. Initialize Path object from input string
 2259 b. Validate path.exists() with detailed error messaging
 2260 c. Verify path.is_dir() for directory type confirmation
 2261 d. Return path.resolve() for absolute path generation
 2262 e. Handle access permissions and directory structure
 2263
 2264 3. FAP separation command construction:
 2265 a. Build command array: ["fap", "separate", input_dir,
 2266 output_dir, "--overwrite", "--recursive"]
 2267 b. Configure input and output directories (same location)
 2268 c. Enable overwrite mode for existing file replacement
 2269 d. Enable recursive processing for subdirectory traversal

```
2268     e. Filter empty arguments and validate command structure
2269
2270 4. Machine learning-based separation execution:
2271   a. Initialize subprocess.Popen with pipe configuration:
2272     - stdout=PIPE, stderr=PIPE for output monitoring
2273     - bufsize=1 for real-time line-buffered output
2274   b. Create threading.Thread for concurrent output handling
2275   c. Execute source separation algorithms on audio files
2276   d. Monitor processing progress through threaded output
2277
2278 5. Real-time progress monitoring system:
2279   Function _read_output(pipe):
2280   a. Implement line-by-line output reading with iter()
2281   b. Decode subprocess output with UTF-8 and error handling
2282   c. Format output with [FAP] prefix for identification
2283   d. Display real-time separation progress to user
2284   e. Ensure proper pipe closure and resource management
2285
2286 6. Separation result validation and output:
2287   a. Wait for subprocess completion with return_code capture
2288   b. Synchronize stdout_thread and stderr_thread completion
2289   c. Evaluate return_code for separation success/failure
2290   d. Generate separated audio files (vocals, instrumentals)
2291   e. Return processing status with detailed error reporting
2292
2293 Technical Specifications:
2294 - Separation engine: Machine learning-based source separation
2295 - Processing scope: Recursive directory traversal capability
2296 - Output formats: Separated vocal and instrumental tracks
2297 - Quality preservation: High-fidelity separation algorithms
2298 - Real-time monitoring: Threaded progress feedback system
2299
2300 Separation Outputs:
2301 - Vocal tracks: Isolated human voice components
2302 - Instrumental tracks: Background music without vocals
2303 - Original preservation: Source files maintained with --overwrite
2304 - Directory structure: Maintains original file organization
2305
2306 Error Handling:
2307 - Directory validation with specific error messages
2308 - FAP tool availability verification and guidance
2309 - Separation algorithm failure detection and reporting
2310 - Thread synchronization and cleanup procedures
2311 - Resource management for large audio file processing
2312 - Comprehensive exception handling with status codes
```

2311 A.7.8 TRANSCRIBER

2312
2313 The Transcriber agent provides comprehensive audio-to-text transcription capabilities using advanced
2314 speech recognition models. It processes multiple audio files within a directory structure and generates
2315 accurate transcription outputs with real-time progress monitoring and robust error handling.

2316
2317 Listing 19: Transcriber Structure

```
2318 Input: data_directory
2319 Output: processing_status
2320
2321 Algorithm:
2322 1. Input validation and directory processing:
```

2322 a. Validate data_dir parameter using InputSchema
2323 b. Convert input path to resolved Path object
2324 c. Verify directory existence and read permissions
2325 d. Ensure input is directory type for batch processing
2326 e. Generate absolute path for consistent file access
2327
2328 2. Path preprocessing and validation:
2329 Function _process_path(input_path):
2330 a. Initialize Path object from input string
2331 b. Validate path.exists() with FileNotFoundError handling
2332 c. Verify path.is_dir() for directory type confirmation
2333 d. Return path.resolve() for absolute path generation
2334 e. Handle access control and permission validation
2335
2336 3. FunASR model configuration and command construction:
2337 a. Build transcription command: ["fap", "transcribe",
2338 "--model-type", "funasr", "--recursive", audio_dir]
2339 b. Configure FunASR model for high-accuracy transcription
2340 c. Enable recursive processing for subdirectory traversal
2341 d. Set model parameters for optimal speech recognition
2342 e. Prepare command array for subprocess execution
2343
2344 4. Subprocess execution with real-time monitoring:
2345 a. Initialize subprocess.Popen with pipe configuration:
2346 - stdout=PIPE, stderr=PIPE for output capture
2347 - bufsize=1 for line-buffered real-time feedback
2348 b. Create threading.Thread for stdout and stderr handling
2349 c. Execute transcription process with concurrent monitoring
2350 d. Handle multiple audio file processing simultaneously
2351
2352 5. Real-time progress tracking system:
2353 Function _read_output(pipe):
2354 a. Implement iter(pipe.readline, b'') for line iteration
2355 b. Decode subprocess output with UTF-8 encoding
2356 c. Apply error handling with 'replace' for corrupted chars
2357 d. Format output with [FAP] prefix for identification
2358 e. Display transcription progress and file completion status
2359
2360 6. Transcription completion and result handling:
2361 a. Wait for subprocess completion with return_code validation
2362 b. Synchronize stdout_thread and stderr_thread completion
2363 c. Evaluate return_code for transcription success/failure
2364 d. Generate transcription files in same directory as audio
2365 e. Return processing status with detailed error information
2366
2367 Technical Specifications:
2368 - Speech recognition model: FunASR for high-accuracy transcription
2369 - Processing scope: Recursive directory traversal capability
2370 - Output format: Text files (.txt) with same basename as audio
2371 - Model features: Multi-language support and noise robustness
2372 - Real-time feedback: Threaded progress monitoring system
2373
2374 Transcription Outputs:
2375 - Text files: Generated with matching audio file basenames
- Directory structure: Maintains original audio file organization
- Encoding: UTF-8 text files for broad compatibility
- Accuracy: High-quality speech-to-text conversion
- Timestamps: Optional timestamp generation for synchronization

2376 Error Handling:
 2377 - Directory validation with specific error messaging
 2378 - FunASR model availability verification and installation guidance
 2379 - Audio format compatibility checking and conversion recommendations
 2380 - Thread synchronization and resource cleanup procedures
 2381 - Comprehensive exception handling with return code analysis
 2382 - Graceful degradation for unsupported audio formats

2383
 2384
 2385

A.7.9 VOICEGENERATOR

2386 The VoiceGenerator agent provides advanced text-to-speech synthesis capabilities for generating high-
 2387 quality voice audio from scene content. It utilizes the CosyVoice2 model for zero-shot voice cloning
 2388 and produces synchronized audio with precise timestamp tracking for video editing workflows.
 2389

2390 Listing 20: VoiceGenerator Structure

2391 Input: video_scene_path, target_vocal_path
 2392 Output: synthesized_audio_path, timestamp_path
 2393
 2394 Algorithm:
 2395 1. Model initialization and dependency management:
 2396 a. Import CosyVoice2 and load_wav from cosyvoice.cli.cosyvoice
 2397 b. Load CosyVoice2 model from pretrained_models/CosyVoice2-0.5B
 2398 c. Configure model parameters: load_jit=False, load_trt=False, fp16=False
 2399 d. Load target vocal prompt using load_wav at 16kHz sample rate
 2400 e. Initialize prompt_speech_16k for zero-shot voice synthesis
 2401
 2402 2. Scene content processing with multilingual support:
 2403 Function _process_with_timestamps(json_file_path):
 2404 a. Load JSON file with UTF-8 encoding for Chinese text support
 2405 b. Extract content_created field from scene JSON structure
 2406 c. Parse content using '/////\\n' delimiter pattern
 2407 d. Create segment list with sequential IDs and content mapping
 2408 e. Generate clean JSON output with ensure_ascii=False
 2409
 2410 3. Text segmentation for optimal TTS processing:
 2411 Function _split_into_sentences(text, max_length=200):
 2412 a. Handle Chinese punctuation
 2413 b. Split by punctuation markers with '|' separator insertion
 2414 c. Apply length-based chunking for sentences exceeding max_length
 2415 d. Further split by commas and Chinese commas
 2416 e. Maintain semantic coherence while respecting length limits
 2417
 2418 4. Audio generation with timestamp tracking:
 2419 Function _generate_audio_for_segments():
 2420 a. Process each segment through sentence chunking
 2421 b. Apply zero-shot voice synthesis for each text chunk:
 2422 cosyvoice.inference_zero_shot(generator, prompt_text,
 2423 prompt_speech_16k, stream=False)
 2424 c. Concatenate chunk waveforms using torch.cat()
 2425 d. Track cumulative timestamps: current_time += segment_duration
 2426 e. Generate timestamp data structure with precise timing
 2427
 2428 5. Zero-shot voice synthesis pipeline:
 2429 a. Create single_sentence_generator() for individual text chunks
 2430 b. Apply inference with target voice characteristics
 2431 c. Extract tts_speech waveform from audio_data response
 2432 d. Concatenate waveforms maintaining audio continuity
 2433 e. Handle synthesis errors with graceful chunk skipping

2430
 2431 6. Audio consolidation and export:
 2432 Function `_combine_audio_files()`:
 2433 a. Combine all segment waveforms using `torch.cat()`
 2434 b. Export final audio using `torchaudio.save()` with `sample_rate`
 2435 c. Generate `cut_points.json` with UTF-8 encoding
 2436 d. Structure timestamp data: `sentence_data.chunks` with IDs
 2437 e. Clean up temporary segment files after combination
 2438
 2439 Technical Specifications:
 2440 - Voice model: CosyVoice2-0.5B for high-quality synthesis
 2441 - Sample rate: 16kHz for prompt audio, model native for output
 2442 - Text processing: Max 200 characters per chunk for optimal quality
 2443 - Audio format: WAV output with torch audio tensor processing
 2444 - Timestamp precision: Millisecond-level accuracy for video sync
 2445
 2446 Multilingual Features:
 2447 - Chinese text support with proper punctuation handling
 2448 - UTF-8 encoding throughout processing pipeline
 2449 - ensure_ascii=False for character preservation
 2450 - Bilingual punctuation recognition and processing
 2451 - Cultural text formatting considerations
 2452
 2453 Error Handling:
 2454 - Model loading verification with dependency checking
 2455 - File existence validation for scene and prompt inputs
 2456 - Chunk-level error isolation with processing continuation
 2457 - Resource cleanup for temporary files and memory management
 2458 - Comprehensive exception handling with detailed error reporting

2459 A.7.10 TTSINFER

2460 The TTSInfer agent provides sophisticated text-to-speech synthesis for rewritten video content,
 2461 utilizing sliced audio clips as voice references. It implements a three-stage Fish-Speech pipeline
 2462 combining VQGAN encoding, text-to-semantic conversion, and audio generation to produce high-
 2463 quality voice clones matching original speaker characteristics.
 2464

2465 Listing 21: TTSInfer Structure

2466 Input: `audio_path`, `speech_text_path`
 2467 Output: `derivative_audio_directory`
 2468
 2469 Algorithm:
 2470 1. Input processing and text segmentation:
 2471 a. Load rewritten text from `speech_path` with UTF-8 encoding
 2472 b. Split content into paragraph segments using line breaks
 2473 c. Validate `audio_path` and extract slice directory structure
 2474 d. Create derivative directory for synthesized output files
 2475 e. Initialize Fish-Speech model paths and dependencies
 2476
 2477 2. Audio reference preparation and duration analysis:
 2478 Function `get_audio_duration(wav_file_path)`:
 2479 a. Open WAV file using `wave.open()` with `contextlib.closing()`
 2480 b. Calculate duration: `frames / sample_rate` for timing reference
 2481 c. Handle audio format errors with exception management
 2482 d. Return precise duration for reference audio segments
 2483
 2484 3. LAB file processing and audio segment mapping:
 2485 a. Scan slice directory for `.lab` files with numerical sorting

```
2484     b. Load LAB content with UTF-8 encoding for text reference
2485     c. Map corresponding .wav files using Path.with_suffix()
2486     d. Create lab_files_with_content structure with paths and content
2487
2488 4. Duration-based audio combination strategy:
2489   For each text paragraph:
2490     a. Start with corresponding LAB file and WAV audio
2491     b. Calculate combined_duration from audio segments
2492     c. If duration < 1 second, aggregate additional segments:
2493         - Cycle through available LAB files using modulo indexing
2494         - Concatenate LAB content and WAV file references
2495         - Continue until minimum 1-second duration achieved
2496     d. Handle multiple WAV concatenation using scipy.io.wavfile
2497
2498 5. Three-stage Fish-Speech synthesis pipeline:
2499   Stage 1 - VQGAN Encoding:
2500   Command: fish_speech/models/vqgan/inference.py
2501   - Input: combined reference WAV file
2502   - Checkpoint: firefly-gan-vq-fsq-8x1024-21hz-generator.pth
2503   - Output: encoded audio features (.npz tokens)
2504
2505   Stage 2 - Text-to-Semantic Conversion:
2506   Command: fish_speech/models/text2semantic/inference.py
2507   - Parameters: --text (target text), --prompt-text (LAB content)
2508   - Input: prompt tokens from Stage 1
2509   - Checkpoint: fish-speech-1.5 semantic model
2510   - Output: semantic codes (temp/codes_0.npz)
2511
2512   Stage 3 - Audio Generation:
2513   Command: fish_speech/models/vqgan/inference.py
2514   - Input: semantic codes from Stage 2
2515   - Checkpoint: firefly-gan-vq-fsq-8x1024-21hz-generator.pth
2516   - Output: final synthesized WAV file
2517
2518 6. Batch processing and audio consolidation:
2519   a. Process each paragraph through complete synthesis pipeline
2520   b. Monitor subprocess execution with stdout/stderr capture
2521   c. Validate return codes for each pipeline stage
2522   d. Collect generated WAV files in derivative directory
2523   e. Perform final audio concatenation using numpy arrays
2524
2525 7. Final audio merging and cleanup:
2526   a. Sort generated WAV files by numerical stem for proper ordering
2527   b. Load audio data using scipy.io.wavfile.read()
2528   c. Validate consistent sample rates across all segments
2529   d. Concatenate audio arrays using numpy.concatenate()
2530   e. Export final merged audio as derivative/final.wav
2531
2532 Technical Specifications:
2533 - Voice synthesis: Fish-Speech 1.5 with VQGAN-based generation
2534 - Audio processing: WAV format with sample rate consistency
2535 - Text encoding: UTF-8 support for multilingual content
2536 - Duration threshold: Minimum 1-second reference for quality synthesis
2537 - Pipeline stages: Sequential VQGAN --> Text2Semantic --> VQGAN
2538
2539 Error Handling:
2540 - Subprocess return code validation for each synthesis stage
2541 - Audio duration calculation with wave format error handling
2542 - Sample rate consistency validation across concatenated segments
```

- 2538 - Temporary file cleanup with graceful error recovery
- 2539 - Exception isolation for individual paragraph processing failures
- 2540 - Directory creation and permission handling for output paths

2542

2543 A.7.11 TTSREPLACE

2544

2545

2546

2547

2548

The TTSReplace agent provides comprehensive video-audio synchronization for replacing original video audio with derivative synthesized speech segments. It implements precise temporal alignment, video speed adjustment, and seamless audio-video merging to produce coherent rewritten video content while maintaining visual continuity.

2549

Listing 22: TTSReplace Structure

2550

2551

2552

Input: video_path
Output: final_video_path

2553

2554

2555

2556

2557

2558

2559

2560

2561

2562

2563

2564

2565

2566

Algorithm:

1. Input validation and directory structure setup:

- a. Validate video_path existence and file accessibility
- b. Extract slice_dir from video filename without extension
- c. Verify derivative audio directory and metadata.json existence
- d. Create final output directory structure for processed clips
- e. Configure UTF-8 encoding for stdout/stderr handling

2. Metadata processing and clip mapping:

- a. Load metadata.json containing temporal clip information
- b. Extract clip data: file names, start times, end times, durations
- c. Map derivative audio files to corresponding video segments
- d. Validate derivative audio availability for each clip
- e. Initialize processed_files array for concatenation tracking

3. Video clip extraction with temporal precision:

For each metadata clip:

- a. Execute FFmpeg video extraction command:
ffmpeg -y -ss start_time -i video_path -to duration
-c:v libx264 -preset fast -vf scale=iw:ih clip_path
- b. Apply precise temporal trimming using start/end timestamps
- c. Maintain video quality with libx264 fast preset
- d. Preserve original resolution with scale=iw:ih filter

2575

2576

2577

2578

2579

2580

2581

2582

2583

2584

2585

2586

2587

2588

2589

2590

2591

4. Audio duration analysis and speed calculation:

- a. Use FFprobe to determine derivative audio duration:
ffprobe -v error -show_entries format=duration
-of default=noprint_wrappers=1:nokey=1 derivative_audio
- b. Calculate speed adjustment factor: target_duration / clip_duration
- c. Determine video playback speed modification requirements
- d. Handle duration mismatches between audio and video segments

5. Video speed adjustment for audio synchronization:

- a. Apply temporal adjustment using setpts filter:
ffmpeg -y -i clip_path -filter:v setpts=speed_factor*PTS
-an adjusted_path
- b. Remove audio track (-an) to prepare for replacement
- c. Modify video playback speed to match derivative audio duration
- d. Maintain frame rate consistency throughout adjustment

6. Audio-video merging with quality preservation:

- a. Combine adjusted video with derivative audio:
ffmpeg -y -i adjusted_path -i derivative_audio

```

2592     -c:v copy -c:a aac -map 0:v:0 -map 1:a:0
2593     -shortest merged_path
2594     b. Use copy codec for video to avoid quality loss
2595     c. Encode audio to AAC for broad compatibility
2596     d. Apply shortest stream duration for synchronization
2597
2598 7. Batch concatenation and final assembly:
2599     a. Generate filelist.txt with absolute paths for FFmpeg concat
2600     b. Write UTF-8 encoded file list for proper path handling
2601     c. Execute FFmpeg concatenation:
2602         ffmpeg -y -f concat -safe 0 -i filelist.txt
2603         -c copy final_output
2604     d. Use copy codec for lossless segment joining
2605
2606 8. Cleanup and resource management:
2607     a. Iterate through output directory for intermediate file removal
2608     b. Preserve only final.mp4 output while removing temporary clips
2609     c. Handle file deletion errors with graceful error reporting
2610     d. Display cleanup progress and completion status
2611
2612 Technical Specifications:
2613 - Video encoding: H.264 with fast preset for efficiency
2614 - Audio encoding: AAC for universal compatibility
2615 - Temporal precision: Frame-accurate clip extraction and alignment
2616 - Speed adjustment: setpts filter for smooth temporal modification
2617 - Concatenation: FFmpeg concat demuxer for seamless joining
2618
2619 Synchronization Features:
2620 - Dynamic speed adjustment based on audio duration differences
2621 - Precise temporal alignment using FFprobe duration analysis
2622 - Quality preservation through copy codec usage during concatenation
2623 - Frame rate consistency maintenance throughout processing pipeline
2624
2625 Error Handling:
2626 - File existence validation for video, audio, and metadata
2627 - Subprocess execution monitoring with return code checking
2628 - Graceful handling of missing derivative audio segments
2629 - Resource cleanup with error isolation for individual file operations
2630 - UTF-8 encoding configuration for international character support

```

2630 A.7.12 TTSSLICER

2631 The TTSSlicer agent provides intelligent audio segmentation for preparing audio content for text-to-
2632 speech processing and transcription workflows. It implements advanced silence detection algorithms
2633 with RMS-based energy analysis to create optimal audio chunks while maintaining temporal precision
2634 and content coherence.
2635

2636 Listing 23: TTSSlicer Structure

```

2637 Input: audio_path
2638 Output: processing_status, segmented_audio_files, metadata
2639
2640 Algorithm:
2641 1. Configuration and parameter initialization:
2642     a. Set duration constraints: min_duration=6.0s, max_duration=8.0s
2643     b. Configure silence detection: min_silence_duration=0.5s, top_db=-35
2644     c. Set analysis parameters: hop_length=10ms, max_silence_kept=0.3s
2645     d. Initialize merge_short=False for individual segment processing

```

- 2646
- 2647
- 2648
- 2649
- 2650
- 2651
- 2652
- 2653
- 2654
- 2655
- 2656
- 2657
- 2658
- 2659
- 2660
- 2661
- 2662
- 2663
- 2664
- 2665
- 2666
- 2667
- 2668
- 2669
- 2670
- 2671
- 2672
- 2673
- 2674
- 2675
- 2676
- 2677
- 2678
- 2679
- 2680
- 2681
- 2682
- 2683
- 2684
- 2685
- 2686
- 2687
- 2688
- 2689
- 2690
- 2691
- 2692
- 2693
- 2694
- 2695
- 2696
- 2697
- 2698
- 2699
2. Audio loading and preprocessing:
 - a. Load audio using `librosa.load()` with original sample rate preservation
 - b. Handle mono/stereo conversion: expand mono to 2D array format
 - c. Validate audio data integrity and duration constraints
 - d. Create output directory structure based on audio filename
 3. RMS-based silence detection pipeline:

Class `_Slicer` implementation:

 - a. Calculate RMS energy using `librosa.feature.rms()`:
 - `frame_length=win_size`, `hop_length=hop_size`
 - Apply threshold: $10^{(\text{top_db}/20.0)}$ for dB to linear conversion
 - b. Analyze energy profile for silence identification
 - c. Track `silence_start` and `clip_start` positions for segmentation
 4. Intelligent segmentation logic:

For each RMS frame:

 - a. Detect silence regions where `rms < threshold`
 - b. Apply segmentation rules:
 - Leading silence: `silence_start == 0` and `i > max_sil_kept`
 - Middle silence: `i - silence_start >= min_interval` and `i - clip_start >= min_length`
 - c. Calculate optimal cut points using `argmin()` on RMS values
 - d. Generate `sil_tags` array with `(start, end)` silence boundaries
 5. Chunk extraction with temporal precision:

Function `_apply_slice(waveform, begin, end)`:

 - a. Convert frame indices to sample indices: `idx * hop_size`
 - b. Extract audio slice maintaining original channel structure
 - c. Calculate precise timestamps: `start_idx/sr`, `end_idx/sr`
 - d. Return structured chunk with audio data and timing information
 6. Duration-based post-processing:
 - a. Short chunk merging (if `merge_short` enabled):
 - Combine consecutive chunks until `max_duration` reached
 - Maintain temporal continuity through concatenation
 - b. Long chunk subdivision:
 - Split chunks exceeding `max_duration` using `_slice_by_max_duration()`
 - Calculate optimal `chunk_size`: `ceil(total_samples / n_chunks)`
 - Ensure uniform distribution across subdivided segments
 7. Audio export and metadata generation:
 - a. Save each chunk using `soundfile.write()`:
 - Transpose multi-channel audio for proper format
 - Maintain original sample rate and bit depth
 - b. Generate metadata with precise timing:
 - `filename`: sequential numbering (`0000.wav`, `0001.wav`, etc.)
 - `start/end` timestamps rounded to 3 decimal places
 - `duration` calculation: `end - start` with precision
 - c. Export `metadata.json` with UTF-8 encoding for compatibility
- Technical Specifications:
- Silence threshold: -35 dB for robust speech detection
 - Temporal resolution: 10ms hop length for precise timing
 - Segment constraints: 6-8 second optimal duration range
 - RMS analysis: Frame-based energy calculation with `librosa`
 - Output format: WAV files with original sample rate preservation
- Segmentation Features:
- Adaptive silence detection with configurable thresholds

- 2700 - Intelligent cut point selection using minimum RMS values
- 2701 - Leading/trailing silence handling with preservation limits
- 2702 - Long audio subdivision for consistent chunk sizes
- 2703 - Metadata tracking for downstream processing integration

2704 Error Handling:

- 2705 - Audio format validation with librosa compatibility checking
- 2706 - Empty audio detection and graceful handling
- 2707 - Directory creation with permission validation
- 2708 - Sample rate preservation across processing pipeline
- 2709 - Robust file I/O with exception management for large files

2711

2712 A.7.13 SVCANALYZER

2713
2714 The SVCAnalyzer agent provides comprehensive MIDI file analysis for music cover creation work-
2715 flows, extracting detailed musical information including note sequences, timing data, and tempo
2716 variations. It performs intelligent lyrics-to-music alignment by matching character counts with actual
2717 note sequences while handling rest periods and tempo changes.

2718

2719 Listing 24: SVCAnalyzer Structure

2720

Input: midi_file_path, lyrics_file_path

2721

Output: song_name, analysis_results_path

2722

2723 Algorithm:

2724

2724 1. Input validation and file processing:

2725

a. Validate MIDI file extension (.mid) and accessibility

2726

b. Validate lyrics file extension (.txt) and encoding

2727

c. Extract song name from lyrics filename without extension

2728

d. Load lyrics content with UTF-8 encoding for character analysis

2729

2729 2. MIDI file parsing and track analysis:

2730

a. Load MIDI file using mido.MidiFile() for comprehensive parsing

2731

b. Iterate through all tracks to identify musical content

2732

c. Extract track names or assign default identifiers

2733

d. Initialize note tracking structures for temporal analysis

2734

2734 3. Tempo change detection and BPM calculation:

2735

Function get_tempo_changes(mid):

2736

a. Scan all tracks for 'set_tempo' messages

2737

b. Record timestamp and tempo value for each change

2738

c. Sort tempo changes chronologically by time

2739

d. Insert default tempo (500000 microseconds) if none at start

2740

e. Calculate BPM: $60000000 / \text{tempo_microseconds}$

2741

2742 4. Note extraction with temporal precision:

2743

For each MIDI track:

2744

a. Process note_on messages with velocity > 0 for note starts

2745

b. Process note_off messages or note_on with velocity = 0 for ends

2746

c. Calculate note duration: $\text{end_time} - \text{start_time}$ in ticks

2747

d. Detect rest periods: $\text{gaps} > \text{ticks_per_beat} / 8$ threshold

2748

e. Group simultaneous notes occurring within 0.01 time units

2749

2749 5. Time conversion with tempo awareness:

2750

Function ticks_to_seconds():

2751

a. Handle variable tempo throughout song duration

2752

b. Calculate duration for each tempo segment separately

2753

c. Apply formula: $(\text{ticks} * \text{tempo}) / (\text{ticks_per_beat} * 1000000)$

d. Accumulate total duration across tempo changes

```
2754     e. Ensure precise timing for musical synchronization
2755
2756 6. Note name conversion and formatting:
2757   Function note_to_name(note_number):
2758   a. Convert MIDI note numbers to musical notation
2759   b. Calculate octave: note_number // 12 - 1
2760   c. Determine note name using chromatic scale array
2761   d. Format as: "{note}{octave}" (e.g., "C4", "F#3")
2762   e. Handle chord notation for simultaneous notes
2763
2764 7. Lyrics-to-music alignment and validation:
2765   a. Count actual notes excluding rest periods
2766   b. Compare note count with lyrics character count
2767   c. For matching counts:
2768       - Insert "AP" markers for rest periods in lyrics
2769       - Map each character to corresponding note
2770       - Maintain temporal alignment throughout song
2771   d. Generate processed lyrics with rest markers
2772
2773 8. Analysis output generation and export:
2774   a. Create structured JSON output with:
2775       - text: processed lyrics with AP markers
2776       - notes: pipe-separated note sequence
2777       - notes_duration: corresponding timing data
2778       - input_type: "word" for character-level processing
2779   b. Save analysis to /analysis/{song_name}.json
2780   c. Ensure UTF-8 encoding with ensure_ascii=False
2781
2782 Technical Specifications:
2783 - MIDI parsing: mido library for comprehensive format support
2784 - Temporal resolution: Tick-based timing with tempo-aware conversion
2785 - Note grouping: 0.01 time unit threshold for simultaneous detection
2786 - Rest detection: ticks_per_beat / 8 minimum gap threshold
2787 - Precision: 6 decimal places for duration measurements
2788
2789 Musical Features:
2790 - Multi-track analysis with automatic track selection
2791 - Chord detection and simultaneous note handling
2792 - Variable tempo support throughout song duration
2793 - Rest period insertion for natural speech rhythm
2794 - Character-to-note mapping for lyrics synchronization
2795
2796 Error Handling:
2797 - MIDI file format validation and parsing error recovery
2798 - Lyrics encoding detection with UTF-8 fallback
2799 - Track selection validation for meaningful musical content
2800 - Character count mismatch detection and reporting
2801 - Directory creation with permission handling for output files
```

2800 A.7.14 SVCCONVERSION

2802 The SVCCConversion agent provides intelligent audio segmentation and timestamp generation for
2803 music cover workflows, converting adapted lyrics into precisely timed JSON format suitable for
2804 video generation. It handles complex temporal alignment by parsing musical structure markers and
2805 generating seamless transition points for multimedia editing.

2806 Listing 25: SVCCConversion Structure

```
2807 Input: adapted_lyrics_string, midi_analysis_path
```

```

2808 Output: timestamp_json_path
2809
2810 Algorithm:
2811 1. Input validation and data loading:
2812   a. Validate adapted lyrics string and analysis file path
2813   b. Load MIDI analysis JSON with UTF-8 encoding
2814   c. Update analysis data with adapted lyrics content
2815   d. Extract duration array from notes_duration field (pipe-separated)
2816
2817 2. Text parsing and temporal alignment:
2818   Function parse_text_to_segments(text, durations):
2819   a. Initialize timeline array and AP time ranges tracker
2820   b. Process character-by-character with duration mapping:
2821     - Detect "AP" markers for musical phrase boundaries
2822     - Map individual characters to corresponding note durations
2823     - Track cumulative time progression through song
2824   c. Generate timeline entries: ("AP", start, end) or ("CHAR", char, start,
↪ end)
2825
2826 3. Segment boundary detection and grouping:
2827   a. Process timeline to identify lyrical segments between AP markers
2828   b. Group consecutive characters into coherent text segments
2829   c. Calculate segment boundaries:
2830     - start: end of previous AP marker or beginning
2831     - end: start of next AP marker or final character
2832   d. Handle edge cases for segments at song boundaries
2833
2834 4. AP marker processing for instrumental breaks:
2835   a. Extract AP time ranges for instrumental/rest periods
2836   b. Identify extended AP durations exceeding 12-second threshold
2837   c. For long AP segments (duration > 12s):
2838     - Generate intermediate chunk points every 12 seconds
2839     - Create "bgm" content markers for background music
2840     - Ensure smooth transitions for video editing workflows
2841
2842 5. Timestamp generation and chunk creation:
2843   a. Generate text chunks from lyrical segments:
2844     - timestamp: segment end time for completion marking
2845     - content: complete text content of segment
2846     - type: "text" for lyrical content identification
2847   b. Generate AP chunks for instrumental breaks:
2848     - timestamp: calculated break points within long AP sections
2849     - content: "bgm" for background music indication
2850     - type: "ap" for instrumental section identification
2851
2852 6. Temporal sorting and priority handling:
2853   a. Combine text and AP chunks into unified entry array
2854   b. Apply sorting criteria: (timestamp, content_type_priority)
2855     - Primary sort: chronological by timestamp
2856     - Secondary sort: text entries before AP entries at same timestamp
2857   c. Ensure temporal consistency for seamless playback
2858
2859 7. JSON output generation and export:
2860   a. Create structured output format:
2861     sentence_data: {
↪       count: total_chunk_number,
↪       chunks: [
↪         {id: sequential_index, timestamp: precise_time, content: text_or_bgm

```

```

2862     ]
2863   }
2864   b. Round timestamps to 3 decimal places for precision
2865   c. Export to dataset/video_edit/voice_gen/gen_audio_timestamps.json
2866   d. Ensure UTF-8 encoding with ensure_ascii=False
2867
2868 Technical Specifications:
2869 - Temporal precision: 3 decimal place timestamp accuracy
2870 - Chunk segmentation: 12-second maximum for extended instrumental breaks
2871 - Character mapping: One-to-one correspondence with MIDI note durations
2872 - Priority system: Text content prioritized over background music markers
2873 - Format compatibility: JSON structure optimized for video editing workflows
2874
2875 Musical Structure Features:
2876 - AP marker recognition for phrase and verse boundaries
2877 - Instrumental break detection and subdivision
2878 - Character-level timing precision for lip-sync accuracy
2879 - Background music insertion points for extended instrumental sections
2880
2881 Error Handling:
2882 - Duration array validation with length consistency checking
2883 - Timeline boundary validation for segment integrity
2884 - File I/O error management for JSON export operations
2885 - Character encoding preservation for international lyrics
2886 - Directory creation with permission handling for output paths

```

2887

2888 A.7.15 SVCCOVERIST

2889

2890 The SVCCoverist agent provides advanced voice timbre cloning and conversion capabilities for
 2891 singing voice synthesis in music cover production. It utilizes the Seed-VC (Voice Conversion)
 2892 framework to perform source-to-target vocal transformation while maintaining musical characteristics
 2893 and timing precision for professional-quality audio output.

2894

Listing 26: SVCCoverist Structure

2895

```

2896 Input: source_audio_path, target_vocal_path
2897 Output: synthesized_audio_path
2898
2899 Algorithm:
2900 1. Input validation and path processing:
2901   a. Validate source audio and target vocal file accessibility
2902   b. Convert file paths to absolute paths for cross-directory execution
2903   c. Extract source and target filenames without extensions
2904   d. Generate output filename: "{source_name}_{target_name}.wav"
2905
2906 2. Directory structure setup and navigation:
2907   a. Calculate final output directory: ".././final" relative to source
2908   b. Create output directory structure with proper permissions
2909   c. Store original working directory for restoration
2910   d. Navigate to tools/seed-vc for model execution environment
2911
2912 3. Environment configuration for Seed-VC execution:
2913   a. Store original PYTHONPATH environment variable
2914   b. Set PYTHONPATH to seed-vc absolute directory path
2915   c. Configure Python module import resolution for seed-vc
2916   d. Ensure model dependencies and checkpoints accessibility
2917
2918 4. Voice conversion command construction:

```

```
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
```

- a. Build subprocess command array:
[python, "inference.py", "--source", source_path,
"--target", target_path, "--output", output_dir,
"--f0-condition", "True"]
- b. Enable F0 conditioning for pitch characteristic preservation
- c. Configure input/output paths for proper file handling

5. Subprocess execution with encoding management:

- a. Execute inference.py with UTF-8 text capture
- b. Handle UnicodeDecodeError with fallback decoding:
 - Capture output as bytes if UTF-8 fails
 - Apply decode with 'replace' error handling
 - Ensure readable output for debugging and monitoring
- c. Monitor stdout and stderr for process feedback

6. Voice conversion processing pipeline:

- a. Load source audio for vocal content extraction
- b. Load target audio for timbre characteristic analysis
- c. Apply neural voice conversion with F0 conditioning:
 - Preserve pitch patterns from source audio
 - Transfer vocal timbre from target speaker
 - Maintain temporal alignment and musical timing
- d. Generate converted audio with target voice characteristics

7. Output validation and cleanup:

- a. Verify successful command execution (return_code == 0)
- b. Validate output file creation in final directory
- c. Restore original working directory and environment:
 - Change back to original directory
 - Restore original PYTHONPATH or remove if empty
 - Clean up temporary environment modifications
- d. Return synthesized audio path for downstream processing

Technical Specifications:

- Voice conversion model: Seed-VC with F0 conditioning enabled
- Audio preservation: Pitch pattern and timing maintenance
- Output format: WAV files with high-quality audio encoding
- Environment isolation: Sandboxed execution with path restoration
- Error handling: Unicode decoding with graceful fallback

Voice Conversion Features:

- Timbre transfer: Source-to-target vocal characteristic mapping
- F0 conditioning: Pitch pattern preservation during conversion
- Musical timing: Temporal alignment maintenance for singing voice
- Quality preservation: High-fidelity audio output for professional use
- Batch processing: Support for multiple source-target combinations

Error Handling:

- File path validation with absolute path conversion
- Directory creation with permission error management
- Subprocess execution monitoring with return code validation
- Unicode encoding error recovery with 'replace' strategy
- Environment restoration with original state preservation
- Comprehensive exception handling with detailed error reporting

2970 A.7.16 SVCSINGLE

2971

2972 The SVCSingle agent provides comprehensive singing voice synthesis for music cover production,
 2973 converting adapted lyrics into high-quality vocal audio with precise temporal alignment. It utilizes
 2974 the DiffSinger framework for neural singing voice generation and implements sophisticated audio
 2975 processing techniques to ensure accurate timing and seamless segment concatenation.

2976

2977

Listing 27: SVCSingle Structure

2978

Input: adapted_lyrics_string, midi_analysis_path, song_name

2979

Output: synthesized_vocal_audio_path

2980

Algorithm:

2981

1. Input validation and MIDI analysis loading:

2982

a. Validate adapted lyrics string and analysis file path

2983

b. Load MIDI analysis JSON with UTF-8 encoding

2984

c. Extract notes and duration data from analysis structure

2985

d. Prepare input structure with text, notes, and timing information

2986

2987

2. Lyrics segmentation and alignment processing:

2988

Function `_split_single_annotation()`:

2989

a. Split text by "AP" markers for phrase boundary detection

2990

b. Filter empty segments and preserve AP delimiters

2991

c. Convert non-AP segments to character-level arrays

2992

d. Generate aligned `notes_list` and `notes_duration_list`

2993

e. Validate index correspondence between text and musical data

2994

2995

3. Segment creation with minimum duration optimization:

2996

Function `_create_segment_with_min_duration()`:

2997

a. Apply duration threshold: 0.15s minimum per segment

2998

b. For segments below threshold, aggregate consecutive characters:

2999

- Combine text content until minimum 0.5s duration

3000

- Concatenate notes using " | " separator

3001

- Merge duration values for unified segment timing

3002

c. Preserve AP markers as individual segments for rest periods

3003

3004

4. Batch audio generation with DiffSinger integration:

3005

a. Create temporary JSON files for each vocal segment

3006

b. Generate filename format: "{song_name}_part_{start}-{end}.json"

3007

c. Execute `run_diffsinger()` for neural voice synthesis:

3008

- Input: JSON annotation files with text and musical data

3009

- Output: WAV files with synthesized singing voice

3010

d. Handle batch processing for efficient resource utilization

3011

3012

5. Temporal precision audio processing:

3013

Function `_phase_vocoder_stretch()`:

3014

a. Load generated audio using `librosa.load()` at 44.1kHz

3015

b. Calculate stretch factor: `current_duration / target_duration`

3016

c. Apply time stretching with `librosa.effects.time_stretch()`

3017

d. Perform sample-level alignment for precise timing:

3018

- Truncate excess samples if audio too long

3019

- Pad with zeros if audio too short

3020

- Ensure exact sample count matching target duration

3021

6. Audio segment consolidation and timing verification:

3022

a. Initialize `combined_audio` with 44.1kHz `AudioSegment`

3023

b. Process each segment with millisecond-level precision:

- Convert numpy arrays to 16-bit PCM format

- Apply clip normalization: `np.clip(audio, -1.0, 1.0) * 32767`

- Force duration alignment with target millisecond values

```

3024     c. Handle AP segments with precise silence generation
3025     d. Concatenate segments maintaining temporal continuity
3026
3027 7. Quality assurance and output generation:
3028     a. Validate total duration against expected timing
3029     b. Report timing errors and adjustment statistics
3030     c. Export final audio to dataset/mad_svc/cover directory
3031     d. Clean up temporary files and intermediate outputs
3032     e. Return absolute path to synthesized vocal audio
3033
3034 Technical Specifications:
3035 - Synthesis engine: DiffSinger neural singing voice model
3036 - Audio format: 44.1kHz WAV with 16-bit PCM encoding
3037 - Temporal precision: Millisecond-level alignment accuracy
3038 - Segment processing: Character-level with minimum duration optimization
3039 - Quality control: Sample-level timing verification and correction
3040
3041 Audio Processing Features:
3042 - Phase vocoder time stretching for duration matching
3043 - Automatic silence insertion for AP markers
3044 - Dynamic segment aggregation for optimal synthesis quality
3045 - Real-time duration monitoring and error reporting
3046 - Batch processing optimization for multiple segment synthesis
3047
3048 Error Handling:
3049 - JSON file validation with UTF-8 encoding support
3050 - Audio file existence verification after synthesis
3051 - Duration mismatch detection with automatic correction
3052 - Temporary file cleanup with exception safety
3053 - Sample rate consistency validation throughout pipeline
3054 - Comprehensive error reporting for debugging and optimization

```

3054 A.7.17 STANDUPCONVERSION

3056 The StandUpConversion agent provides precise temporal mapping for stand-up comedy audio
3057 segments, converting individual audio files into structured timestamp format suitable for video
3058 generation workflows. It analyzes audio duration and creates synchronized timing data for seamless
3059 multimedia editing and video synchronization.

3061 Listing 28: StandUpConversion Structure

```

3062 Input: segment_directory, metadata_path
3063 Output: timestamp_json_path
3064
3065 Algorithm:
3066 1. Input validation and metadata loading:
3067     a. Validate segment directory and metadata file paths
3068     b. Load metadata JSON containing script analysis results:
3069         - tone information for each segment
3070         - text content without formatting markers
3071         - reaction indicators for audience response integration
3072     c. Initialize timing tracking variables and chunk array
3073
3074 2. Audio file analysis and duration calculation:
3075     For each metadata entry:
3076     a. Construct audio file path: {seg_dir}/{index}.wav
3077     b. Load audio data using soundfile.read():
3078         - Extract audio samples and sample rate
3079         - Handle various audio formats with sf compatibility

```

```
3078     c. Calculate precise duration: len(audio_samples) / sample_rate
3079     d. Track cumulative timing for sequential segments
3080
3081 3. Timestamp generation with cumulative timing:
3082     a. Maintain current_time tracker for running total
3083     b. For each segment:
3084         - Calculate end_time = current_time + segment_duration
3085         - Round timestamp to 3 decimal places for precision
3086         - Update current_time for next segment calculation
3087     c. Handle timing gaps and overlaps with precise calculation
3088
3089 4. Chunk data structure creation:
3090     a. Generate structured chunk entries:
3091         {
3092             id: sequential_index (1-based),
3093             timestamp: cumulative_end_time,
3094             content: original_text_content
3095         }
3096     b. Preserve original text content from metadata analysis
3097     c. Maintain sequential ordering for proper video synchronization
3098
3099 5. Error handling and file validation:
3100     a. Validate audio file existence for each expected segment
3101     b. Handle corrupted or missing audio files gracefully:
3102         - Log specific file errors with detailed messages
3103         - Continue processing remaining segments
3104         - Maintain timing accuracy despite missing segments
3105     c. Apply robust audio format compatibility checking
3106
3107 6. JSON output generation and export:
3108     a. Structure final timestamp data:
3109         sentence_data: {
3110             count: total_valid_chunks,
3111             chunks: [chunk_array_with_timing_data]
3112         }
3113     b. Export to timestamps.json in metadata directory
3114     c. Ensure UTF-8 encoding with ensure_ascii=False
3115     d. Apply proper JSON formatting with 2-space indentation
3116
3117 7. Timing validation and quality assurance:
3118     a. Verify sequential timestamp progression
3119     b. Check for negative durations or timing anomalies
3120     c. Validate total duration against expected performance length
3121     d. Report timing statistics for quality control
3122
3123 Technical Specifications:
3124 - Audio analysis: soundfile library for precise duration calculation
3125 - Timing precision: 3 decimal places for millisecond-level accuracy
3126 - Format compatibility: Multi-format audio support through soundfile
3127 - Data structure: JSON with nested sentence_data organization
3128 - Error tolerance: Graceful handling of missing or corrupted segments
3129
3130 Timestamp Features:
3131 - Cumulative timing calculation for proper sequence alignment
3132 - Sample-rate aware duration computation for accuracy
3133 - Sequential ID assignment for video editing reference
3134 - Content preservation from original script analysis
3135 - Robust error handling for production workflow reliability
```

- 3132 Video Synchronization Preparation:
- 3133 - Compatible format for VideoConversion agent integration
 - 3134 - Precise timing data for frame-accurate video editing
 - 3135 - Content mapping for automated video scene generation
 - 3136 - Seamless transition support for multimedia workflows
 - 3137 - Quality assurance validation for downstream processing

- 3138 Error Handling:
- 3139 - Audio file existence verification with detailed error reporting
 - 3140 - Sample rate consistency checking across segments
 - 3141 - Metadata correlation validation with audio file availability
 - 3142 - JSON export error management with fallback procedures
 - 3143 - Comprehensive exception handling for production stability

3145

3146 A.7.18 CROSSTALKCONVERSION

3147

3148 The CrossTalkConversion agent provides precise temporal mapping for crosstalk audio segments,
 3149 converting dual-performer dialogue into structured timestamp format suitable for video genera-
 3150 tion workflows. It analyzes audio duration and creates synchronized timing data with performer
 3151 identification for seamless multimedia editing and video synchronization.

3152

Listing 29: CrossTalkConversion Structure

3153

3154 Input: segment_directory, metadata_path
 3155 Output: timestamp_json_path

3156

Algorithm:

3157

- 3158 1. Input validation and metadata loading:
 - 3159 a. Validate segment directory and metadata file paths
 - 3160 b. Load metadata JSON containing crosstalk analysis results:
 - 3161 - role information for each dialogue segment
 - 3162 - tone data for delivery style reference
 - 3163 - text content without formatting markers
 - 3164 - optional reaction indicators for audience response
 - 3165 c. Initialize timing tracking variables and chunk array
- 3166 2. Audio file analysis and duration calculation:

3167 For each metadata entry:

 - 3168 a. Construct audio file path: {seg_dir}/{index}.wav
 - 3169 b. Load audio data using soundfile.read():
 - 3170 - Extract audio samples and sample rate information
 - 3171 - Handle various audio formats with soundfile compatibility
 - 3172 c. Calculate precise duration: len(audio_samples) / sample_rate
 - 3173 d. Track cumulative timing for sequential dialogue segments
- 3174 3. Timestamp generation with cumulative timing:
 - 3175 a. Maintain current_time tracker for running total duration
 - 3176 b. For each dialogue segment:
 - 3177 - Calculate end_time = current_time + segment_duration
 - 3178 - Round timestamp to 3 decimal places for precision
 - 3179 - Update current_time for next segment calculation
 - 3180 c. Handle timing gaps and overlaps with precise calculation
- 3181 4. Content formatting with performer identification:
 - 3182 a. Generate performer-tagged content structure:
 - 3183 "[{performer_role}] {dialogue_text}"
 - 3184 b. Preserve original role information from metadata:
 - 3185 - Funny guy (dou_gen) performer identification
 - Setup guy (peng_gen) performer identification

3186 c. Maintain dialogue attribution for video scene generation
3187
3188 5. Chunk data structure creation:
3189 a. Generate structured chunk entries:
3190 {
3191 id: sequential_index (1-based),
3192 timestamp: cumulative_end_time,
3193 content: "[role] dialogue_text"
3194 }
3195 b. Preserve performer role tags for video editing reference
3196 c. Maintain sequential ordering for proper dialogue flow
3197
3198 6. Error handling and file validation:
3199 a. Validate audio file existence for each expected segment
3200 b. Handle corrupted or missing audio files gracefully:
3201 - Log specific file errors with detailed messages
3202 - Continue processing remaining dialogue segments
3203 - Maintain timing accuracy despite missing segments
3204 c. Apply robust audio format compatibility checking
3205
3206 7. JSON output generation and export:
3207 a. Structure final timestamp data:
3208 sentence_data: {
3209 count: total_valid_chunks,
3210 chunks: [chunk_array_with_timing_and_roles]
3211 }
3212 b. Export to timestamps.json in metadata directory
3213 c. Ensure UTF-8 encoding with ensure_ascii=False
3214 d. Apply proper JSON formatting with 2-space indentation
3215
3216 8. Timing validation and quality assurance:
3217 a. Verify sequential timestamp progression across dialogue
3218 b. Check for negative durations or timing anomalies
3219 c. Validate total duration against expected performance length
3220 d. Report timing statistics and performer distribution
3221
3222 Technical Specifications:
3223 - Audio analysis: soundfile library for precise duration calculation
3224 - Timing precision: 3 decimal places for millisecond-level accuracy
3225 - Format compatibility: Multi-format audio support through soundfile
3226 - Data structure: JSON with nested sentence_data organization
3227 - Role preservation: Performer identification tags in content field
3228
3229 Crosstalk-Specific Features:
3230 - Dual-performer role tracking for video scene generation
3231 - Performer-tagged content for automated video editing
3232 - Sequential dialogue timing for natural conversation flow
3233 - Error tolerance for production workflow reliability
3234 - Metadata correlation with original script analysis
3235
3236 Video Synchronization Preparation:
3237 - Compatible format for VideoConversion agent integration
3238 - Performer identification for automated scene switching
3239 - Precise timing data for frame-accurate video editing
3240 - Content mapping for dual-performer video generation
3241 - Seamless transition support for dialogue-based multimedia
3242
3243 Error Handling:
3244 - Audio file existence verification with detailed error reporting

- 3240 - Sample rate consistency checking across dialogue segments
- 3241 - Metadata correlation validation with audio file availability
- 3242 - JSON export error management with fallback procedures
- 3243 - Comprehensive exception handling for production stability
- 3244 - Performer role validation for consistent content formatting

3247 A.8 KNOWLEDGE ANALYSIS AGENTS

3248 Knowledge analysis agents are specialized components that process and synthesize multimedia
 3249 content to generate coherent narratives and creative outputs. These agents combine video content
 3250 extraction, narrative summarization, and rhythm-aware storyboard generation to create comprehensive
 3251 multimedia knowledge representations. The knowledge analysis pipeline incorporates semantic
 3252 understanding, temporal synchronization, and creative synthesis to transform raw content into
 3253 structured creative materials suitable for video production workflows.

3255 A.8.1 RHYTHMCONTENTGENERATOR WITH GLOBAL-AWARE MECHANISM

3256 The RhythmContentGenerator extracts video segment content and creates scene-focused narrative
 3257 summaries that incorporate user creative requirements. The Global-Aware Mechanism enhances
 3258 video retrieval by refining raw user input into captions-aware queries through a two-stage process:
 3259 first building and analyzing a comprehensive caption bank from video content, then integrating this
 3260 with rhythm analysis to generate fine-grained, contextually grounded storyboard-based subqueries.

3262 Listing 30: RhythmContentGenerator Structure

```

3263 Input: user_requirements, rhythm_analysis_directory, video_segments_path
3264 Output: video_scene_path, video_summary_path
3265
3266 Algorithm:
3267 1. Content extraction and caption bank construction:
3268   a. Load video segments from kv_store_video_segments.json
3269   b. Extract content from all video segments across videos
3270   c. Transform raw captions to numbered video segments format:
3271     'Caption:' > 'Video Segments {id}:'
3272   d. Build comprehensive caption bank with sequential numbering
3273
3274 2. Video content contextualization:
3275   System: "You are a creative beat sync video producer who is good at
3276   write scenes from ground truth video segments, strictly following
3277   the user's requirements."
3278
3279   Video context prompt template:
3280   "Here is the visual description of all video segments from source video:
3281
3282   {video_summary}
3283
3284   Each video segment represents a scene from the ground truth video.
3285   Later, you'll need to use these video segments content to write
3286   storyboards."
3287
3288 3. Rhythm analysis integration:
3289   a. Parse rhythm_points.json for beat count extraction
3290   b. Load rhythm visualization from rhythm_detection.png
3291   c. Generate rhythm reference context:
3292     "Background Music Visualization with Rhythm Points for Reference:
3293     - Plot shows musical intensity and rhythm patterns over time
     - Peaks represent high-energy moments suitable for impactful scenes
     - Valleys indicate calmer segments for transitional content
     - Use visualization to guide scene pacing and emotional intensity"
```

```

3294
3295 4. Global-aware storyboard generation with conversation state:
3296   Conversation state management:
3297   - Initialize conversation tracking array
3298   - Maintain context across multi-turn interactions
3299   - Store user and assistant messages for continuity
3300
3301   Storyboard creation prompt with global awareness:
3302   "Build rhythm-synchronized video storyboards from ground truth video
3303   segments content, aligning with user's requirements
3304
3305   {rhythm_reference}
3306
3307   Total Scenes Required: {sections_num}
3308   User's creative requests (high priority): '{user_idea}'
3309
3310   Video Storyboards Guidelines:
3311   1. Scene Structure: Begin each with /////, number 1 to {sections_num}
3312   2. Visual Requirements: Detailed character appearances, rich motion
3313   descriptions, no dialogue required
3314   3. Rhythm Integration: Match scene intensity with visualization patterns
3315   4. Content Rules: Max two sentences per scene, focus on visual elements,
3316   maintain narrative flow, base on previous grounded video segments
3317
3318   Format Output: /////\n[Scene description]\n\n///// \n[Scene description]"
3319
3320 5. API interaction with retry mechanism:
3321   @tenacity.retry configuration:
3322   - wait=exponential(multiplier=1, min=2, max=60)
3323   - stop=after_attempt(5)
3324   - before_sleep=logging with attempt number
3325
3326   Conversation-aware API calls:
3327   def _call_gpt_api(user_prompt, temperature=0.7):
3328       self.conversation.append({"role": "user", "content": user_prompt})
3329       if len(self.conversation) > 1:
3330           initial_message = "In previous conversation, you said: " +
3331               assistant_messages[-1] + "\n\n"
3332           combined_prompt = initial_message + user_prompt
3333           response = gpt(model=self.model, system=self.system_message,
3334               user=combined_prompt)
3335           assistant_response = response.choices[0].message.content
3336           self.conversation.append({"role": "assistant", "content":
3337   ↪ assistant_response})
3338
3339 6. Output generation and validation:
3340   a. Create structured JSON with user_idea, video_summary, segment_scene
3341   b. Save storyboard to video_scene.json with UTF-8 encoding
3342   c. Generate preview with truncation for large content
3343   d. Return file paths for downstream video processing
3344
3345   Error Handling:
3346   - Tenacity retry logic for robust API interactions
3347   - Fallback mechanisms for missing rhythm analysis files
3348   - Graceful degradation when video content unavailable
3349   - Boolean conversion for string-based configuration parameters
3350   - Default value assignment for missing beat count data
3351
3352

```

3348 A.8.2 NEWSCONTENTGENERATOR

3349

3350 The NewsContentGenerator creates news summary content by implementing a dual-agent pipeline
 3351 that processes reference video transcripts and adapts them according to user requirements and
 3352 presentation styles. The algorithm employs a presenter agent for content adaptation and a judge
 3353 agent for structural formatting.

3354

Listing 31: NewsContentGenerator Structure

3355

```
3356 Input: user_requirements, news_presentation_style, video_directory
3357 Output: video_scene_json
```

3358

3359 Algorithm:

3360

- 3361 1. Transcript discovery and loading:
 - 3362 a. Scan video directory for .lab files with multiple encoding support
 - 3363 b. Parse .lab file format: [start_time end_time transcription]
 - 3364 c. Extract transcription text while preserving content integrity
 - 3365 d. Handle encoding fallbacks: utf-8, gb18030, gbk, gb2312, cp1252

3366

- 3367 2. Presentation style processing:
 - 3368 a. Load presentation methodology from file or direct string input
 - 3369 b. Parse formatting guidelines and style requirements
 - 3370 c. Prepare template for content adaptation pipeline

3371

- 3372 3. Presenter agent processing with embedded prompt:

3373 System: "You are an experienced expert in news writing skit review
 3374 copy. Pay special attention to user's words count requirements."

3375

3376 User prompt template:
 3377 "Create skit narration copy, strictly following user's ideas and
 3378 presentation methods."

3379

3380 User's idea: '{user_idea}'
 3381 Grounded text content: {content}
 3382 Follow this presentation method: {present_content}

3383

- 3384 Requirements:
- 3385 1. Format: Remove section numbers, min 11 words per sentence
 - 3386 2. Content: Use original dialogues, focus on plot elements
 - 3387 3. Language: Third-person, convert numbers to words, separate
 3388 abbreviations (ChatGPT --> Chat GPT)"

3389

- 3390 4. Judge agent formatting with embedded prompt:

3391 System: "You are a content formatting specialist with expertise
 3392 in following guidelines"

3393

3394 User prompt template:
 3395 "Format content with requirements:
 3396 - Remove all commas
 3397 - Start with ////\n
 3398 - Chunk each period with \n\n////\n
 3399 - Keep original content, separate sentences

3400

3401 Example: ////\nGood morning everyone nice to meet you again.
 \n\n////\nThe weather is very nice today."

3402

- 3403 5. Visual scene generation:
 - 3404 a. Process formatted content for scene extraction
 - 3405 b. Generate English visual-scene keywords and descriptions
 - 3406 c. Apply scene translation with embedded prompt:

3402 "Deduce visual-scene keywords, each section some scene keywords
3403 (proper nouns: iPhone 16, SWE Arena Benchmark)
3404 Keep same paragraph separators, max 1 sentence per section"
3405
3406 6. Content validation and output:
3407 a. Count content sections using ///// markers
3408 b. Generate structured JSON with user_idea, content_created,
3409 segment_scene
3410 c. Save to video_scene.json with UTF-8 encoding
3411 d. Provide section statistics and preview
3412
3412 Error Handling:
3413 - Multi-encoding file reading with graceful fallbacks
3414 - API retry logic with exponential backoff (5 attempts)
3415 - Content truncation for oversized inputs (15K character limit)
3416 - Fallback formatting for judger agent failures
3417
3418

3419 A.8.3 NEWS PRESENTATION STYLE SYSTEM PROMPT

3420 This subsection outlines the standardized presentation methodology for video event overview content
3421 generation, ensuring consistent narrative structure and professional delivery across all multimedia
3422 production workflows.
3423

3424 Listing 32: Video Event Overview Presentation Style

3425 Core Requirements:
3426 1. Third-person narrative perspective throughout content
3427 2. Factual accuracy: avoid fabricated specifications (e.g., iPhone 15, not
3428 ↪ iPhone 15.0)
3429 3. Proper subject attribution: understand main characters (Apple released
3430 ↪ iPhone 4,
3431 not "iPhone 4 released new product")
3432 4. No greeting/farewell formalities (avoid "good morning," "thank you for
3433 ↪ watching")
3434 5. Focus: comprehensive video event overview generation
3435
3435 OPENING SECTION (First Three Sentences):
3436 Objective: Create compelling hook with core concept presentation
3437 - Immediately establish central narrative tension
3438 - Highlight primary innovation or conflict
3439 - Engage audience with compelling opening statement
3440
3441 Example Structure:
3442 "Meta just announced revolutionary virtual reality equipment. It can visualize
3443 immersive environments in real time that would require fastest GPUs several
3444 ↪ hours
3445 to render. This breakthrough addresses accelerated processing challenges
3446 ↪ called [specific technology]."
3447
3447 Alternative Opening Approaches:
3448 - Situational context establishment
3449 - Hypothetical scenario presentation
3450 - Generational impact framing
3451 - Direct technical revelation
3452
3453 CONTENT BODY SECTION (Primary Narrative):
3454 Structure: Timeline-based plot advancement with clear progression
3455 - Unfold events chronologically with logical sequence
- Maintain clear narrative thread throughout

- 3456 - Continuously advance plot without stagnation
- 3457 - Avoid generalized event summaries or character abstractions
- 3458 - Incorporate key character dialogue appropriately
- 3459 - Use accessible language: "really," "and," "how," "if," "because," "but"

3460

3461 Prohibited Elements:

- 3462 - Analytical reflections or thematic summaries
- 3463 - Macro-level content summarization
- 3464 - Theme or character development overviews
- 3465 - Generalized terminology: "The event shows how...", "showing the...",
- 3466 "episodes...", "reveals...", "demonstrates..."
- 3467 - Direct mention of "event" within body text
- 3468 - Literary analysis or interpretive commentary

3469

3470 Required Elements:

- 3471 - Tight pacing with smooth scene transitions
- 3472 - Constant plot advancement and momentum
- 3473 - Strong character development and interaction
- 3474 - Factual content delivery without speculation

3475

3476 CLOSING SECTION (Final Sentences):

3477 Objective: Thematic elevation without content summarization

- 3478 - Refine deeper meaning and significance
- 3479 - Sublimate central themes naturally
- 3480 - Avoid high-level content recapping
- 3481 - Connect to broader implications or impact
- 3482 - Maintain narrative flow to conclusion

3483

3484 Quality Assurance Standards:

- 3485 - Grammatical accuracy and professional language use
- 3486 - Consistent third-person perspective maintenance
- 3487 - Factual verification of all technical specifications
- 3488 - Character attribution accuracy verification
- 3489 - Prohibition compliance monitoring for banned phrases
- 3490 - Timeline coherence and logical progression validation

3491

3492

3493 A.8.4 COMMENTARYCONTENTGENERATOR

3494

3495 The CommentaryContentGenerator creates commentary content from text source materials with specialized formatting for video presentations. It employs a dual-agent architecture to transform source texts into structured video-ready narratives while maintaining user-specified creative requirements.

3496

3497

Listing 33: CommentaryContentGenerator Structure

3498

3499

3500

3501

3502

3503

3504

3505

3506

3507

3508

3509

Input: user_requirements, source_text_path, commentary_presentation_style
Output: video_scene_json

Algorithm:

1. Source text processing:
 - a. Load source text with multi-encoding support (utf-8, gb18030, gbk, etc.)
 - b. Validate content size and apply truncation if exceeding 30K characters
 - c. Generate text statistics for processing optimization
 - d. Handle encoding fallbacks with binary reading approach
2. Presenter agent processing with embedded prompts:

System: "You are an experienced expert in writing review copy.
Pay special attention to user's words count requirements."

User prompt template:

```
3510 "Create narration copy, strictly following user's ideas and
3511 presentation methods.
3512
3513 User idea: '{user_idea}'
3514 Grounded text content: {content}
3515 Follow this presentation method: {present_content}
3516
3517 Requirements:
3518 1. Format: Response in source language, remove chapter numbers,
3519 max 3 commas per sentence
3520 2. Content: Follow word count requirements, use original dialogues,
3521 focus on plot elements
3522 3. Language: Clear narrative flow, no user requirement mentions"
3523
3524 3. Judger agent formatting with embedded prompts:
3525 System: "You are a content formatting specialist with expertise
3526 in following guidelines"
3527
3528 User prompt template:
3529 "Format content with requirements:
3530 - Start each sentence with ////
3531 - Remove chapter numbers and punctuation after segmentation
3532 - Keep original content, separate each sentence
3533 - Purpose: sentence-level segmentation
3534
3535 Example: ////\nsentence one \n\n////\nsentence two"
3536
3537 4. Visual scene description generation:
3538 a. Process formatted content for scene inference
3539 b. Apply scene description generation with embedded prompt:
3540 "Convert to English visual-scene descriptions:
3541 - Keep //// markers unchanged
3542 - Max 1 sentence per section
3543 - Replace low-quality descriptions with conflict scenes
3544 - Describe character appearances when names mentioned
3545 - Generate high-ignition story moments"
3546
3547 5. Content validation and structuring:
3548 a. Count content sections using //// markers
3549 b. Validate scene coherence and visual consistency
3550 c. Generate structured JSON output:
3551 - reqs: user requirements
3552 - content_created: formatted narrative content
3553 - segment_scene: visual scene descriptions
3554 d. Save to video_scene.json with UTF-8 encoding
3555
3556 6. Quality assurance and reporting:
3557 a. Display section statistics and preview
3558 b. Validate content structure integrity
3559 c. Report processing completion status
3560
3561 Error Handling:
3562 - Multi-encoding file reading with graceful degradation
3563 - Content truncation for oversized inputs (30K limit)
3564 - API retry logic with exponential backoff (5 attempts)
3565 - Fallback formatting for agent processing failures
3566 - Binary reading fallback for encoding issues
```

3564 A.8.5 COMMENTARY PRESENTATION STYLE SYSTEM PROMPT

3565
3566 This subsection provides comprehensive guidelines for creating engaging narrative content with
3567 proper structure, pacing, and audience engagement techniques.
3568

3569 Listing 34: Content Writing Methodology
3570

3571 OPENING SECTION: First three sentences with compelling attraction, quickly
3572 ↪ highlighting core elements
3573
3574 MAIN CONTENT SECTION (Majority of content, maintain plot progression in latter
3575 ↪ half, strictly adhere to word count requirements):
3576 - Unfold according to plot timeline with clear main thread
3577 - Smooth chapter transitions, continuous plot advancement
3578 - Avoid generalizing events and characters
3579 - Appropriately use key character dialogue
3580 - Tight pacing, smooth scene transitions, continuous plot development,
3581 ↪ distinct character portrayal
3582 - Maximum 3 commas per sentence
3583
3584 PROHIBITED ELEMENTS:
3585 - Literary analysis, reflection, and thematic summaries
3586 - Macro-level generalizations about story themes or character development
3587 - Forbidden phrases: "The story shows...", "demonstrates...", "plot...", "
3588 ↪ reveals...", "someone's story..."
3589 - Cannot mention "story" in main content
3590
3591 OPENING SUMMARY SECTION:
3592 First three sentences provide content overview, followed by 2-3 introductory
3593 ↪ sentences before narrative begins
3594
3595 Example 1: Situational and Immersive Openings
3596 Identity Immersion:
3597 "His name is Wang Xiaoming. You think he's just an ordinary college student?
3598 ↪ No. He's a time traveler from the future, shouldering the mission to change
3599 ↪ history."
3600
3601 Knowledge-Guided:
3602 "Do you know what the world's most dangerous job is? Turns out it's
3603 ↪ maintenance engineers at Antarctic research stations, dancing with death
3604 ↪ daily in minus 70-degree environments."
3605
3606 Hypothetical Scenarios:
3607 "If someone gave you 100 million dollars, but the condition was never using
3608 ↪ social media again, would you accept? For modern people, this might be the
3609 ↪ cruelest multiple choice question."
3610
3611 Role-Playing:
3612 "If you were a newly graduated medical student suddenly facing an
3613 ↪ unprecedented pandemic outbreak, what would you do? This was exactly Dr. Li
3614 ↪ 's situation in early 2020."
3615
3616 Time-Travel Hypothetical:
3617 "If you traveled back to the Titanic in 1912, knowing the ship would sink in
3618 ↪ four days, how would you save over 2000 lives without being labeled insane
3619 ↪ ?"
3620
3621 Example 2: Theme-First Openings
3622 Real Event Revelation:

3618 "This isn't fiction; this is real. An engineer was trapped in Mount Everest's
3619 ↪ death zone for 72 hours straight-no oxygen, no supplies. Doctors said he
3620 ↪ couldn't survive. However, a miracle happened."
3621

3622 Ability Paradox:

3623 "If one day you could suddenly see others' death times, this wouldn't be a
3624 ↪ superpower but torture. Julian possesses this 'curse,' watching his lover's
3625 ↪ countdown decrease daily while remaining powerless."
3626

3627 Expectation Reversal:

3628 "He's two meters tall, muscular, with piercing eyes. Is this a professional
3629 ↪ boxer? No, he's a world-famous children's book author whose gentle stories
3630 ↪ have accompanied generations."
3631

3632 Extreme Situations:

3633 "Minus 40 degrees, no food, pitch black, with only aurora as light source. How
3634 ↪ does one survive 30 days in Arctic wilderness? This was extreme explorer
3635 ↪ Zhang San's challenge, and his choices completely changed survival theory."
3636

3637 Suspenseful Questions:

3638 "What could make a billionaire abandon all wealth to live as a hermit in deep
3639 ↪ mountains for 20 years? Today, we reveal the shocking truth behind this
3640 ↪ Silicon Valley legend's decision."
3641

3642 MAIN CONTENT SECTION (Majority of content, maintain plot progression in latter
3643 ↪ half):

3644 Unfold according to plot timeline with clear main thread, ensuring continuous
3645 ↪ advancement in latter sections
3646

3647 PROHIBITED ELEMENTS:

- 3648 - Literary analysis, reflection, and thematic summaries
- 3649 - Macro-generalizations about story themes or character development
- 3650 - Forbidden phrases: "The story shows...", "demonstrates...", "plot...", "
3651 ↪ reveals...", "someone's story..."
- 3652 - Cannot mention "story" in main content
3653

3654 REQUIRED TRANSITION WORDS:

3655 Use common transitional words: "unexpectedly," "suddenly," "turns out," "but,"
3656 ↪ "however," "yet," "as a result," "finally," "until," "if," "while," "
3657 ↪ indeed," "discovered," "only," "surprisingly," "afterward," "exactly," "not
3658 ↪ only," "nevertheless," "little did they know," "moreover," "of course," "
3659 ↪ because," "therefore," etc.
3660

3661 NARRATIVE TECHNIQUES:

- 3662 - Dense use of transition words with high frequency of "but," "yet," "however"
- 3663 - Unexpected twists for dramatic effect using "turns out," "little did they
3664 ↪ know"
- 3665 - Plot advancement after each transition bringing new developments
- 3666 - Concise action descriptions with short sentences for complex actions
- 3667 - Sensory descriptive words enhancing immersion
- 3668 - Frequent emotional adjectives
- 3669 - Emphatic words highlighting characteristics
- 3670 - Contrasting words creating dramatic tension
- 3671 - Narrative connective words for smooth scene transitions
- Conversational tone with natural, approachable language
- Short sentences for rapid pacing at key moments
- Repetitive sentence structures for emphasis
- Rhetorical devices including metaphors and hyperbole
- Strategic suspense placement

3672 - Dramatic contrast throughout
 3673
 3674 CLOSING SECTION (Final sentence):
 3675 Extract deeper meaning, elevate themes without mentioning "story ending" or "
 3676 ↪ story beginning" to avoid breaking audience immersion
 3677

3678
 3679 A.8.6 TTSWRITER

3680 The TTSWriter agent provides intelligent text rewriting capabilities for creating derivative video
 3681 content based on sliced audio transcripts. It employs a dual-LLM approach using Claude for creative
 3682 content generation and DeepSeek for precise text extraction, enabling coherent content transformation
 3683 while maintaining original speech patterns and timing constraints.
 3684

3685 Listing 35: TTSWriter Structure

3686 Input: user_requirements, audio_path
 3687 Output: rewritten_speech_path
 3688
 3689 Algorithm:
 3690 1. Input validation and file path resolution:
 3691 a. Validate user requirements and audio file path
 3692 b. Generate lab_path: audio_filename + '.lab' extension
 3693 c. Extract slice_lab_dir from audio path without extension
 3694 d. Verify transcription file availability and accessibility
 3695
 3696 2. Transcript aggregation and segmentation analysis:
 3697 a. Scan slice_lab_dir for .lab files with sorted ordering
 3698 b. Load each lab file content with UTF-8 encoding
 3699 c. Aggregate slice_lab array with individual segment transcripts
 3700 d. Load complete original transcript from main .lab file
 3701 e. Generate structured output format for LLM processing
 3702
 3703 3. Creative content generation with Claude LLM:
 3704 System context: "Parody text recreation expert, specializing in generating
 3705 text suitable for new scenarios"
 3706
 3707 User prompt structure:
 3708 "Perform creative recreation based on the following points and user
 3709 requirements:
 3710 1. Ensure smooth text flow between slices while recreating each slice
 3711 2. Imitate the language style and sentence structure of each slice,
 3712 only replace specific content
 3713 3. When replacing original text vocabulary, length variation should
 3714 not exceed two characters
 3715
 3716 User requirements: {user_requirements}
 3717 Original text: {original_transcript}
 3718 Output format: {structured_format}"
 3719
 3720 4. Structured content generation workflow:
 3721 a. Format template creation for each audio slice:
 3722 "{index}. Original slice content: {original_content}\n Recreation:"
 3723 b. Apply user requirements to content transformation guidelines
 3724 c. Maintain linguistic style and sentence structure consistency
 3725 d. Enforce character length constraints (+2 characters maximum)
 e. Generate comprehensive rewritten content with original context
 3726
 3727 5. Text extraction and refinement with DeepSeek LLM:
 System context: "Text extraction expert"

3726

3727

User prompt structure:

3728

"Extract the ****recreation**** content of each slice and output line by line

3729

Output format:

3730

Recreation of slice 1

3731

Recreation of slice 2

3732

3733

Text to be extracted: {generated_content}"

3734

3735

6. Dual-stage output processing and file generation:

3736

a. Save raw generation output to 'raw_speech.txt'

3737

b. Apply extraction LLM for clean content isolation

3738

c. Generate final 'speech.txt' with line-separated segments

3739

d. Ensure UTF-8 encoding for multilingual content support

3740

e. Return speech_path for downstream TTS processing

3741

Content Transformation Guidelines:

3742

- Linguistic style preservation: Maintain original speech patterns

3743

- Structural consistency: Preserve sentence construction and rhythm

3744

- Length constraints: Maximum +-2 character variation per replacement

3745

- Semantic adaptation: Transform content while preserving meaning flow

3746

- Temporal alignment: Ensure rewritten segments fit original timing

3747

Quality Assurance Features:

3748

- Dual-LLM validation for content accuracy and extraction precision

3749

- Structured format enforcement for consistent processing

3750

- Character-level length monitoring for TTS compatibility

3751

- Cultural context preservation for appropriate content adaptation

3752

- Error handling for LLM API failures with graceful degradation

3753

3754

3755

Error Handling:

3756

- Lab file existence validation with detailed error reporting

3757

- UTF-8 encoding enforcement for international character support

3758

- LLM API failure detection with exception management

3759

- File I/O error handling for transcript processing

3760

- Directory structure validation for slice organization

3761

3762

3763

A.8.7 SVCADAPTER

3764

The SVCAdapter agent provides sophisticated lyrics adaptation capabilities for music cover creation,

3765

maintaining original melody structure while enabling creative content transformation. It employs

3766

a multi-stage approach using dual-LLM processing for high-quality lyrical recreation with precise

3767

character count alignment and structural preservation.

3768

Listing 36: SVCAdapter Structure

3769

Input: user_requirements, midi_analysis_path, song_name

3770

Output: adapted_lyrics_string

3771

Algorithm:

3772

1. Input validation and MIDI analysis processing:

3773

a. Validate user requirements and file path accessibility

3774

b. Load MIDI analysis JSON with UTF-8 encoding

3775

c. Extract original lyrics text from analysis data

3776

d. Initialize song name for output file generation

3777

3778

2. Lyrics structure parsing and segmentation:

3779

Function parse_lyrics_structure(lyrics):

a. Split lyrics by 'AP' delimiter markers for phrase separation

3780 b. Create structure array with "LYRICS" placeholders and "AP" markers
3781 c. Extract lyrics_parts array containing actual lyrical content
3782 d. Maintain temporal alignment between structure and content
3783

3784 3. Template generation for controlled adaptation:
3785 Function generate_lyrics_template(lyrics_parts):
3786 a. Create numbered segments with original lyrics content
3787 b. Calculate character count constraints for each segment
3788 c. Generate structured format:
3789 "{index}. Original lyrics segment: {content}
3790 Character limit: {count}
3791 Recreation:"
3792 d. Ensure precise character count preservation

3793 4. Full lyrics generation with Claude LLM:
3794 System context: "Professional lyrics adaptation AI for high-quality
3795 ↪ recreation"
3796

3797 User prompt structure:
3798 "Perform recreation based on following points and user requirements:
3799 1. Strictly follow character count limits for each recreation segment
3800 2. Focus on rhyme and rhythm while maintaining narrative flow
3801 3. Ensure complete semantic integrity in lyrical content
3802 4. Use reasonable word combinations and inter-segment rhyming

3803 User requirements: {user_requirements}
3804 Original lyrics: {original_lyrics}
3805 Output format: {template}"
3806

3807 5. Content extraction and refinement with DeepSeek LLM:
3808 System context: "Lyrics extraction expert"
3809

3810 User prompt structure:
3811 "Extract the **recreation** content from each lyrics segment:
3812 Output format:
3813 Recreation of segment 1
3814 Recreation of segment 2

3815 Text to extract: {generated_lyrics}"
3816

3817 6. Character count alignment and quality assurance:
3818 Function align_extract_parts():
3819 a. Compare extract_parts length with lyrics_parts requirements
3820 b. For mismatched character counts (max 5 retry attempts):
3821 - Generate context-aware alignment prompts with surrounding segments
3822 - Apply Claude LLM for precise character count correction
3823 - Maintain semantic coherence with previous/next segments
3824 c. Apply fallback padding/truncation if alignment fails:
3825 - Add Chinese character "la" characters for insufficient length
3826 - Truncate excess characters while preserving meaning

3827 7. Structure reconstruction and output generation:
3828 a. Merge aligned lyrics parts with original structure markers
3829 b. Replace "LYRICS" placeholders with adapted content segments
3830 c. Preserve "AP" delimiters for proper phrase separation
3831 d. Generate multiple output formats:
3832 - raw_lyrics.txt: initial generation output
3833 - lyrics.txt: line-separated adapted segments
 - script.txt: complete restructured lyrics

```

3834     - {song_name}_cover.json: updated MIDI analysis with new lyrics
3835
3836 Technical Specifications:
3837 - Character preservation: Exact count matching for musical timing
3838 - Structural integrity: AP delimiter preservation for phrase boundaries
3839 - Multi-format output: TXT and JSON formats for different use cases
3840 - Quality assurance: Iterative refinement with context-aware correction
3841 - Semantic coherence: Meaning preservation across segment boundaries
3842
3843 Adaptation Features:
3844 - Rhyme scheme preservation with enhanced creativity
3845 - Rhythmic pattern maintenance for musical compatibility
3846 - Context-aware segment generation with surrounding lyrical flow
3847 - Automated character count correction with intelligent padding
3848 - Multi-stage validation for lyrical quality and structural accuracy
3849
3850 Error Handling:
3851 - JSON parsing validation with encoding error management
3852 - Character count mismatch detection with automated correction
3853 - LLM API failure handling with graceful degradation
3854 - File I/O error management for multiple output formats
3855 - Structure length validation with detailed error reporting

```

3856 A.8.8 STANDUPADAPTER

3857
3858 The StandUpAdapter agent provides sophisticated comedy script adaptation capabilities for trans-
3859 forming reference content into structured stand-up comedy performances. It employs Claude LLM
3860 for creative content generation while maintaining professional comedy formatting standards including
3861 tone markers and atmosphere cues for enhanced delivery guidance.

3862 Listing 37: StandUpAdapter Structure

```

3864 Input: user_requirements, reference_script_path
3865 Output: segmented_standup_script
3866
3867 Algorithm:
3868 1. Input validation and reference script loading:
3869     a. Validate user requirements and script file path
3870     b. Load reference script content with UTF-8 encoding
3871     c. Extract data directory path for output file placement
3872     d. Strip whitespace and prepare content for adaptation
3873
3874 2. Stand-up comedy format specification:
3875     a. Define mandatory tone markers:
3876         - [Natural]: conversational delivery style
3877         - [Confused]: questioning or bewildered tone
3878         - [Empathetic]: understanding or sympathetic delivery
3879         - [Exclamatory]: energetic or surprised expression
3880     b. Configure atmosphere cues for audience interaction:
3881         - [Laughter]: comedic punchline moments
3882         - [Cheers]: celebratory or triumphant points
3883
3884 3. Creative adaptation with Claude LLM:
3885     System context: "Professional stand-up comedy adaptation specialist"
3886
3887     User prompt structure:
3888     "Adapt the following reference script into stand-up comedy format.
3889
3890     Content to adapt: {reference_script}

```

3888 Additional requirements: {user_requirements}
3889
3890 Format specifications:
3891 1. Each line must begin with tone markers: [Natural] [Confused]
3892 [Empathetic] [Exclamatory]
3893 2. Add atmosphere cues [Laughter] or [Cheers] at key moments
3894 3. Keep each line independent with consistent structure
3895
3896 Important notes:
3897 - NO titles, introductions, or conclusions
3898 - Preserve core humor while localizing cultural references
3899 - Incorporate English stand-up linguistic features and rhythm
3900 - Use atmosphere cues sparingly (3-4 total)
3901 - Generate 3-5 minute performance script"
3902
3903 4. Comedy structure optimization:
3904 a. Preserve original humor essence while adapting format
3905 b. Localize cultural references for target audience
3906 c. Incorporate stand-up comedy linguistic patterns:
3907 - Timing-based humor with pause indicators
3908 - Observational comedy structure
3909 - Callback references and running gags
3910 d. Apply rhythm and pacing appropriate for live performance
3911
3912 5. Tone marker integration and line formatting:
3913 a. Ensure every line begins with appropriate tone marker
3914 b. Match tone markers to content emotional context
3915 c. Structure independent lines for performance flexibility:
3916 "[Tone]content...\n[Tone]content...[Atmosphere]"
3917 d. Maintain comedic flow between individual segments
3918
3919 6. Atmosphere cue placement and timing:
3920 a. Identify punchline moments for [Laughter] placement
3921 b. Recognize triumph or celebration points for [Cheers]
3922 c. Apply strategic placement for maximum comedic impact:
3923 - Immediately after dialogue delivery
3924 - At peak comedic tension release points
3925 d. Limit to 3-4 cues total to avoid oversaturation
3926
3927 7. Output generation and file management:
3928 a. Generate complete stand-up script with title
3929 b. Format output structure:
3930 "# title\n[Tone]content...\n[Tone]content...\n..."
3931 c. Save to 'stand-up.txt' in reference script directory
3932 d. Return script string for immediate use or further processing
3933
3934 Technical Specifications:
3935 - Content adaptation: Reference-based creative transformation
3936 - Format compliance: Strict tone marker and atmosphere cue requirements
3937 - Performance duration: 3-5 minute script generation target
3938 - Linguistic adaptation: English stand-up comedy style integration
3939 - Cultural localization: Reference adaptation for target audience
3940
3941 Comedy Adaptation Features:
- Humor preservation with format transformation
- Cultural reference localization for audience relevance
- Professional stand-up delivery guidance through tone markers
- Audience interaction planning with atmosphere cues
- Independent line structure for performance flexibility

3942
3943
3944
3945
3946
3947
3948
3949
3950
3951
3952
3953
3954
3955
3956
3957
3958
3959
3960
3961
3962
3963
3964
3965
3966
3967
3968
3969
3970
3971
3972
3973
3974
3975
3976
3977
3978
3979
3980
3981
3982
3983
3984
3985
3986
3987
3988
3989
3990
3991
3992
3993
3994
3995

Error Handling:

- Reference script file validation with encoding verification
- Claude LLM API error detection with graceful failure handling
- Output file creation with directory permission management
- UTF-8 encoding preservation for international character support
- Exception isolation with detailed error reporting for debugging

A.8.9 STANDUPSYNTH

The StandUpSynth agent provides comprehensive audio synthesis for stand-up comedy performances, combining advanced text-to-speech generation with audience reaction integration. It employs CosyVoice2 for high-quality speech synthesis and implements intelligent script parsing to create realistic comedy performances with appropriate vocal delivery and audience responses.

Listing 38: StandUpSynth Structure

Input: segmented_script, target_vocal_directory, reaction_directory
Output: merged_audio_path, segment_directory, metadata_path

Algorithm:

1. Input validation and environment setup:
 - a. Validate script content and directory paths
 - b. Convert target_vocal_dir and reaction_dir to absolute paths
 - c. Navigate to tools/CosyVoice directory for model execution
 - d. Initialize CosyVoice2 model with pretrained weights:
 - Model: CosyVoice2-0.5B for high-quality synthesis
 - Configuration: load_jit=False, load_trt=False, fp16=False
2. Script parsing and segment analysis:
 - a. Split script by newlines and filter empty lines
 - b. Skip first line (title) and process content lines
 - c. Create segment directory structure for individual audio files
 - d. Initialize counters and result arrays for processing tracking
3. Intelligent script analysis with DeepSeek LLM:

System context: "Stand-up comedy segment analyzer"

User prompt structure:
"Analyze tone, text content, and atmosphere marker:
{comedy_line}"

Output strictly in JSON format:

 1. tone: 'Natural', 'Empathetic', 'Confused', or 'Exclamatory'
 2. text: segment content without markers
 3. reaction: 'Laughter' or 'Cheers' (only if atmosphere marker present)
 4. Strict reliance on existing markers, no interpretation
 5. NO extra characters or explanations"
4. JSON response processing and validation:
 - a. Strip markdown code block markers (```json, ```)
 - b. Parse JSON response with error handling
 - c. Extract tone, text, and optional reaction fields
 - d. Validate tone values against supported categories
 - e. Handle malformed JSON with graceful error reporting
5. Voice synthesis with tone-specific prompts:
 - a. Load corresponding prompt files:
 - Prompt text: {tone}.lab file for speech context

3996 - Prompt audio: {tone}.wav file for voice characteristics
3997 b. Apply CosyVoice2 zero-shot synthesis:
3998 - Input: parsed text content
3999 - Context: tone-specific prompt text and audio
4000 - Output: high-quality speech audio at model sample rate
4001 c. Save individual segments: {segment_index}.wav
4002
4003 6. Audience reaction integration:
4004 a. Check for reaction field in parsed JSON
4005 b. Load corresponding reaction audio: {reaction}.wav
4006 c. Combine speech and reaction using AudioSegment:
4007 - Load original synthesized speech
4008 - Load reaction audio (laughter/cheers)
4009 - Concatenate: original_audio + reaction_audio
4010 - Export combined audio overwriting original segment
4011
4012 7. Audio consolidation and final merge:
4013 Function merge_audio_files():
4014 a. Initialize silent AudioSegment for accumulation
4015 b. Iterate through numbered segment files (0.wav to cnt.wav)
4016 c. Load and concatenate each segment with error handling
4017 d. Export final merged audio to ../final/stand_up.wav
4018 e. Return absolute path for downstream processing
4019
4020 8. Metadata generation and cleanup:
4021 a. Compile processing results into structured JSON
4022 b. Save metadata: tone, text, and reaction information per segment
4023 c. Export to stand-up.json in target vocal directory
4024 d. Restore original working directory after processing
4025 e. Return comprehensive output paths and metadata
4026
4027 Technical Specifications:
4028 - Speech synthesis: CosyVoice2-0.5B with zero-shot voice cloning
4029 - Audio processing: PyDub for segment manipulation and concatenation
4030 - Script parsing: DeepSeek LLM for intelligent content analysis
4031 - Output format: WAV files with model native sample rate
4032 - Metadata: JSON structure with segment-level processing information
4033
4034 Comedy Performance Features:
4035 - Tone-specific voice synthesis for natural delivery variation
4036 - Audience reaction integration for realistic performance atmosphere
4037 - Segment-based processing for precise timing control
4038 - Automatic tone detection from script markers
4039 - Professional audio quality with seamless segment transitions
4040
4041 Error Handling:
4042 - CosyVoice2 model loading validation with detailed error reporting
4043 - JSON parsing with markdown cleanup and format validation
4044 - Audio file loading verification for segments and reactions
4045 - Graceful handling of missing reaction files or invalid tones
4046 - Comprehensive exception management with processing continuation
4047 - Working directory restoration for environment consistency

4045
4046
4047
4048
4049

A.8.10 CROSSTALKADAPTER

The CrossTalkAdapter agent provides specialized script adaptation capabilities for transforming reference content into traditional Chinese crosstalk (xiangsheng) dialogue format. It employs Claude LLM

4050 for cultural localization and dialogue structure optimization while maintaining authentic crosstalk per-
 4051 formance characteristics including role-specific delivery patterns and traditional interactive elements.

4052 Listing 39: CrossTalkAdapter Structure

```

4054 Input: user_requirements, reference_script_path, dou_gen_directory,
4055 ↪ peng_gen_directory
4056 Output: segmented_crosstalk_script
4057
4058 Algorithm:
4059 1. Input validation and role configuration:
4060   a. Validate user requirements and reference script file path
4061   b. Extract performer directories for role-specific voice synthesis
4062   c. Generate role names from directory basenames:
4063      - dou_gen_name: comic lead performer (funny guy)
4064      - peng_gen_name: straight man performer (setup guy)
4065   d. Load reference script content with UTF-8 encoding
4066
4067 2. Traditional crosstalk role definition:
4068   a. Configure dou_gen (funny guy) characteristics:
4069      - Comic lead role delivering main jokes
4070      - Drives narrative progression and humor development
4071      - Primary responsibility for punchline delivery
4072   b. Configure peng_gen (setup guy) characteristics:
4073      - Straight man role providing reactions and setup
4074      - Creates opportunities for comic lead responses
4075      - Maintains dialogue rhythm and audience engagement
4076
4077 3. Cultural adaptation with Claude LLM:
4078   System context: "Professional cross talk (xiang sheng) adaptation
4079   ↪ specialist"
4080
4081   User prompt structure:
4082   "Adapt English stand-up comedy material into authentic traditional
4083   Chinese crosstalk dialogue format.
4084
4085   Material to adapt: {reference_script}
4086
4087   Crosstalk roles:
4088   - {dou_gen_name}: Comic lead (funny guy), delivers main jokes and drives
4089   ↪ narrative
4090   - {peng_gen_name}: Straight man (setup guy), reacts and plays off comic
4091   ↪ lead
4092
4093   Format Requirements:
4094   1. Each performer's lines on separate lines starting with their name
4095   2. Begin each line with tone marker: [Natural] [Confused] [Emphatic]
4096   3. Same tone should not appear consecutively for more than two lines
4097
4098   Additional requirements: {user_requirements}"
4099
4100 4. Dialogue structure optimization:
4101   a. Implement traditional crosstalk format requirements:
4102      - Line-by-line role alternation for dynamic interaction
4103      - Tone marker integration: [Natural], [Confused], [Emphatic]
4104      - Consecutive tone limitation to prevent monotony
4105   b. Preserve cultural authenticity through:
4106      - Traditional crosstalk speech patterns and rhythm
4107      - Common xiangsheng phrases and interactive elements
4108      - Cultural reference localization for Chinese audience
  
```

```

4104
4105 5. Performance format standardization:
4106   a. Apply consistent dialogue structure:
4107      "[tone] Role_name: dialogue_content"
4108   b. Ensure proper role distribution between performers:
4109      - Balanced dialogue allocation
4110      - Appropriate setup-punchline timing
4111      - Natural conversational flow maintenance
4112   c. Integrate traditional crosstalk elements:
4113      - Call-and-response patterns
4114      - Misunderstanding-based humor
4115      - Cultural wordplay and linguistic jokes
4116
4117 6. Content localization and cultural adaptation:
4118   a. Transform English humor into Chinese comedy traditions
4119   b. Adapt cultural references for Chinese audience familiarity
4120   c. Incorporate traditional crosstalk performance elements:
4121      - Verbal sparring between performers
4122      - Audience-directed asides and commentary
4123      - Regional dialect influences and speech patterns
4124   d. Maintain original humor essence while ensuring cultural relevance
4125
4126 7. Output generation and file management:
4127   a. Generate complete crosstalk script with title header
4128   b. Format output structure: "Title\n[tone] Role: content...\n..."
4129   c. Save to 'cross-talk.txt' in reference script directory
4130   d. Return script string for immediate use or further processing
4131   e. Ensure UTF-8 encoding preservation for Chinese character support
4132
4133 Technical Specifications:
4134 - Content adaptation: Reference-based cultural transformation
4135 - Role management: Dual-performer dialogue optimization
4136 - Format compliance: Traditional crosstalk structure requirements
4137 - Cultural localization: English-to-Chinese humor translation
4138 - Performance guidance: Tone markers for delivery optimization
4139
4140 Crosstalk Adaptation Features:
4141 - Authentic xiangsheng dialogue patterns and timing
4142 - Cultural reference localization for target audience
4143 - Traditional performer role dynamics (funny guy/setup guy)
4144 - Balanced dialogue distribution for engaging performance
4145 - Professional comedy structure with setup-punchline optimization
4146
4147 Error Handling:
4148 - Reference script file validation with encoding verification
4149 - Claude LLM API error detection with graceful failure handling
4150 - Output file creation with directory permission management
4151 - UTF-8 encoding preservation for Chinese character integrity
4152 - Exception isolation with detailed error reporting for debugging workflow

```

4151 A.8.11 CROSSTALKSYNTH

4153 The CrossTalkSynth agent provides comprehensive audio synthesis for crosstalk performances,
4154 combining dual-performer voice generation with intelligent dialogue analysis. It employs CosyVoice2
4155 for high-quality speech synthesis and DeepSeek LLM for script parsing to create authentic crosstalk
4156 performances with role-specific vocal characteristics and appropriate delivery styles.

4157

Listing 40: CrossTalkSynth Structure

4158
4159
4160
4161
4162
4163
4164
4165
4166
4167
4168
4169
4170
4171
4172
4173
4174
4175
4176
4177
4178
4179
4180
4181
4182
4183
4184
4185
4186
4187
4188
4189
4190
4191
4192
4193
4194
4195
4196
4197
4198
4199
4200
4201
4202
4203
4204
4205
4206
4207
4208
4209
4210
4211

Input: segmented_script, funny_guy_directory, setup_guy_directory
Output: merged_audio_path, segment_directory, metadata_path

Algorithm:

1. Input validation and role setup:
 - a. Validate script content and performer directory paths
 - b. Convert directories to absolute paths for cross-platform compatibility
 - c. Extract performer names from directory basenames:
 - dou_gen_name: funny guy performer (comic lead)
 - peng_gen_name: setup guy performer (straight man)
 - d. Initialize CosyVoice2 model in tools/CosyVoice environment
2. CosyVoice2 model initialization:
 - a. Navigate to CosyVoice working directory for proper execution
 - b. Load pretrained model: CosyVoice2-0.5B with configuration:
 - load_jit=False, load_trt=False, fp16=False
 - High-quality synthesis settings for crosstalk dialogue
 - c. Handle model loading errors with graceful failure reporting
3. Script parsing and line-by-line analysis:
 - a. Split script by newlines and filter empty content
 - b. Skip title line (first line) and process dialogue content
 - c. Create segment directory for individual audio file storage
 - d. Initialize counters and result arrays for processing tracking
4. Intelligent dialogue analysis with DeepSeek LLM:

System context: "Crosstalk dialogue line analyzer"

User prompt structure:
"Analyze crosstalk dialogue line for performer role, tone, text content:
{dialogue_line}"

Output JSON format with STRICT rules:

 1. role: either {funny_guy_name} or {setup_guy_name}
 2. tone: 'Natural', 'Emphatic', or 'Confused'
 3. text: dialogue content without formatting
 4. reaction: 'Laughter' or 'Cheers' (only if markers present)
 5. NO extra characters before/after JSON

Example: {'role': '{performer}', 'tone': 'Natural', 'text': '...'}
5. JSON response processing and validation:
 - a. Strip markdown code block markers (```json, ```)
 - b. Parse JSON response with comprehensive error handling
 - c. Extract role, tone, text, and optional reaction fields
 - d. Validate performer names against configured roles
 - e. Handle malformed JSON with graceful error recovery
6. Role-specific voice synthesis:
 - a. Load performer-specific prompt files:
 - Prompt text: {data_dir}/{role}/{tone}.lab
 - Prompt audio: {data_dir}/{role}/{tone}.wav at 16kHz
 - b. Apply CosyVoice2 zero-shot inference:
 - Input: extracted dialogue text
 - Context: role and tone-specific prompts
 - Output: high-quality speech synthesis
 - c. Save individual segments: {segment_index}.wav

- 4212 7. Audio consolidation and performance assembly:
 4213 Function `merge_audio_files()`:
 4214 a. Initialize silent `AudioSegment` for sequential concatenation
 4215 b. Load numbered audio segments (`0.wav` to `cnt.wav`)
 4216 c. Concatenate segments maintaining dialogue flow:
 4217 - Preserve timing between performer exchanges
 4218 - Maintain audio quality throughout concatenation
 4219 d. Export final performance: `../final/cross_talk.wav`
- 4220 8. Metadata generation and workflow completion:
 4221 a. Compile processing results into structured JSON metadata
 4222 b. Include role, tone, text, and reaction information per segment
 4223 c. Export metadata to `cross-talk.json` in data directory
 4224 d. Restore original working directory after processing
 4225 e. Return comprehensive output paths for downstream processing

4226 Technical Specifications:

- 4227 - Speech synthesis: `CosyVoice2-0.5B` with zero-shot voice cloning
 4228 - Audio processing: `PyDub` for segment manipulation and merging
 4229 - Script analysis: `DeepSeek LLM` for intelligent dialogue parsing
 4230 - Role management: Dual-performer voice synthesis with tone variation
 4231 - Output format: WAV files with model native sample rate

4232 Crosstalk Performance Features:

- 4233 - Dual-performer role-specific voice synthesis
 4234 - Tone-aware delivery variation (Natural, Emphatic, Confused)
 4235 - Authentic crosstalk dialogue timing and rhythm
 4236 - Professional audio quality with seamless performer transitions
 4237 - Intelligent script parsing for automated role and tone detection

4238 Error Handling:

- 4239 - `CosyVoice2` model initialization validation with detailed error reporting
 4240 - JSON parsing with markdown cleanup and format validation
 4241 - Audio file loading verification for role-specific prompts
 4242 - Graceful handling of missing prompt files or invalid roles
 4243 - Comprehensive exception management with processing continuation
 4244 - Working directory restoration for environment consistency

4245

4246 A.8.12 VIDEOCONVERSION

4247

4248 The `VideoConversion` agent transforms audio content with JSON timestamps into comprehensive
 4249 visual scene descriptions optimized for video generation workflows. It processes timestamped audio
 4250 segments and converts narrative content into actionable visual scenes through intelligent content
 4251 analysis.
 4252
 4253

4254

Listing 41: `VideoConversion` Structure

4255

4256 Input: `timestamp_path`
 4257 Output: `video_scene_path`

4258

Algorithm:

4259

1. Timestamp file processing:
 4260 a. Load JSON timestamp data from `cut_points.json`
 4261 b. Validate file structure and content availability
 4262 c. Extract content segments from `sentence_data.chunks` array
 4263 d. Handle file format validation and error reporting

4264

2. Content formatting and preparation:
 4265 a. Extract content from each timestamped chunk

```

4266     b. Format content with ///// separators between segments
4267     c. Create unified content string: "//////\n" + join(contents)
4268     d. Validate segment count and content integrity
4269
4270 3. LLM-based scene description generation with embedded prompts:
4271   System prompt: "You are a visual scene descriptor. Follow exact
4272   requirements:
4273   - Keep number of ///// marks unchanged
4274   - Deduce English visual-scene descriptions only
4275   - Keep same sentence separators and spacing
4276   - Max 1 sentence per scene section
4277   - Don't directly translate sentences
4278   - Describe character appearances when names mentioned"
4279
4280   User prompt template:
4281   "Content to process: {content}"
4282
4283   Example Input: ////\n[Emily] and [Jackson] stood together, ocean
4284   breeze ruffling hair, surrounded by vastness of ocean.
4285   \n\n////\nLeader increased Xiao Wang's business freedom.
4286
4287   Example Output: ////\nRed hair girl Emily and brown hair boy Jackson
4288   standing on sunset seaside with wind-blown hair
4289   \n\n////\nWhite t-shirt young employees in office environment"
4289
4290 4. Scene description processing:
4291   a. Apply LLM processing with DeepSeek model
4292   b. Generate visual scene keywords and descriptions
4293   c. Maintain temporal alignment with original timestamps
4294   d. Preserve character appearance details and environmental context
4294
4295 5. Output structuring and validation:
4296   a. Create structured JSON with segment_scene field
4297   b. Preserve original content_created for reference
4298   c. Validate scene coherence and visual consistency
4299   d. Save to video_scene.json in scene_output directory
4299
4300 6. Quality assurance and reporting:
4301   a. Verify segment count matches input timestamps
4302   b. Display processing statistics and completion status
4303   c. Handle error cases with detailed error messaging
4304   d. Return structured output for downstream processing
4305
4306 Error Handling:
4307 - Timestamp file existence validation
4308 - JSON structure validation with KeyError handling
4309 - Empty content segment detection and reporting
4310 - Directory creation with proper permissions
4311 - LLM response parsing with fallback mechanisms

```

4313 A.8.13 VIDEOQA

4314 The VideoQA agent provides comprehensive video content analysis through automated transcription
4315 and interactive question-answering capabilities. It processes multiple videos simultaneously, generates
4316 combined transcripts, and enables users to query video content through natural language interactions.
4317

4318 Listing 42: VideoQA Structure with System Prompt

```
4319 Input: video_directory, output_path, save_history_flag
```

```
4320 Output: transcript_path, qa_session_results
4321
4322 Algorithm:
4323 1. Initialize video processing environment:
4324   a. Load Whisper model for speech recognition
4325   b. Configure device (CUDA/CPU) and memory allocation
4326   c. Set up processing directories and file paths
4327
4328 2. Video discovery and validation:
4329   a. Scan directory for supported video formats
4330   b. Validate file integrity and accessibility
4331   c. Sort files for consistent processing order
4332
4333 3. Batch transcription pipeline:
4334   For each video file:
4335     a. Load video using Whisper pipeline
4336     b. Apply automatic speech recognition:
4337        - Chunk audio into 30-second segments
4338        - Process with batch size optimization
4339        - Extract timestamps and text content
4340     c. Format transcript with video identification
4341     d. Append to combined transcript buffer
4342
4343 4. Transcript consolidation:
4344   a. Merge individual transcripts with video markers
4345   b. Save combined transcript to designated path
4346   c. Generate processing statistics and metadata
4347
4348 5. Interactive Q&A session initialization:
4349   a. Load combined transcript for query processing
4350   b. Truncate content if exceeding API limits (50K chars)
4351   c. Initialize conversation history tracking
4352
4353 6. Q&A processing with embedded prompt:
4354   While user input != 'quit':
4355     a. Receive user question via command line interface
4356     b. Process question through QA agent with prompt:
4357
4358       System Message: "You are a helpful assistant who answers questions
4359       about video content based strictly on the provided transcripts from
4360       multiple videos."
4361
4362       User Prompt Template:
4363       "Answer the user's question carefully based only on the information
4364       contained in the video transcripts. The transcript contains content
4365       from multiple videos, each marked with '=== VIDEO: filename ==='.
4366       If the answer cannot be found in the transcripts, state that clearly.
4367
4368       User's question: '{user_question}'
4369
4370       Video transcripts content: {content}
4371
4372       Requirements:
4373       1. Be concise and direct in your answer
4374       2. Only use information from the transcripts
4375       3. If the answer is not in the transcripts, say 'I don't have enough
4376       information from the videos to answer this question'
4377       4. When referencing information, mention which video file it came
4378       from if relevant
```

- 4374 5. Don't make up information not present in the transcripts"
 4375
 4376 c. Query LLM with constructed prompt and transcript context
 4377 d. Generate factual answer based on video content
 4378 e. Display answer and log interaction
 4379 f. Update conversation history with timestamp
 4380
 4381 7. Session finalization:
 4382 a. Save Q&A history to persistent storage
 4383 b. Generate session summary statistics
 4384 c. Return processing results and file paths
 4385
 4386 8. Error handling and recovery:
 4387 - Fallback to existing transcripts on failure
 4388 - Graceful session termination on interruption
 4389 - Comprehensive logging for debugging

4391 A.8.14 VIDEOSUMMARIZATIONGENERATOR

4392 The VideoSummarizationGenerator creates comprehensive video summarization content by pro-
 4393 cessing video files through automatic speech recognition and applying user-specified presentation
 4394 styles. It combines advanced transcription capabilities with intelligent content adaptation to generate
 4395 structured video summaries.
 4396

4397 Listing 43: VideoSummarizationGenerator Structure

4398 Input: user_idea, video_directory, presentation_style_path, output_path
 4399 Output: content_output, status, processed_videos, transcript_source
 4400
 4401 Algorithm:
 4402 1. Whisper model initialization with retry logic:
 4403 a. Load AutoModelForSpeechSeq2Seq from whisper-large-v3-turbo
 4404 b. Configure device detection (CUDA/CPU) and dtype optimization
 4405 c. Initialize AutoProcessor with tokenizer and feature extractor
 4406 d. Create ASR pipeline with optimized parameters:
 4407 - max_new_tokens=128, chunk_length_s=30, batch_size=16
 4408 - generate_kwargs={"language": "en", "task": "transcribe"}
 4409
 4410 2. Content source processing:
 4411 a. Analyze video_dir input type (file/directory/text)
 4412 b. For video files:
 4413 - Scan supported extensions (.mp4, .avi, .mov, .mkv, etc.)
 4414 - Execute batch transcription with progress tracking
 4415 - Combine transcripts with video identification markers
 4416 c. For text files:
 4417 - Load with multi-encoding support and validation
 4418 - Handle binary reading fallback for encoding issues
 4419
 4420 3. Batch video transcription pipeline:
 4421 For each video file:
 4422 a. Apply Whisper ASR with retry logic (3 attempts)
 4423 b. Process with chunk-based approach for large files
 4424 c. Generate combined transcript with video markers:
 4425 "=== VIDEO: filename.mp4 ===\n[transcript_content]"
 4426 d. Handle transcription errors with graceful degradation
 4427
 4. Presenter agent processing with embedded prompts:
 System: "You are an experienced expert in writing transcripts
 summarization. Pay special attention to user's words count

```

4428 requirements."
4429
4430 User prompt template:
4431 "Create content summarization, strictly following user's ideas and
4432 presentation methods, answer using user's idea language.
4433
4434 User's idea: '{user_idea}'
4435 Grounded text content: {content}
4436 Follow this presentation method: {present_content}
4437
4438 Requirements:
4439 1. Format: Less point forms
4440 2. Content: Follow word count requirements, use original dialogues
4441 3. Language: Third-person perspective, clear narrative flow"
4442
4443 5. Content adaptation and optimization:
4444 a. Load presentation style from file or direct string input
4445 b. Apply content truncation for oversized inputs (15K limit)
4446 c. Generate statistics: word count, character count, source type
4447 d. Process through presenter agent with comprehensive retry logic
4448
4449 6. Output generation and validation:
4450 a. Create structured content output with user_idea and content_created
4451 b. Save to specified output path as plain text
4452 c. Generate processing metadata:
4453     - processed_videos: list of successfully transcribed files
4454     - transcript_source: source type classification
4455     - status: success/error with detailed error information
4456
4457 Error Handling:
4458 - Comprehensive retry logic for all major operations (3-5 attempts)
4459 - Exponential backoff with configurable wait times
4460 - Graceful fallback for transcription failures
4461 - Multi-encoding file reading with robust error recovery
4462 - Device fallback for CUDA unavailability
4463 - Path validation and directory creation

```

4463 A.9 PROMPT OF INTENT ANALYSIS

4464
4465 To mitigate the memory burden on LLMs and achieve more precise intent-tool alignment, we
4466 conduct intent analysis of user requirements and provide a detailed prompt as presented in Listing 44.
4467 Additionally, VideoAgent adaptively refines its intent analysis based on the feedback provided during
4468 the self-reflection phase, with the corresponding prompt shown in Listing 45.

4469 Listing 44: System Prompt of Intent Analysis

```

4470
4471 You are an intent analyst.
4472 I will provide a set of candidate intents and a user's requirements.
4473 Please select as many relevant candidate intents as possible that match the
4474 ↪ user's requirements. Consider factors such as keywords, audio types (speech
4475 ↪ , song, music, etc.), and other relevant dimensions.
4476
4477 Candidate intents:
4478 {intents}
4479
4480 User requirements:
4481 {reqs}
4482
4483 Please Only output pure List format:

```

```
4482 ['intent1', 'intent2', ...]
4483
4484 Note! Don't output any analysis and explanations!
4485
```

Listing 45: System Prompt of Intent Analysis

```
4488 You are an intent analyst.
4489 Previous analysis attempt failed with the following reflection:
4490 {reflection}
4491
4492 Previous selected intents:
4493 {previous_intents}
4494
4495 Please re-analyze the user's requirements with the candidate intents below,
4496 ⇨ considering the reflection.
4497 Select as many relevant candidate intents as possible.
4498
4499 Candidate intents:
4500 {intents}
4501
4502 User requirements:
4503 {reqs}
4504
4505 Please Only output pure List format:
4506 ['intent1', 'intent2', ...]
4507
4508 Note! Don't output any analysis and explanations!
4509
```

4510 A.10 GRAPH CONSTRUCTION RULES

4511 We provide a comprehensive prompt for the agent graph, which guides the collaboration among
4512 tooluse agents and also allows for a high degree of flexibility in constructing agent workflows, not
4513 being limited to predefined pipelines, as presented in Listing 46.
4514

Listing 46: System Prompt of Agent Graph

```
4516 You are an Agent Graph Designer. I will provide User Requirement and
4517 ⇨ Registered Agents (Name, Description and Parameters information):
4518
4519 Your task is to:
4520 1. Judge Feasibility:
4521   - Evaluate implementation feasibility of User Requirements
4522   - Output: "Feasible" or "Infeasible" (strictly one of these)
4523 2. Design Executable Agent Graph:
4524   - Format: List
4525   - Agent Graph shall contain metadata for each Agent Node including:
4526     * name: (string)
4527     * inputs: (list of input parameter objects with):
4528     * parameter: input parameter name
4529     * description: brief parameter description
4530     * outputs: (list of output parameter objects with):
4531     * parameter: output parameter name
4532     * description: brief parameter description
4533     * links: (list of dictionaries) where each dictionary specifies:
4534       - key of dictionaries: target agent name
4535       - value of dictionaries: target agent's input parameter name that
4536 ⇨ this output connects to
4537 3. Generate Agent Chain:
```

4536 - Format: List
 4537 - Generate the Agent Chain based on the description of the Agent and the
 4538 ↪ sequential information contained in the designed Agent Graph
 4539 4. Generate User Input Graph
 4540 - Format: List
 4541 - Parameter nodes with no in-degree (no incoming edges) are uniformly
 4542 ↪ considered to require user input.
 4543 - Parameter nodes with no in-degree may have different names but share the
 4544 ↪ same user input, meaning a single user input parameter can point to
 4545 ↪ multiple such nodes.
 4546 - Parameter nodes with no in-degree that are linked to user input should be
 4547 ↪ represented in the format ****AgentName.input_parameter****
 4548 - Generate the User Input Graph based on the Agent descriptions and
 4549 ↪ parameter passing information in the designed Agent Graph.
 4549 5. Output Reasoning:
 4550 - If Feasible, Provide concise reasoning (<200 words) explaining the entire
 4551 ↪ workflow logic
 4552 - If Infeasible, Specify exact failure reasons (<200 words)
 4553
 4554 In addition to the above formatting requirements, please also note the
 4555 ↪ following:
 4556 1. For each element of ****outputs**** in each Agent Node:
 4557 - Ensure that the ****links**** in the ****outputs**** point to an input parameter
 4558 ↪ that actually exists in the next Agent Node.
 4559 - The output parameter name does not need to match the input parameter name
 4560 ↪ in the next Agent Node.
 4561 - Ensure the output parameter's description and type match the input
 4562 ↪ parameter requirements of the next Agent Node. For example, a file path
 4563 ↪ output cannot be passed to a directory path input.
 4563 2. Final JSON Output Format Specification:
 4564 {
 4565 "Feasibility": "Feasible" or "Infeasible",
 4566 "Agent Graph": ...,
 4567 "Agent Chain": ...,
 4568 "User Input Graph": ...,
 4569 "Reasoning": ...
 4570 }
 4570 Strictly follow JSON output format!
 4571

4572 A.11 TWO-STEP REVIEWER

4573 During the self-reflection phase, we propose a two-step reviewer strategy to maximize the accuracy
 4574 of the evaluation. The corresponding prompts can be found in Listing 47 and Listing 48 respectively.

4575 Listing 47: System Prompt for Initial Evaluation

4576 You are an agent graph validation system.
 4577 I will provide:
 4578 1. User Requirement
 4579 2. Registered agent metadata
 4580 3. Candidate Agent Graph
 4581 4. An Agent Chain derived from the Candidate Agent Graph
 4582 5. Required User Inputs
 4583
 4584 User Requirements:
 4585 {reqs}
 4586
 4587 Metadata of registered agents:
 4588 {tools}

```

4590
4591 Task: Evaluate the candidate agent graph:
4592
4593 Candidate Agent Graph:
4594 {agent_graph}
4595
4596 Agent Chain:
4597 {agent_chain}
4598
4599 Required User Inputs:
4600 {user_inputs}
4601
4602 Evaluation Criteria:
4603 1. Based on the Metadata of registered agents and parameter passing in the
4604 ↪ Agent Graph, determine from multiple aspects whether the user requirements
4605 ↪ can be fulfilled:
4606     - Execution sequence of agents in the Agent Graph
4607     - For parameter nodes with no incoming edges, they are uniformly considered
4608 ↪ as user inputs, but it is necessary to determine whether they should be
4609 ↪ provided by the user or by the parent agent
4610     - Validate that the necessary output parameters are correctly routed to the
4611 ↪ intended agent and the expected input parameters.
4612 2. There should be no functionally redundant agents (e.g., repeatedly adding
4613 ↪ audio tracks to a video).
4614 3. For vaguely mentioned requirements in user needs, lenient evaluation is
4615 ↪ acceptable. For example, if the user requests audio quality improvement, it
4616 ↪ 's sufficient as long as at least one relevant agent in the graph meets
4617 ↪ this requirement.
4618
4619 Please Only output pure JSON format:
4620 {{
4621 "Result": '0' if correct else '1',
4622 "Reasoning": Concisely state the key reasons why a score of '0' or '1' was
4623 ↪ assigned (<100 words).
4624 }}

```

Listing 48: System Prompt for Secondary Evaluation

```

4625 You are an agent graph reflection system.
4626 I will provide:
4627 1. User Requirement
4628 2. Registered agent metadata
4629 3. Candidate Agent Graph
4630 4. An Agent Chain and User Input Graph derived from the Candidate Agent Graph
4631 5. Previous validation result
4632
4633 User Requirements:
4634 {reqs}
4635
4636 Metadata of registered agents:
4637 {tools}
4638
4639 Task: Evaluate the candidate agent graph:
4640
4641 Candidate Agent Graph:
4642 {agent_graph}
4643
4644 Agent Chain:
4645 {agent_chain}

```

4644

4645

Required User Input Graph:

4646

```
{user_inputs}
```

4647

4648

Previous validation result:

4649

```
{judge_res}
```

4650

4651

Reflection Task:

4652

1. If the previous validation result is '0', please reflect on whether there

4653

↪ were any overlooked aspects based on the ****Evaluation Criteria**** and the

4654

↪ reasoning behind the previous validation result.

4655

2. If the previous validation result is '1', please reflect on whether the

4656

↪ reasoning behind the previous validation result was correct.

4657

Evaluation Criteria:

4658

1. Based on the Metadata of registered agents and parameter passing in the

4659

↪ Agent Graph, determine from multiple aspects whether the user requirements

4660

↪ can be fulfilled:

4661

- Execution sequence of agents in the Agent Graph

4662

- For parameter nodes with no incoming edges, they are uniformly considered

4663

↪ as user inputs, but it is necessary to determine whether they should be

4664

↪ provided by the user or by the previous agent

4665

- Validate that the necessary output parameters are correctly routed to the

4666

↪ intended agent and the expected input parameters.

4667

- Validate that the output parameters' description and type match the input

4668

↪ requirements of the next agent.

4669

- Not all output parameters are necessarily mapped to the input

4670

↪ requirements of the next agent. Redundant output parameters may exist, but

4671

↪ they should not interfere with the fulfillment of user requirements.

4672

2. There should be no functionally redundant agents (e.g., repeatedly adding

4673

↪ audio tracks to a video).

4674

3. For vaguely mentioned requirements in user needs, lenient evaluation is

4675

↪ acceptable. For example, if the user requests audio quality improvement, it

4676

↪ 's sufficient as long as at least one relevant agent in the graph meets

4677

↪ this requirement.

4678

Please Only output pure JSON format:

4679

```
{
```

4680

```
  "Result": '0' if correct else '1',
```

4681

```
  "Reasoning": Concisely state the key reasons why a score of '0' or '1' was
```

4682

```
  ↪ assigned (<100 words).
```

4683

```
}
```

4684

4685

```
}
```

4686

4687

A.12 CASES OF MULTI-MODAL AGENT WORKFLOW

4688

4689

4690

4691

4692

4693

4694

4695

4696

4697

4698

4699

4700

To demonstrate VideoAgent's high sensitivity to user requirements and its generalization capability, we also provide several constructed agent workflow cases. VideoAgent can be extensively customized to tailor agent workflows according to user needs. For example, as shown in Listing 49, it performs loudness normalization on audio. Users can further extend this functionality with additional audio processing tasks, such as vocal separation or adding background music. Moreover, VideoAgent supports customized audio and video creation, allowing users to freely propose diverse creative requirements, including voice cloning, song generation, and music video production.

Listing 49: User Requirement of Audio Preprocessing

4695

User Requirement: I need to extract the audio from a video file, improve its

4696

↪ sound quality by normalizing the loudness, and then mix it with some

4697

↪ background music. The result should be a clean audio track that balances

4698

↪ the original audio with background music.

4698

4699

4700

4701

Listing 50: User Requirement of Audio Preprocessing

4702

4703

4704

4705

4706

4707

4708

4709

4710

4711

4712

4713

4714

4715

4716

4717

4718

4719

4720

4721

4722

4723

4724

4725

4726

4727

4728

4729

4730

4731

4732

4733

4734

4735

```

{
  "Feasibility": "Feasible",
  "Intent List": ["Loudness Normalization", "Add BGM", "Audio Optimization"],
  "Agent Graph": [
    {"node": "LoudnessNormalizer", "inputs": [{"name": "data_dir", "
    ↪ description": "Directory of audio files to be normalized"}], "outputs": [{"
    ↪ name": "status", "description": "Execution status (success/error)", "links
    ↪ ": []}]},
    {"node": "Mixer", "inputs": [{"name": "bgm_path", "description": "Path to
    ↪ the BGM file"}, {"name": "audio_path", "description": "Audio to be mixed
    ↪ into the background music"}], "outputs": [{"name": "audio_path", "
    ↪ description": "Path to the synthesized audio", "links": []}]},
  ],
  "Agent Chain": ["LoudnessNormalizer", "Mixer"],
  "User Input Graph": [
    {"node": "raw_audio_directory", "description": "Directory containing the
    ↪ user's raw audio recording that needs to be normalized", "links": [{"
    ↪ LoudnessNormalizer": "data_dir"}]},
    {"node": "normalized_audio_path", "description": "Path to the normalized
    ↪ audio file after processing by LoudnessNormalizer", "links": [{"Mixer": "
    ↪ audio_path"}]},
    {"node": "background_music_path", "description": "File path to the
    ↪ background music to be mixed with the voice recording", "links": [{"Mixer":
    ↪ "bgm_path"}]}
  ],
  "Reasoning": "This workflow normalizes the user's raw audio recording and then
  ↪ mixes it with background music. First, LoudnessNormalizer processes the
  ↪ raw audio files to ensure consistent volume levels. Since
  ↪ LoudnessNormalizer doesn't output file paths (only a status), we require
  ↪ the user to provide the path to the normalized audio file after processing.
  ↪ Then, the Mixer agent combines this normalized audio with background music
  ↪ , ensuring the music doesn't overpower the voice. The workflow addresses
  ↪ the requirement to add background music to a recording while maintaining
  ↪ proper audio balance."
}

```

4736

Listing 51: User Requirement of Storytelling Video

4737

4738

4739

4740

4741

4742

4743

Listing 52: Case of Storytelling Video

4744

4745

4746

4747

4748

4749

4750

4751

```

{
  "Feasibility": "Feasible",
  "Intent List": ["Text-to-Speech", "Commentary", "Add BGM", "Video Edit"],
  "Agent Graph": [
    {"node": "CommentaryContentGenerator", "inputs": [{"name": "reqs", "
    ↪ description": "User's idea for the commentary video including word count
    ↪ requirements"}, {"name": "source_text", "description": "File path to the
    ↪ novel source text"}, {"name": "comm_present_style", "description": "File
    ↪ path to commentary presentation style for content generation"}], "outputs":
    ↪ [{"name": "video_scene_path", "description": "File path storing scene

```

```

4752 ↪ semantics for video storyboard sound synthesis.", "links": [{"
4753 ↪ VoiceGenerator": "video_scene_path"}, {"VideoSearcher": "video_scene_path
4754 ↪ "}]}}],
4755   {"node": "VoiceGenerator", "inputs": [{"name": "video_scene_path", "
4756 ↪ description": "Path to a custom scene JSON file"}, {"name": "
4757 ↪ target_vocal_path", "description": "Path to the target timbre for voice
4758 ↪ generation"}], "outputs": [{"name": "audio_path", "description": "Path to
4759 ↪ the synthesized audio", "links": [{"Mixer": "audio_path"}]}, {"name": "
4760 ↪ timestamp_path", "description": "Path to video frame timestamp", "links":
4761 ↪ [{"VideoEditor": "timestamp_path"}]}]},
4762   {"node": "Mixer", "inputs": [{"name": "bgm_path", "description": "Path to
4763 ↪ the BGM file"}, {"name": "audio_path", "description": "Audio to be mixed
4764 ↪ into the background music"}], "outputs": [{"name": "audio_path", "
4765 ↪ description": "Path to the synthesized audio", "links": [{"VideoEditor": "
4766 ↪ audio_path"}]}]},
4767   {"node": "VideoPreloader", "inputs": [{"name": "video_dir", "description":
4768 ↪ "Directory containing the source MP4 video files to be processed"}], "
4769 ↪ outputs": [{"name": "status", "description": "Execution status (success/
4770 ↪ error)", "links": []}],
4771   {"node": "VideoSearcher", "inputs": [{"name": "video_scene_path", "
4772 ↪ description": "File path storing scene semantics for video storyboard sound
4773 ↪ synthesis."}], "outputs": [{"name": "status", "description": "Execution
4774 ↪ status (success/error)", "links": []}],
4775   {"node": "VideoEditor", "inputs": [{"name": "video_dir", "description": "
4776 ↪ Directory containing source video files"}, {"name": "audio_path", "
4777 ↪ description": "Path to the audio"}, {"name": "timestamp_path", "description
4778 ↪ ": "JSON File path used to store and load the timestamp of the end of each
4779 ↪ video segment"}], "outputs": [{"name": "video_path", "description": "Path
4780 ↪ to the generated video file", "links": []}]
4781 ],
4782 "Agent Chain": ["CommentaryContentGenerator", "VoiceGenerator", "Mixer", "
4783 ↪ VideoPreloader", "VideoSearcher", "VideoEditor"],
4784 "User Input Graph": [
4785   {"node": "commentary_requirements", "description": "User's specific
4786 ↪ requirements for the commentary video including word count and style", "
4787 ↪ links": [{"CommentaryContentGenerator": "reqs"}]},
4788   {"node": "novel_file", "description": "File path to the novel text that
4789 ↪ will be the subject of commentary", "links": [{"CommentaryContentGenerator
4790 ↪ ": "source_text"}]},
4791   {"node": "commentary_style", "description": "File path to the desired
4792 ↪ presentation style for the commentary", "links": [{"
4793 ↪ CommentaryContentGenerator": "comm_present_style"}]},
4794   {"node": "voice_timbre", "description": "Path to the target vocal timbre
4795 ↪ file for professional voiceover", "links": [{"VoiceGenerator": "
4796 ↪ target_vocal_path"}]},
4797   {"node": "background_music", "description": "Path to the background music
4798 ↪ file", "links": [{"Mixer": "bgm_path"}]},
4799   {"node": "visuals_directory", "description": "Directory containing source
4800 ↪ video files for visual content", "links": [{"VideoPreloader": "video_dir"},
4801 ↪ {"VideoEditor": "video_dir"}]}
4802 ],
4803 "Reasoning": "The workflow begins with generating commentary content about the
4804 ↪ novel based on user requirements. This content is structured with scene
4805 ↪ semantics. The VoiceGenerator then creates professional voiceover audio and
4806 ↪ provides timestamps for visual syncing. In parallel, VideoPreloader
4807 ↪ prepares video files from the user's library while VideoSearcher identifies
4808 ↪ appropriate visual clips based on the commentary content. The voiceover is
4809 ↪ mixed with background music by the Mixer agent. Finally, VideoEditor
4810 ↪ combines the mixed audio track with matching visuals using the timestamps

```

```

4806 ↪ to synchronize content, producing a complete commentary video about the
4807 ↪ novel with professional narration, appropriate background music, and
4808 ↪ relevant visuals."
4809 }
4810

```

A.13 MLLM/VLM FOR VIDEO EDITING EVALUATION

This subsection presents the standardized system prompt methodology for MLLM integration in video editing benchmark tasks, focusing on precise clip-to-caption matching and automated video content analysis workflows.

Listing 53: MLLM Video Clip Matching Prompt with Frame Integration

```

4818 Prompt Construction with Frame Integration:
4819
4820 1. Text prompt construction:
4821 prompt = f"""
4822 Here is a video broken down into {len(chunks)} clips of 3 seconds each.
4823 Each clip shows multiple frames from that time segment.
4824 The clips are structured as follows:
4825 """
4826
4827 for chunk in chunks:
4828     prompt += f"Clip {chunk['chunk_idx']}\n"
4829
4830 prompt += f"""
4831 Identify {num_periods} specific clips that best correspond to the following
4832 ↪ captions:
4833 """
4834
4835 for i, caption in enumerate(shuffled_captions):
4836     prompt += f"{i+1}. {caption}\n"
4837
4838 prompt += f"""
4839 Response with a JSON structure containing the clip numbers that best match
4840 ↪ each caption.
4841
4842 Rules:
4843 1. Each clip_id should be a number (not a string)
4844 2. Include exactly {num_periods} selections
4845 3. Choose from the clip numbers shown (0 to {len(chunks)-1})
4846 4. Keep the reason concise (1-2 sentences)
4847 5. Ensure your output is valid JSON - no trailing commas, proper quotes, etc.
4848 6. Do not include the ```json prefix or ``` suffix around your response
4849 7. Do not answer anything unrelated
4850
4851 Return only the JSON object with no additional text.
4852
4853 Output Example:
4854 {{
4855 "selections": [
4856     {{
4857         "caption": "[caption text]",
4858         "clip_id": [clip number],
4859         "reason": "[brief description of what is seen in this clip]"
4860     }},
4861     ...
4862 ]
4863 }}

```

```
4860 """
4861
4862 2. Frame integration into multimodal message structure:
4863 messages = [{"role": "user", "content": []}]
4864
4865 # Add text prompt as first content element
4866 messages[0]["content"].append({"type": "text", "text": prompt})
4867
4868 # Sequential frame insertion for visual analysis
4869 for img_base64 in frames:
4870     messages[0]["content"].append({
4871         "type": "image_url",
4872         "image_url": {
4873             "url": f"data:image/jpeg;base64,{img_base64}"
4874         }
4875     })
4876
4877
4878
4879
4880
4881
4882
4883
4884
4885
4886
4887
4888
4889
4890
4891
4892
4893
4894
4895
4896
4897
4898
4899
4900
4901
4902
4903
4904
4905
4906
4907
4908
4909
4910
4911
4912
4913
```