

Supplementary Materials

1. S2I+I2V vs S2V



Figure 9. Comparison of subject-to-image-to-video [9] and subject-to-video (ours).

As mentioned in the main text, combining subject-to-image (S2I) and image-to-video (I2V) can achieve similar effects to subject-to-video (S2V), but there are some difficult limitations. Firstly, existing methods [9, 14, 21] for generating subject-consistent images or ID-consistent images still exhibit noticeable artificial artifacts, and there is significant room for improvement in the dimension of subject consistency. Equally important, I2V cannot ensure consistency of the subject during motion. As illustrated in Figure 9, when inputting a reference portrait, S2I first generates a reference image for the initial frame of I2V. If the initial frame includes a back view or occlusions, I2V may “imagine” a false ID during the process of removing the occlusion, leading to a failure in maintaining consistency.

2. Copy-paste problem

In the field of video generation, the copy-paste issue is particularly prominent, manifesting as the leakage of image content into the generated video. Some methods sample keyframes from a video and use them as image conditions to reconstruct the video. However, this approach allows the model to employ shortcut learning strategies, simplifying the content understanding process. Figure 10 shows examples of the copy-paste issue, sampling from the initial, middle, and final frames: In the first row, the girl’s expression



Figure 10. Intuitive cases of copy-paste problems. The red font in the text prompt does not function as intended.

remains unchanged, ignoring the text prompt. In the second row, the cartoon character’s movements remain stiff and identical to the reference. The third row illustrates a common case where the generated video is too similar to I2V, diminishing the effectiveness of scene-related text and reducing content diversity. To address this, we focus on constructing cross-video multi-subject pairings, ensuring subjects match in content while allowing for non-rigid deformations and changes in color distribution, thereby avoiding the copy-paste problem.

3. Ablation study supplement

Multi-subject confusion issue. When multiple reference subjects are input simultaneously, appearance confusion may occur. Our solution aligns text descriptions with video subjects during training, ensuring distinct descriptions for each subject. During inference, a rephraser adjusts the input text prompts to align with the training data format. For example, in the first row of Figure 12, the original prompt “A family of three is having a meal at the table” caused confusion. The rephrased prompt “a woman in black, a young girl in white, and an elderly man in a suit eating together at the table” resolved this issue. In the second row of Figure 12, the original prompt “a girl in casual clothes walking by the beach” failed to match the reference. The rephrased prompt “a girl in a white T-shirt and jeans walking by the beach”

	w/o text-image alignment	w/ text-image alignment
Success rate	65%	95%

Table 4. Success rate of multi-subject generation with and without text-image alignment.

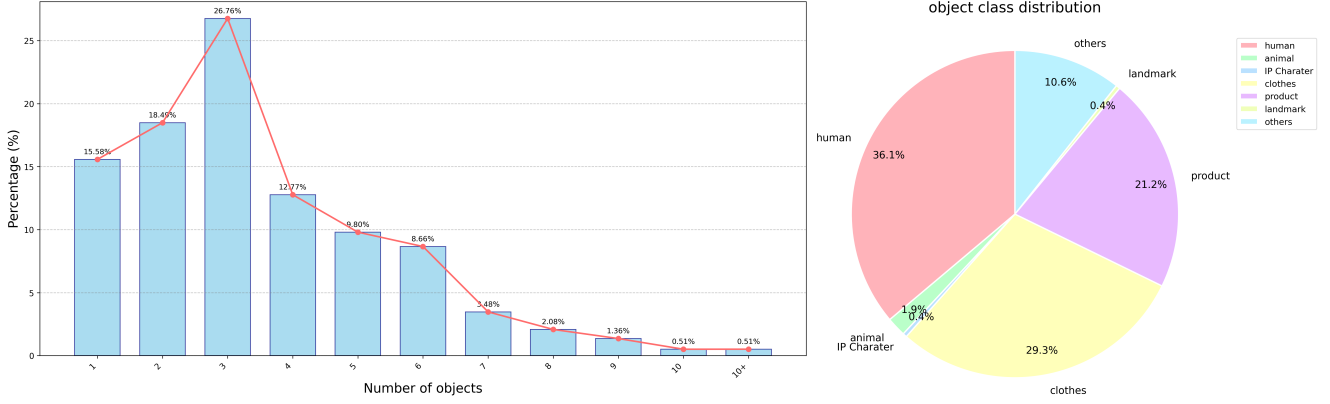


Figure 11. Distribution of object frequencies and class.



Figure 12. Examples of multi-subject confusion: On the left are the multi-subject reference images, while the right columns present the cases of confusion and the successful cases after improvement.

successfully matched the reference. Quantitative analysis, shown in Table 4, indicates a significant increase in the success rate of subject-consistent generation with this method. Aligning image and text is crucial for multi-subject generation tasks. This approach, which requires no additional complex data structures or model designs, significantly optimizes the multi-subject confusion problem.

4. Data pipeline for face ID

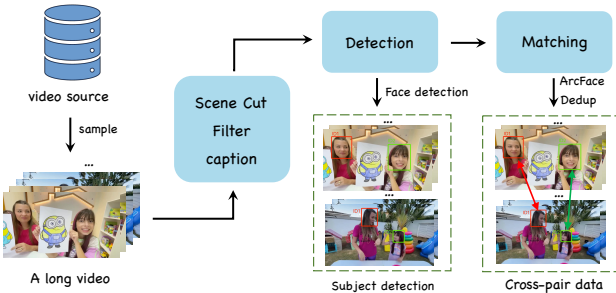


Figure 13. Facial data processing pipeline for constructing ID cross-pair

To enhance facial ID consistency, we developed an additional data pipeline for processing facial data. As shown in Figure 13, the facial data pipeline reuses the scene segmentation, video filtering, and annotation steps from the general subject pipeline. During the detection stage, we use an internal facial detection tool to identify each face in the video reference frames and calibrate it with the VLM [1] results from the captions using IOU (Intersection Over Union). In the matching stage, we calculate facial similarity using Arcface [8] features and add a deduplication operator [25] to further calibrate the recognition results.

5. Data distribution

Distribution of video object quantities. We sample three frames at [0.05, 0.5, 0.95] of the video timeline and perform object detection on these frames. We filter out objects that meet the following criteria: (1) objects that are small in size or occupy a small proportion of the frame; (2) objects with a high degree of overlap with other objects; and (3) incomplete objects judged by the VLM [1]. The final distribution of the number of objects per video is shown in the table on the left side of Figure 11.

Distribution of video object types. We use LLM [37] to classify the noun fields in all captions into the following categories: human, animal, clothes, product, landmark, IP character, and others. The distribution is shown in the accompanying Figure 11, with human, clothes, and product categories accounting for the majority.

6. Model architecture

The architecture of the *Phantom* model is shown in Figure 14, which supplements the missing details in the main text. As illustrated, it integrates the VAE and CLIP encoders to process reference images, while the text encoder handles captions. The encoded features are combined with added noise and processed through multiple MMDiT blocks, resulting in the final output. This design ensures a bal-

Acknowledgments We would like to express our gratitude to the Bytedance-Seed team for their support. Special thanks to Haoyuan Guo, Zhibei Ma, Sen Wang and Lu Jiang for their assistance with the model and data. In addition, we are also very grateful to Siying Chen, Qingyang Li, and Wei Han for their help with the evaluation.

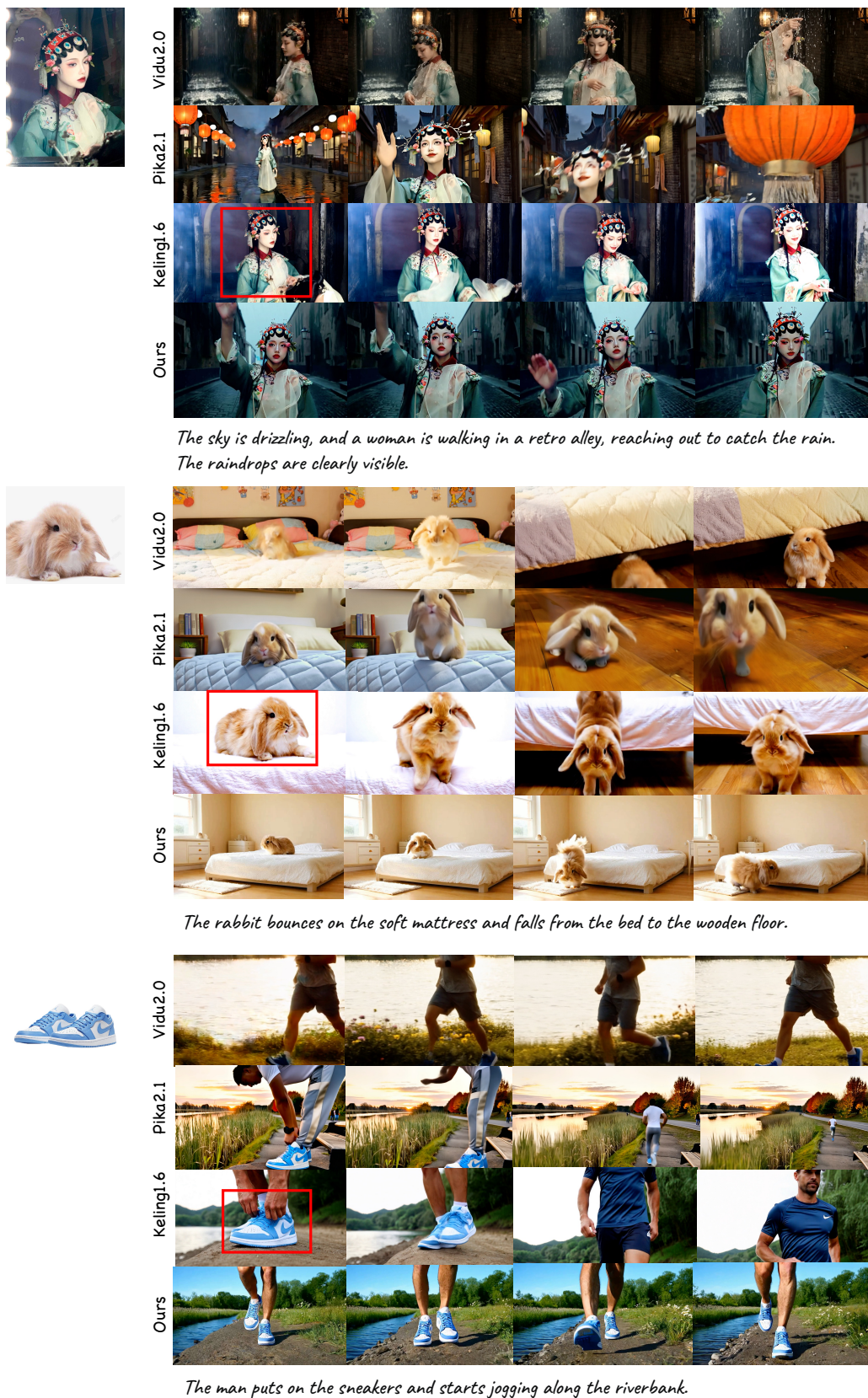
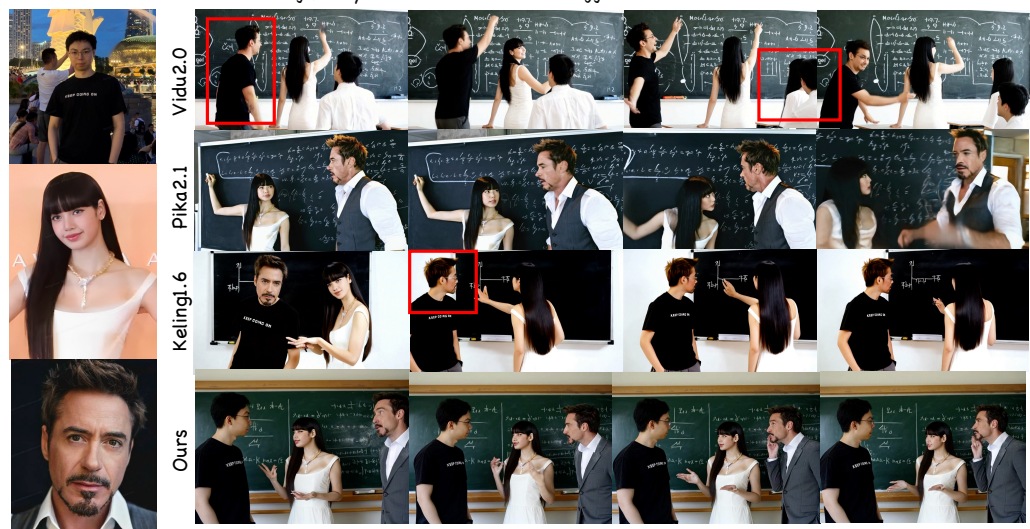


Figure 15. Comparative results of single reference subject-to-video generation.



On the stage, they both shook hands and hugged.



They were fiercely discussing the solution to a math problem in front of the blackboard, pointing and pointing at the blackboard.



A character walking on the moon.

Figure 16. Comparative results of multi-reference subject-to-video generation.