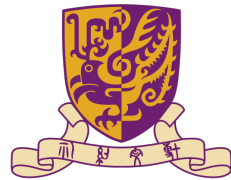


OmniShow: Unifying Multimodal Conditions for Human-Object Interaction Video Generation

Donghao Zhou^{1,*}, Guisheng Liu^{2,*}, Hao Yang², Jiatong Li^{2,†}, Jingyu Lin³, Xiaohu Huang⁴, Yichen Liu²,
Xin Gao², Cunjian Chen³, Shilei Wen^{2,§}, Chi-Wing Fu¹, Pheng-Ann Heng^{1,§}

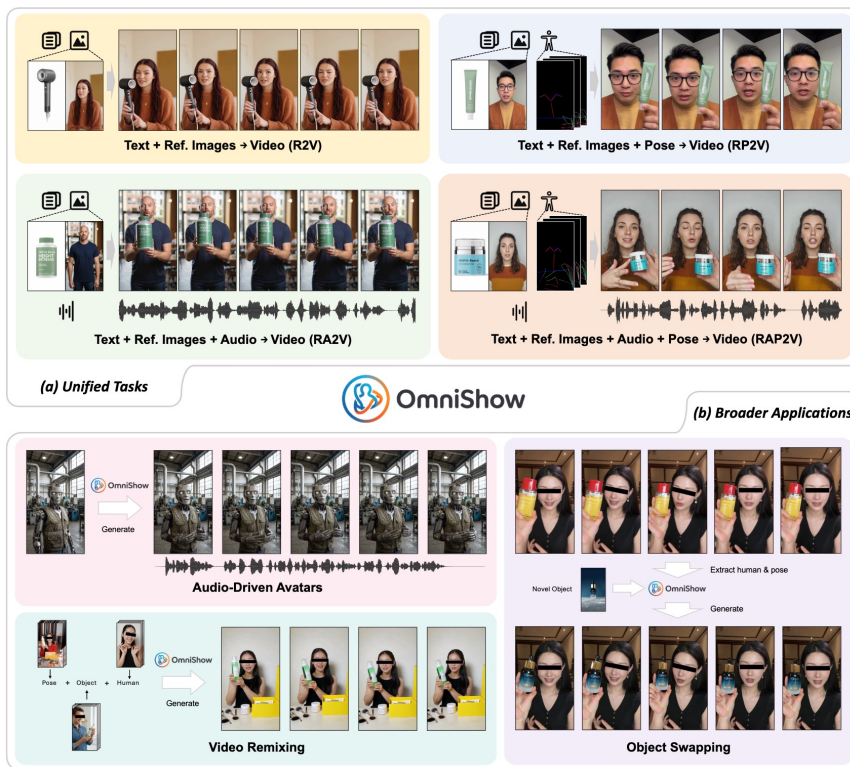
¹The Chinese University of Hong Kong, ²ByteDance, ³Monash University, ⁴The University of Hong Kong

*(*Equal Contribution, †Project Lead, §Corresponding Author)*



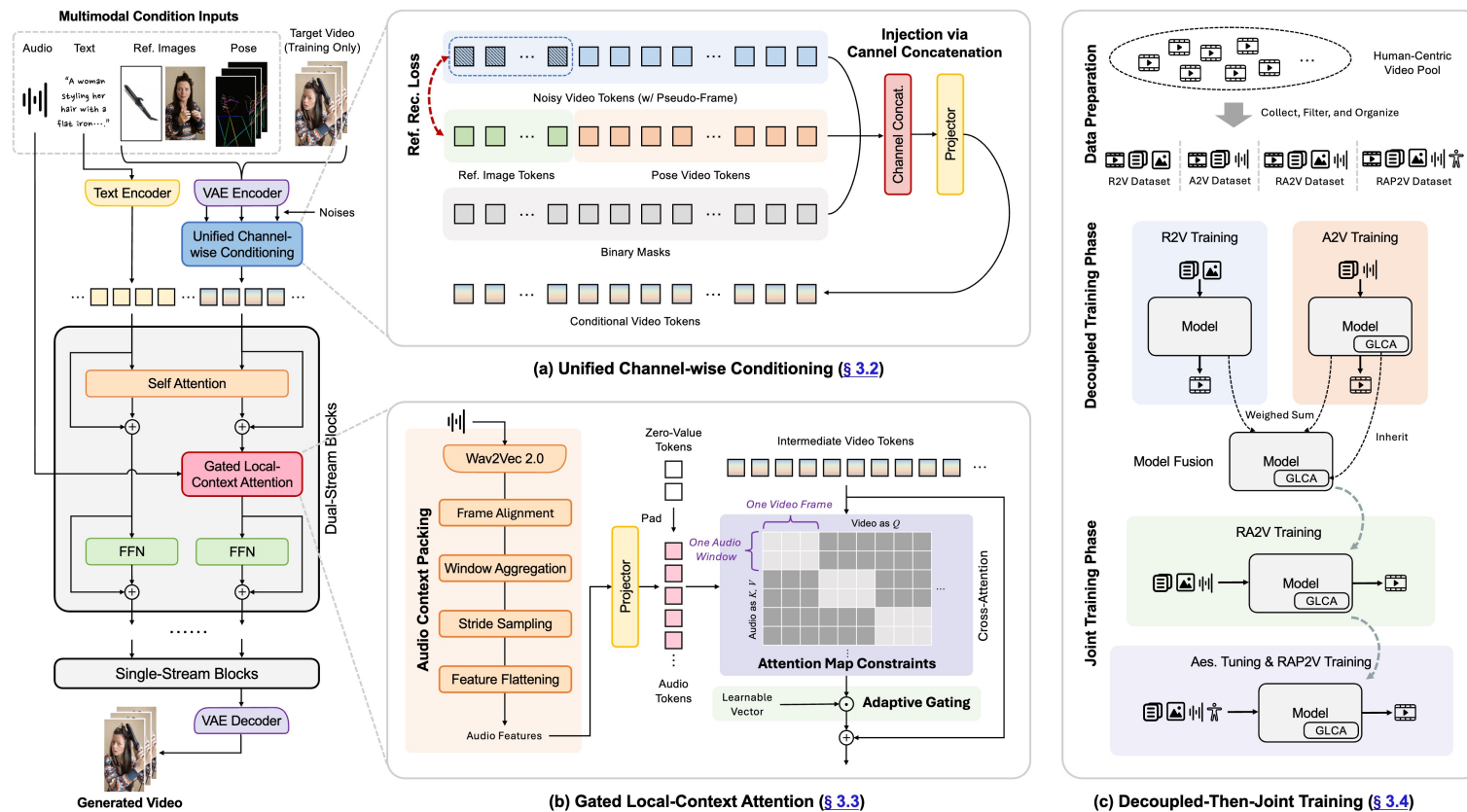
Overview

We introduce **OmniShow**, the **first all-in-one-model** that unifies **text**, **reference image**, **audio**, and **pose** conditions for **Human-Object Interaction Video Generation (HOIVG)**, which also supports broader real-world applications.



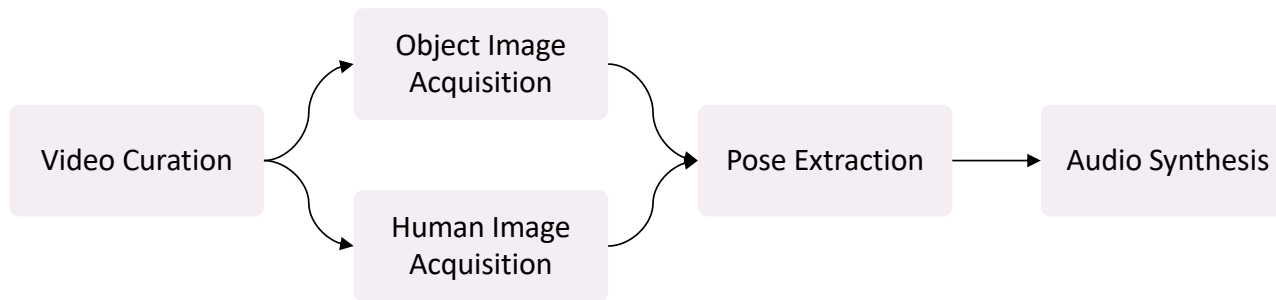
Pipeline

We introduce **OmniShow**, the **first all-in-one-model** that unifies **text, reference image, audio, and pose** conditions for **Human-Object Interaction Video Generation (HOIVG)**, which also supports broader real-world applications.



Methodology

Existing benchmarks often focus on limited-modality control (e.g., text+pose or text+images). We construct **HOIVG-Bench** to bridge this gap by providing an evaluation suite comprising 135 **carefully curated samples** and **dedicated metrics**.



Text Alignment:

- TA from VideoReward

Reference Consistency:

- Face Sim from OpenS2V
- NexusScore from OpenS2V

Pose Accuracy:

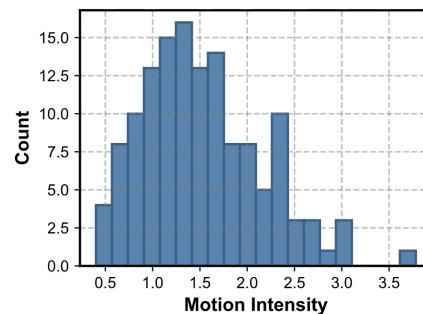
- AKD based on DWPose
- PCK based on DWPose




Audio-Visual Synchronization:

- Sync-C based on SyncNet
- Sync-D based on SyncNet

Video Quality:

- AES from VBench
- IQA from VBench
- VQ from VideoReward
- MQ from VideoReward



Text	Reference Image	Pose	Audio
<p>The man dressed in a grey long-sleeve shirt holds a black cylindrical container labeled HOIVG-Bench, standing indoors in a cozy room with wooden furniture in the background. The product is small-size, easily held in one hand. He gestures expressively with his free hand while speaking directly to the camera, highlighting the product's modular design and ergonomic support.</p>			

Experiments

We compare our OmniShow with existing state-of-the-art methods on **HOIVG-Bench**, showing that OmniShow can achieve **superior overall performance** across **diverse multimodal task settings**.

Method	Text Align.	Reference Consistency		Audio-Visual Sync.		Pose Accuracy		Video Quality			
	TA \uparrow	FaceSim \uparrow	NexusScore \uparrow	Sync-C \uparrow	Sync-D \downarrow	AKD \downarrow	PCK \uparrow	AES \uparrow	IQA \uparrow	VQ \uparrow	MQ \uparrow
<i>Text+Reference-to-Video (R2V)</i>											
HunyuanCustom [21]	7.523	0.440	0.359	-	-	-	-	0.452	0.697	10.11	5.286
HuMo-1.7B [4]	7.087	0.647	0.333	-	-	-	-	0.441	0.723	9.76	3.406
HuMo-17B [4]	7.949	0.843	0.346	-	-	-	-	0.448	0.726	9.97	3.685
VACE [31]	<u>8.413</u>	0.759	<u>0.368</u>	-	-	-	-	0.457	0.722	10.72	5.442
Phantom-1.3B [42]	8.342	0.708	0.351	-	-	-	-	<u>0.459</u>	0.722	10.90	<u>5.637</u>
Phantom-14B [42]	8.609	0.876	0.366	-	-	-	-	0.449	0.741	<u>10.93</u>	5.517
OMNISHOW (Ours)	7.746	<u>0.874</u>	0.389	-	-	-	-	0.468	<u>0.740</u>	11.12	5.885
<i>Text+Reference+Audio-to-Video (RA2V)</i>											
HunyuanCustom [21]	7.289	0.457	<u>0.350</u>	6.072	10.08	-	-	<u>0.439</u>	0.715	9.15	3.658
HuMo-1.7B [4]	7.489	0.575	0.329	7.234	9.117	-	-	0.428	0.731	9.97	4.182
HuMo-17B [4]	8.146	<u>0.805</u>	0.344	<u>8.013</u>	<u>8.316</u>	-	-	0.439	<u>0.739</u>	<u>10.27</u>	<u>4.269</u>
OMNISHOW (Ours)	<u>8.093</u>	0.810	0.369	8.612	7.608	-	-	0.465	0.742	10.86	5.554
<i>Text+Reference+Pose-to-Video (RP2V)</i>											
AnchorCrafter [64]	2.669	0.404	0.215	-	-	0.229	0.176	0.499	0.673	8.95	4.241
VACE [31]	7.690	0.600	<u>0.352</u>	-	-	<u>0.206</u>	<u>0.336</u>	<u>0.450</u>	<u>0.712</u>	<u>10.14</u>	5.393
OMNISHOW (Ours)	<u>6.526</u>	<u>0.474</u>	0.418	-	-	0.174	0.460	0.447	0.722	10.28	<u>4.937</u>

Quantitative Comparison

Experiments

OmniShow can achieve **industry-graded quality** while achieve **consistent controllability** across **multiple tasks**. Note that OmniShow is the **first model** that supports four multimodal inputs at one time.



Qualitative Comparison

Experiments

The model derived from the A2V training, denoted as **OmniShow-A2V**, also achieve **top-tier performance of audio-driven avatars**.

Method	*IQA \uparrow	*AES \uparrow	Sync-C \uparrow	Sync-D \downarrow
FantasyTalking [58]	2.11	1.12	1.11	12.88
HunyuanVideo-Avatar [7]	1.76	1.18	4.89	9.37
Hallo3 [10]	2.31	<u>1.48</u>	4.26	10.22
MultiTalk [34]	2.07	1.30	<u>6.34</u>	8.47
OmniAvatar [18]	2.16	1.31	5.40	9.13
OMNISHOW-A2V (Ours)	<u>2.26</u>	1.51	6.49	<u>8.97</u>

Quantitative Comparison on the EMTD Benchmark

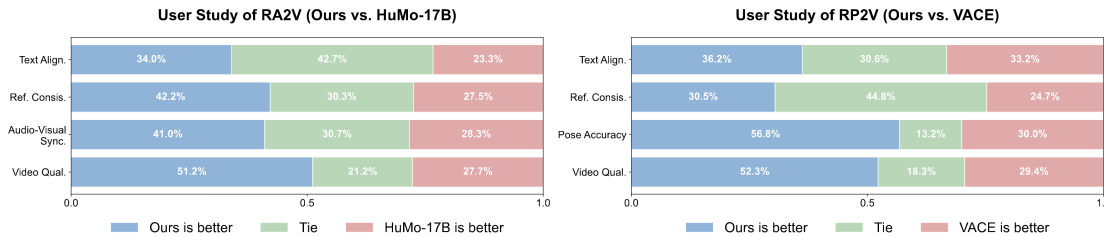
Compared with a representative cascaded baseline (VACE + LatentSync), OmniShow achieve **best results across all evaluation metrics**, avoiding the artifacts and blur often introduced by isolated lip-synchronization process.

Method	Text Align.	Reference Consistency		Audio-Visual Sync.		Pose Accuracy		Video Quality			
	TA \uparrow	FaceSim \uparrow	NexusScore \uparrow	Sync-C \uparrow	Sync-D \downarrow	AKD \downarrow	PCK \uparrow	AES \uparrow	IQA \uparrow	VQ \uparrow	MQ \uparrow
Cascaded Baseline	6.885	0.591	0.341	7.016	7.823	0.198	0.340	0.417	0.709	10.05	3.911
OMNISHOW (Ours)	7.134	0.645	0.353	7.699	7.674	0.172	0.478	0.424	0.725	11.06	5.880



Qualitative Comparison with VACE + LatentSync

Judged by human evaluators, our OmniShow achieves superior performance in overall quality, showing great alignment with **human preference**.



Qualitative Comparison with VACE + LatentSync

We conducted detailed ablation studies for our key techniques and design choices, to show the **effectiveness** and **robustness** of our method.

Variant	FaceSim \uparrow	NexusScore \uparrow	AES \uparrow	Variant	Sync-C \uparrow	Sync-D \downarrow	AES \uparrow	Variant	NexusScore \uparrow	Sync-D \downarrow	AES \uparrow
Token-Concat	0.601	0.344	0.466	w/o Audio Context	8.872	7.878	0.533	Multi-Stage (R2V \rightarrow RA2V)	0.360	13.23	0.473
w/o Ref. Rec. Loss	0.678	0.352	0.466	w/o Attn.-Map Constraints	2.201	13.01	0.545	Multi-Stage (A2V \rightarrow RA2V)	0.342	7.38	0.456
Ours	0.707	0.353	0.471	w/o Adaptive Gating	8.872	7.819	0.529	Ours	0.364	8.14	0.474
				Ours	9.023	7.419	0.540				

RoPE Strategy	FaceSim \uparrow	NexusScore \uparrow	AES \uparrow	Context Setup	Sync-C \uparrow	Sync-D \downarrow	AES \uparrow
Spatiotemporal Shift	0.279	0.339	0.456	Context Window = 11	7.020	9.588	0.527
Native (Ours)	0.707	0.353	0.471	Context Window = 5 (Ours)	9.023	7.419	0.540

Quantitative Ablation Studies

Thank you very much for your listening!