



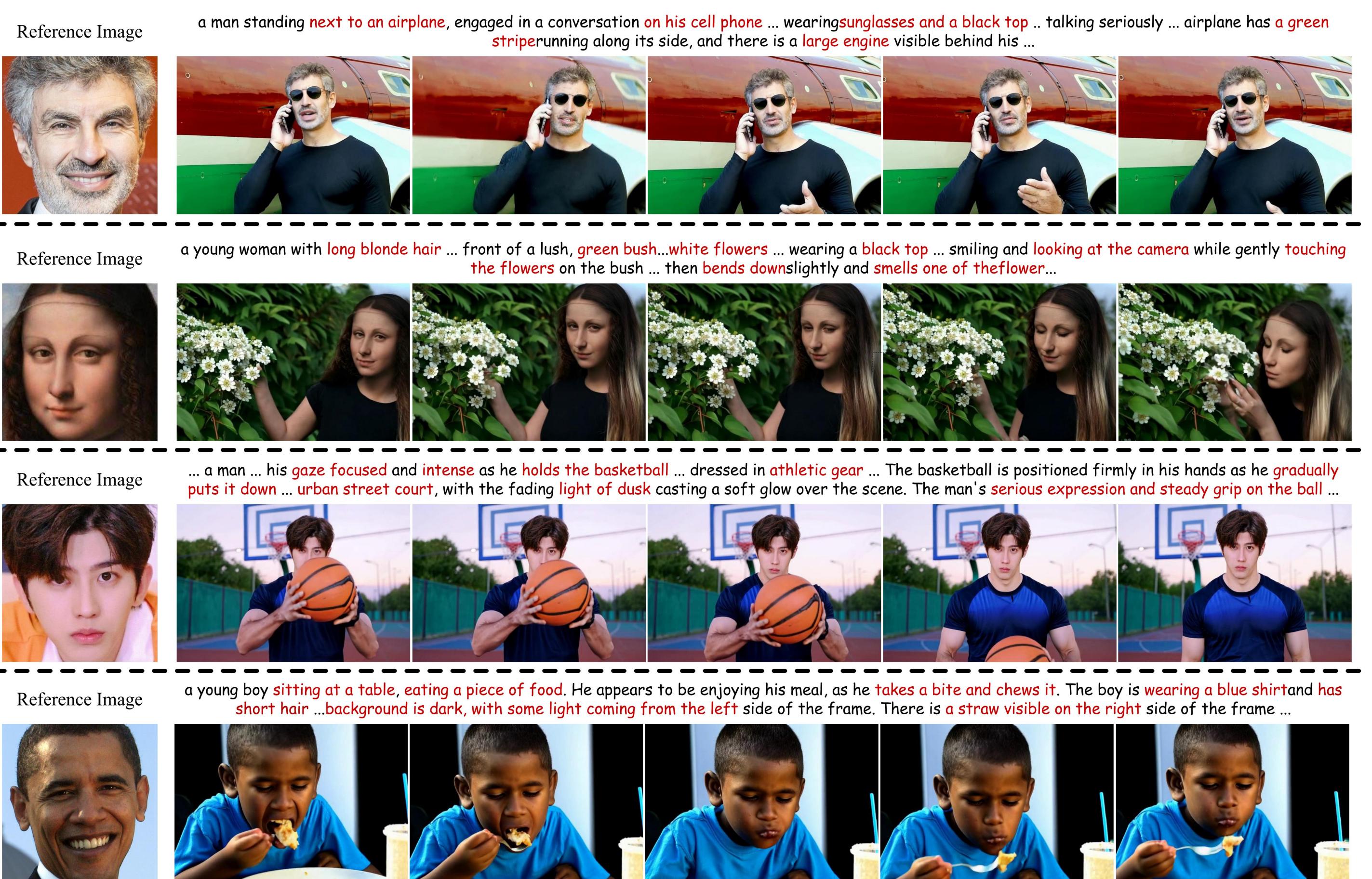
# Identity-Preserving Text-to-Video Generation by Frequency Decomposition

Shanghai Yuan, Jinfa Huang, Xianyi He, Yunyang Ge, Yujun Shi, Liuhan Chen, Jiebo Luo, Li Yuan



## Contributions

- A tuning-free pipeline without tedious case-by-case finetuning.
- A frequency-aware heuristic identity-preserving Diffusion Transformer (DiT) -based control scheme.
- A novel data pipeline to process and construct a high quality identity preservation video dataset



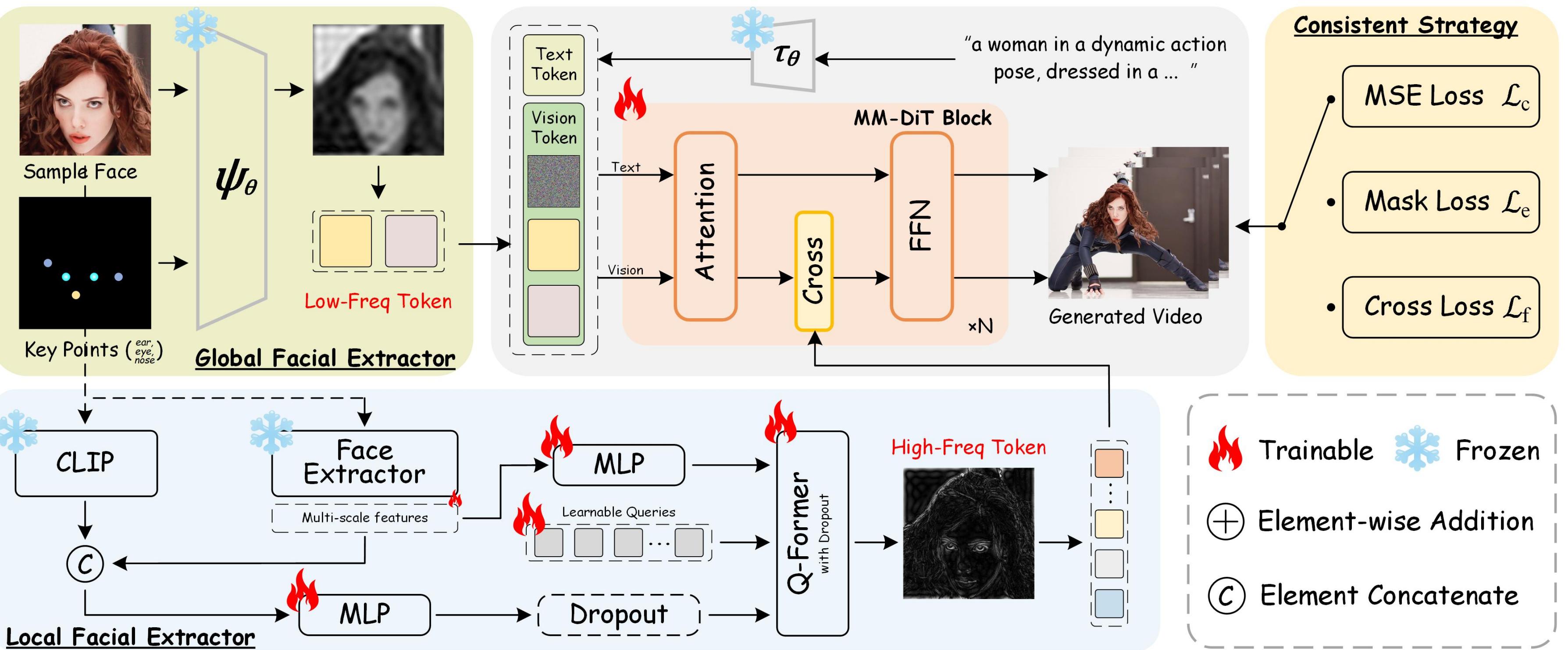
## Key Findings

These findings are not merely synthesized from and validated by existing diffusion and ViT literature, but with additional support by our experiments.

- Finding 1.** Low-frequency features are essential for pixel-level prediction tasks in diffusion models, as they ease model training. *This can solve the problem of DiT being difficult to train.*
- Finding 2.** Transformers have limited perception of high-frequency information, which is important for preserving facial features. *Inspires us to design local extractor to enhance performance.*

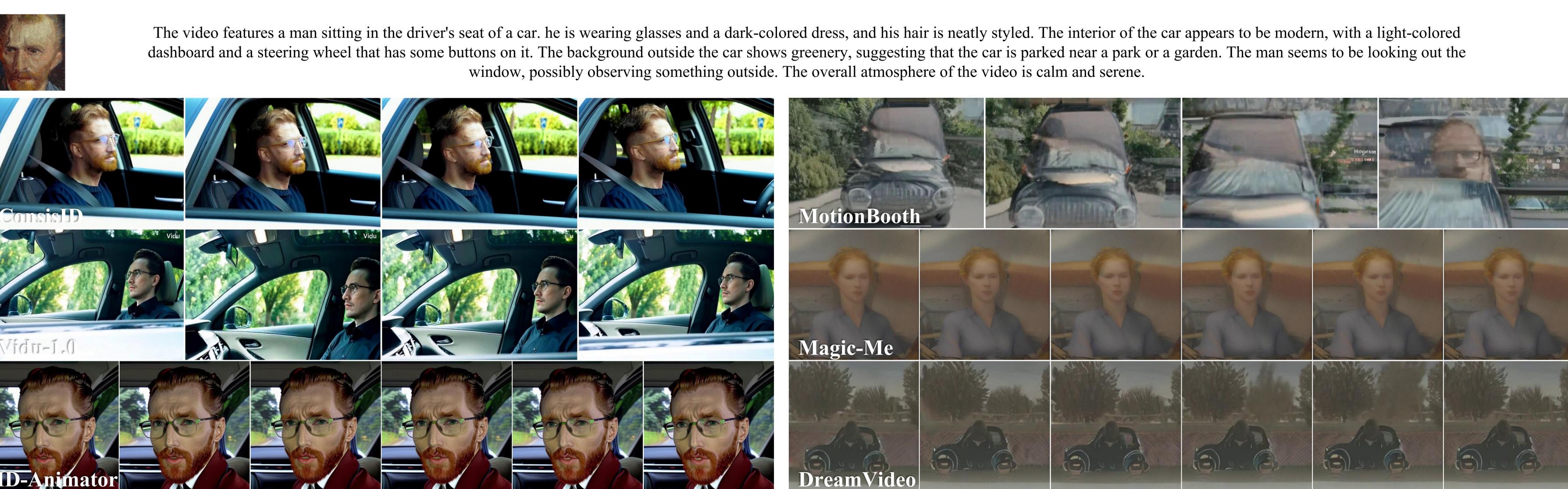
## Framework

- CosisID is designed based on both high-frequency and low-frequency perspectives and is complemented by the Consistency Training Strategy.



## Comparison

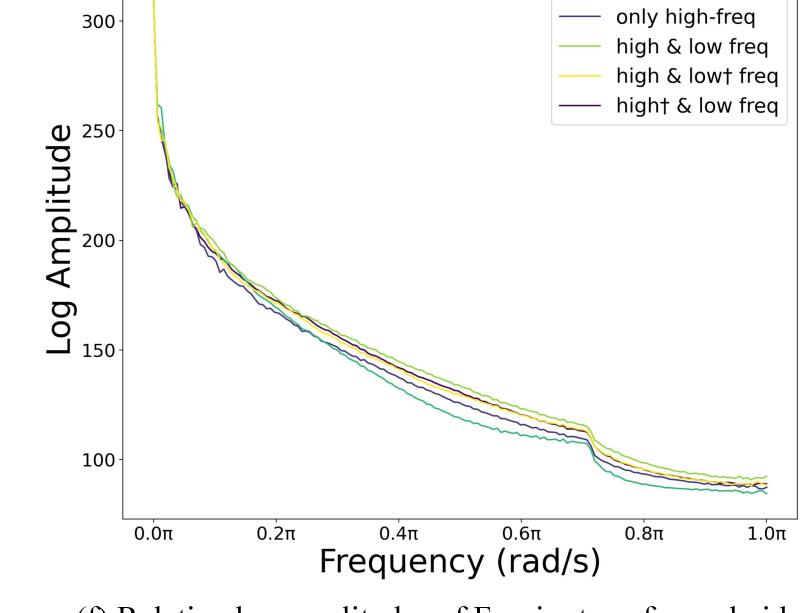
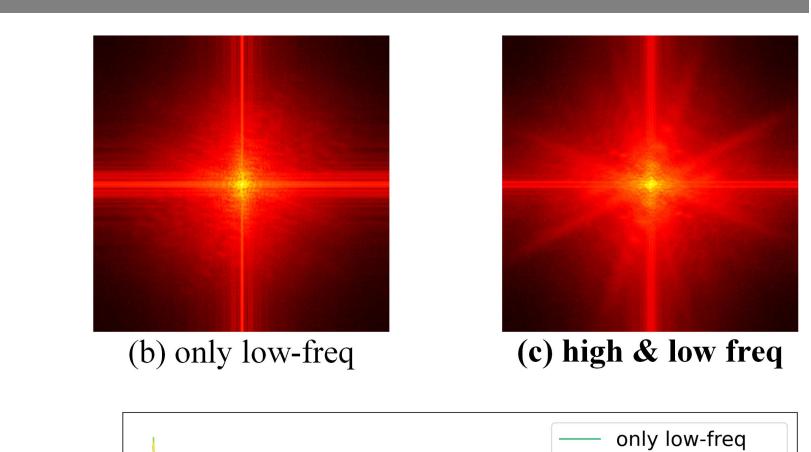
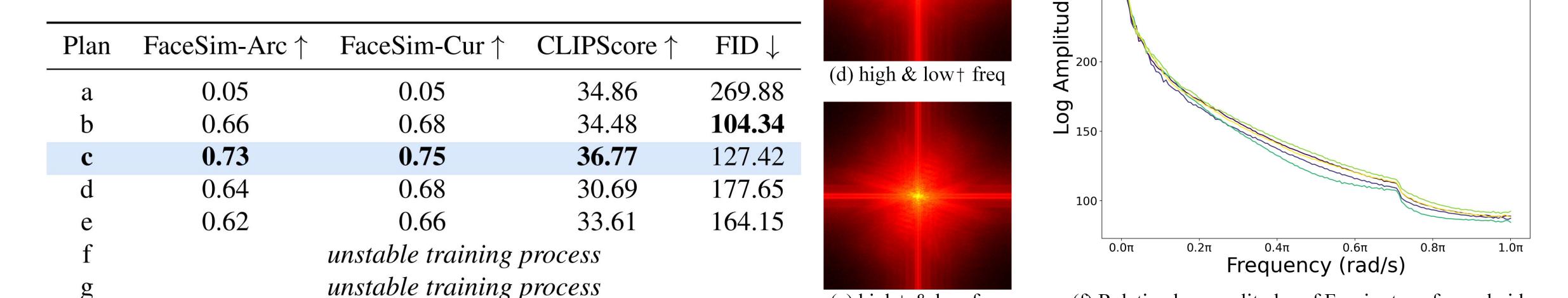
- We surpasses existing tuning-free/based, open/closed-source and I2V methods.



	FaceSim-Arc ↑	FaceSim-Cur ↑	CLIPScore ↑	FID ↓
Vidu 1.5 [6]	0.36	0.39	32.87	215.42
<b>CosisID</b>	<b>0.52</b>	<b>0.54</b>	<b>33.08</b>	<b>163.68</b>
DreamVideo [72]	0.03	0.03	26.03	237.91
MotionBooth [73]	0.05	0.06	24.42	287.90
Magic-Me [46]	0.09	0.10	23.14	237.35
<b>CosisID</b>	<b>0.46</b>	<b>0.47</b>	<b>27.45</b>	<b>181.97</b>

## Validate Key Findings

- Training instability when low-frequency features are omitted.
- Injecting high-frequency feature can lead to improvements.



- CosisID can be zero-shot transferred to animation generation.

