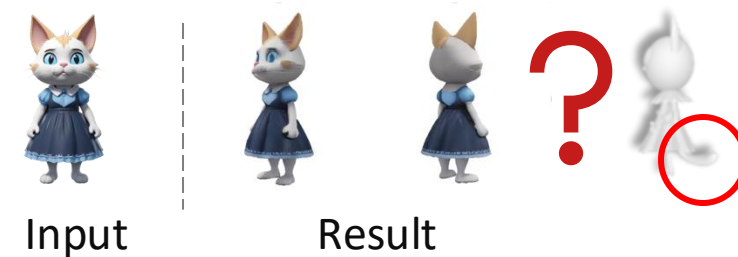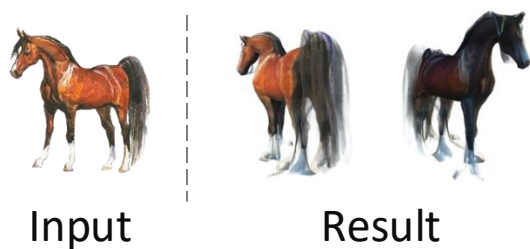# Phidias: A Generative Model for Creating 3D Content from Text, Image, and 3D Conditions with Reference-Augmented Diffusion

*Zhenwei Wang*[*1], *Tengfei Wang*[*2], *Zexin He*[3], *Gerhard Hancke*[1], *Ziwei Liu*[4], *Rynson W.H. Lau*[1]*

[*]*Joint First Author*,  [1]*City University of Hong Kong*,  [2]*Shanghai AI Lab*,
[3]*Chinese University of Hong Kong*,  [4]*S-Lab, Nanyang Technological University*

**Problems of existing image to 3D generative models:**

- **Generation quality**
- **Generalization Ability**
- **Controllability**



Input        Result

Input        Result
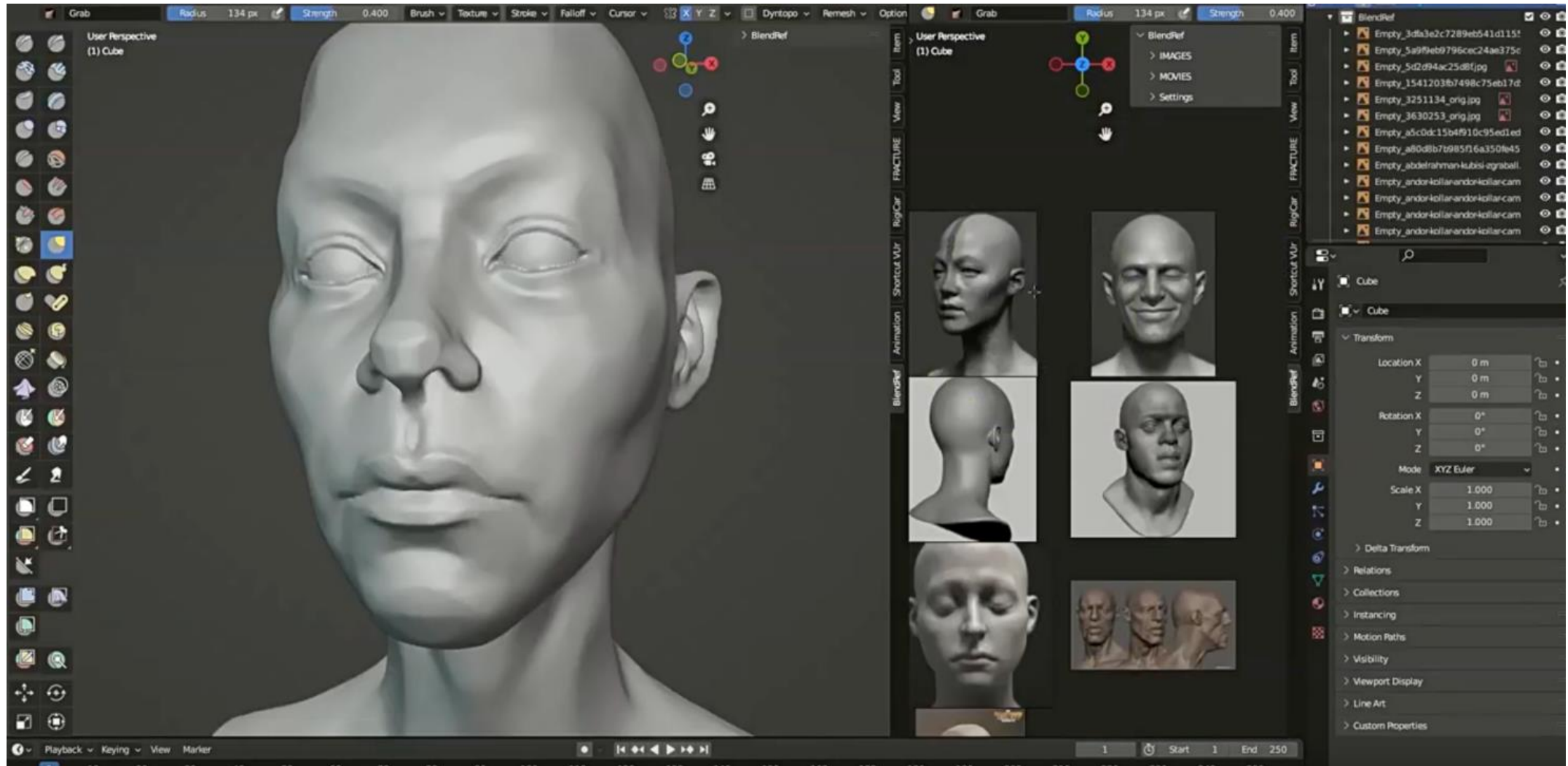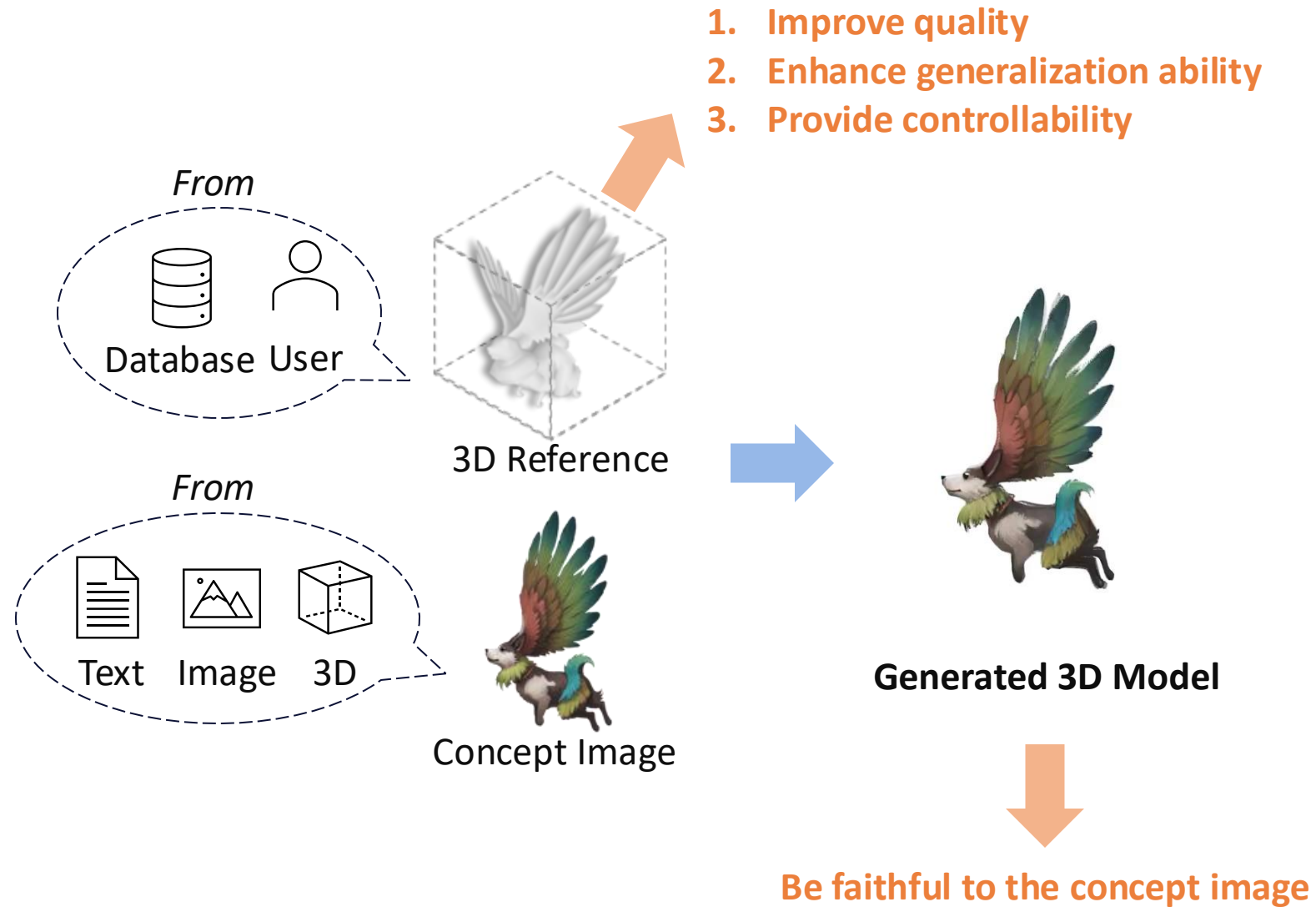
Input        Result

**Success of Retrieval augmented generation (RAG) for language and image**

**What about RAG for 3D?**
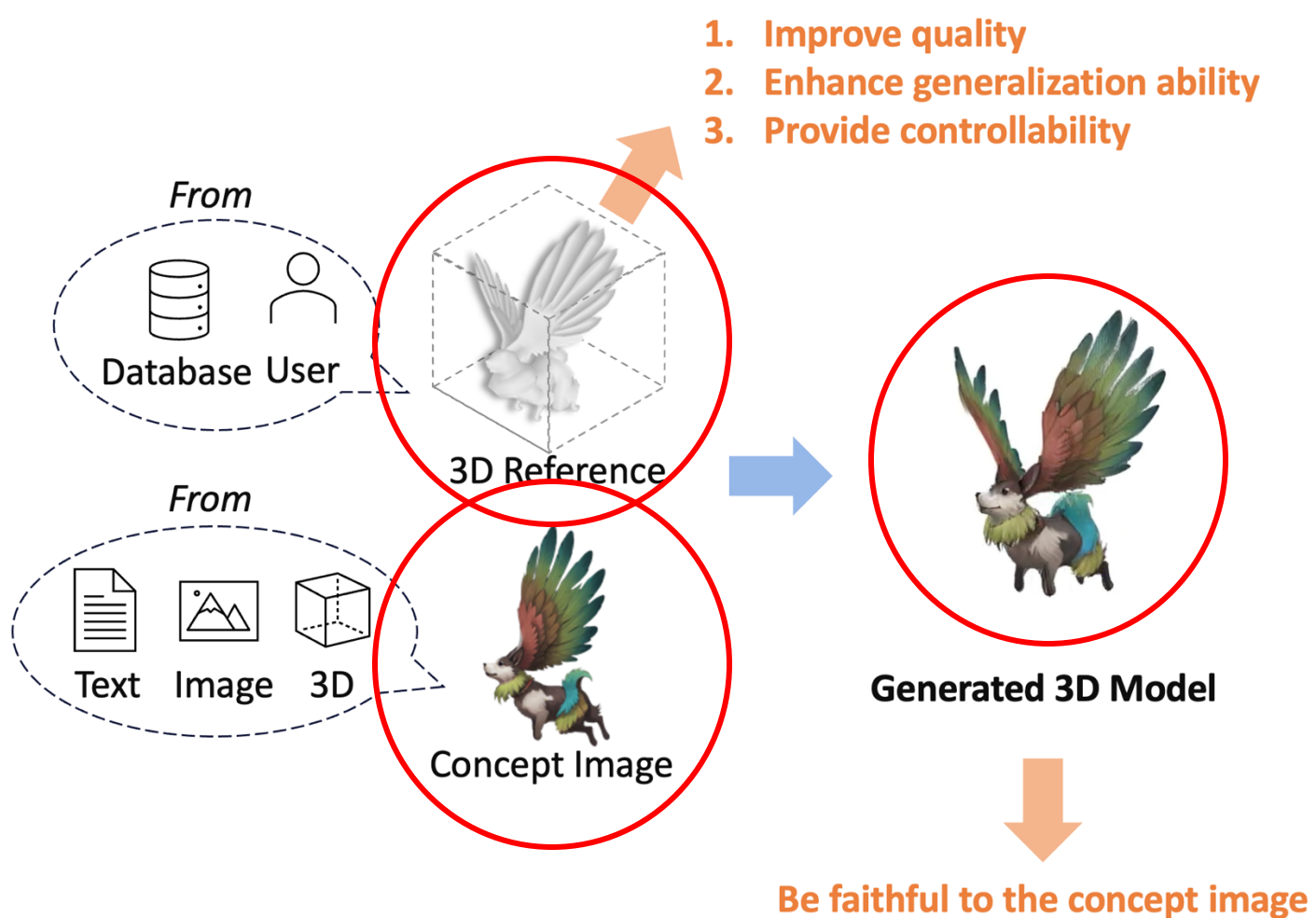
Image to 3D

Image Input — Retrieve — 3D Reference — Generated 3D Model

Image Input — Retrieve — 3D Reference — Generated 3D Model

Text to 3D

"a young anime-style girl with blue hair"

Text Input — Generated Image — Retrieve — 3D Reference — Generated 3D Model — 3D Reference — Generated 3D Model

3D to 3D

3D Input — Image variations — Self-Reference — 3D Reference — Generated 3D Variation 1 — Generated 3D Variation 2

*From*

Database  User

3D Reference

*From*

Text  Image  3D

Concept Image

**Input**

Dynamic Reference Routing

Meta-ControlNet

Adaptive Control Signal

Pretrained Multi-View Diffusion

Generated Multi-View Images

**Reference-Augmented Multi-View Generation**
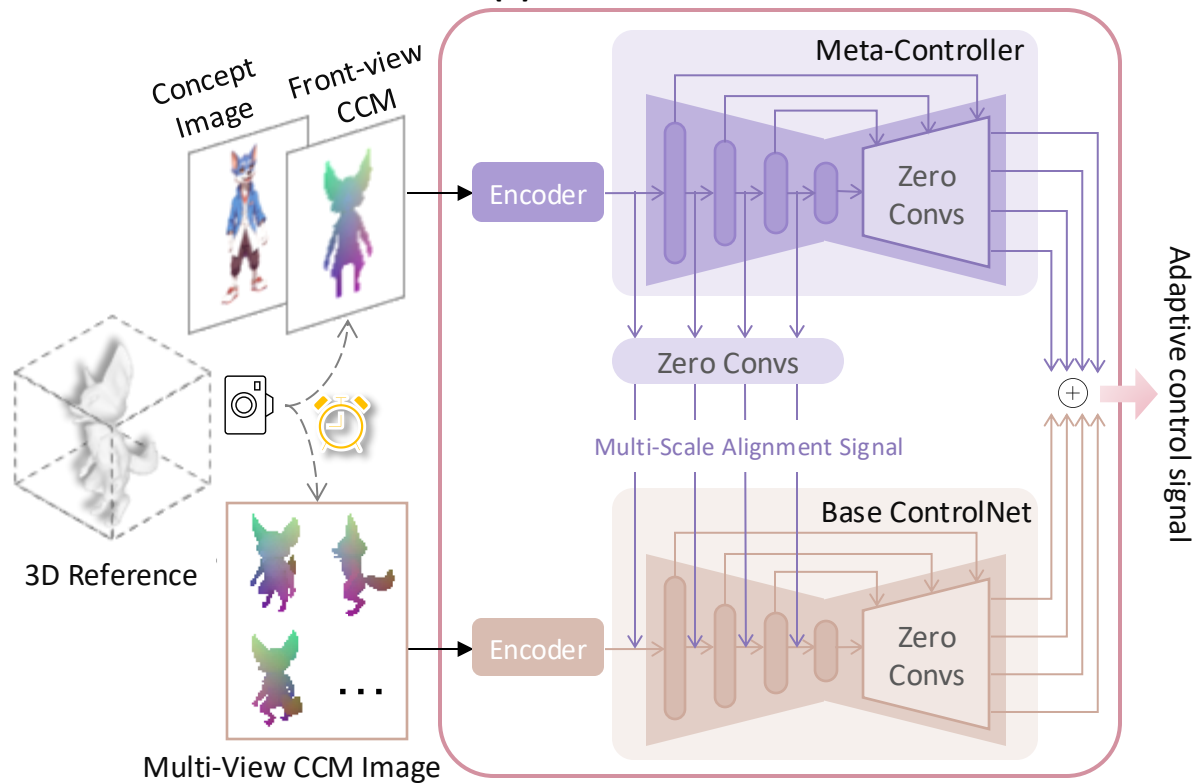
Sparse-View 3D Reconstruction

**Generated 3D Model**
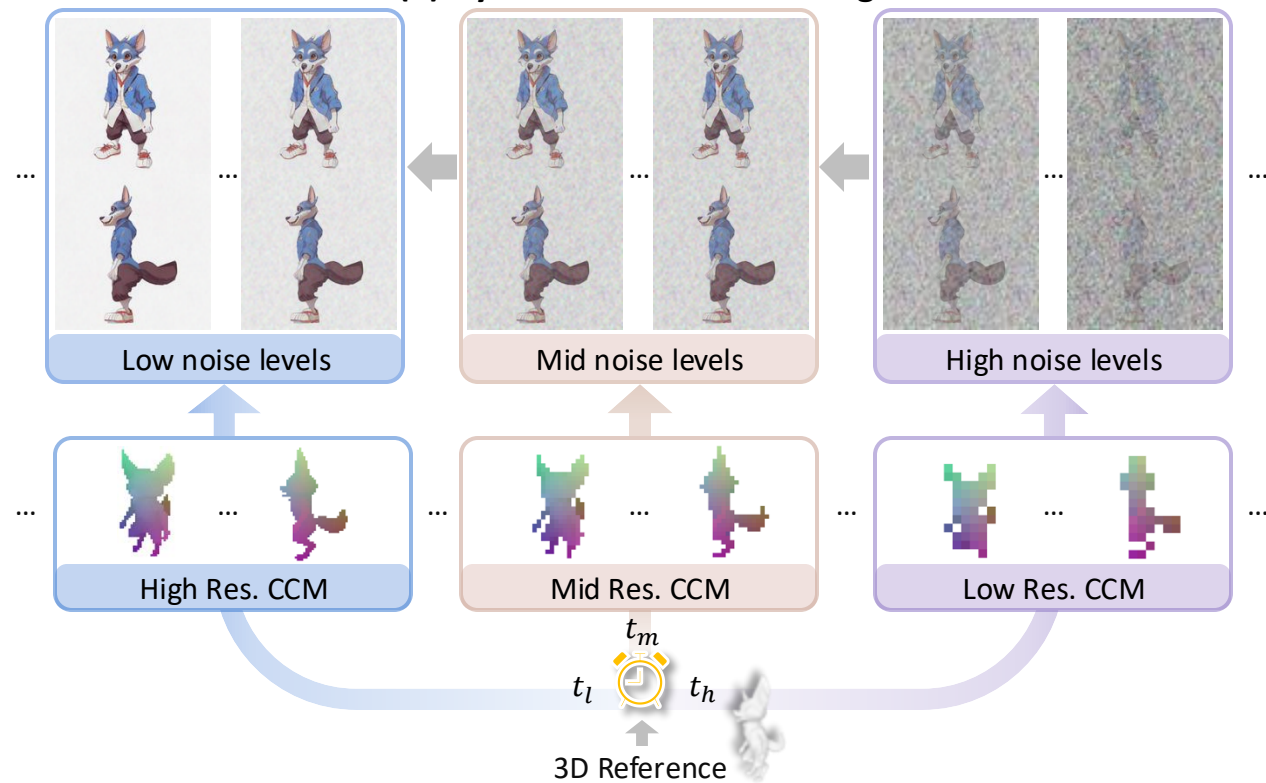
**(a) Meta-ControlNet**

**(b) Dynamic Reference Routing**
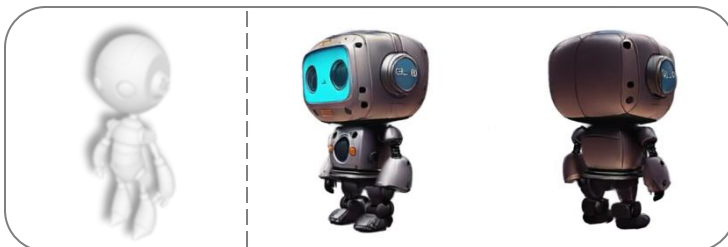
Image Input     3D Reference     Generated 3D Model     Image Input     3D Reference     Generated 3D Model
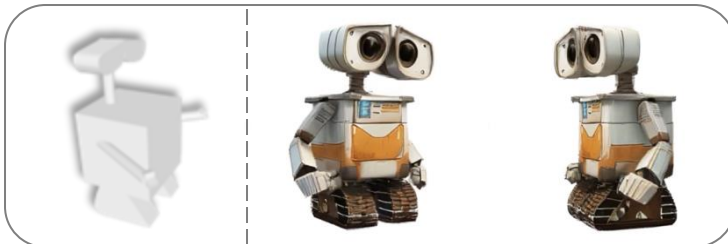
Image Input    3D Reference    Generated 3D Model    Image Input    3D Reference    Generated 3D Model

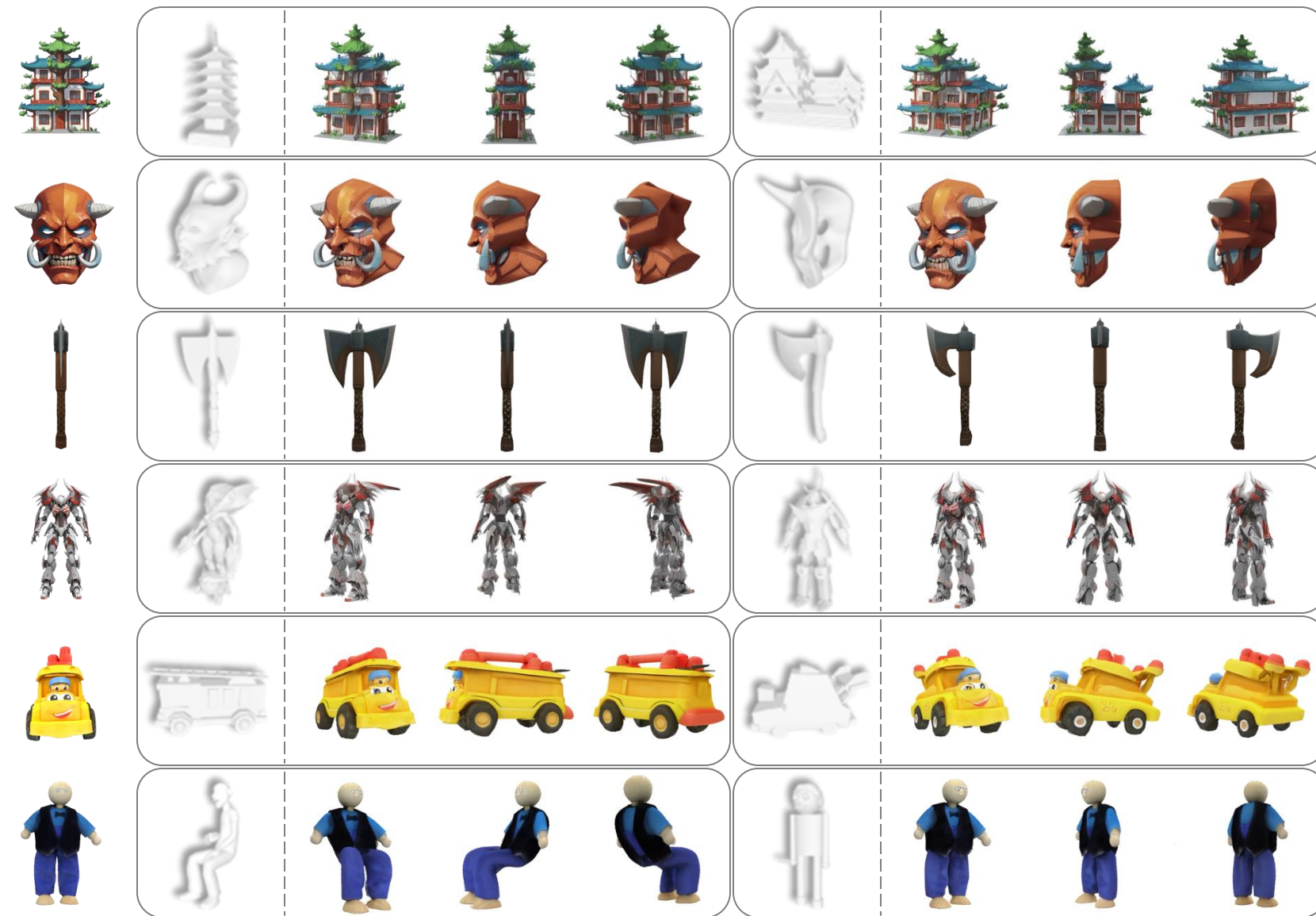Input Image | Retrieved 3D Reference 1 | Generated Model 1 | Retrieved 3D Reference 2 | Generated Model 2
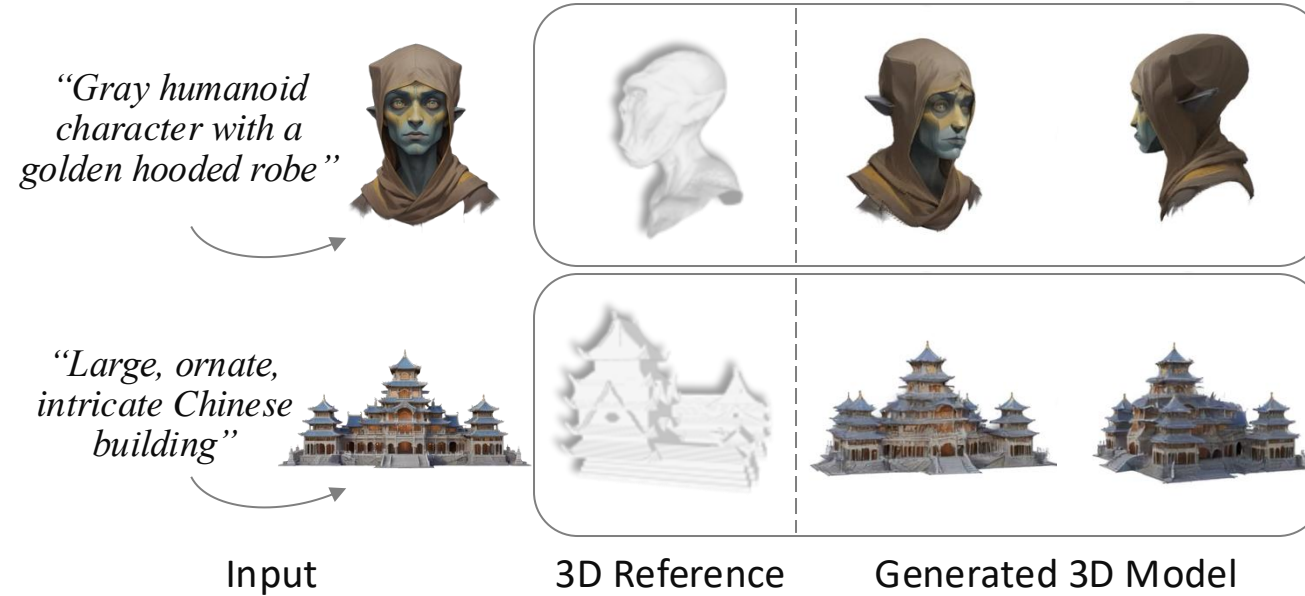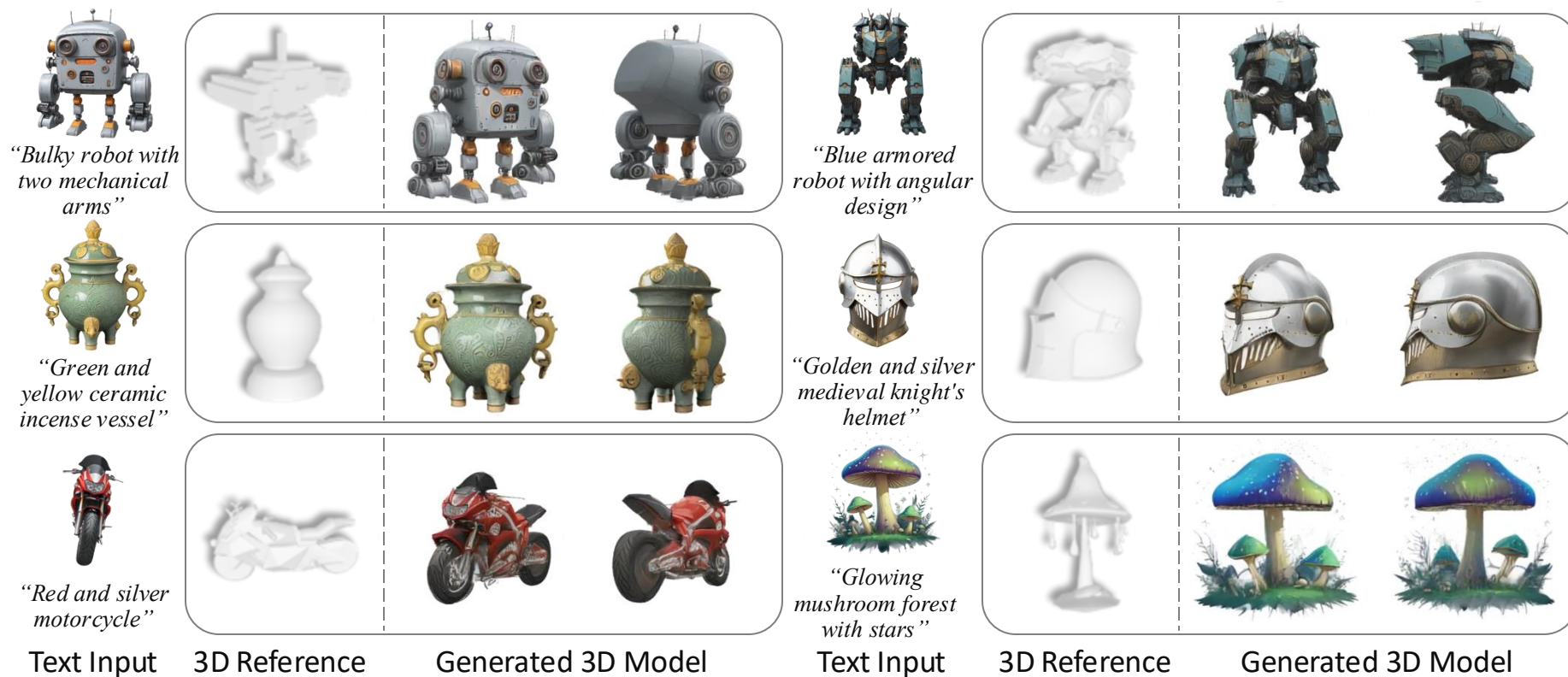
- **Retrieval-Augmented Text-to-3D Generation**

- **Theme-Aware 3D-to-3D Generation**

- **Interactive 3D Generation with Coarse Guidance**

- **High-Fidelity 3D Completion**

- **Retrieval-Augmented Text-to-3D Generation**



"Gray humanoid character with a golden hooded robe"

"Large, ornate, intricate Chinese building"

Input      3D Reference      Generated 3D Model

- **Retrieval-Augmented Text-to-3D Generation**



*"Bulky robot with two mechanical arms"*

*"Green and yellow ceramic incense vessel"*

*"Red and silver motorcycle"*

*"Blue armored robot with angular design"*

*"Golden and silver medieval knight's helmet"*

*"Glowing mushroom forest with stars"*

Text Input    3D Reference    Generated 3D Model    Text Input    3D Reference    Generated 3D Model

- **Theme-Aware 3D-to-3D Generation**



| 3D Input | Self-Reference | Generated 3D Variation 1 | Generated 3D Variation 2 |

- **Theme-Aware 3D-to-3D Generation**



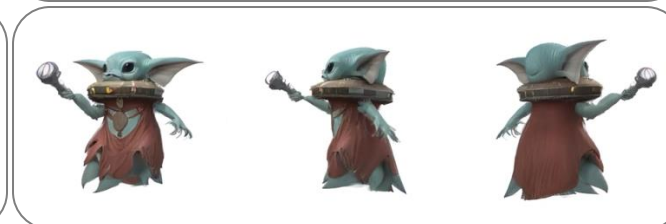3D Input    Self-Reference        Generated 3D Variation 1        Generated 3D Variation 2

- **Theme-Aware 3D-to-3D Generation**



3D Input    Self-Reference      Generated 3D Variation 1      Generated 3D Variation 2

- **Interactive 3D Generation with Coarse Guidance**
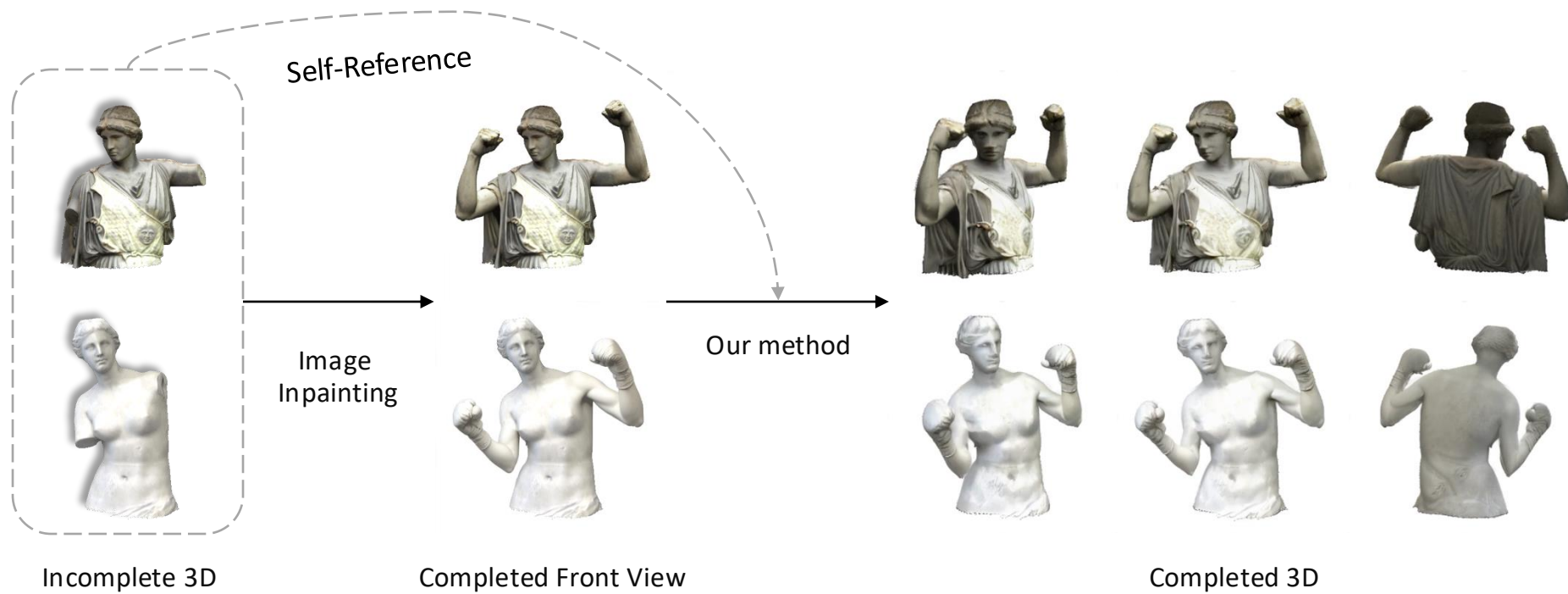


Input Image

Coarse shape

Generated 3D
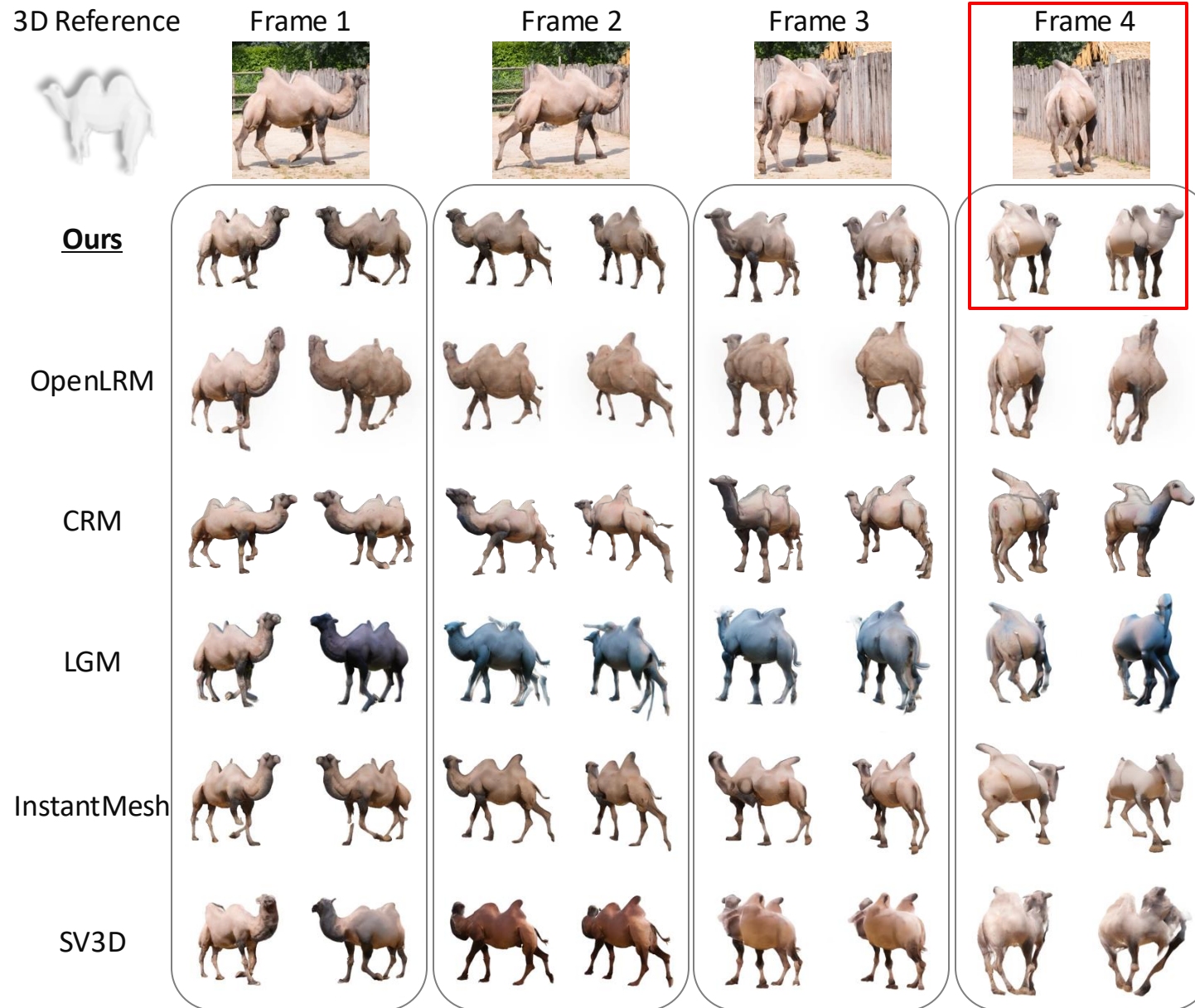
Coarse shape

Generated 3D

- **High-Fidelity 3D Completion**



Incomplete 3D       Completed Front View       Completed 3D

No 3D Reference

Random 3D Reference

GT 3D Reference

Input Image

Top-1 Retrieval

Top-3 Retrieval

Top-5 Retrieval

**(a) Angle deviation between input image and 3D reference**

| Input Image | 3D Ref. CCM | Generated 3D Model |

(30°, 20°)   (90°, −10°)

(150°, 20°)   (210°, −10°)

**(b) Semantic-aligned but structural-misaligned 3D reference**

| Input Image | 3D Ref. CCM | Generated 3D Model |

(30°, 20°)   (90°, −10°)

(150°, 20°)   (210°, −10°)