



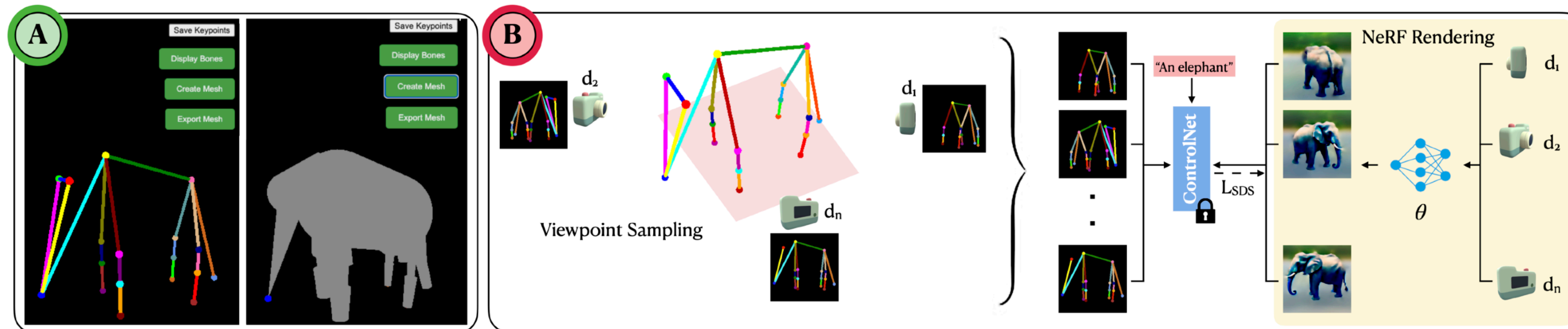
Inconsistencies in 3D animal generation for existing methods

- Existing text-to-3D methods utilizing text-to-image diffusion models generate animals with inaccurate geometry and anatomy owing to disconnected views during training.
- Training a diffusion model with camera parameters as input has shown to alleviate this problem, however requires training on 3D object datasets which are also limited.
- Animal generation using 3DMMs either 1) employ 3D scanning, thus limited in representation capability, or 2) utilize image sets and produce low-detail 3D meshes.

Generating anatomically consistent text-to-3D animals

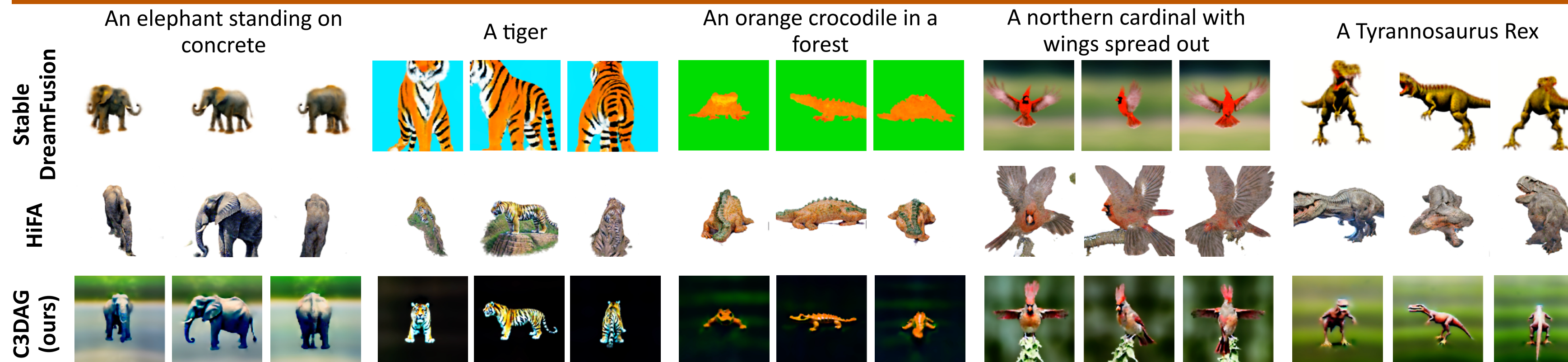
- We use 2D views of a 3D pose prior from various cameras to control image generation, achieving **multi-view 3D consistency** and relieving the need to train on 3D datasets.
- Controlled pose allows for generation of 3D animals in specific poses which are hard to describe only via text thus enhancing imagination-to-generation alignment.
- We showcase the benefits of our method compared to SOTA text-to-3D generation models such as HiFA[1], Stable-DreamFusion[2], and 3DMM based methods such as 3DFauna[3].

We create a webUI based 3D pose editor and shape initializer that uses simple geometric constructs such as ellipsoids, cylinders, and cones to generate initial shapes.



We pre-train a NeRF using the initial shape, then fine-tune it with SDS loss from our tetrapod-pose guided ControlNet, using 2D views of the 3D pose as control signals.

Results



Subsequent Work

YOU DREAM : Generating Anatomically Controllable Consistent Text-to-3D Animals

- automatic 3D pose generation using multi-agent LLMs
- enhanced 3D pose control for generating imaginary creatures
- improved 3D generation quality



References

- [1] <https://github.com/JunzheJosephZhu/HiFA>
- [2] <https://github.com/ashawkey/stable-dreamfusion>
- [3] https://huggingface.co/spaces/Kyle-Liz/3DFauna_demo