

1 Appendix

2 A Implementation Details

3 In the training process of the NeRF [Mildenhall *et al.*, 2021]
4 setting, we set rendering resolution to 128×128 , and batch
5 size to 1. We apply our random multi-view render system to
6 capture a combined image with four sub-images with rotation
7 angle α set to 90° . We use AdamW optimizer [Kingma and
8 Ba, 2014] with learning rate 1×10^{-2} and 1×10^{-3} for geom-
9 etry and background modeling. The background is replaced
10 with random colors with 80% of chance. In the DMTet [Shen
11 *et al.*, 2021] setting, most of the parameters stay the same, but
12 in the self-boost stage, we increase the resolution to 512×512
13 for a better result. The initialization stage of 3D Gaussian
14 Splatting [Kerbl *et al.*, 2023] is somehow different from the
15 other two methods as they use hash-grid while 3D Gaussian
16 Splatting is able to initialize from point cloud representation
17 directly. The rendering resolution is also 512×512 .

18 We apply the CFG trick and negative prompts following the
19 example from MVDream [Shi *et al.*, 2023], further append
20 prompt “, 3d asset” or “, multi-view of the 3d asset” to get a
21 more consistent result.

22 B Simply Combination Ablation Study

23 Our BoostDream method does not just simply combine the
24 feed-forward approach with the SDS-based method. To fur-
25 ther test the benefits of applying our multi-view based strat-
26 egy, we also design an ablation study using DreamFusion
27 [Poole *et al.*, 2022] with the same initialization stage as our
28 method. We use the results from Shap-E [Jun and Nichol,
29 2023] in the initialization stage and use the same prompt
30 text as input to optimize the NeRF representation with Deep-
31 Floyd [StabilityAI, 2023]. The results of the original Dream-
32 Fusion, the DreamFusion with initialization stage, and our
33 BoostDream-NeRF are shown in Figure 1. We can see in
34 the first row even with the proper initialization, DreamFusion
35 still suffers from the Janus problem and has coarse results
36 compared to our BoostDream results.

37 C Control Condition Ablation Study

38 We also test our method with different multi-view control
39 conditions replacing the normal map. We choose canny edge
40 [Canny, 1986] and depth map [Ranftl *et al.*, 2020] as guid-
41 ance obtained through the same multi-view render system as
42 normal map. The results are shown in the Figure 2. Canny
43 edge just contains the edge information of the 3D asset. Intui-
44 tively, it is unsuitable as a control condition when generating
45 high-quality 3D assets. The results also illustrate this point:
46 when using canny edge as the control condition, the 3D asset
47 suffers from incomplete generation. Especially in the second
48 row, the bear turns out to be unnatural and has strange colors.
49 Instead of canny edge using edge information to guide the re-
50 finement process, the depth map utilizes depth information,
51 leading to complete generation results. However, we find that
52 the generated results are less detailed when the control con-
53 dition is depth map. This can be explained by the fact that
54 minor details information is not prominent in depth map but
55 salient in normal map [Zhang *et al.*, 2023]. We can further

validate this idea with the last column, the generated 3D as-
56 sets are high-quality and with more details when under the
57 guidance of normal map.

58 D Result on Different 3D Representations

59 This section supplements the comparison experiment in Sec-
60 tion 4.3. We implement our BoostDream on other differen-
61 tiative representations, including DMTet [Shen *et al.*, 2021]
62 and 3D Gaussian Splatting [Kerbl *et al.*, 2023]. The re-
63 sults are shown in Figure 3, illustrating the generality of our
64 method in generating high-quality assets using different dif-
65 ferential 3D representations.

66 References

- [Canny, 1986] John Canny. A computational approach to
67 edge detection. *IEEE Transactions on pattern analysis and
68 machine intelligence*, (6):679–698, 1986.
- [Jun and Nichol, 2023] Heewoo Jun and Alex Nichol. Shap-
69 e: Generating conditional 3d implicit functions. *arXiv
70 preprint arXiv:2305.02463*, 2023.
- [Kerbl *et al.*, 2023] Bernhard Kerbl, Georgios Kopanas,
71 Thomas Leimkühler, and George Drettakis. 3d gaus-
72 sian splatting for real-time radiance field rendering. *ACM
73 Transactions on Graphics*, 42(4), July 2023.
- [Kingma and Ba, 2014] Diederik P Kingma and Jimmy Ba.
74 Adam: A method for stochastic optimization. *arXiv
75 preprint arXiv:1412.6980*, 2014.
- [Mildenhall *et al.*, 2021] Ben Mildenhall, Pratul P Srinivasan,
76 Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi,
77 and Ren Ng. Nerf: Representing scenes as
78 neural radiance fields for view synthesis. *Communications
79 of the ACM*, 65(1):99–106, 2021.
- [Poole *et al.*, 2022] Ben Poole, Ajay Jain, Jonathan T Bar-
80 ron, and Ben Mildenhall. Dreamfusion: Text-to-3d using
81 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- [Ranftl *et al.*, 2020] René Ranftl, Katrin Lasinger, David
82 Hafner, Konrad Schindler, and Vladlen Koltun. Towards
83 robust monocular depth estimation: Mixing datasets for
84 zero-shot cross-dataset transfer. *IEEE transactions on pat-
85 tern analysis and machine intelligence*, 44(3):1623–1637,
86 2020.
- [Shen *et al.*, 2021] Tianchang Shen, Jun Gao, Kangxue Yin,
87 Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahe-
88 dra: a hybrid representation for high-resolution 3d shape
89 synthesis. In *Advances in Neural Information Processing
90 Systems (NeurIPS)*, 2021.
- [Shi *et al.*, 2023] Yichun Shi, Peng Wang, Jianglong Ye,
91 Mai Long, Kejie Li, and Xiao Yang. Mvdream:
92 Multi-view diffusion for 3d generation. *arXiv preprint
93 arXiv:2308.16512*, 2023.
- [StabilityAI, 2023] StabilityAI. Deepfloyd. [https://
94 huggingface.co/DeepFloyd](https://huggingface.co/DeepFloyd), 2023.
- [Zhang *et al.*, 2023] Lvmin Zhang, Anyi Rao, and Maneesh
95 Agrawala. Adding conditional control to text-to-image dif-
96 fusion models, 2023.

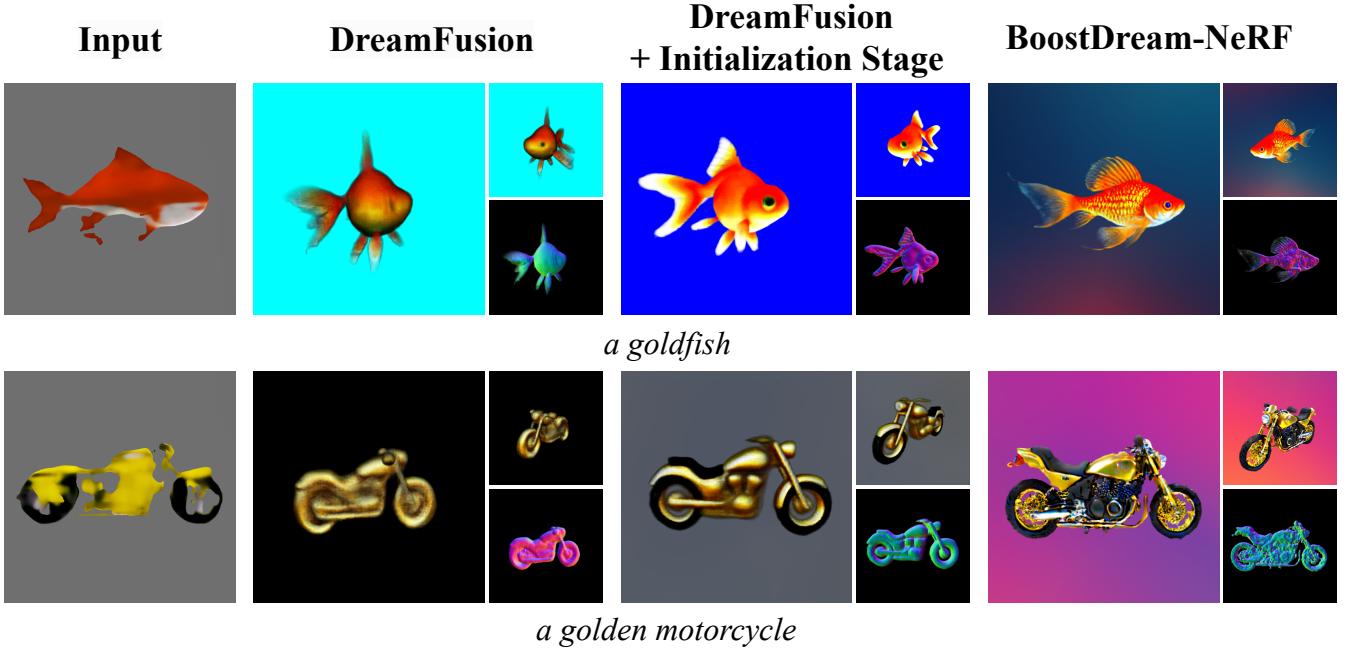


Figure 1: Simply Combination Ablation Study. The first column is the input coarse model generated by Shap-E [Jun and Nichol, 2023], while the next three columns are the results for the original DreamFusion [Poole *et al.*, 2022], the DreamFusion with initialization stage, and our BoostDream-NeRF, respectively.

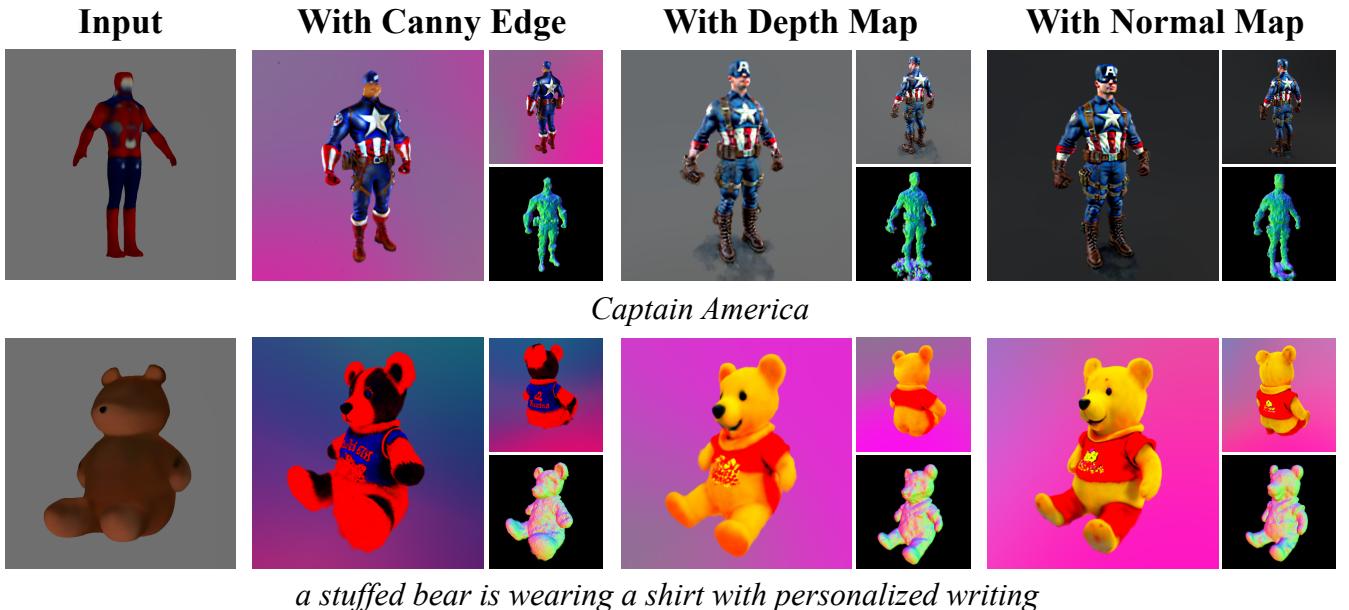
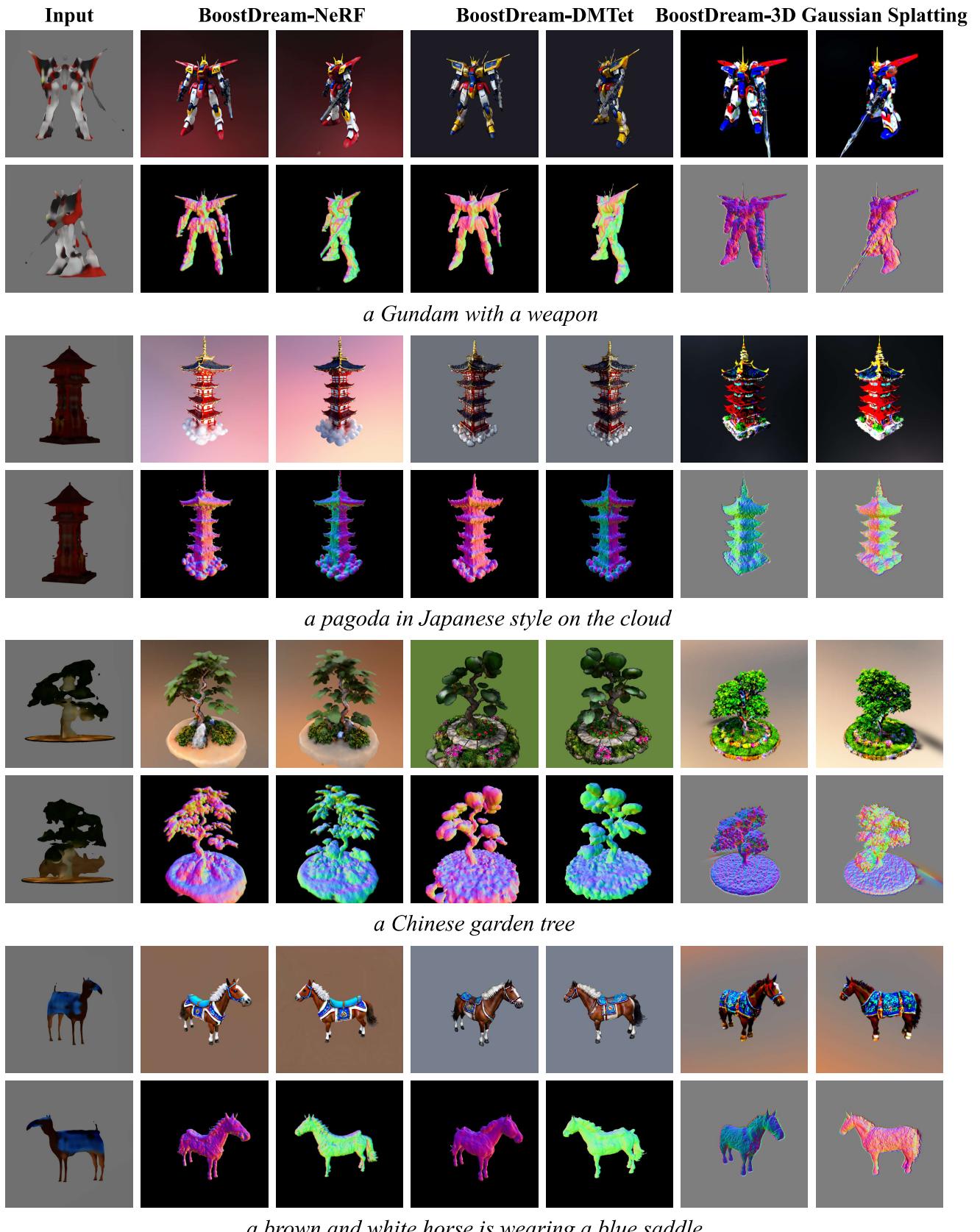


Figure 2: Control Condition Ablation Study. The first column is the input coarse model generated by Shap-E [Jun and Nichol, 2023], while all other columns are the output of our BoostDream method with different control conditions.



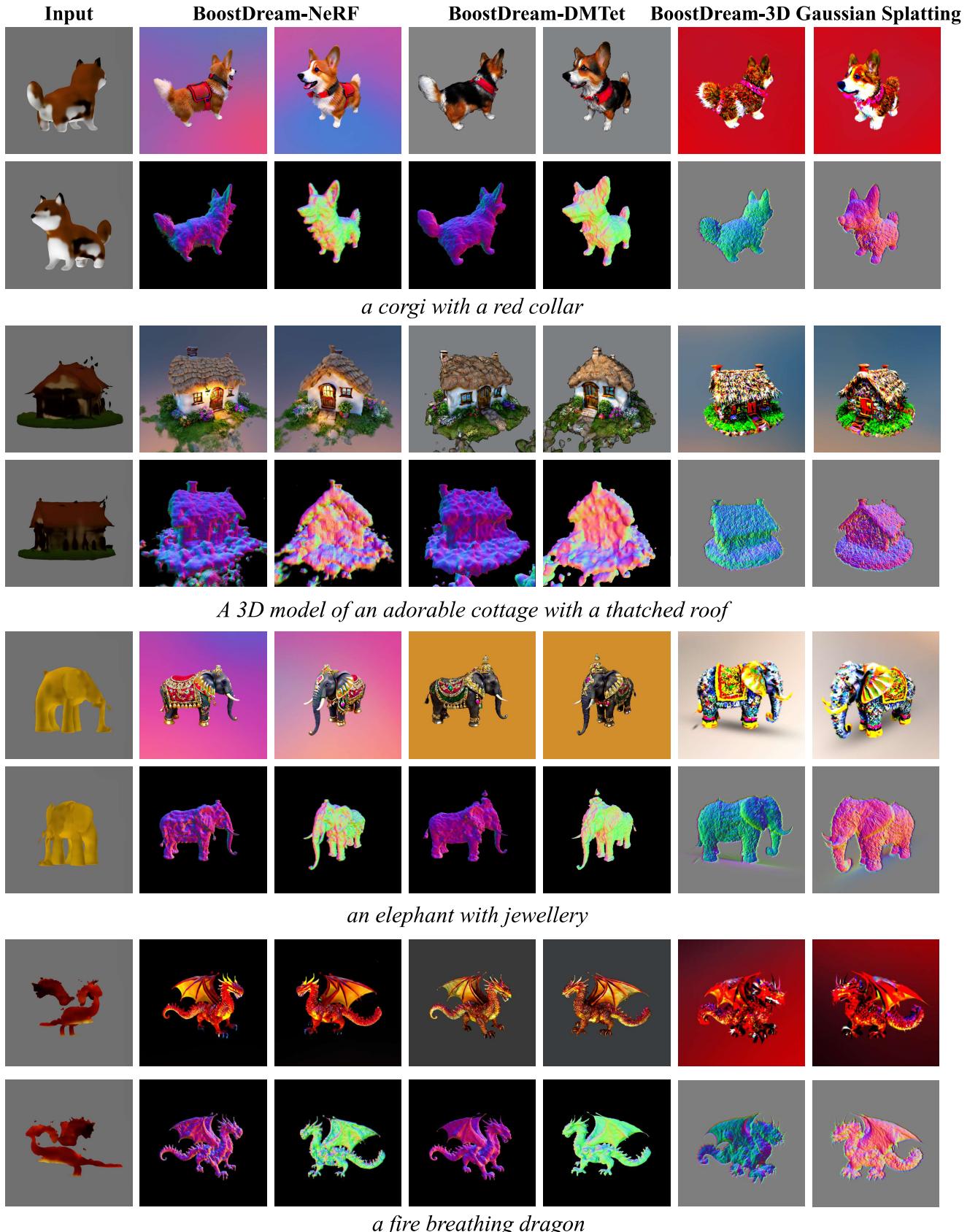


Figure 3: Result on Different 3D Representations. The first column is the input coarse model generated by Shap-E [Jun and Nichol, 2023], while the next three columns are the results of our BoostDream method implemented with NeRF [Mildenhall *et al.*, 2021], DMTet [Shen *et al.*, 2021] and 3D Gaussian Splatting [Kerbl *et al.*, 2023], respectively.