

CTR3D: Cross-view Token Reduction for Dense Multi-view Generation

Supplementary Material

846 6. Discussion and limitation

847 In this paper, we propose a method for generating dense-
848 view images from a single input image. Our approach en-
849 ables the formation of images and normals from 12 distinct
850 views. Notably, CTR3D represents the first work capa-
851 ble of producing 12 views of images and normals simul-
852 taneously with nearly the same memory usage and running
853 speed compared with existing method.

854 **Performance justification.** While full attention theoret-
855 ically yields the highest quality, its quadratic complexity
856 severely limits the view number. Our method overcomes
857 this bottleneck, enabling significantly more views (12 vs.
858 6). This provides richer geometric information, leading to
859 more accurate and complete reconstruction (e.g., Figure 7).

860 **Multi-view generation from text-generated images.** We
861 performed multi-view generation from text-generated im-
862 ages and evaluated CLIP scores on renderings of the re-
863 constructed 3D meshes. As shown in Tab. 5, our method
864 achieved the best scores compared to baselines. Fig. 7
865 further illustrates our method’s effectiveness, showing how
866 generating more views (including non-horizontal) results in
867 more complete meshes. Additionally, we also compared our
868 method with the recent baseline method LaRa [5] which
869 uses Zero123++ [48] to generate multi-view images from
870 text-generated images, and then reconstructs these multi-
871 view images into 3D Gaussians. As shown in Tab. 5 and
872 Fig.8, our method achieves better novel-view generation
873 than LaRa.

874 **Compared with existing token reduction methods.** Al-
875 though numerous methods have been proposed to accel-
876 erate token reduction in transformers [18, 19, 23], we
877 found that these approaches are generally unsuitable for
878 multi-view generation. This is primarily because, in ad-
879 dition to accelerating reduction, it is essential to restore
880 the number of tokens to ensure high-fidelity image gen-
881 eration. In our experiments, we not only compared our
882 method with ToMe but also implemented other learnable
883 reduction techniques [46, 69] that incorporate reduction-
884 recuperation mechanisms. However, we observed that these
885 methods typically incur significant overhead when restoring
886 the number of tokens. This issue is particularly pronounced
887 in high-resolution multi-view scenarios, where the number
888 of tokens increases exponentially. Consequently, the train-
889 ing costs associated with these methods are prohibitively
890 high, and our computational resources are insufficient to
891 support their training. As a result, we have excluded com-
892 parisons with these methods from our experimental results.

893 **Hyperparameter analysis.** Note that our current choice of

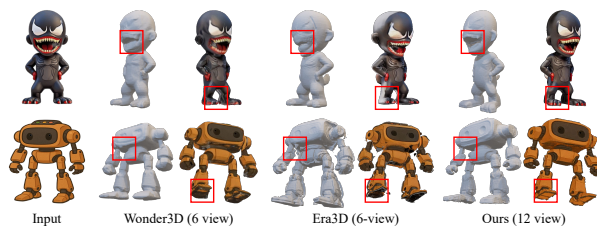


Figure 7. Qualitative comparison of text-to-image-to-multi-view.

hyperparameter K (which is [4096, 1536, 384, 96]) in our
designed model is empirical, balancing efficiency and per-
formance, which is effective in all of our experiments. To
provide a more detailed analysis, we re-evaluated the speed
of our Cross-view Token Reduced Attention layer under dif-
ferent hyperparameter settings. Tab. 6 illustrates how differ-
ent K values affect runtime across varying number of views
and generation image resolutions.

Compared with existing native 3D generation methods.

In this paper, we mainly focus on efficient high-resolution
dense multi-view generation. We present a comparative
analysis of novel view synthesis against baseline methods,
specifically Wonder3D [39] and Era3D [25]. Additionally,
we highlight the primary application of our method: re-
construction. Existing literature suggests that an increased
number of views contributes to improved reconstruction
outcomes, a claim supported by our findings on the GSO
dataset. In this study, we focus on demonstrating the results
of integrating our CTR3D with NeuS2 [56] (the same re-
construction method used in Wonder3D and Era3D) for the
reconstruction task. The reconstruction comparisons under-
score the significance of utilizing multiple views for effec-
tive reconstruction and illustrate the capability of CTR3D
in generating multi-view consistent images and normals.
Additionally, in Fig. 8, we compare the reconstruction of
our generated multi-view images using NeRF versus NeuS,
which indicates that NeuS presents higher capability for
representing details.

It is important to note that we employed an earlier re-
construction method to facilitate quantitative comparisons
with previous works, which may result in our reconstruction
quality appearing inferior to that of recent native 3D gen-
eration methods [27, 59, 67]. The primary objective of this
paper is to introduce a lightweight token reduction method
that enhances the number of generated views based on the
stable diffusion model, rather than to conduct a quality com-
parison of final reconstructions or generated geometries.

Furthermore, our method has the potential to com-



Figure 8. Qualitative comparison with LaRa.

Metric	LaRa	wonder3d	era3d	ours
CLIP Score \uparrow	0.7609	0.8206	0.8359	0.8425

Table 5. Quantitative comparison of text-to-image-to-multi-view.

K	(8192,3072,768,192)		(4096,1536,384,96)		(2048,768,192,48)	
(resolution, view)	(512,6)	(512,12)	(512,6)	(512,12)	(512,6)	(512,12)
Times (ms)	2.137	3.262	1.856	3.041	1.736	2.820

Table 6. Running time under different K settings.

932 plement recent native 3D generation techniques, such as
 933 Clay [67] and CraftsMan [27]. However, the lack of pub-
 934 licly available code and models for multi-view image in-
 935 puts prevents us from obtaining their generated results. It
 936 is also worth noting that these methods currently support
 937 only 4 input views. While they do not yet accommodate
 938 12-view inputs, they have demonstrated the importance and
 939 effectiveness of multi-view image inputs. With ongoing ad-
 940 vancements, we anticipate the emergence of native 3D gen-
 941 eration methods capable of supporting more input views in
 942 the future. At that time, the gap between generation and
 943 reconstruction is likely to narrow, producing results more
 944 aligned with the input images and further highlighting the
 945 value of our CTR3D approach.

946 **Limitations.** Due to constrained computational resources,
 947 we only conducted experiments generating 12 views. In the
 948 future, we aim to explore scenarios with denser viewpoints
 949 or higher resolutions. Additionally, our approach may occa-
 950 sionally produce smooth textures in the generated images.
 951 To address this, we plan to incorporate a refinement stage or
 952 add a lightweight detail-preserving branch in future work.

953 7. More visualized results

954 In this section, we present more testing results compared
 955 with Wonder3D [39] and Era3D [25]. As shown in Fig-
 956 ure 9, our results present better completeness and quality.
 957 We argue that this improvement mainly comes from dense
 958 views and our light-weight token reduction design. The
 959 dense views provide more overlap region for reducing the
 960 complexity of reconstruction, while the light-weight token
 961 reduction promises the cross-view attention can work on
 962 dense views as well as preserving the multi-view consis-
 963 tency.

964 We also provide the generated 12 views and recon-
 965 structed mesh in Figure 10.

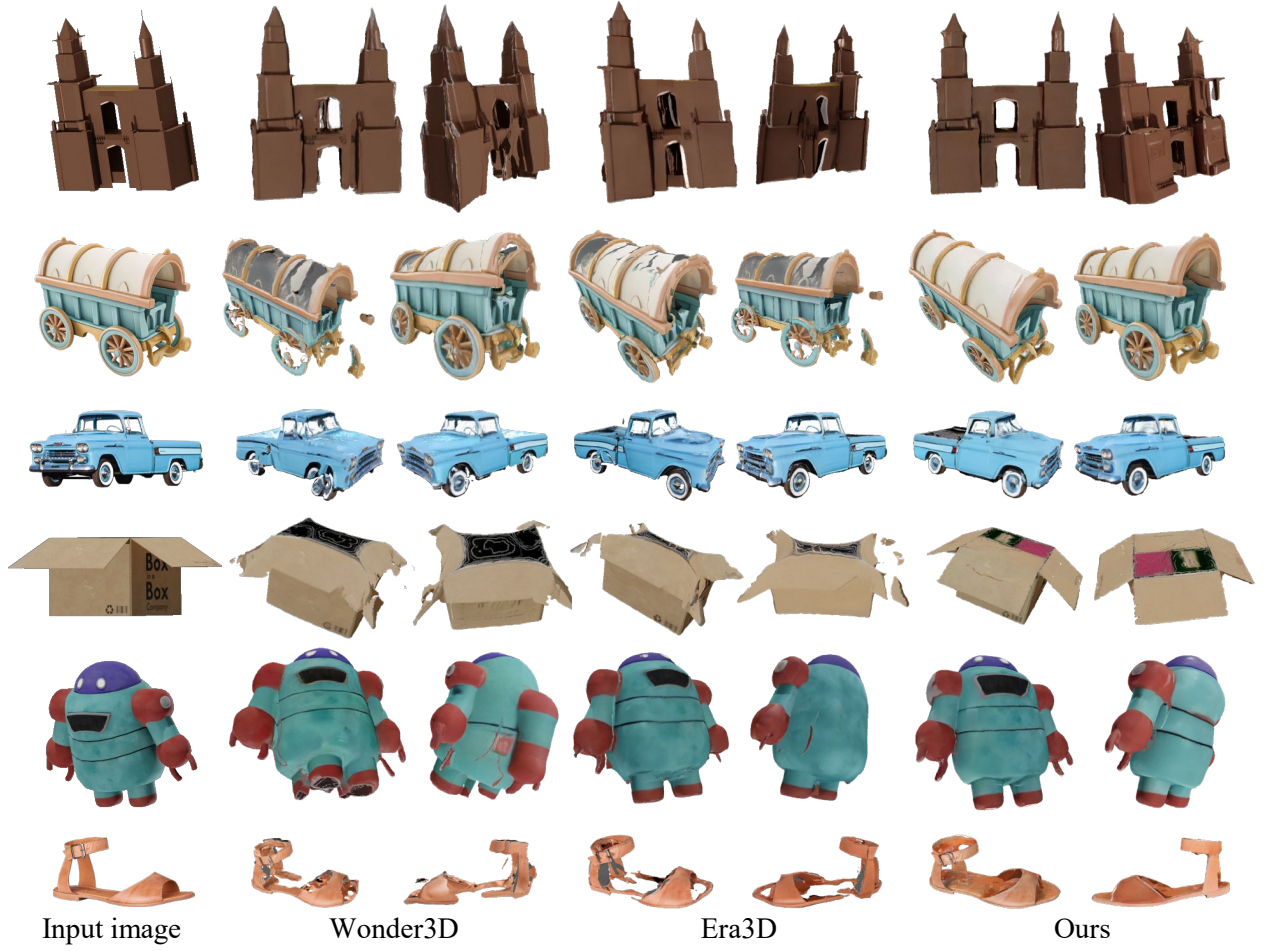


Figure 9. The visualized comparison results.



Figure 10. The visualized results generated from our CTR3D.