

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053



Figure 1: Panorama in a out of training-distribution setting.

## A APPENDIX

### A.1 OUT-OF-DISTRIBUTION EXAMPLES



Figure 2: Out-of-distribution examples generated by Cubediff. On the right are artistic generations, on the left fantasy sceneries like Alice in Wonderland.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

---

## A.2 TEXT PROMPT EXAMPLES

The text prompts for the qualitative examples are the following:

### City skyline

- A glittering cityscape at night, densely packed high-rise buildings illuminated against a dark sky. The buildings are various heights and designs, reflecting light onto the calm water in the foreground. A hill or mountain range forms a dark backdrop behind the city.
- Dark, silhouetted mountains stretching across the entire view, with the faintest glow of city lights on the horizon. The sky is a deep, nighttime blue.
- The side profile of a mountainous region, dark and slightly textured, extending from the foreground to the distant horizon. The city lights are visible in the distance to the right, providing a subtle contrast against the mountain's blackness.
- Similar to the right. A darker, less detailed section of the mountain range appears to the left.
- An expansive, dark night sky with subtle cloud formations. The upper portions of the city's tallest buildings are faintly visible as a horizontal line against the darkness.
- The dark, calm water of the bay reflecting the city lights, showing subtle ripples and disturbances. The reflection is most intense closest to the buildings and gradually fades into darkness.

### Green house

- A dense, overgrown greenhouse teems with lush greenery, creating a vibrant, jungle-like scene. Large ferns, broad-leafed tropical plants, and hanging baskets overflowing with foliage dominate the view. Large rocks, blanketed in thick, *\*vibrant\** green moss, suggest a path, nearly disappearing beneath the surrounding plants. The glass roof, tinged with green, casts a soft, diffused light.
- The right side of the greenhouse is a solid wall of plant life. Large, vibrant green leaves and hanging plants create a dense, tropical atmosphere. The rocks, almost entirely smothered in *\*bright\** green moss and overflowing greenery, offer only the faintest hint of a path. Wild, untamed growth dominates the scene.
- The greenhouse extends into the distance, an endless expanse of green fading into shadow. The rocks, now completely obscured by plants and a thick carpet of *\*emerald\** moss, suggest a path swallowed by the jungle. The sheer volume of greenery creates a sense of depth and wildness. The glass roof is barely visible.
- To the left, a tall palm tree rises amidst the dense foliage. Large ferns and other leafy plants create a solid wall of green around it. The rocks, barely discernible beneath the thick layer of *\*brilliant\** green moss and creeping vines, continue on this side, obscuring the path almost entirely.
- Looking up, the glass panes of the greenhouse roof, tinted green with clinging vines and moss, are mostly obscured by the dense canopy. Hanging plants, heavy with ferns and leafy vines, cascade downwards, creating a lush, verdant ceiling. The metal framework is almost entirely hidden.
- Looking down, the large rocks are almost completely hidden beneath a thick carpet of *\*luminescent\**, almost *\*glowing\** green moss, fallen leaves, vines, and other plant debris. Only the edges of the rocks peek through the dense greenery, making the path nearly invisible. The texture of the moss appears incredibly soft and velvety, a seamless blanket of vibrant green.

### Living room

- A cozy, modern living room with large windows allowing natural light to flood in. The room is furnished with a soft, gray sofa facing the windows and a wooden coffee table in

- 
- 108 the center. A bookshelf filled with plants and books stands against the far wall. A fireplace  
109 in the corner crackles softly, casting a warm glow across the room.
- 110 • A side view of the living room from the right, showing the side of the gray sofa facing  
111 toward the large windows. The coffee table is positioned in front of the sofa. To the left,  
112 the wall-mounted television is visible above the fireplace, and the wooden floor stretches  
113 across the room. A small side table with a lamp sits next to the sofa.
  - 114 • A view of the back wall of the room, where the bookshelf is the primary focus. The large  
115 windows let in a soft light, but the sofa and coffee table are out of sight from this angle.  
116 The fireplace is visible on the right side of the room, softly illuminating the space, and the  
117 television mounted above it is partially visible.
  - 118 • A side view from the left side of the living room, showing the bookshelf along the far wall,  
119 and the curtains gently swaying in front of the large windows. The coffee table sits on a  
120 stylish rug, but the sofa itself is mostly out of view, hidden from this angle. The soft glow  
121 of the fireplace on the far side of the room adds warmth to the scene.
  - 122 • Looking up at the ceiling, the room features modern recessed lighting, casting a soft, even  
123 glow across the space. Wooden beams accent the edges of the ceiling, adding a rustic touch.  
124 The tops of the windows and the moving curtains are visible from this angle, as the diffuse  
125 light from outside fills the room.
  - 126 • Looking down at the floor, you see polished wooden floorboards and a stylish area rug that  
127 lies under the coffee table. The legs of the sofa are visible at the edge of the rug, and a few  
128 books and a small potted plant sit on the coffee table. The contrast between the rug and the  
129 wooden floor gives the room a balanced, warm feel.

### 131 **Snow landscape**

- 132 • A view of a frozen river winding through a snowy forest. The river is partially frozen, with  
133 snow-covered banks on either side and dark, still water visible in the center. Bare, leafless  
134 trees line the edges of the river, their branches covered in snow. The overcast sky casts a  
135 soft, cold light over the scene, creating a peaceful and serene winter atmosphere.
- 136 • A side view of the snowy forest, showing the frozen river cutting through the landscape.  
137 Bare trees with snow-covered branches stretch across the scene, and the forest extends into  
138 the distance, with the winding river creating a natural divide. The snowy ground and trees  
139 give the area a quiet, isolated feel, and the soft light from the sky casts long shadows across  
140 the snow.
- 141 • A view from behind the river, where the water flows into the distance, disappearing into the  
142 snow-covered forest. The winding shape of the river is prominent, with the snow forming  
143 smooth, white edges along the banks. The bare trees on either side create a tunnel-like  
144 effect as they stretch over the river. The air feels crisp and cold, and the entire landscape is  
145 blanketed in a thick layer of snow.
- 146 • A left side view of the river, with snow-covered branches hanging over the water. The  
147 frozen river winds through the snowy forest, with the bare trees standing tall on either side.  
148 The snow is thick and untouched, creating a pristine winter scene. The soft light of the  
149 overcast sky adds a calm, cold atmosphere to the landscape.
- 150 • An aerial view of the frozen river, showing its winding path through the snow-covered  
151 forest. The dark water contrasts sharply with the white snow, and the leafless trees form a  
152 web of branches stretching across the landscape. The river curves gently through the scene,  
153 and the blanket of snow gives the forest a peaceful, quiet feel.
- 154 • Looking down at the frozen river from above, the snow forms a thick blanket along the  
155 banks, with patches of dark water visible in the center. The ground is covered in snow, with  
156 bare tree branches reaching over the river. The scene feels still and cold, with no signs of  
157 movement, and the snow seems to absorb all sound, creating a peaceful, quiet atmosphere.

158 For completeness, we provide the text prompts used for the qualitative comparison on the Laval  
159 Indoor dataset. On the left, we depict the input image and on the right we provide the text prompt.  
160  
161

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215



A concrete stairwell with orange railings leads up to a yellow door with a number 2 on it.

A bedroom with a window overlooking a snowy forest, a bed, a desk, and a dresser.

### A.3 PERSPECTIVE IMAGES FOR EVALUATION

Please note that all perceptual/text alignment metrics require another network to be computed. However, as these networks are trained with perspective images alone, the metrics would not give meaningful results when computed on panoramas. To circumvent this problem, we instead we render 10 random perspective images with a FOV of  $90^\circ$  for each panorama and use those for metric computation. Notice that we do not sample with an elevation of less than  $-45^\circ$  or more than  $45^\circ$  as other works such as MVDiffusion do not generate full  $360^\circ$  panoramas.

### A.4 USER STUDY

As described in the main paper we perform a two-alternative forced choice (2AFC) considering all competitors and variants of our method. In Figure 3, we show the percentage of wins against all 1258 draws and corresponding confidence intervals. This indicates ours methods performs significantly better in the user study.

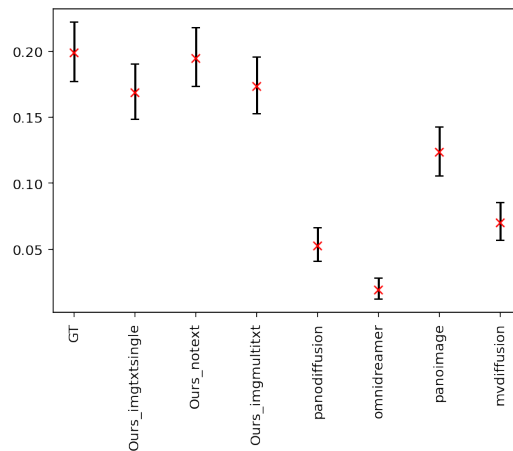


Figure 3: **Results of User study.** In this figure, we show the percentage of wins against all draws including the confidence interval

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

---

## A.5 MORE QUALITATIVE RESULTS

We provide more qualitative results. In Figures 4 and 5, we show the input modalities in green and the generated faces and panorama. We see that our model aligns well to the given text prompts.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323



Figure 4: Generated panoramas with multiple text prompts and image condition

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377



Figure 5: Generated panoramas with multiple text prompts and image condition

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

### A.6 MORE RESULTS OF UNSYNCHRONIZED GROUPNORM



Figure 6: Additional results of predictions with unsynchronized Group Norm.



432 A.7 MORE RESULTS OF NON-OVERLAPPING PREDICTIONS  
433

434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485



Figure 7: Additional results of non-overlapping predictions.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

### A.8 INDIVIDUAL FACE OVERLAPS FROM QUALITATIVE COMPARISON

Here, we depict the individual faces generated by our three CubeDiff methods used in the equirectangular panoramas in Figure 4 in the main paper. We show both uncropped and cropped faces as requested by reviewers. Additionally, we show the ground truth panoramas, the individual textual face descriptions. The corresponding input conditioning image is always the first (front) image of the individual faces (and thus equal for all models).

#### Ours<sub>img</sub>

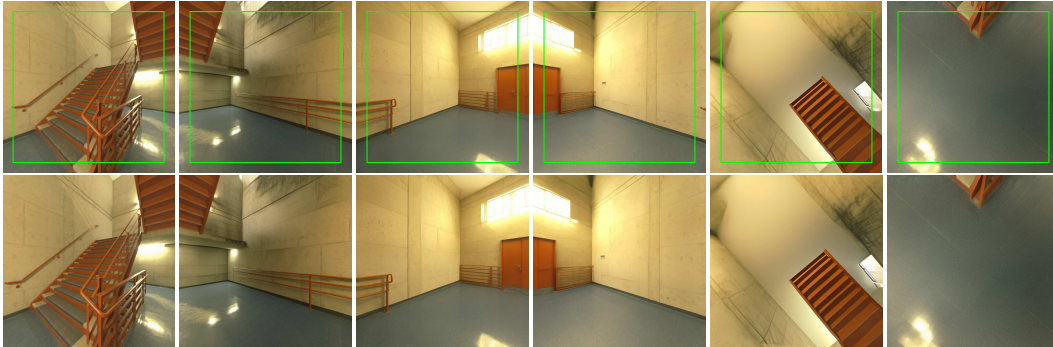


Figure 8: Our generated faces with input conditioning image only. Top row shows the uncropped faces, bottom row shows the cropped faces.



Figure 9: Our generated faces with input conditioning image. Top row shows the uncropped faces, bottom row shows the cropped faces.

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

**Ours<sub>img+txt</sub>**

- A concrete stairwell with orange railings leads up to a yellow door with a number 2 on it.

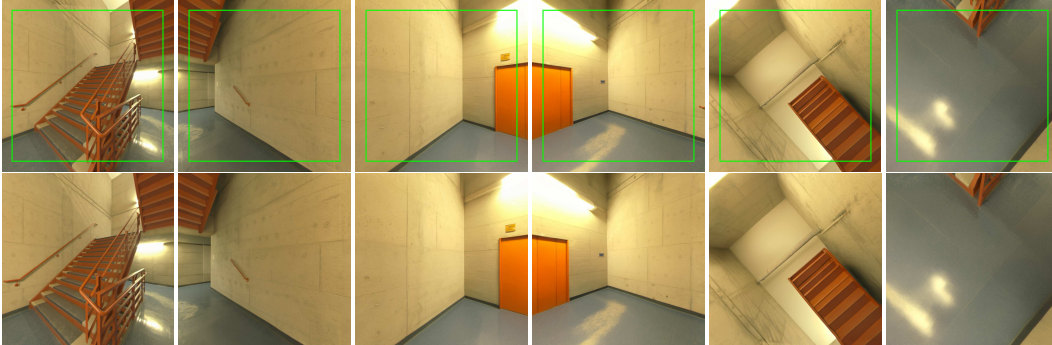


Figure 10: Our generated faces with single caption input. Top row shows the uncropped faces, bottom row shows the cropped faces.

- A bedroom with a window overlooking a snowy forest, a bed, a desk, and a dresser.



Figure 11: Our generated faces with single caption input. Top row shows the uncropped faces, bottom row shows the cropped faces.

594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

**Ours<sub>img+multitxt</sub>**

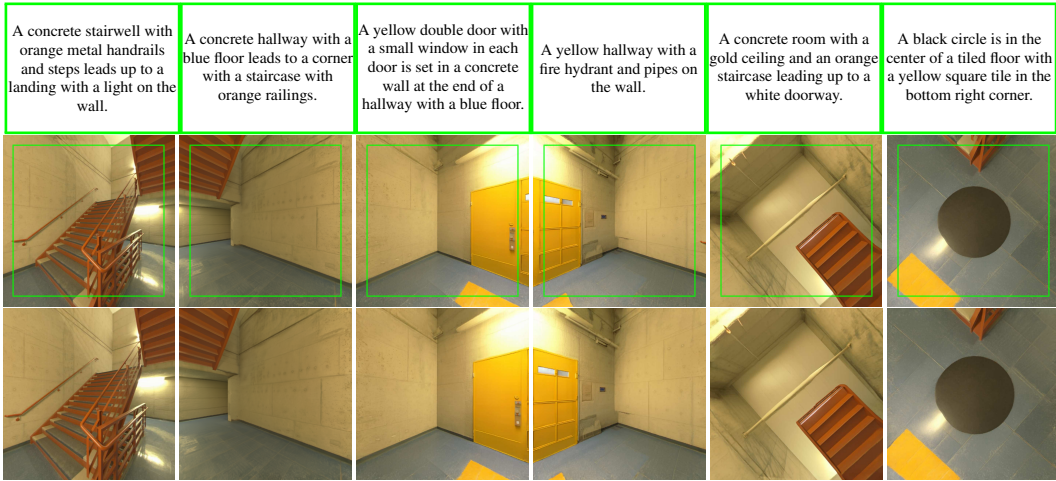


Figure 12: Our generated faces with multi caption input. Top row shows the uncropped faces, bottom row shows the cropped faces.



Figure 13: Our generated faces with multi caption input. Top row shows the uncropped faces, bottom row shows the cropped faces.

**Ground truth**



Figure 14: Ground truth panoramas from the Laval Indoor dataset.

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

### A.9 VAE RECONSTRUCTIONS

Below, we present pairs of images: the ground truth perspective images with a 95° field of view (FoV) and their corresponding reconstructed images. The reconstructed images are produced by passing the ground truth images through the encoder of our VAE and then decoding the resulting latent representations using the decoder of the same VAE.

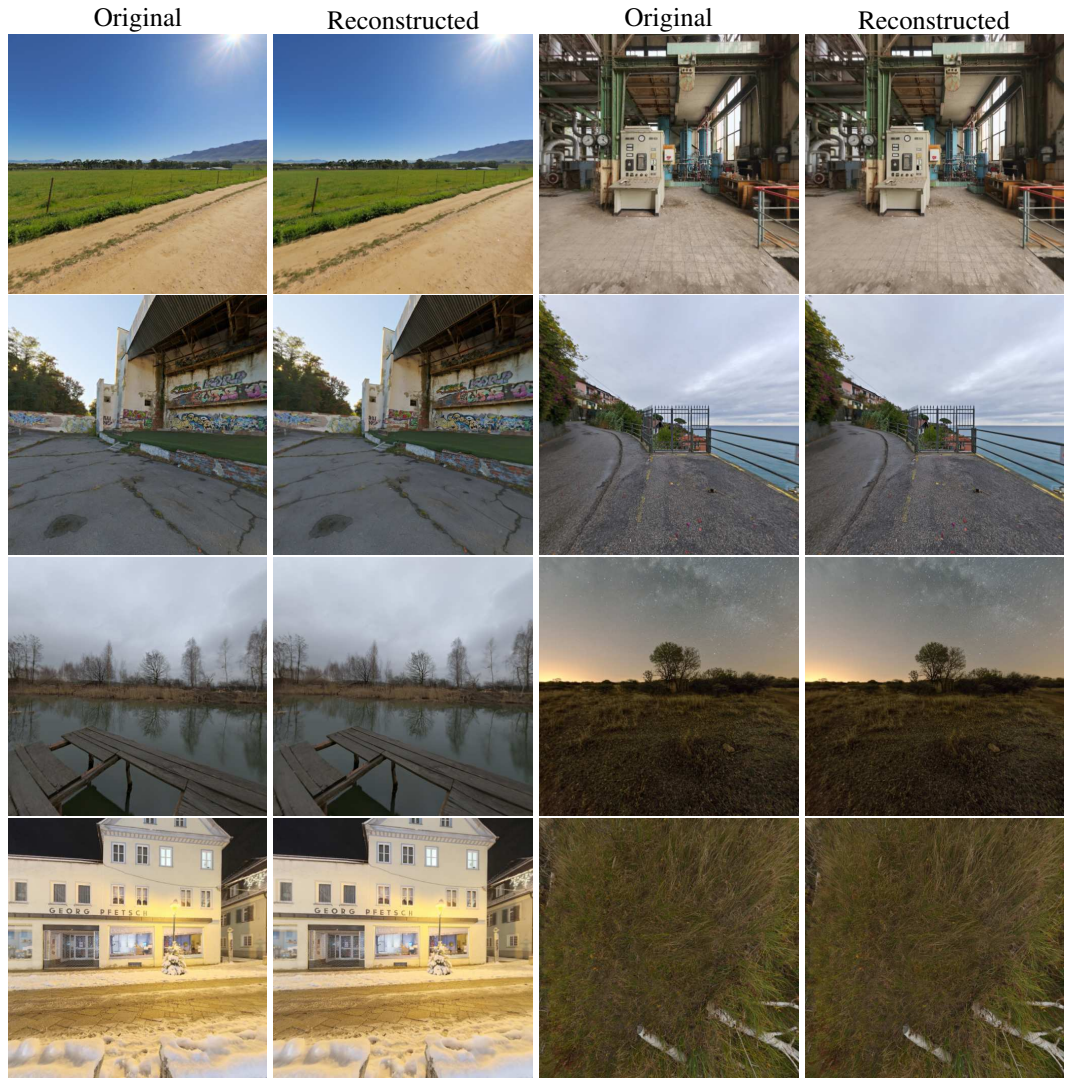


Figure 15: Examples of ground truth and encoded-decoded perspective images with a FoV of 95° using our VAE. The VAE is capable of reconstructing perspective images with out loss of quality.

A.10 DETAILED ARCHITECTURE

We illustrate our latent diffusion model in Figure 16. The model’s input is a concatenation of encoded latents, positional encodings, and an input mask indicating the conditioning image. To generate the initial latents, the input image is encoded with a VAE, while Gaussian noise is sampled for the other five faces. The VAE architecture is identical to that of Stable Diffusion’s VAE, with one modification: all GroupNorm layers are replaced with synchronized GroupNorms, where normalization is computed across both the spatial and frame dimensions.

The combined input is downsampled three times to a resolution of  $B \times 6 \times 32 \times 32$ . The first and last blocks of the model exclude attention layers and operate independently on each face. Once the final layer is computed, the output is processed through the synchronized decoder. Notably, except for the GroupNorm layers, all computations are performed per face, with no awareness of the overall panorama structure.

Despite this simplicity, our approach outperforms existing methods, as demonstrated by our results.

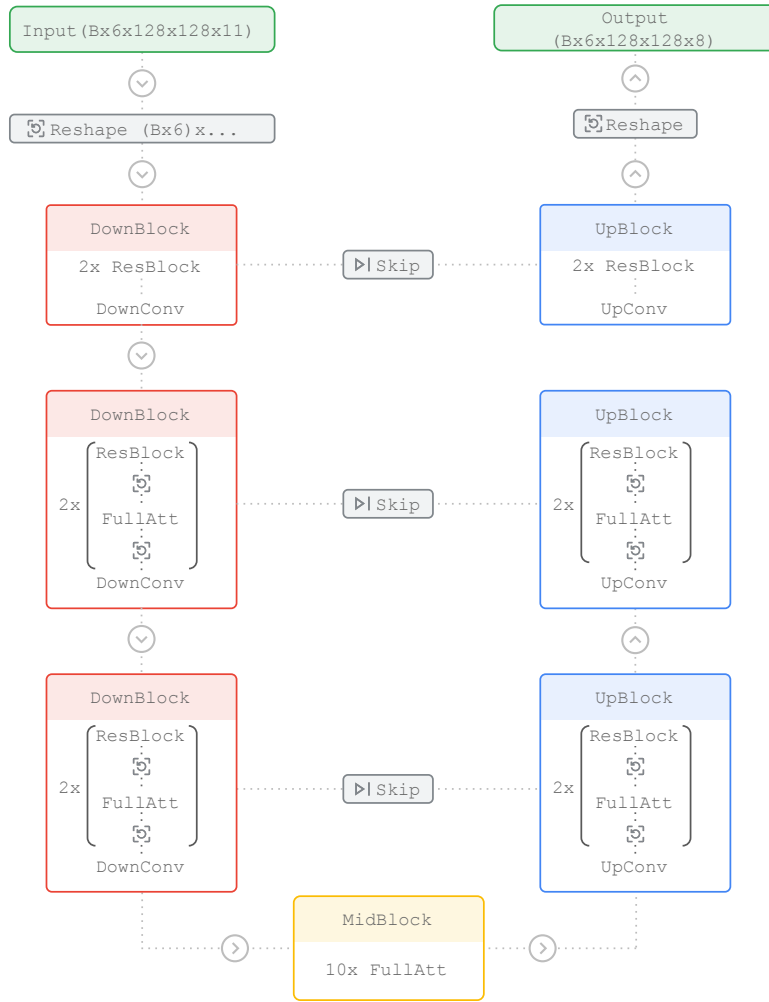


Figure 16: Illustration of our latent diffusion model.

A.11 ABLATIONS ON PANORAMIC DATA

To evaluate the impact of dataset size on the performance of our method, we conducted an ablation study by training CubeDiff on three subsets of panoramic data: a tiny dataset containing approximately 700 panoramas from the Polyhaven dataset, a medium dataset of about 20,000 panoramas from the Structured3D dataset (the same dataset PanoDiffusion used and comparable in size to MVDiffusion), and a full dataset with over 40,000 panoramas. The results demonstrate that CubeDiff performs robustly across all settings. Even the tiny model, trained on only 700 panoramas, achieves competitive results, while the medium model closely matches the performance of the full model and significantly outperforms baseline methods in most metrics. Qualitative results further confirm the ability of the tiny and medium models to generate visually consistent and high-quality panoramas, demonstrating CubeDiff’s robustness even with constrained data. These findings indicate that the superior performance of CubeDiff stems not only from data volume but also from the strength of the cubemap representation and its compatibility with pretrained latent diffusion models.

	LAVAL Indoor			SUN360		
	FID ↓	KID ( $\times 10^2$ ) ↓	Clip-FID ↓	FID ↓	KID ( $\times 10^2$ ) ↓	Clip-FID ↓
Text2Light	28.3	1.45	11.5	60.1	4.31	31.3
PanFusion	41.7	2.85	19.8	30.0	1.42	7.8
OmniDreamer	71.0	5.17	23.9	92.3	8.89	51.7
PanoDiffusion	58.6	4.08	26.6	52.9	3.51	28.9
Diffusion360	33.1	2.07	16.9	45.4	3.73	18.5
MVDiffusion	25.7	1.11	13.5	50.9	3.71	15.4
<b>Ours<sub>tiny</sub></b>	27.3	1.05	8.8	41.7	2.99	14.7
<b>Ours<sub>medium</sub></b>	13.8	0.66	8.5	23.9	1.28	10.7
<b>Ours<sub>full</sub></b>	10.0	0.35	4.1	24.1	1.33	7.0

Table 1: **Quantitative Ablation on the Laval Indoor and SUN360 dataset.** We train a model (Ours<sub>tiny</sub>) on a tiny dataset and another model (Ours<sub>medium</sub>) on a medium dataset



Figure 17: Qualitative results of the ablated models. Both the tiny and the medium model are able to generate consistent panoramas.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863



Figure 18: Example of an OOD generation of our tiny (top row) model, medium model (second row) and MVDiffusion (last row).