

000

001 **Disambiguating Monocular Depth Estimation**

002 **with a Single Transient: Supplemental**

003 **Information**

004

005

006 Anonymous ECCV submission

007

008 Paper ID 3668

009

010

011

012

013

014 **1 Ablation study on number of SID bins**

015

016

017 We conducted an ablation study on the effect of the number of SID bins [2] on

018 both runtime and RMSE. We performed this analysis using SPAD data with

019 a signal-to-background (SBR) of 100, simulated on the test set of NYU Depth

020 v2. We used DenseDepth [1] for our MDE CNN. Only the histogram matching

021 portion was timed, not the CNN nor the denoising pipeline.

022

023

024

# of sid bins	RMSE	Approx. Time/image (sec)
70	0.351	0.24
140	0.346	0.63
210	0.345	1.12
280	0.345	1.84

030 Fig. 1: Effect of number of SID bins on RMSE and runtime. The marginal improvement

031 in RMSE is offset by the increase in runtime as the number of bins grows.

032

033

034

035

036 **2 Ablation study on effect of reflectance estimation**

037

038

039 We conducted an ablation study on whether the use of a reflectance estimate

040 has an impact on the runtime and quality of the solution. We performed this

041 analysis using SPAD data with a signal-to-background (SBR) of 100, simulated

042 on the test set of NYU Depth v2 and using DenseDepth [1] for our MDE CNN.

043 Only the histogram matching portion was timed, not the CNN nor the denoising

044 pipeline.

045

046

047

048

049

050

051

052

053

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

	Intensity-weighted histogram	Intensity-aware pixel movement	Avg. RMSE	Time per image (sec)
Yes	Yes	0.346	4.6	
	No	0.346	0.6	
No	Yes	0.444	4.7	
	No	0.444	0.6	

Fig. 2: Effect of reflectance modeling on RMSE and runtime. When the SPAD is simulated with the reflectance info but no reflectance estimate is used to generate a weighted histogram from the CNN depth map, the results are significantly worse. Furthermore, once the pixel movement matrix has been computed, the pixel movement procedure need not take into account the weights of the pixels being moved, since doing so provides no improvement and can take appreciably longer than a vectorized implementation that does not take pixel weights into account.

3 Pseudocode, pixel shifting, and dither artifacts

We give pseudocode for our algorithm here. In the first part of our algorithm we compute the pixel shifting matrix mapping the histogram h_s (computed from the initial depth map and reflectance estimate) to h_t (computed from the captured transient).

Algorithm 1 Find Pixel Movement

```

procedure FINDPIXELMOVEMENT( $h_s$  of length  $M$ ,  $h_t$  of length  $N$ )
    Initialize  $T$  as an  $M \times N$  array of zeros.
    for  $m$  in  $1, \dots, M$  do
        for  $n$  in  $1, \dots, N$  do
             $p_s \leftarrow \sum_{i=1}^{n-1} T[m, i]$ 
             $p_t \leftarrow \sum_{i=1}^{m-1} T[i, n]$ 
             $T[m, n] \leftarrow \min(h_s[m] - p_s, h_t[n] - p_t)$ 
        end for
    end for
    return  $T$ 
end procedure

```

Given this pixel movement matrix T , we apply the appropriate movements to the initial depth map I . The pixels of the image I take depth bin values in $\{0, \dots, K - 1\}$.

Algorithm 2 Move Pixels

```

090
091 procedure MOVEPIXELS(input image I size M × N, pixel movement matrix T of
092 size K × K)
093     for k in  $0, \dots, K - 1$  do
094          $p[k, :] \leftarrow T[k, :] / \sum_{i=1}^K T[k, i]$ 
095     end for
096     for m in  $1, \dots, M$  do
097         for n in  $1, \dots, N$  do
098             Sample  $k'$  according to  $p[I[m, n], :]$ .
099              $I[m, n] \leftarrow k'$ .
100        end for
101    end for
102    return I
103 end procedure
104

```

Because the pixel shifting process in Algorithm 2 contains a sampling step, it is possible for *dither artifacts* to appear in the output image *I*, as shown in figure 6. Specifically, when there are multiple possible output depth bins for a given input depth bin, and a large region of equal depth in the input image, the randomness in the pixel shifting algorithm will distribute the pixels of large, equal-depth region in the input across the multiple possible output depth bins in a random fashion.

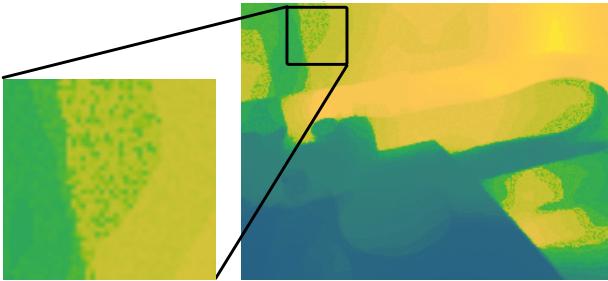


Fig. 3: Example of dither artifacts. Sometimes, when our histogram matching is applied to images with large regions of similar depths, dither artifacts will occur.

4 Nominal values for SBR calculation

Here we give the values used in the equation for received power, repeated here for convenience:

$$P_R = P_T \cdot \rho \cdot \frac{A_{rec}}{\pi R^2} \cdot \eta. \quad (1)$$

The object in question is a vertical, planar, perfectly Lambertian surface. The following table gives the values used for this calculation.

Symbol	Description	Nominal Value
ρ	Albedo of lambertian surface	0.3
P_T	Total irradiance at wavelength (W/m^2)	0.026
R	Distance to surface (m)	3
A_{rec}	Area of detector (m^2)	1.96×10^{-9}
η	Quantum efficiency of detector	0.3
P_R	Received power at detector (W)	1.62×10^{-13}

Fig. 4

Once P_R is determined, we compute the number of photons using the laser wavelength $\lambda = 532$ nm as

$$N = \frac{P_R \lambda}{hc} \quad (2)$$

where $h \approx 6.62610^{-34}$ is Planck's constant and $c \approx 3 \times 10^8$ is the speed of light. Now using the fact that our laser runs at 10 MHz, we get the number of photons per pulse as 0.043, which puts us in the low-flux regime.

5 Comparison of diffused vs. scanned imaging

In our experiments, we capture measurements by scanning the scene with a single-pixel SPAD detector whose optical path is aligned with a laser. This arrangement allows us to capture a reference "ground truth" depth map for quantitative validation of our method. To emulate measurements captured using a system where the laser and detector are diffused over the scene, we digitally sum the measurements to obtain a single transient.

In order to verify that digital summation of scanned measurements yields results that are similar to those captured by a diffused laser and detector, we capture an example scene using a modified hardware prototype in both scanned and diffused modes. The hardware prototype (shown in Fig. ??) consists of a more powerful laser (Katana 05HP, 532 nm) which we operate at approximately 25 mW output power with two single-pixel SPAD detectors. One SPAD is aligned with the optical path of the laser, and the other SPAD is operated without a lens to integrate light from the entire scene. Both SPADs are fitted with a 10 nm bandpass filter centered at 532 nm, which reduces the amount of integrated ambient light. We attach a holographic diffuser (Thorlabs ED1-S50) to the laser output in order to diffuse light onto the scene or alternatively remove the diffuser and use a pair of scanning mirrors to scan the scene.

The modified hardware setup is used to capture an example scene in both scanned and diffused modes, and the resulting transients are used to refine an initial depth estimate from the Kinect RGB image. We illustrate the results of this procedure in Fig. 5. The reconstructions from the scanned and diffused measurements are similar in reconstruction quality and also show similar quantitative

improvement in terms or error over the initial depth estimate. The unnormalized photon counts are also shown in Fig. 5, and we note that the counts show similar trends. The number of recorded photons in these experiments is shown in Table 1. In both cases, the rate of detected photons is far less (<5%) than the number of emitted laser pulses, and so we conclude that the measurements are captured in the low-flux regime where pileup effects are negligible. We attribute most of the differences between the scanned and diffused transients to the baseline between the positions of the diffused and scanned SPADs (see Fig. ??).

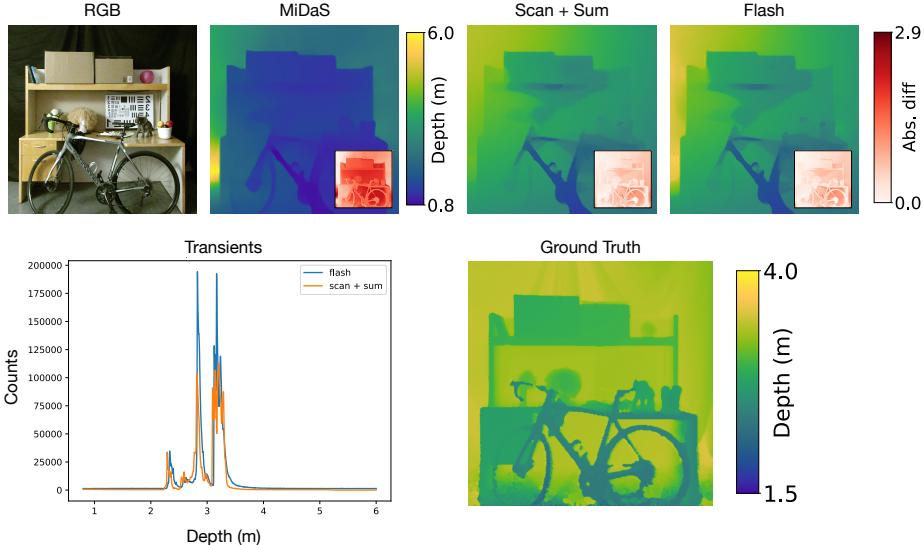


Fig. 5: Comparing Scan + Sum and Diffused SPAD. The transients shown are equal-time (and thus equal energy), and display qualitative similarity in the absence of pileup.

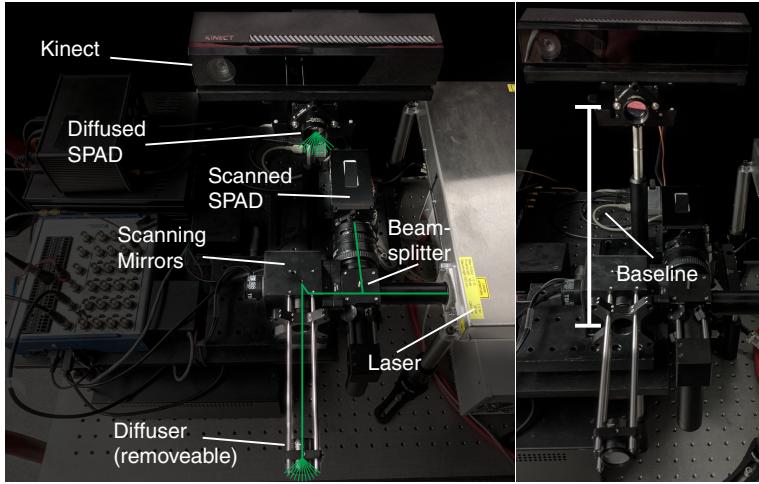


Fig. 6: Modified hardware setup. The setup is used to compare scanned and diffused measurements and employs two SPAD detectors and two laser configurations. In the first configuration, the scene is illuminated by sending the laser light through a holographic diffuser and a lensless SPAD integrates light from the entire scene. In the second, the SPAD is aligned with the optical path of the laser and the scene is scanned using a pair of scanning mirrors. The baseline between the two SPADs (right) results in some observed differences in the recorded transients.

Experiment	Detected Photons	Laser Pulses	Detection Rate
Scanned	1.4×10^7	6×10^8	2.3%
Diffused	2.4×10^7	6×10^8	4.0%

Table 1: **Recorded photons for diffused vs. scanned scene.** In each capture mode, scanned or diffused, the number of detected photons does not exceed 5% of the number of emitted laser pulses, placing the capture within the low-flux regime where pileup effects are negligible.

6 Additional results on NYU Depth v2

Figures 7–15 show additional results for our method on the NYU Depth v2 dataset when the depth estimate is initialized with the DenseDepth [1] (Figures 7–9), DORN [2] (Figures 10–12) and MiDaS [3] (Figures 13 – 15) monocular depth estimators.

We compare the output of the network z_0 , the median-rescaled network output (where the depth map z_0 is scaled pixel-wise by a scalar $\frac{\text{median}(z_{GT})}{\text{median}(z_0)}$, z_{GT}

270 being the ground truth depth map), the network output matched to the ground
271 truth depth histogram, and the output of our histogram matching method under
272 a signal-to-background ratio (SBR) of 100. We use the luminance of the RGB
273 image as our reflectance map for both SPAD simulation and histogram match-
274 ing. We show absolute difference maps and also give the root-mean-square error
275 (RMSE) for each example.

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

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

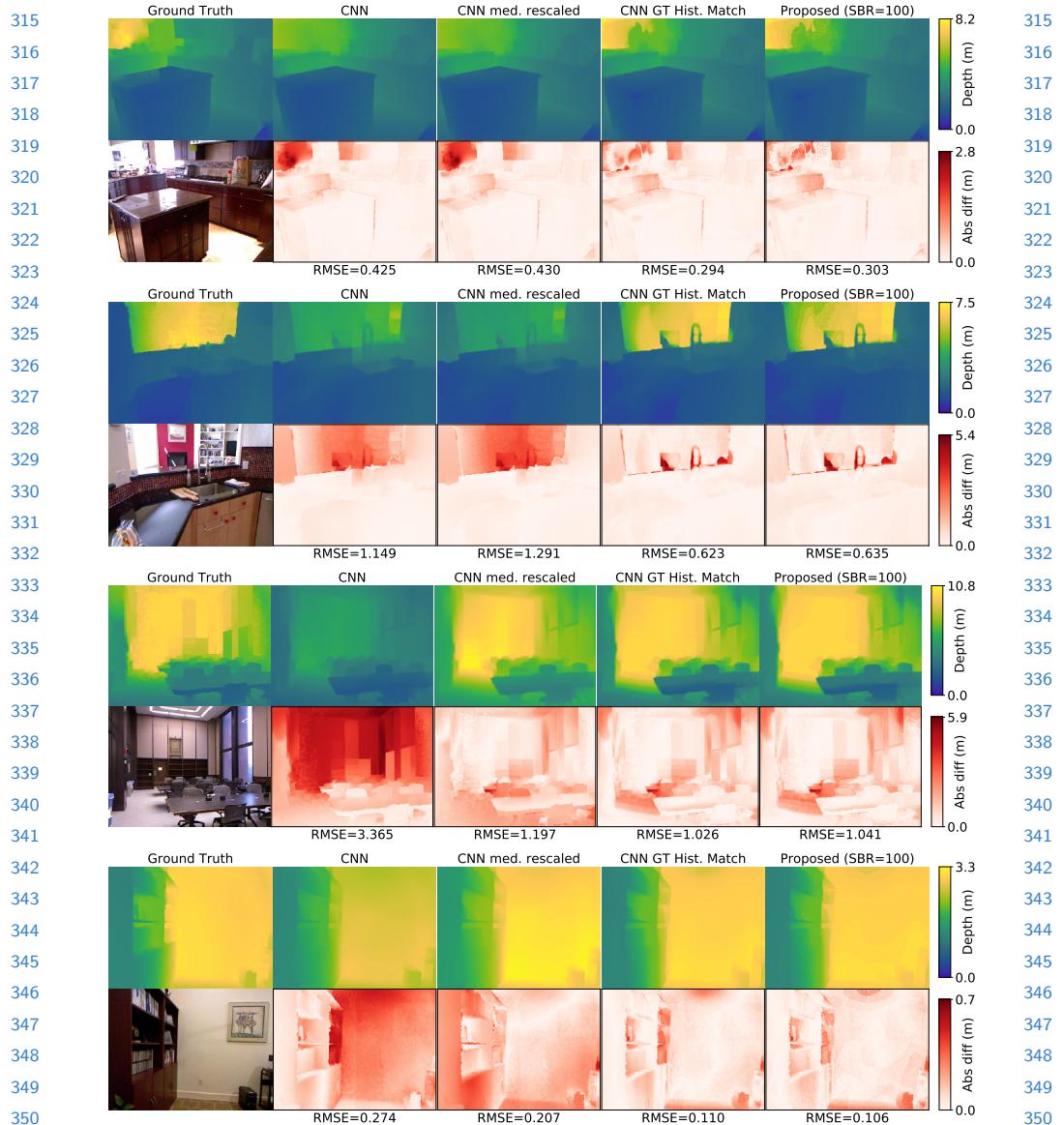


Fig. 7: Results with DenseDepth as the monocular depth estimator. Our method is able to scale and shift the depth maps to mitigate gross errors in depth scaling.

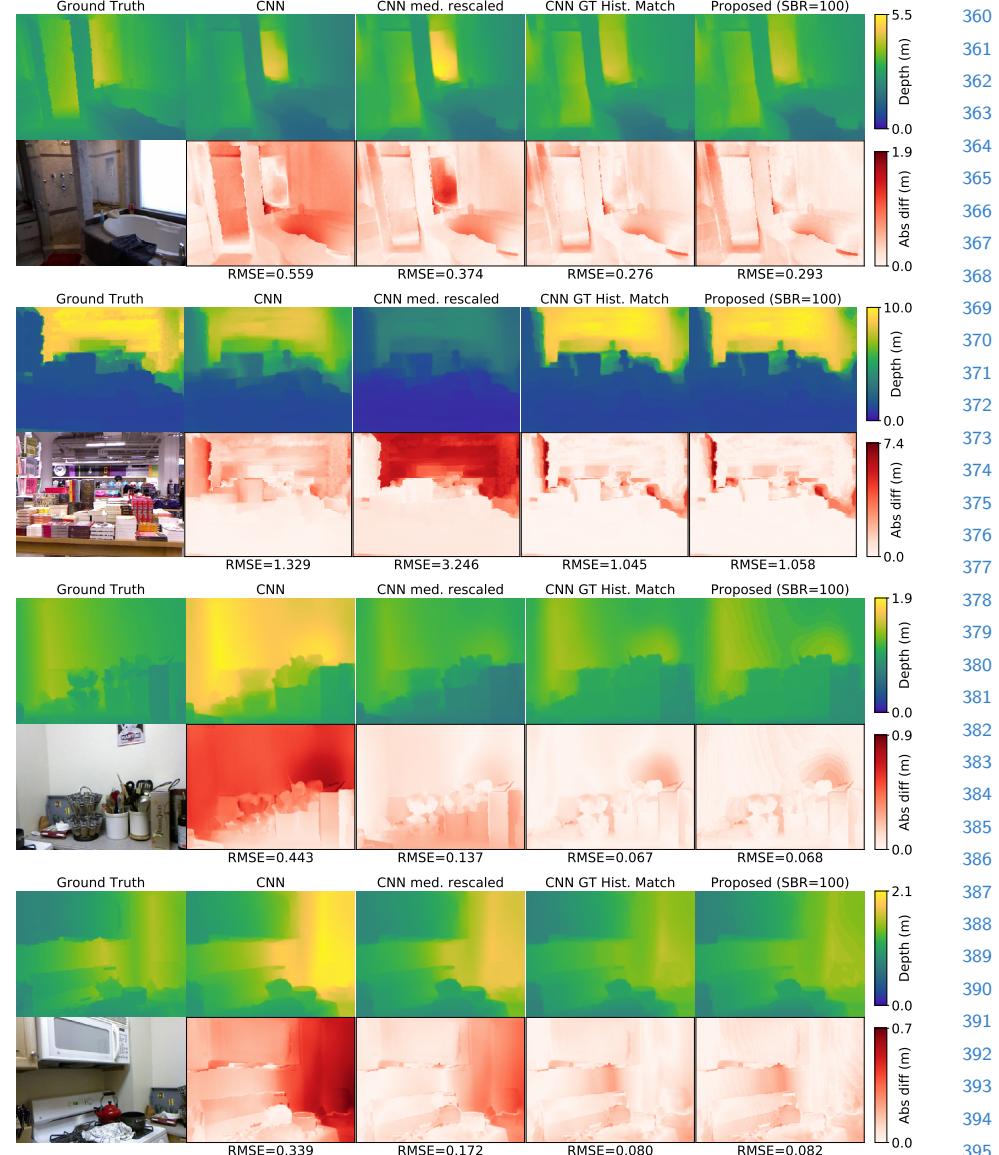


Fig. 8: Results with DenseDepth as the monocular depth estimator. Our method is able to scale and shift the depth maps to mitigate gross errors in depth scaling.

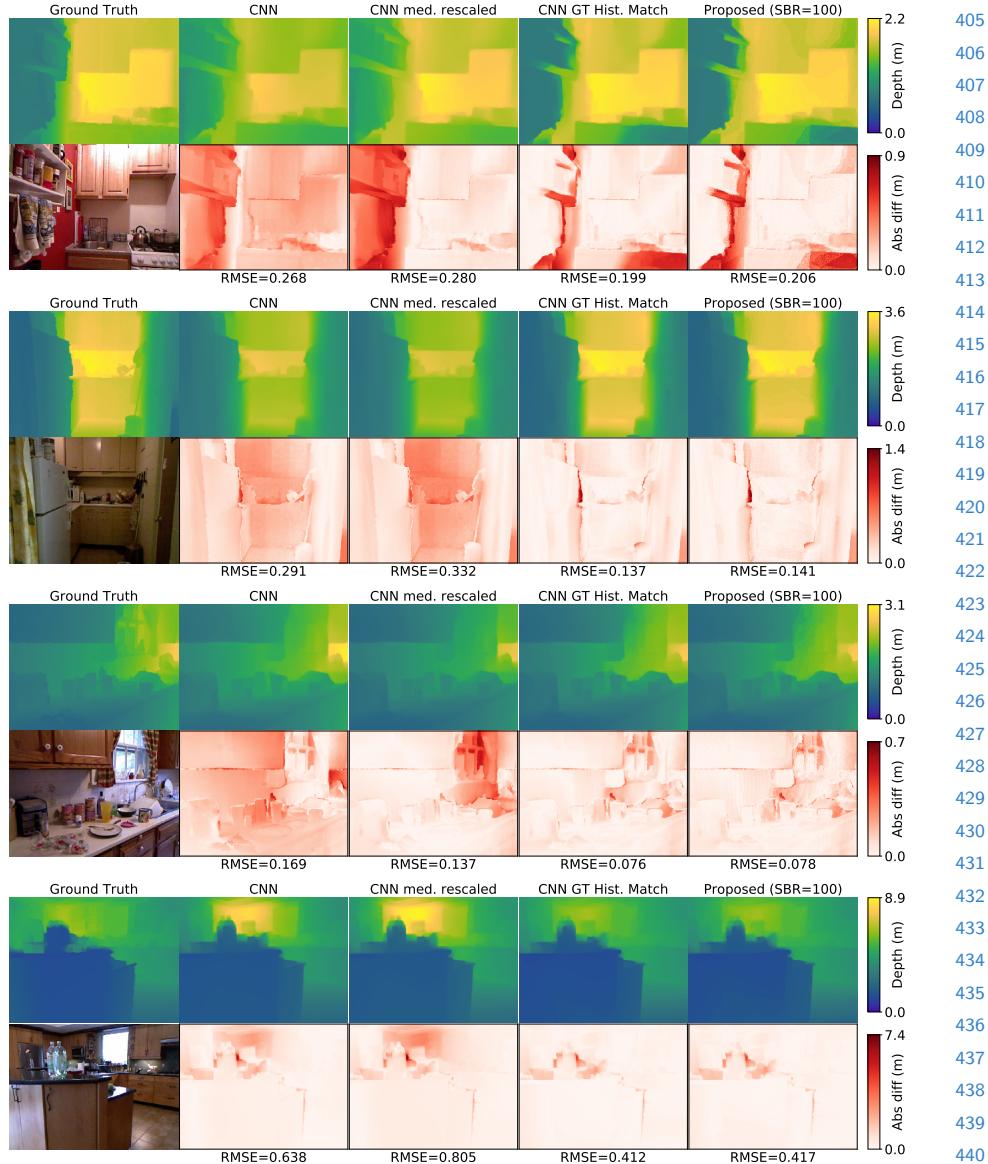


Fig. 9: Results with DenseDepth as the monocular depth estimator. Our method is able to scale and shift the depth maps to mitigate gross errors in depth scaling.

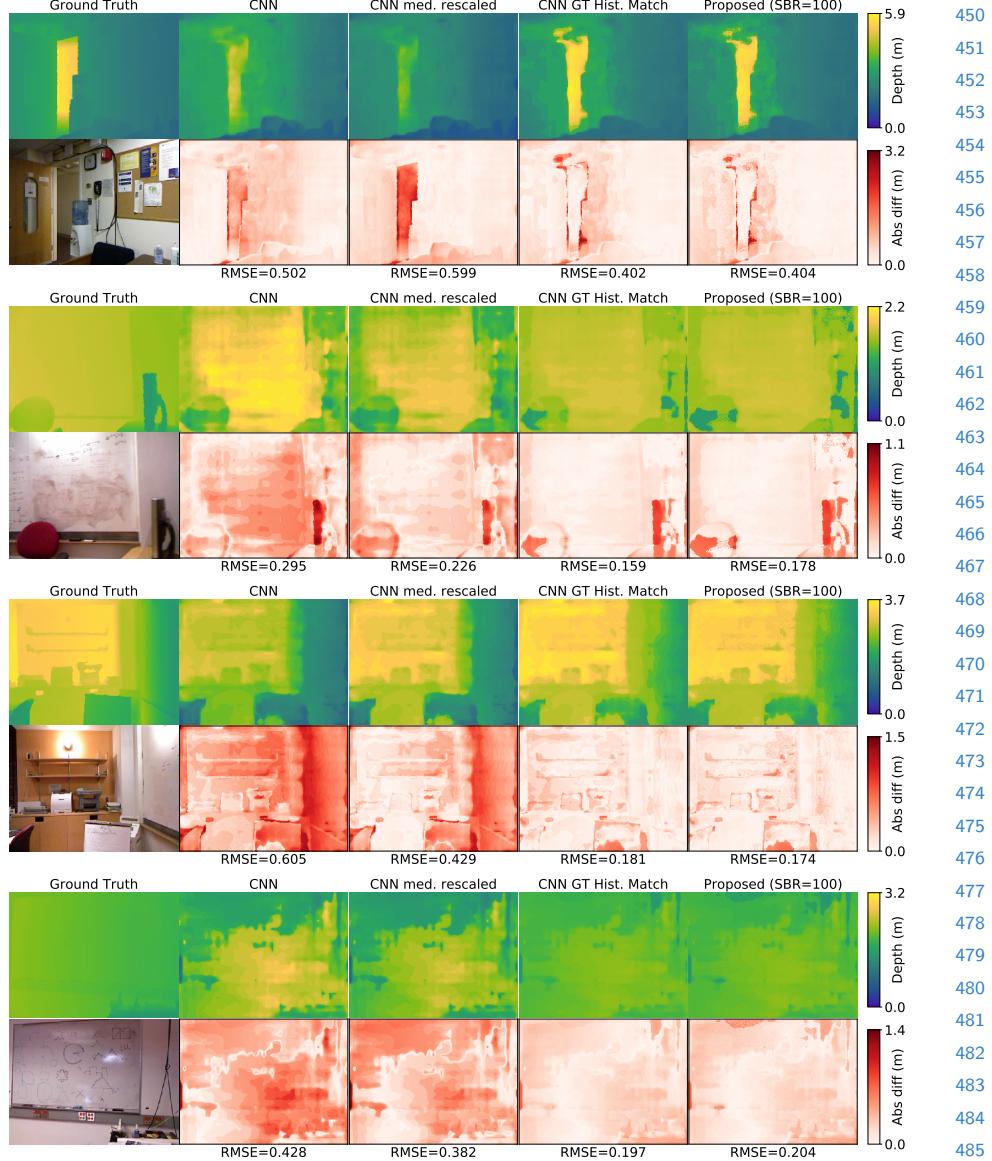


Fig. 10: Results with DORN as the monocular depth estimator. Our method is able to scale and shift the depth maps to mitigate gross errors in depth scaling.

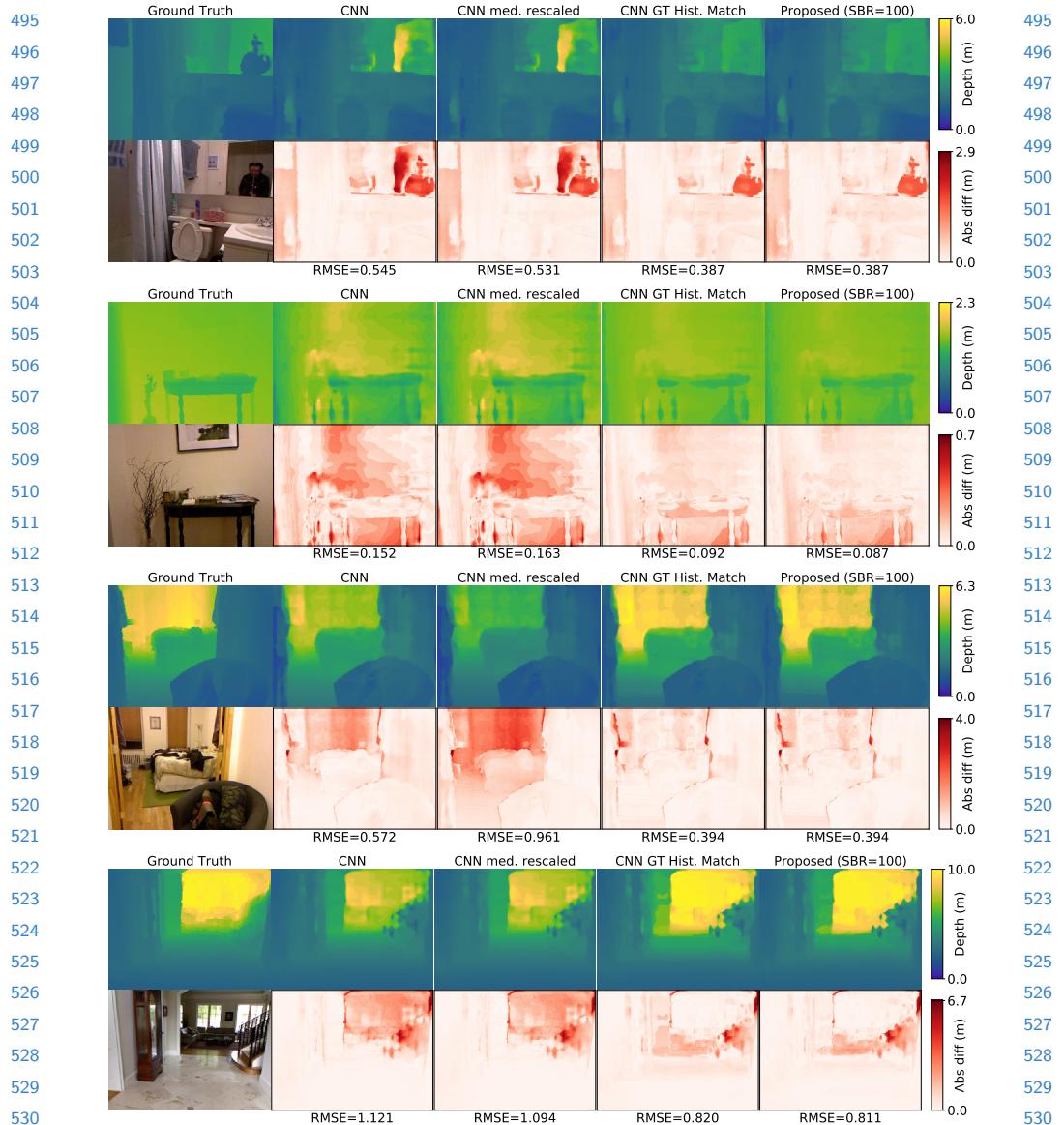
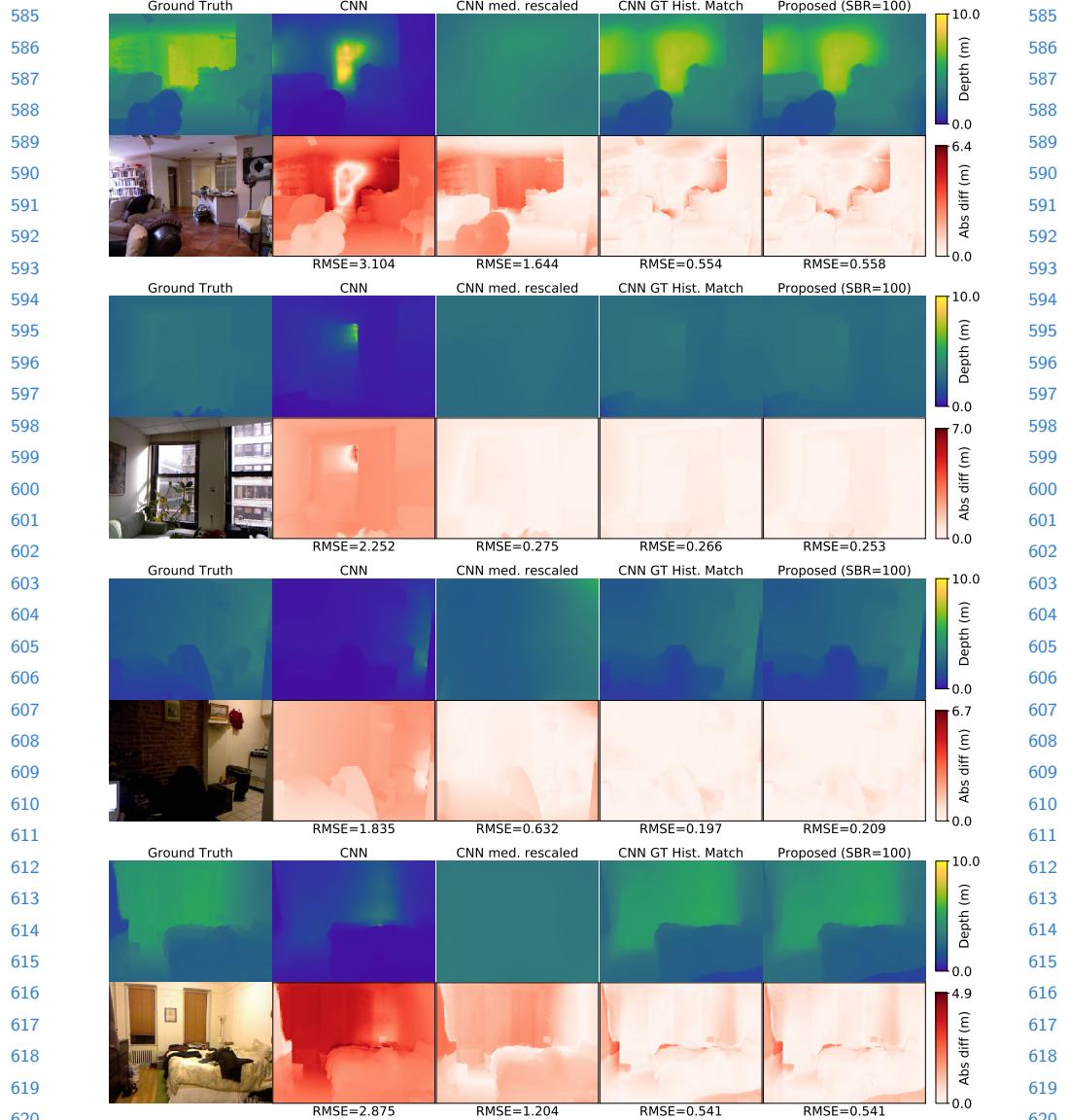


Fig. 11: Results with DORN as the monocular depth estimator. Our method is able to scale and shift the depth maps to mitigate gross errors in depth scaling.



577 Fig. 12: Results with DORN as the monocular depth estimator. Our method is able to
 578 scale and shift the depth maps to mitigate gross errors in depth scaling.



621 Fig. 13: Results with MiDaS as the monocular depth estimator. Our method is able to
 622 scale and shift the depth maps to mitigate gross errors in depth scaling.
 623
 624
 625
 626
 627
 628
 629

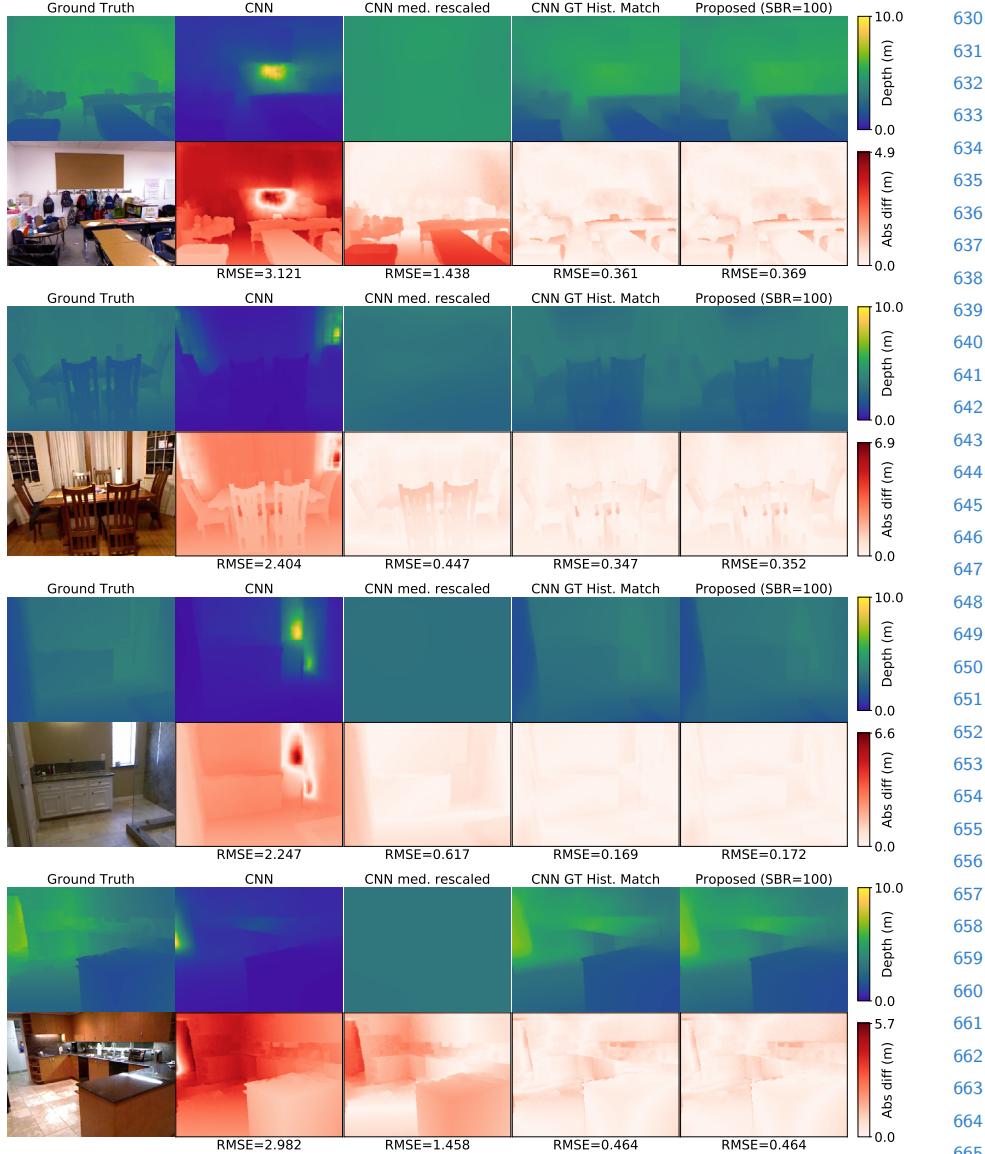


Fig. 14: Results with MiDaS as the monocular depth estimator. Our method is able to scale and shift the depth maps to mitigate gross errors in depth scaling.

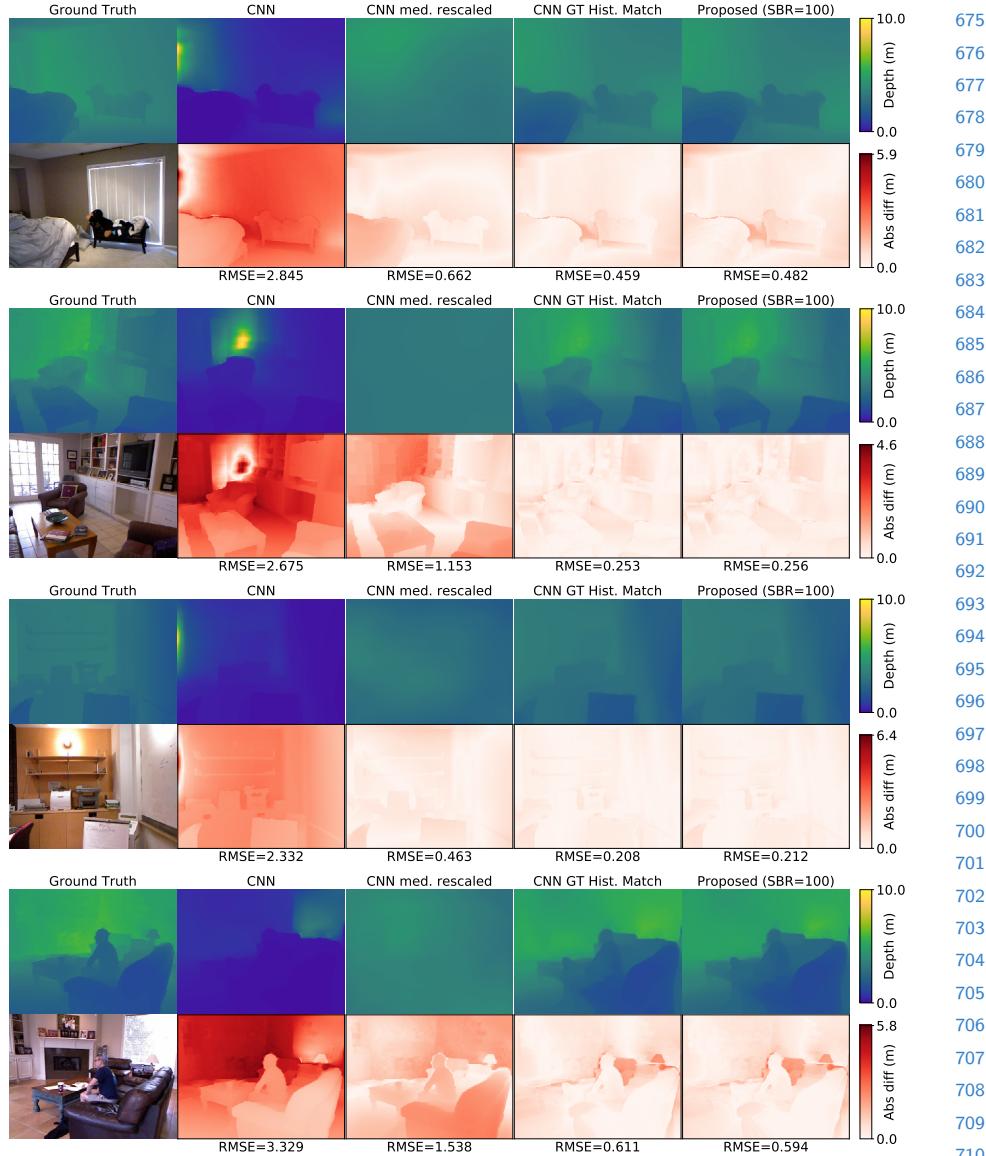


Fig. 15: Results with MiDaS as the monocular depth estimator. Our method is able to scale and shift the depth maps to mitigate gross errors in depth scaling.

7 Additional results for hardware prototype

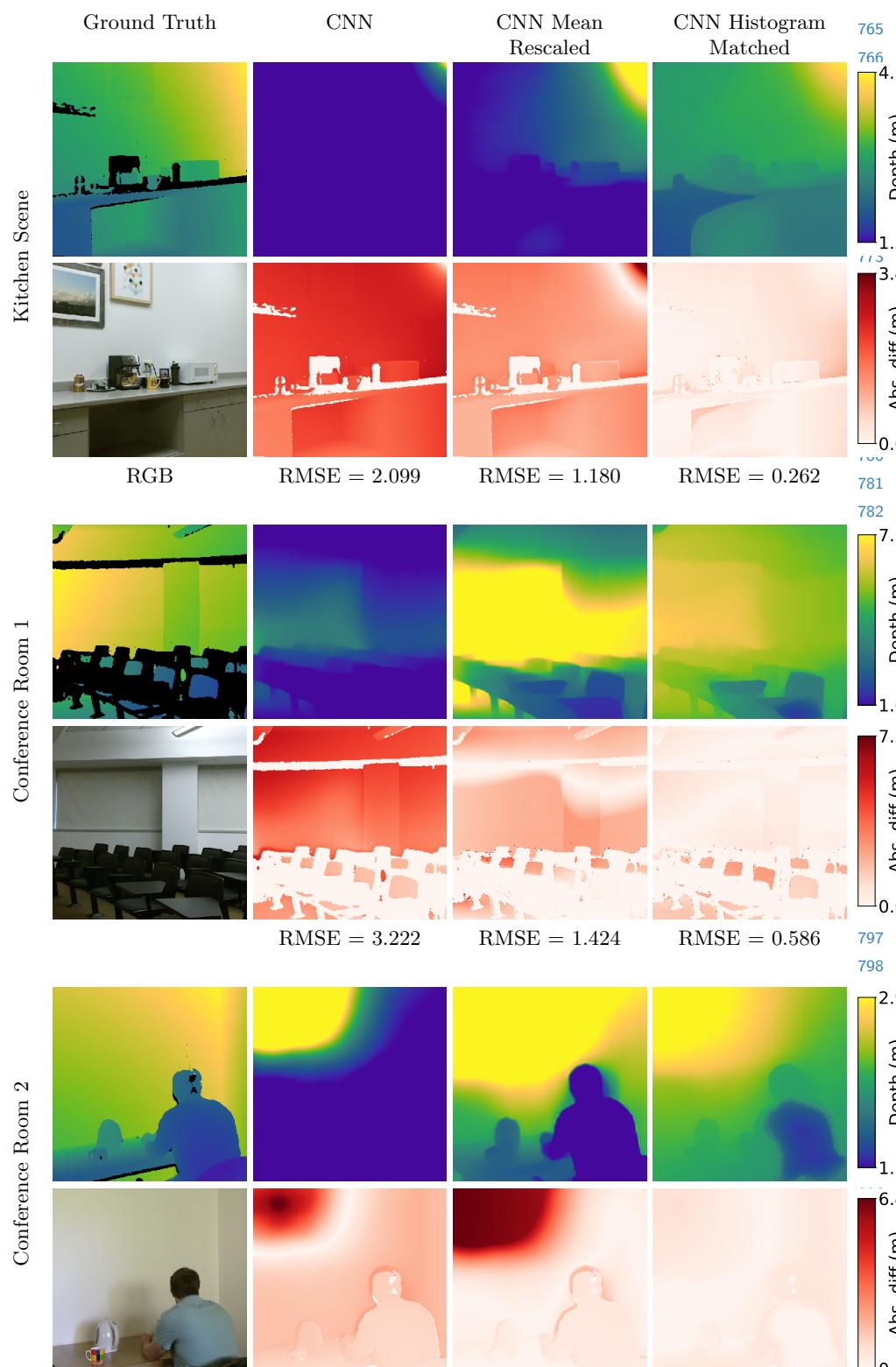
Figures 16–24 show all the captured results when the depth estimate is initialized with the MiDaS [3] (Figures 16–18), DenseDepth (Figures 19–21), and DORN (Figures 22–24). We compare the output of the network z_0 , the mean-rescaled

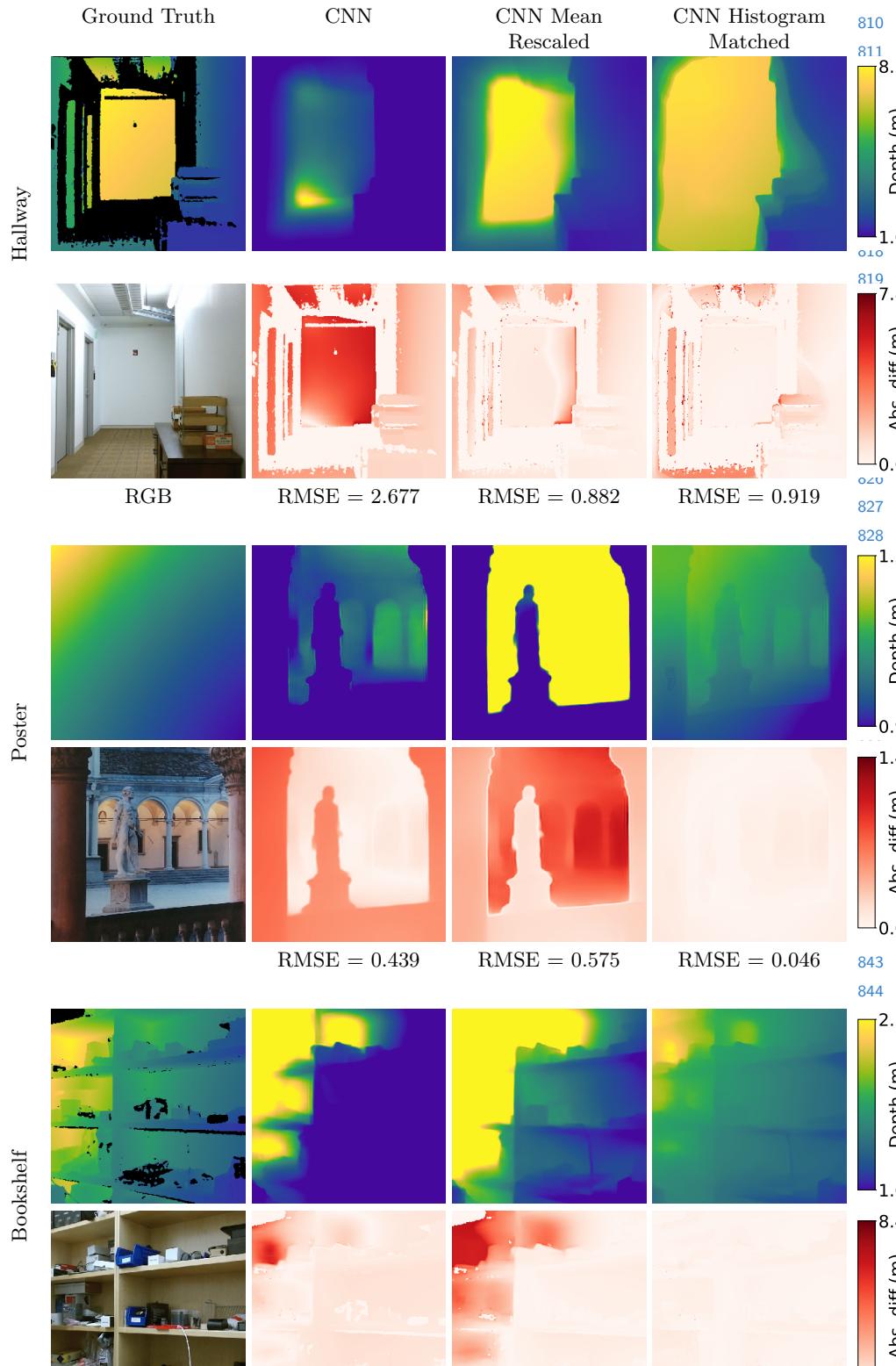
720 network output where the depth map z_0 has been scaled pixel-wise by the scalar
721 $\frac{\text{median}(h_{\text{target}})}{\text{median}(z_0)}$ (h_{target} is the processed SPAD transient), and the output of our
722 method. As our laser is red, we use the R channel of the RGB image as our
723 reflectance map. We show absolute difference maps and also give the root-mean-
724 square-error (RMSE) for each example.

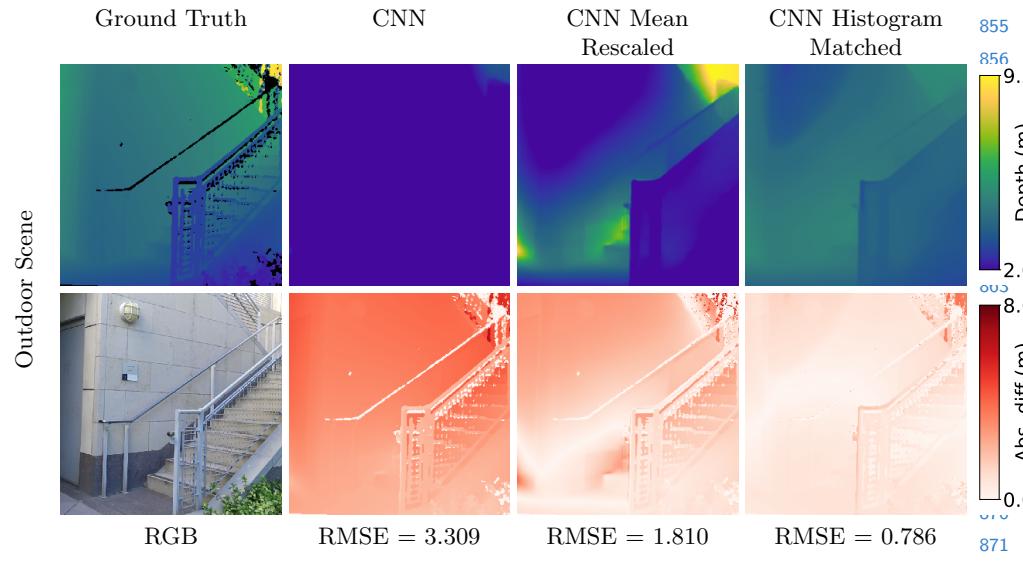
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759

760 Black pixels in the ground truth depth correspond to locations where our
761 scanner was unable to produce an accurate depth estimate (this can occur for a
762 variety of reasons including dark albedo and surface specularity). These pixels
763 are masked off and not used in the RMSE calculation, and appear as an absolute
764 difference of 0 in the difference maps.

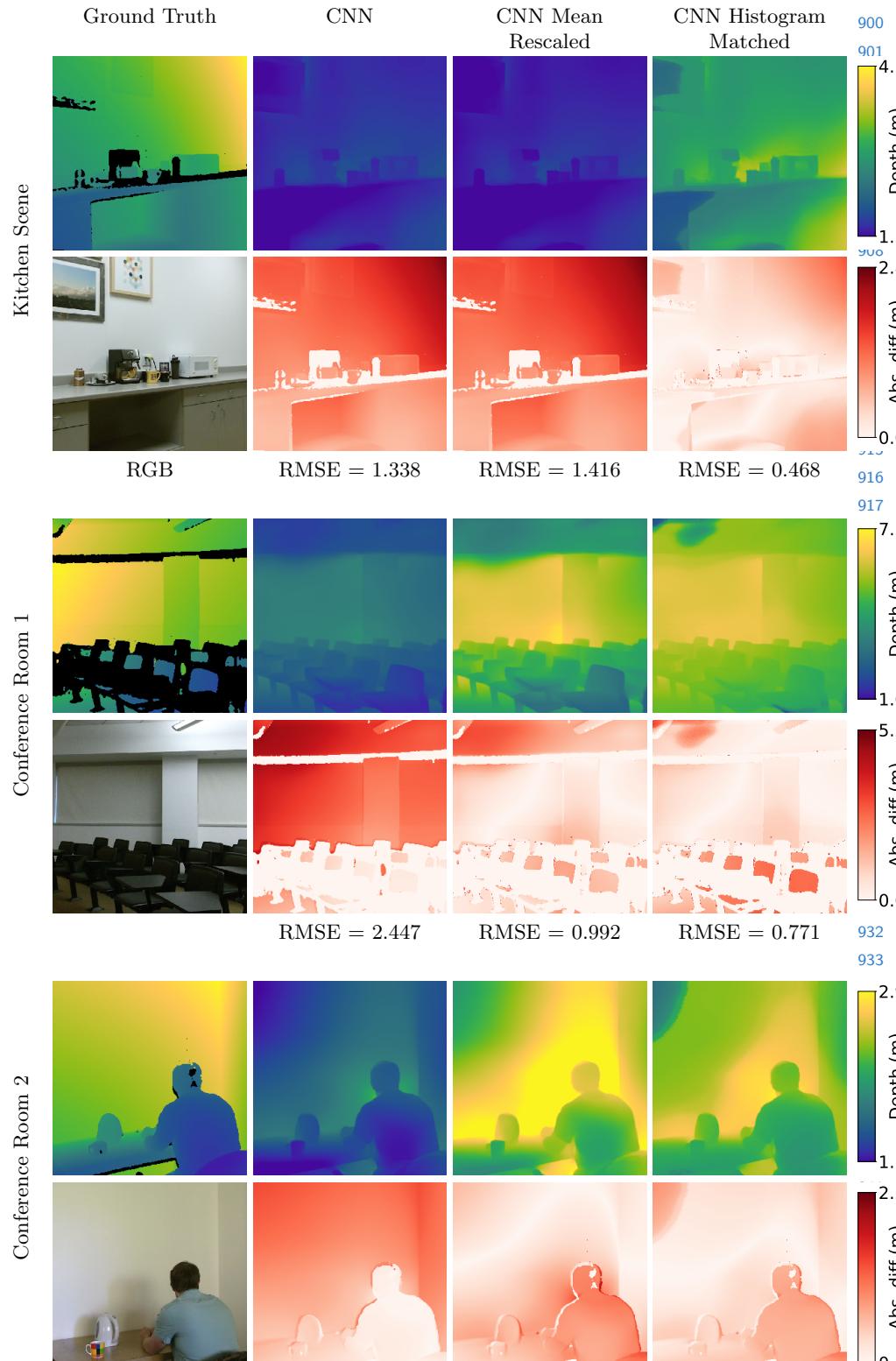
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764

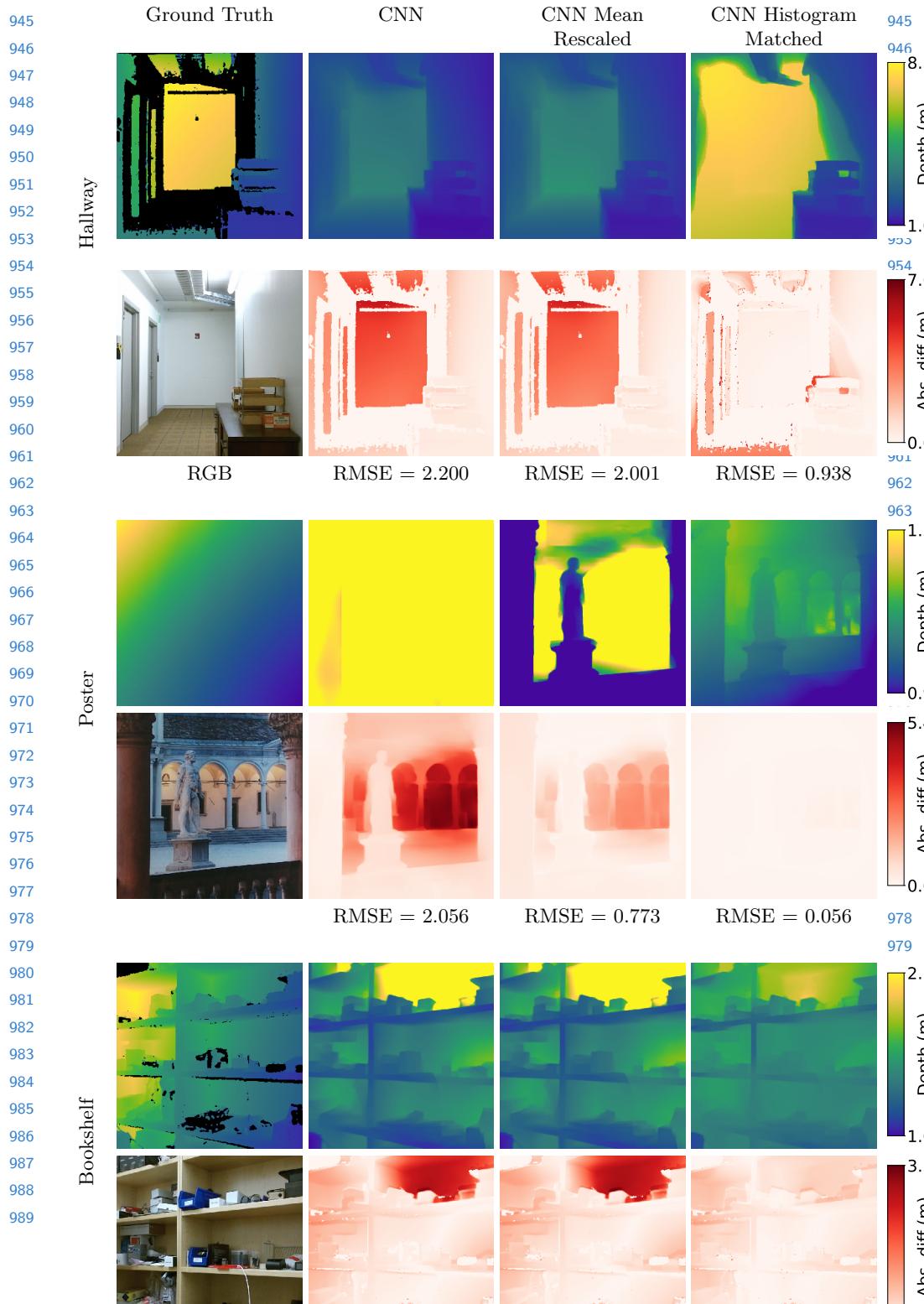






872 Fig. 18: Captured results initialized using the MiDaS CNN on an outdoor scene. Second
 873 row shows absolute difference between above estimates and ground truth. MiDaS
 874 does not output metric depth, so the CNN depth map is scaled to be in the range
 875 (0.494, 11.094) by default. However, MiDaS does produce accurate ordinal depth,
 876 leading to stronger performance of our histogram matching compared to other methods.





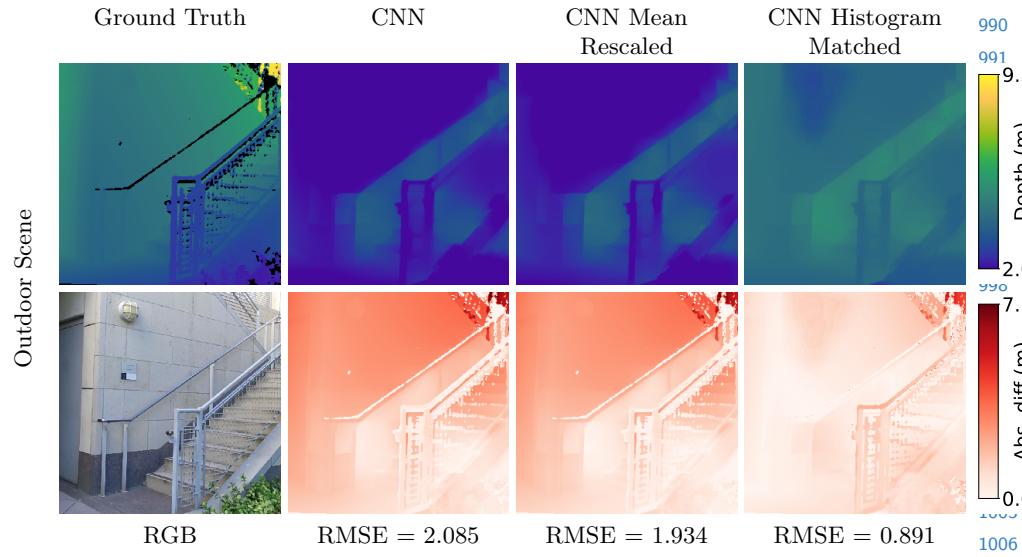
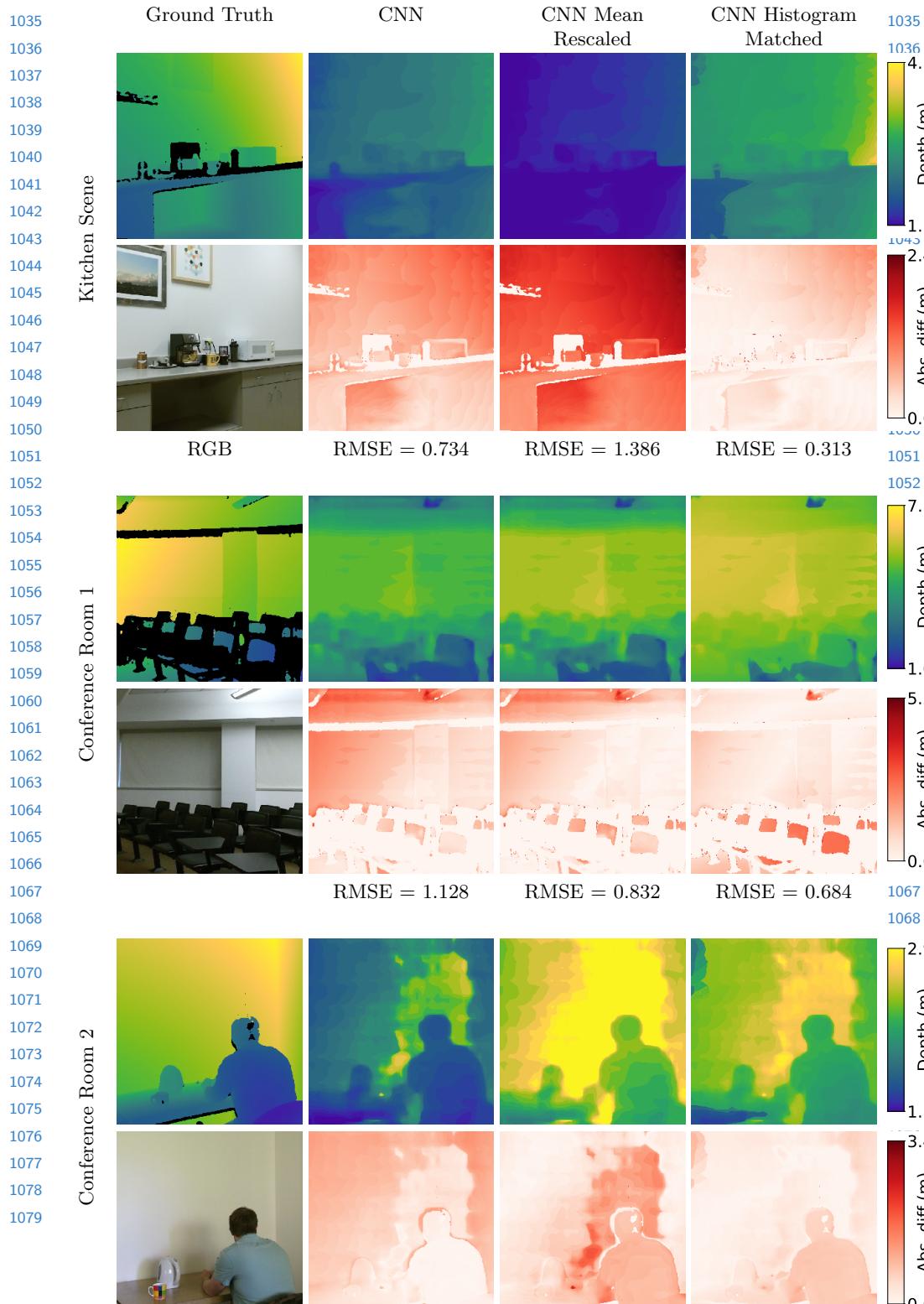
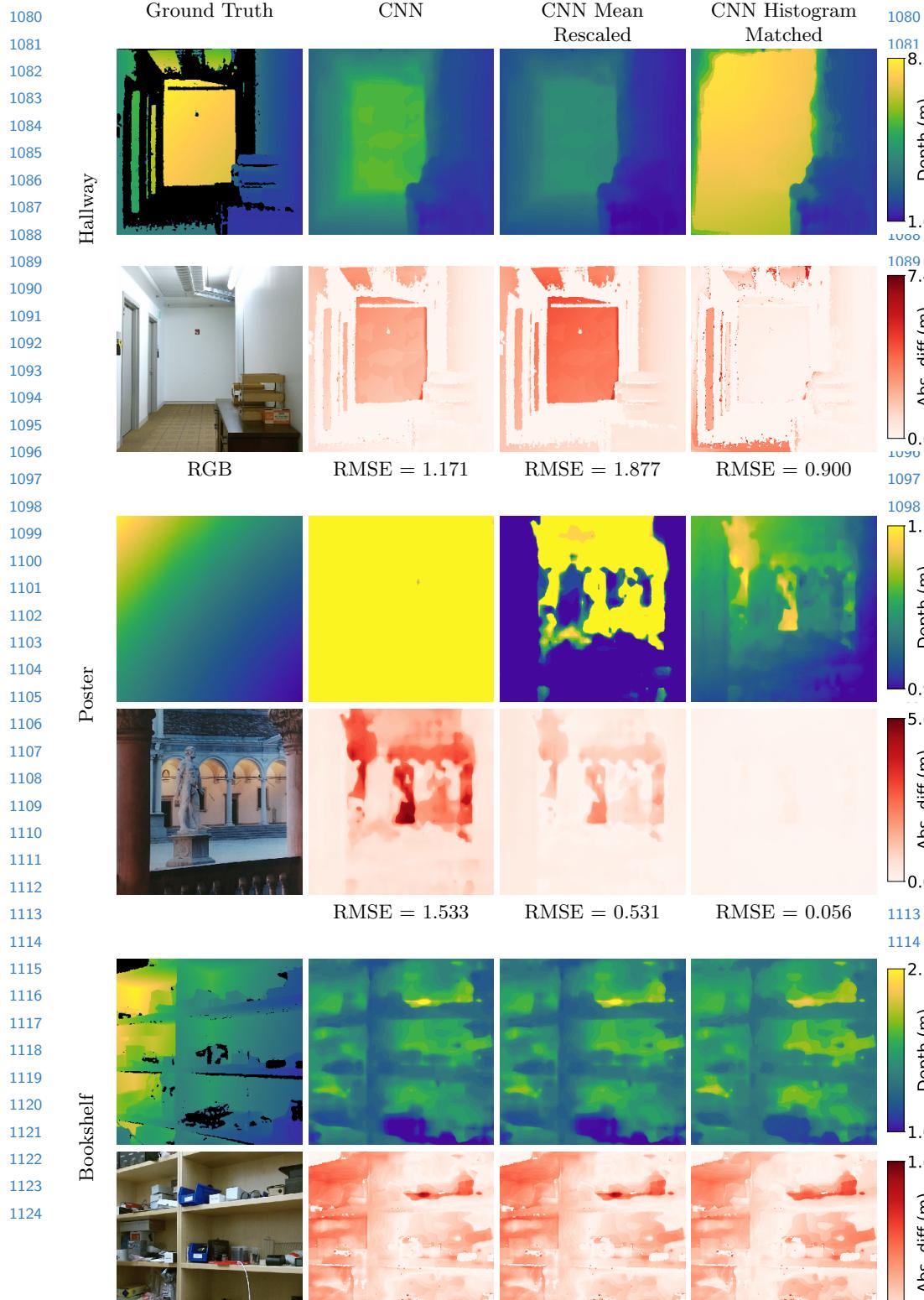


Fig. 21: Captured results initialized using the DenseDepth CNN on an outdoor scene. Second row shows absolute difference between above estimates and ground truth.





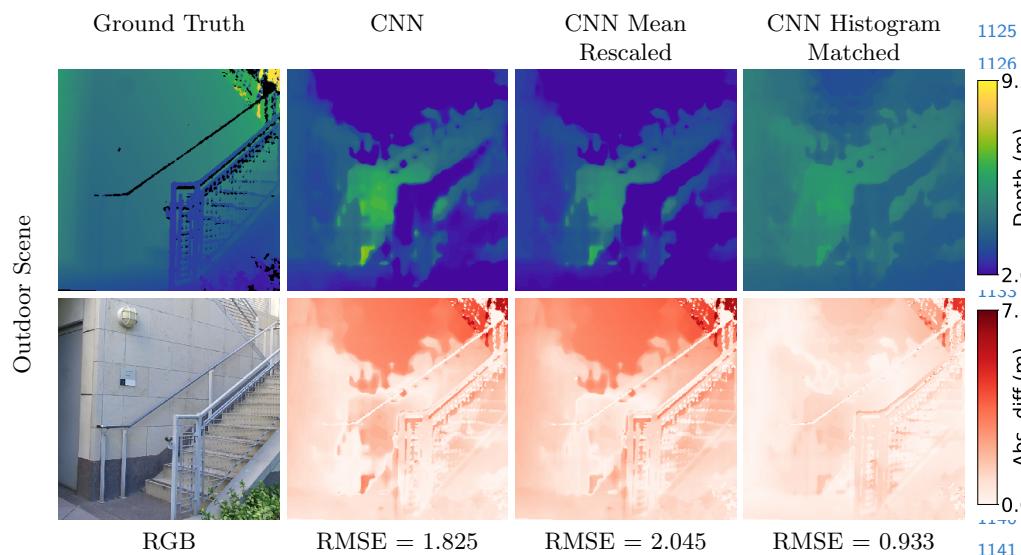


Fig. 24: Captured results initialized using the DORN CNN on an outdoor scene. Second row shows absolute difference between above estimates and ground truth.

1170 References

- 1172 1. Alhashim, I., Wonka, P.: High quality monocular depth estimation via transfer learn-
1173 ing. arXiv:1812.11941v2 (2018)
- 1174 2. Fu, H., Gong, M., Wang, C., Batmanghelich, K., Tao, D.: Deep ordinal regression
1175 network for monocular depth estimation. In: Proc. CVPR (2018)
- 1176 3. Lasinger, K., Ranftl, R., Schindler, K., Koltun, V.: Towards robust monocular depth
1177 estimation: Mixing datasets for zero-shot cross-dataset transfer. arXiv:1907.01341
(2019)