

# Supplemental Material: Deep Novel View Synthesis from Unstructured Input

Anonymous ECCV submission

Paper ID 3314

## 1 Introduction

In this supplementary document, we present additional experimental results in Section 2. This includes quantitative and qualitative results for our own recordings, as well as additional qualitative results for the Tanks and Temples sequences and for DTU view interpolation and extrapolation. In Section 3, we provide details on our recurrent mapping and blending network.

## 2 Additional Results

For an additional evaluation of novel view synthesis from unstructured input, we recorded new challenging sequences. To get unstructured input for the source images, we first recorded for each scene a video sequence with a duration of one too two minutes. In addition, we recorded for each scene up to two additional videos that serve as ground-truth targets for novel view synthesis and are not included in the source input.

We summarize the quantitative results on those recordings in Table 1. Note that all results are slightly worse than the Tanks and Temples results presented in the main paper. This is an indication for the challenging setup provided in our recordings. However, we also see that the results of our method are clearly better than those of the state-of-the-art methods. In Figures 1, 2, 3, 4, 5, and 6 we show qualitative results of all methods on a subset of our own recordings.

We also present additional qualitative results for our evaluations on the Tanks and Temples sequences in Figures 7, 8, 9, 10, 11, 12, 13, and 14 and for the DTU view interpolation and extrapolation in Figures 15, 16, 17, 18, 19, 20, 21, 22, and 23.

Table 1: Quantitative results on our own recordings. Mean over all scenes.

	↓LPIPS	↑SSIM	↑PSNR
EVS [2]	14.96	40.74	12.82
LLFF [3]	17.36	30.45	10.52
NPBG [1]	9.85	70.65	18.31
Our	<b>7.60</b>	<b>77.42</b>	<b>19.86</b>

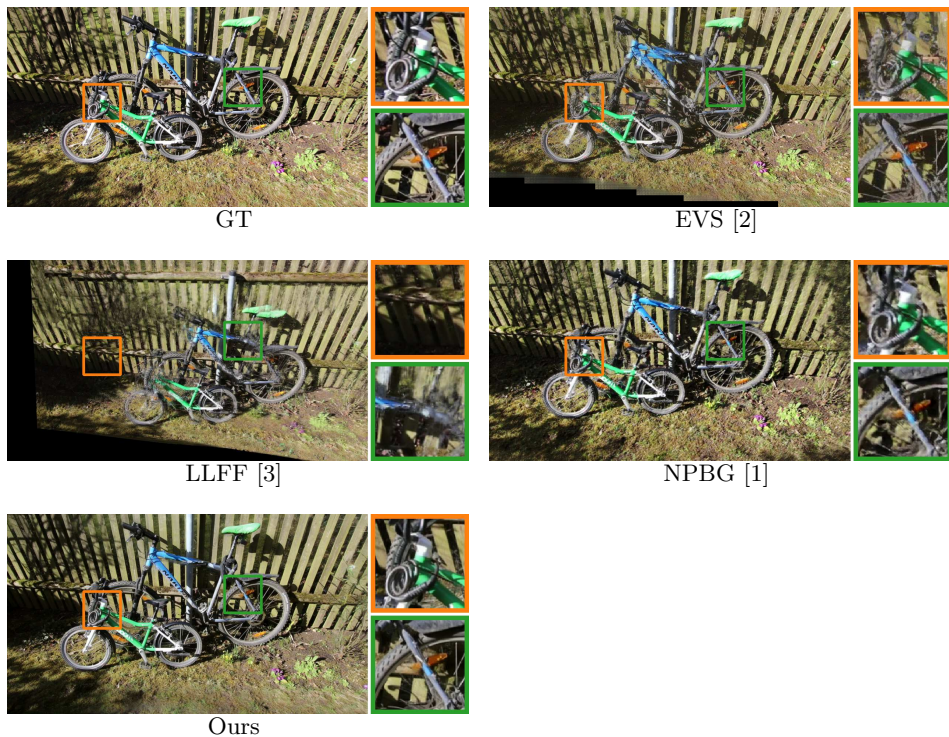


Fig. 1: Own recording *Bike*, view 6

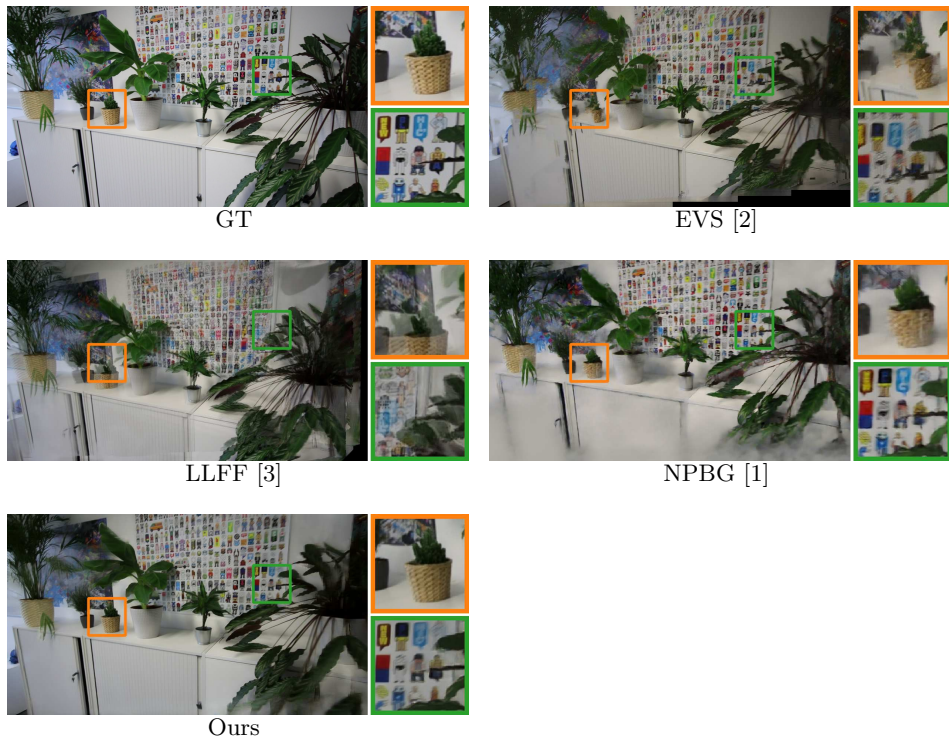


Fig. 2: Own recording *Flowers*, view 127

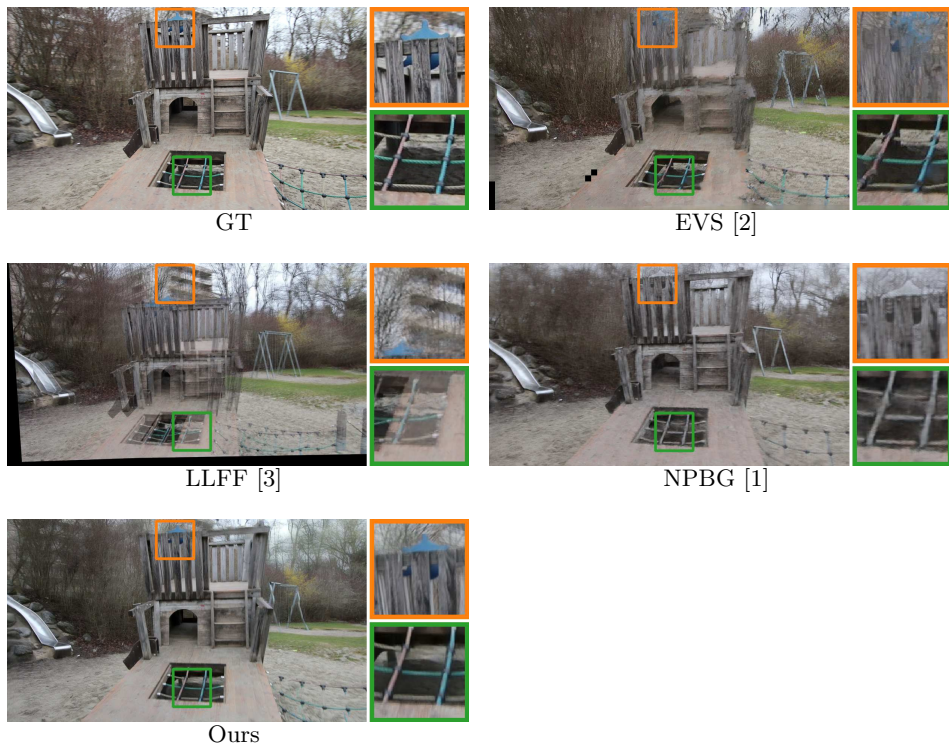


Fig. 3: Own recording *Pirate*, view 139



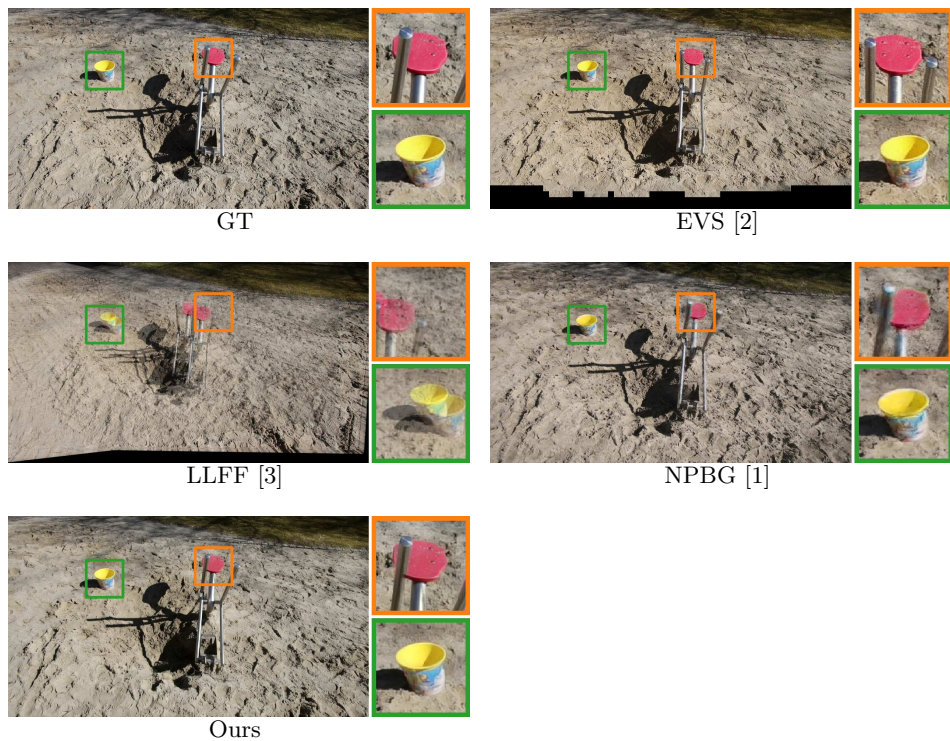


Fig. 4: Own recording *Playground*, view 43

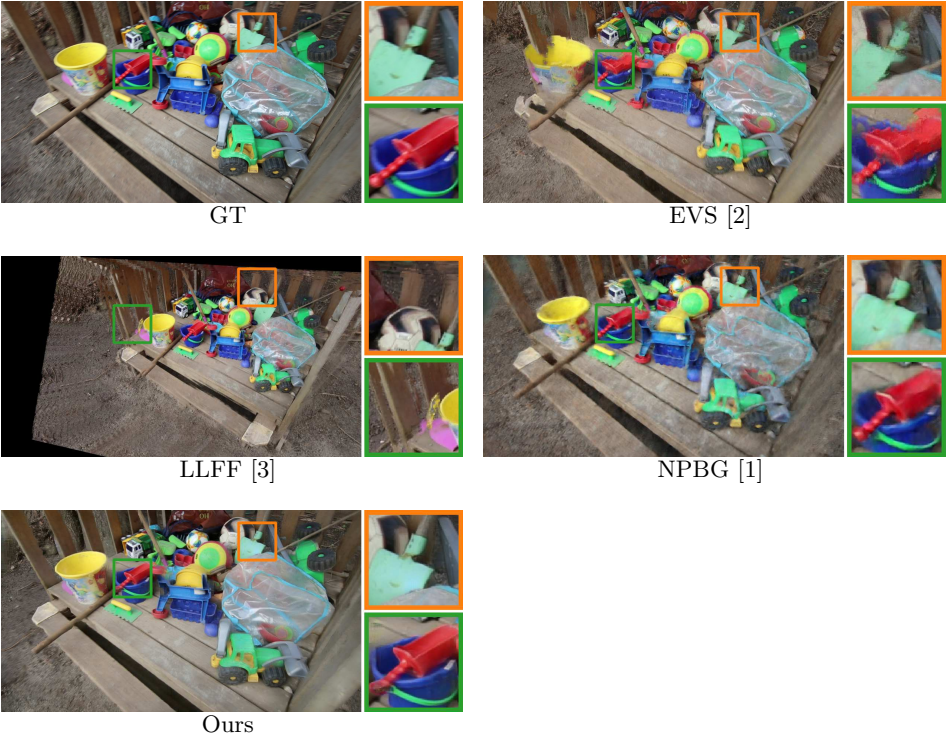


Fig. 5: Own recording *Sandbox*, view 102

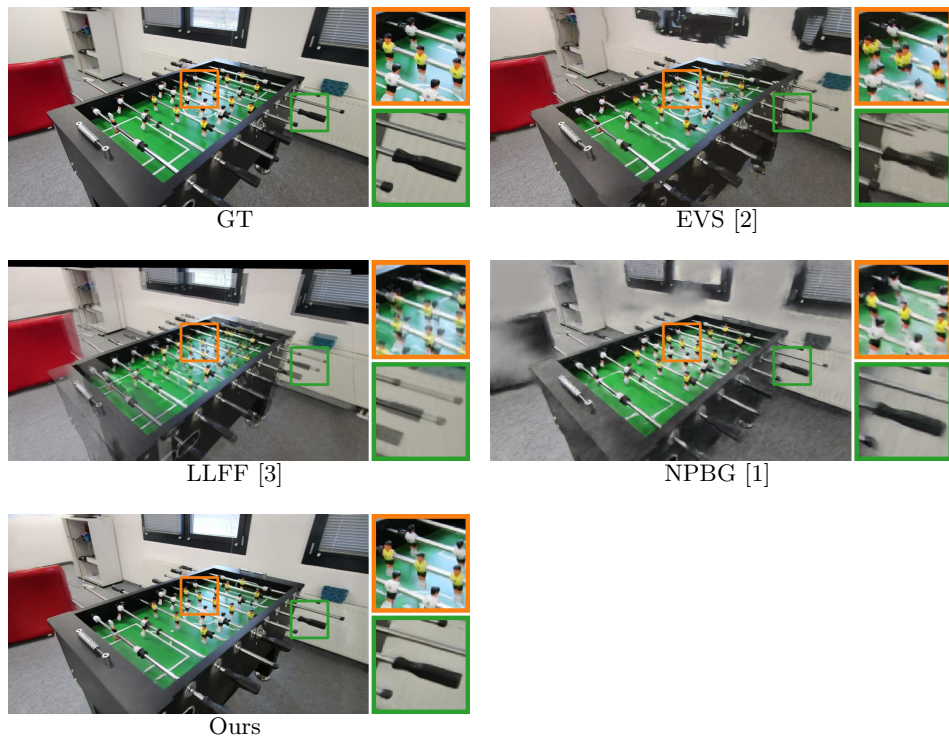


Fig. 6: Own recording *Soccertable*, view 127

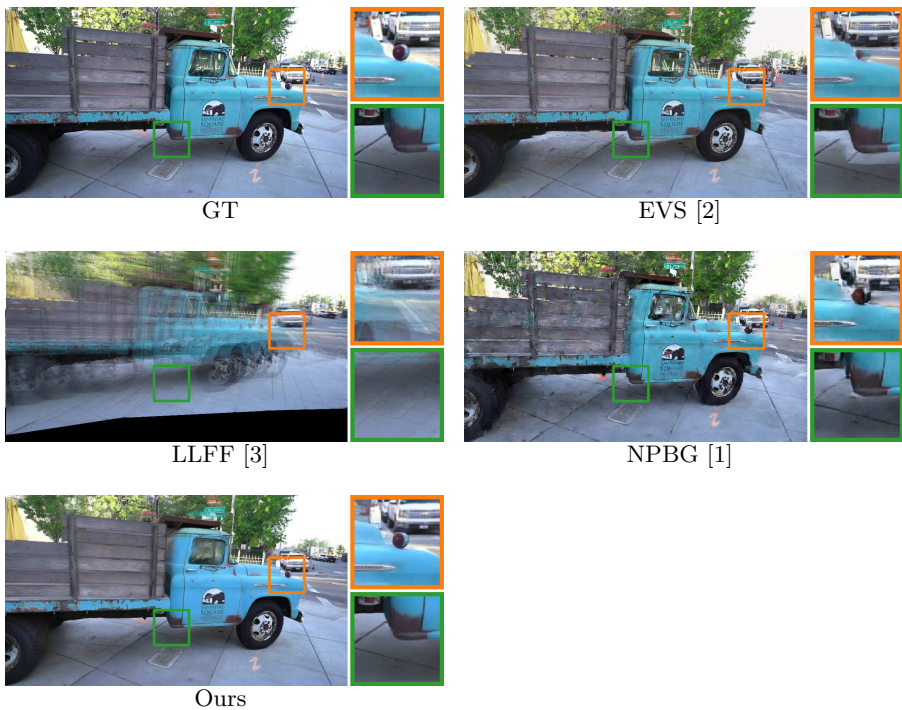


Fig. 7: Tanks and Temples *Truck*, view 0



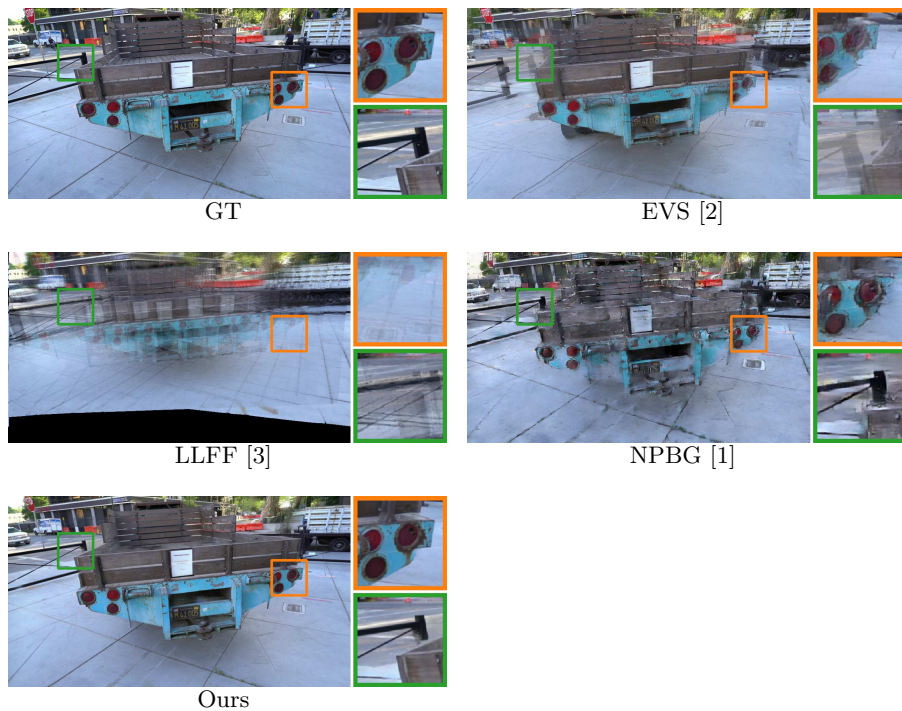


Fig. 8: Tanks and Temples *Truck*, view 16



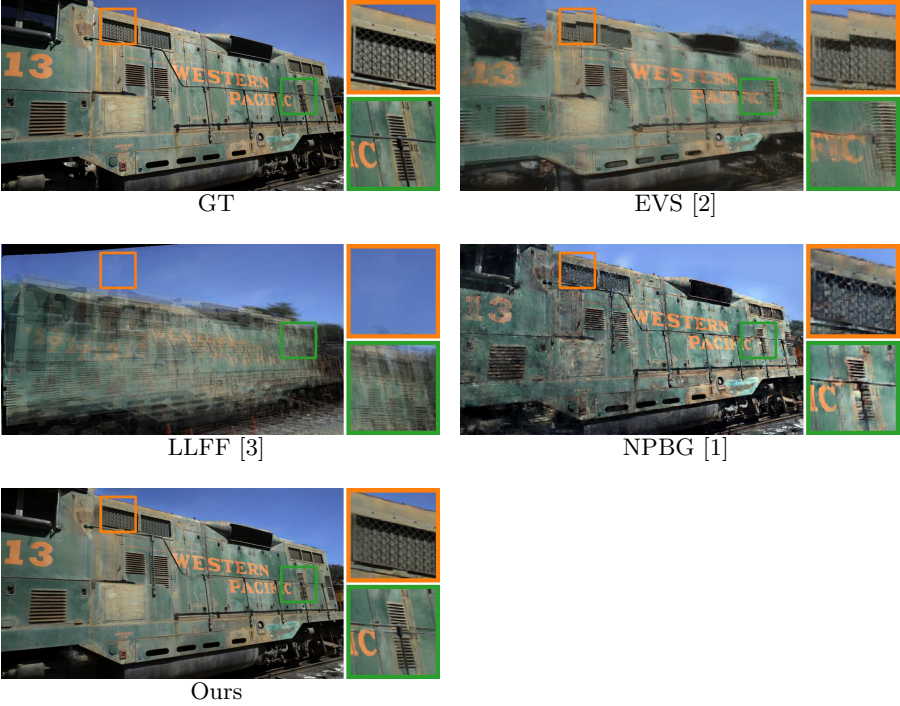


Fig. 9: Tanks and Temples *Train*, view 4

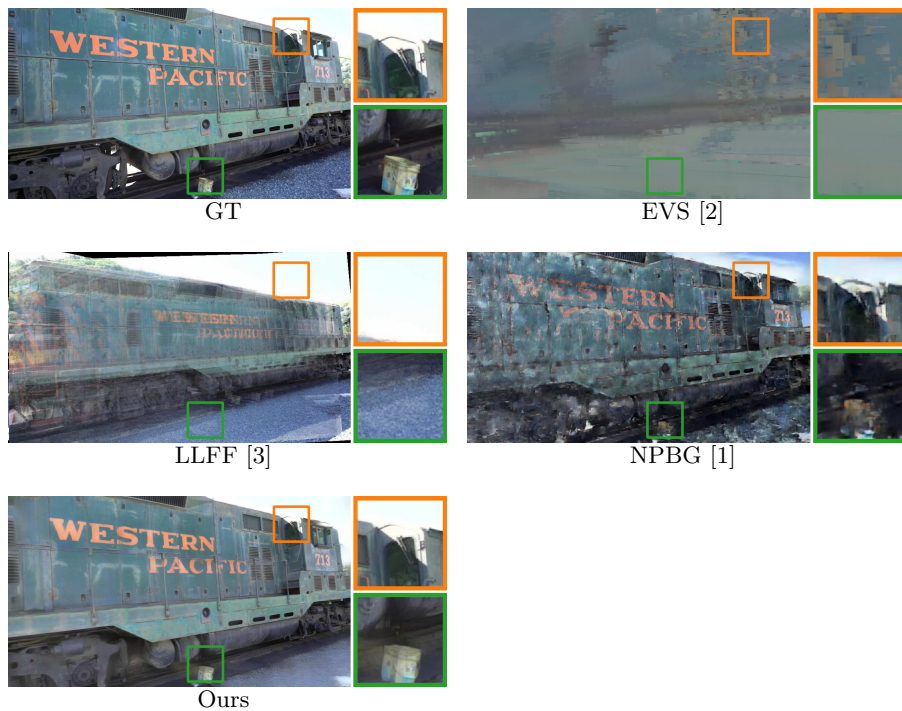


Fig. 10: Tanks and Temples *Train*, view 28

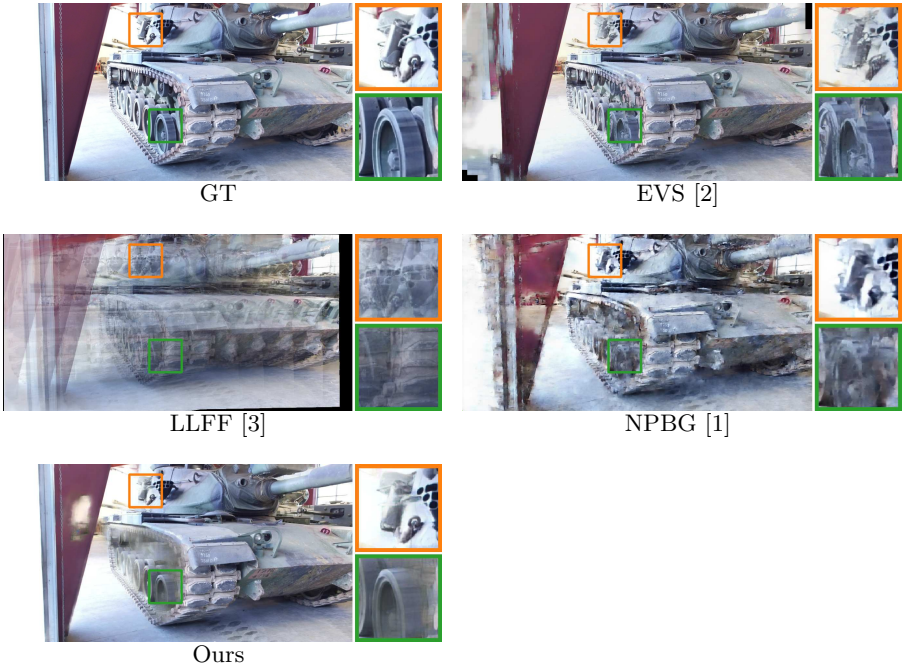


Fig. 11: Tanks and Temples *M60*, view 0

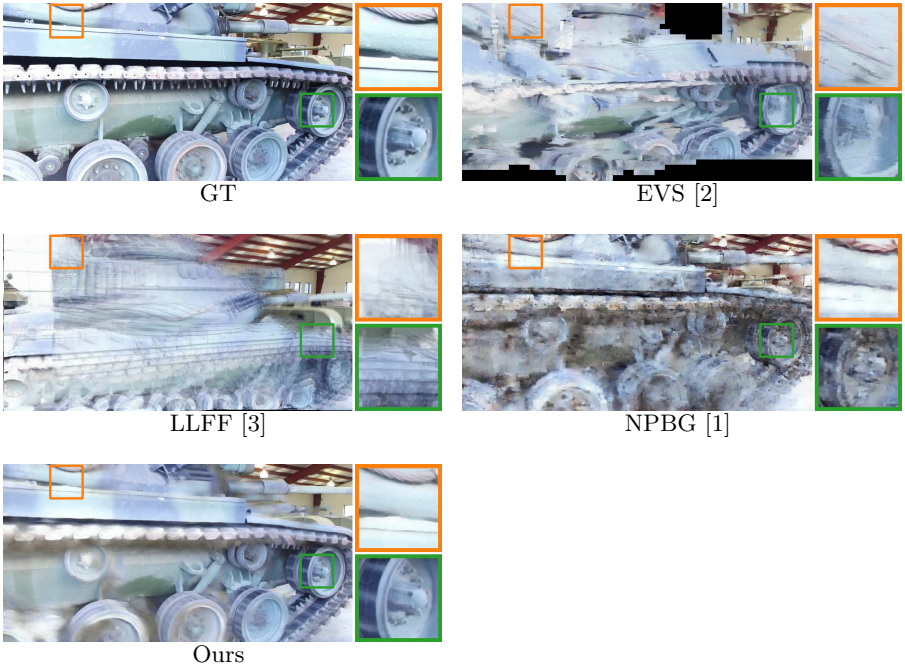


Fig. 12: Tanks and Temples *M60*, view 19

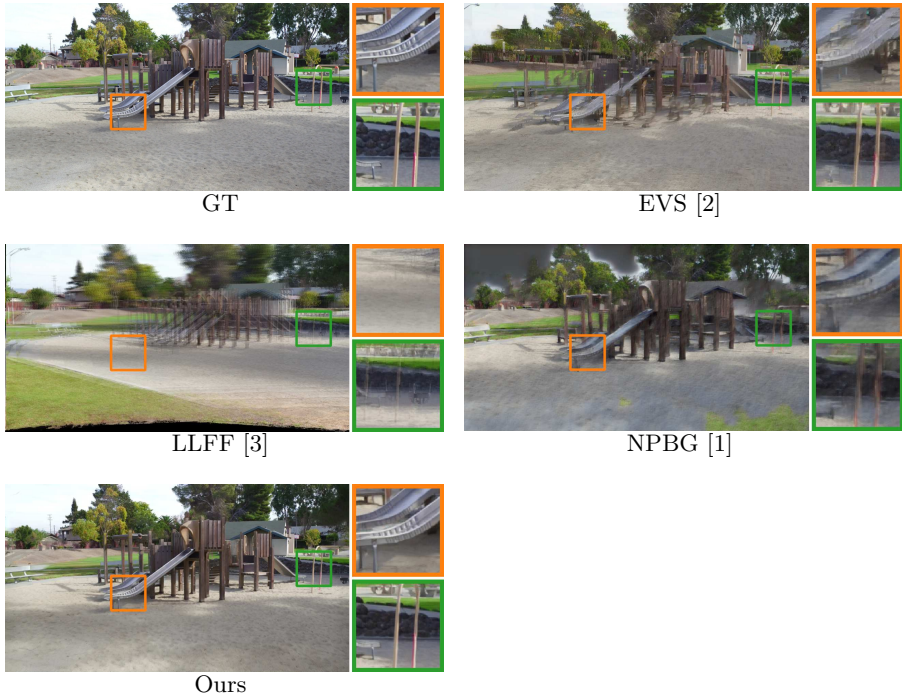


Fig. 13: Tanks and Temples *Playground*, view 10



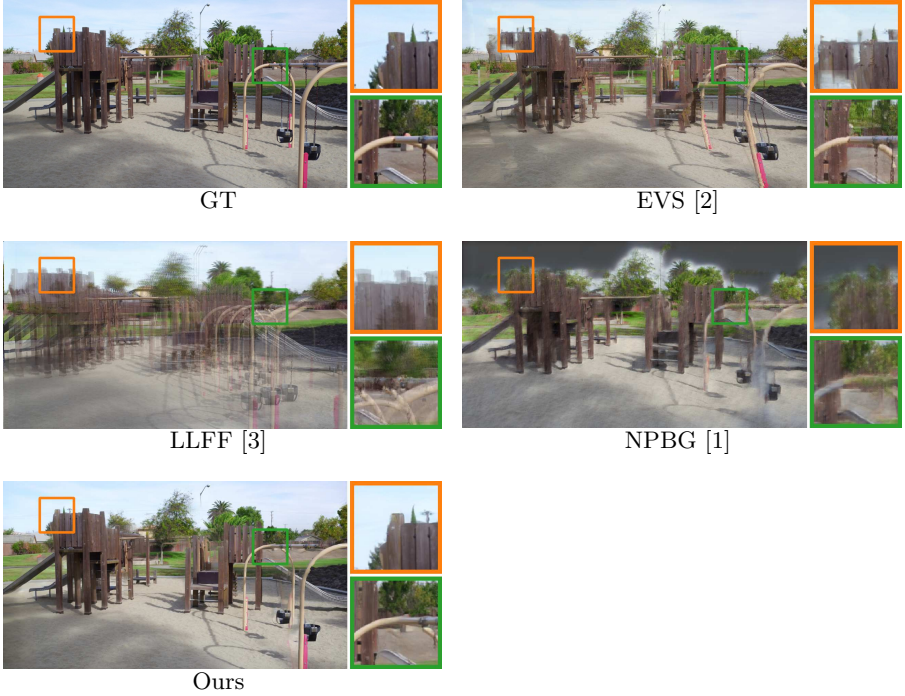


Fig. 14: Tanks and Temples *Playground*, view 31

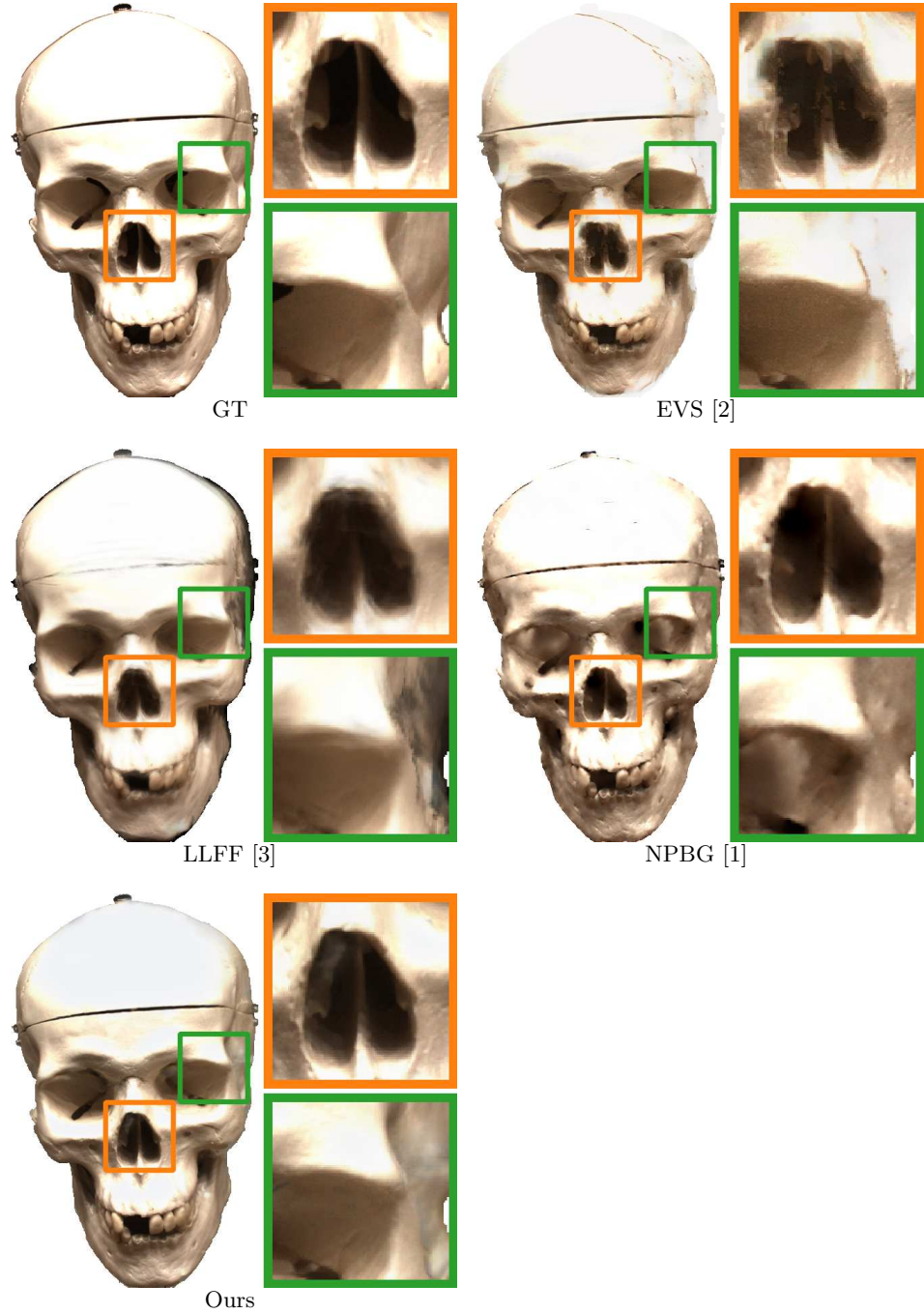


Fig. 15: DTU set 65, view 0, interpolation

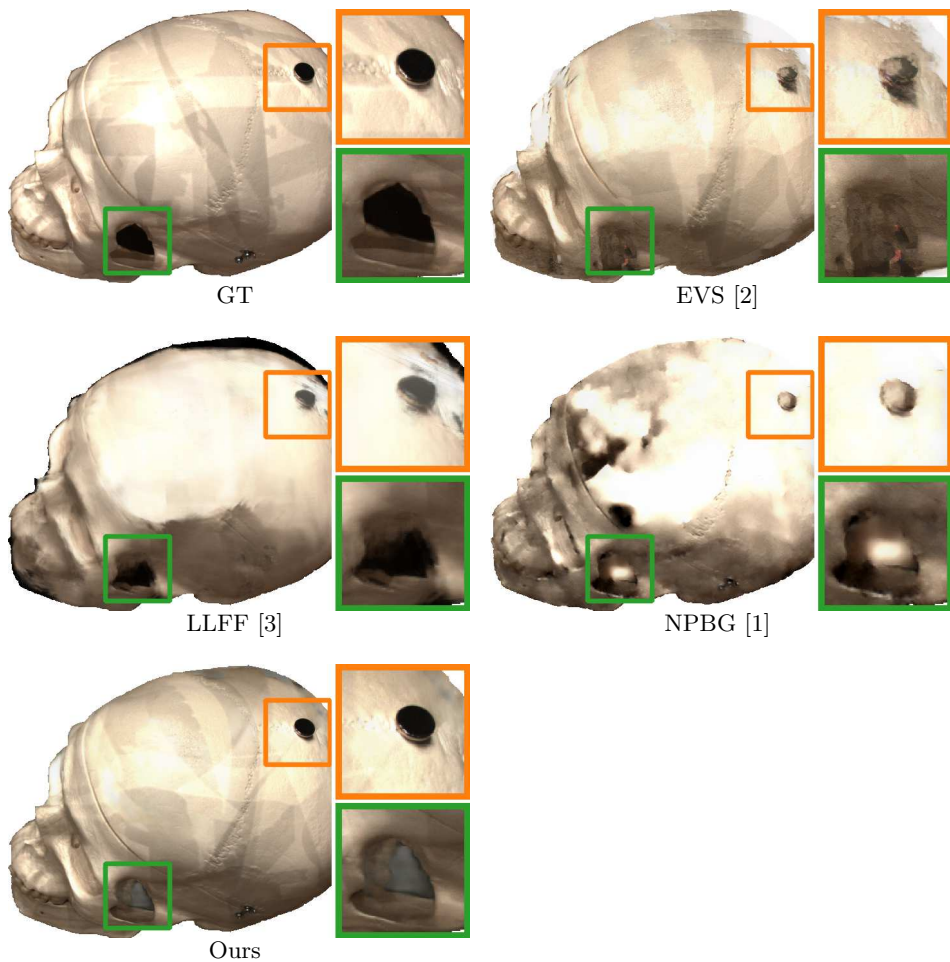


Fig. 16: DTU set 65, view 1, extrapolation

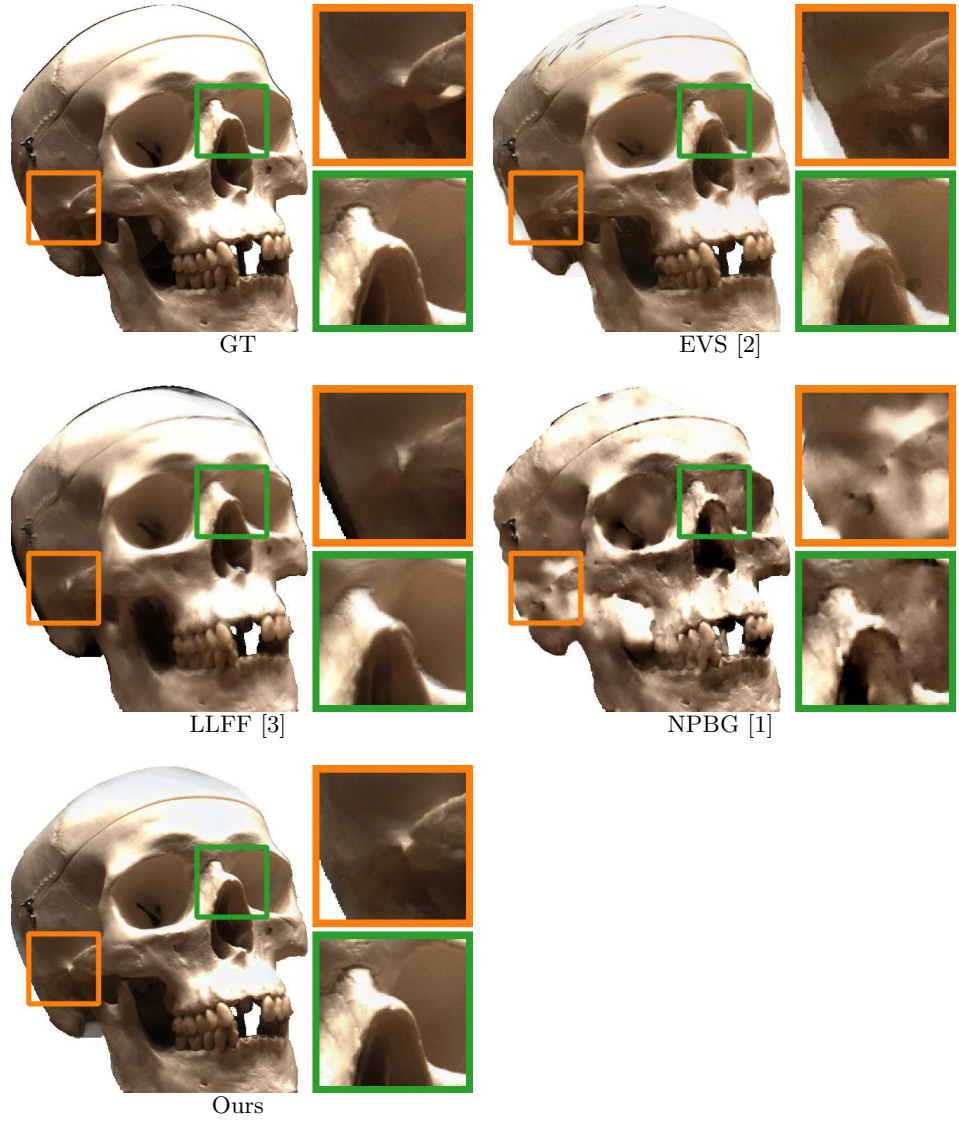


Fig. 17: DTU set 65, view 3, extrapolation



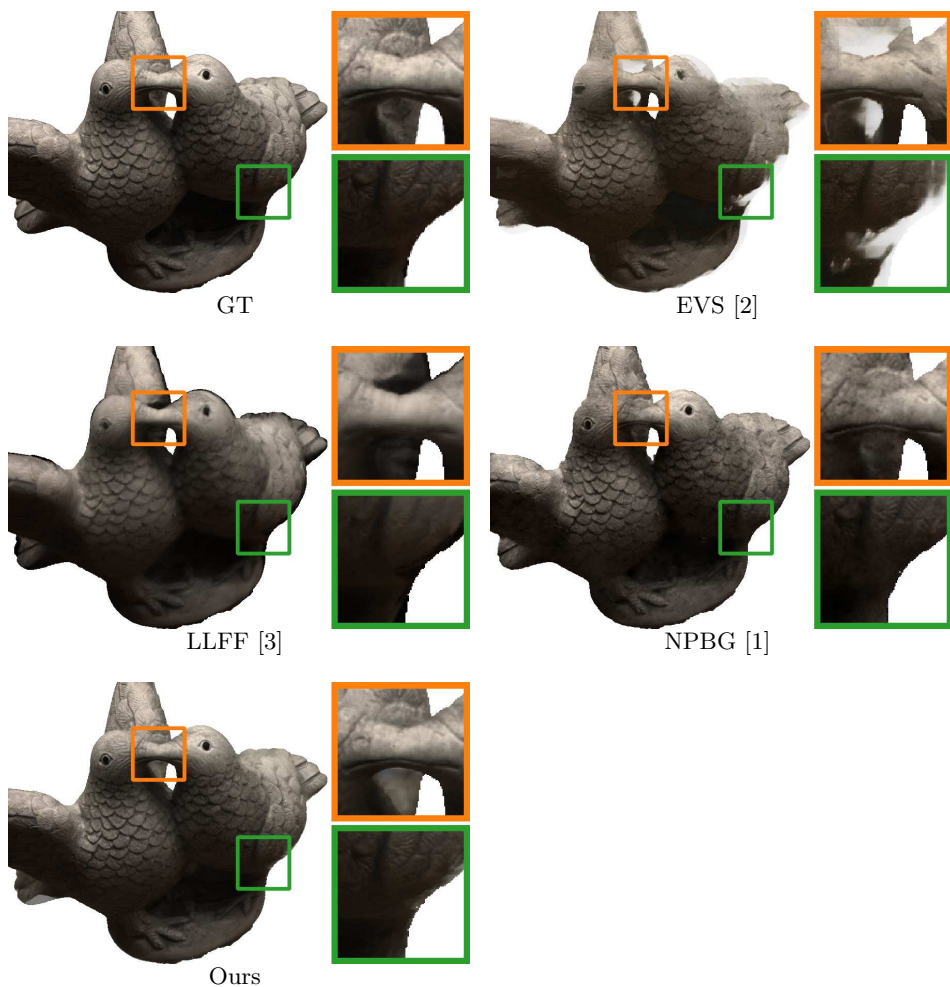


Fig. 18: DTU set 106, view 0, interpolation



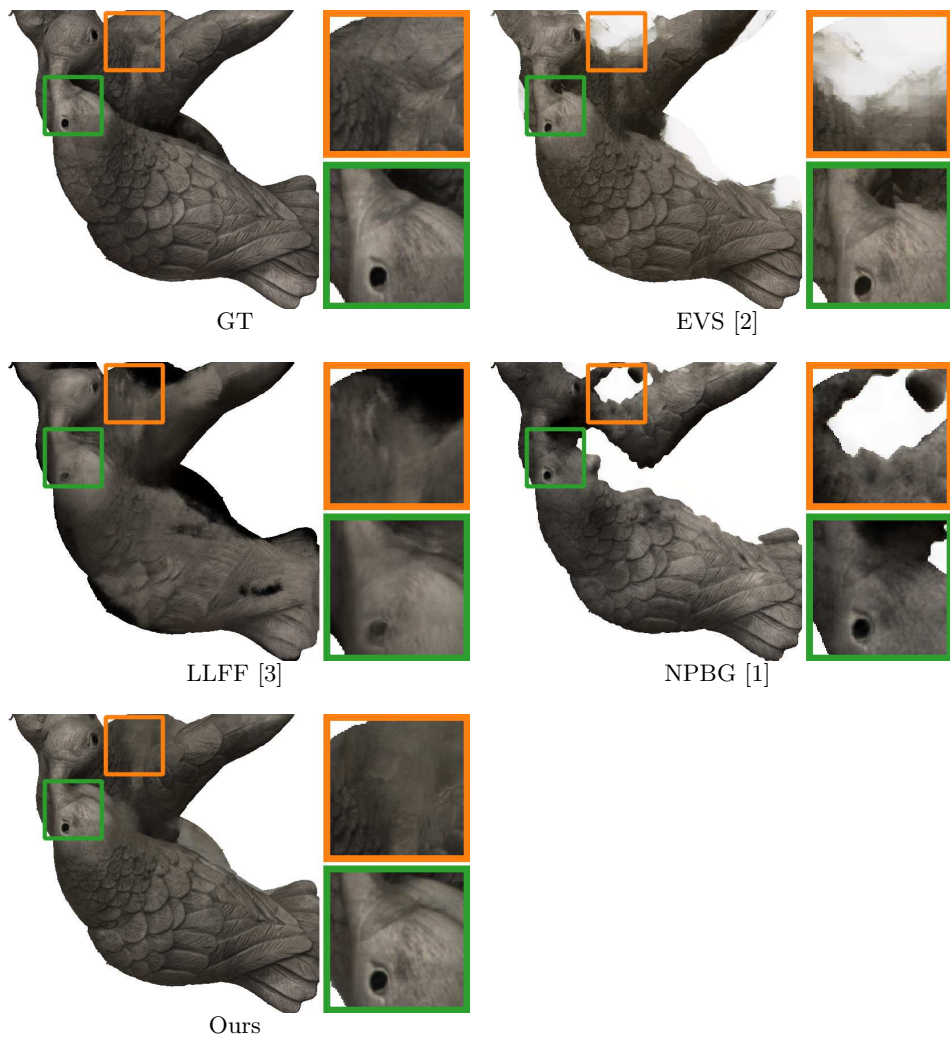


Fig. 19: DTU set 106, view 1, extrapolation

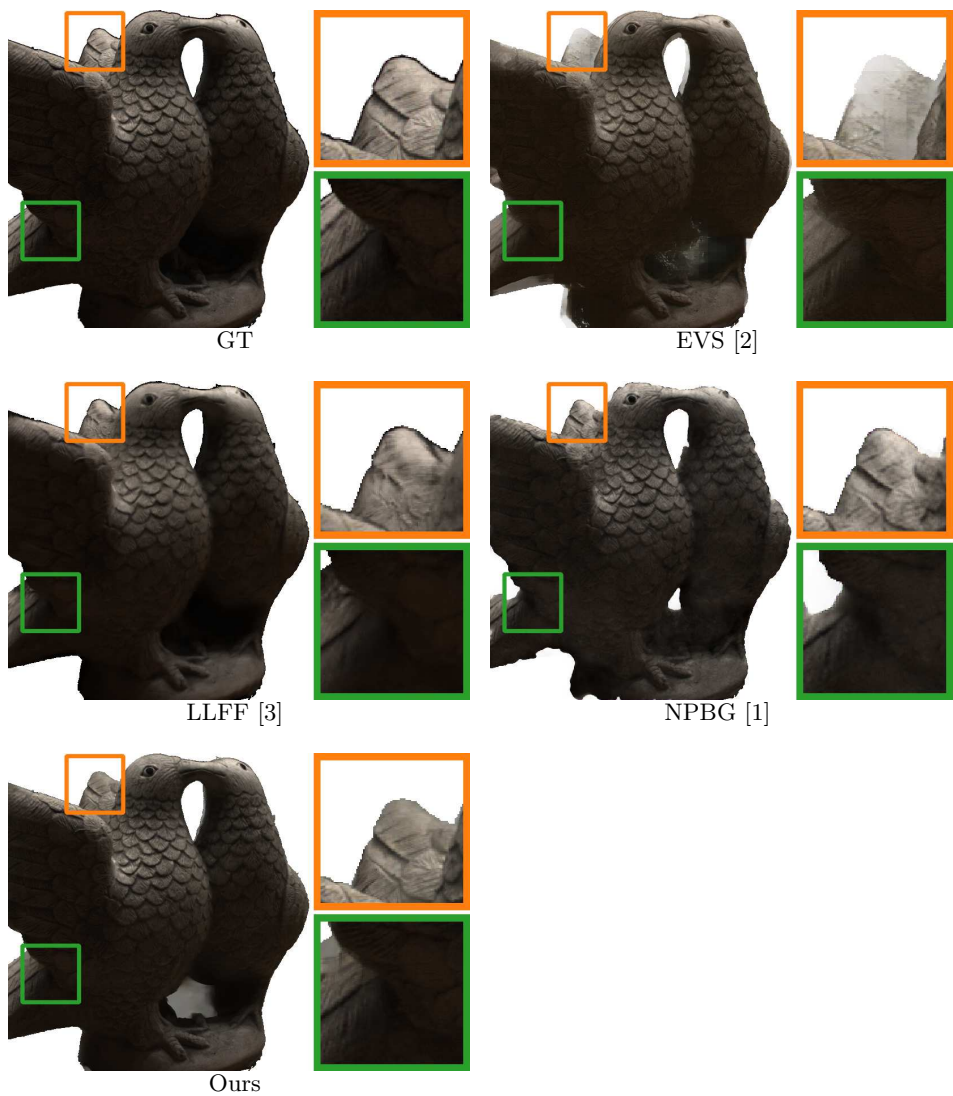


Fig. 20: DTU set 106, view 3, extrapolation



Fig. 21: DTU set 118, view 0, interpolation

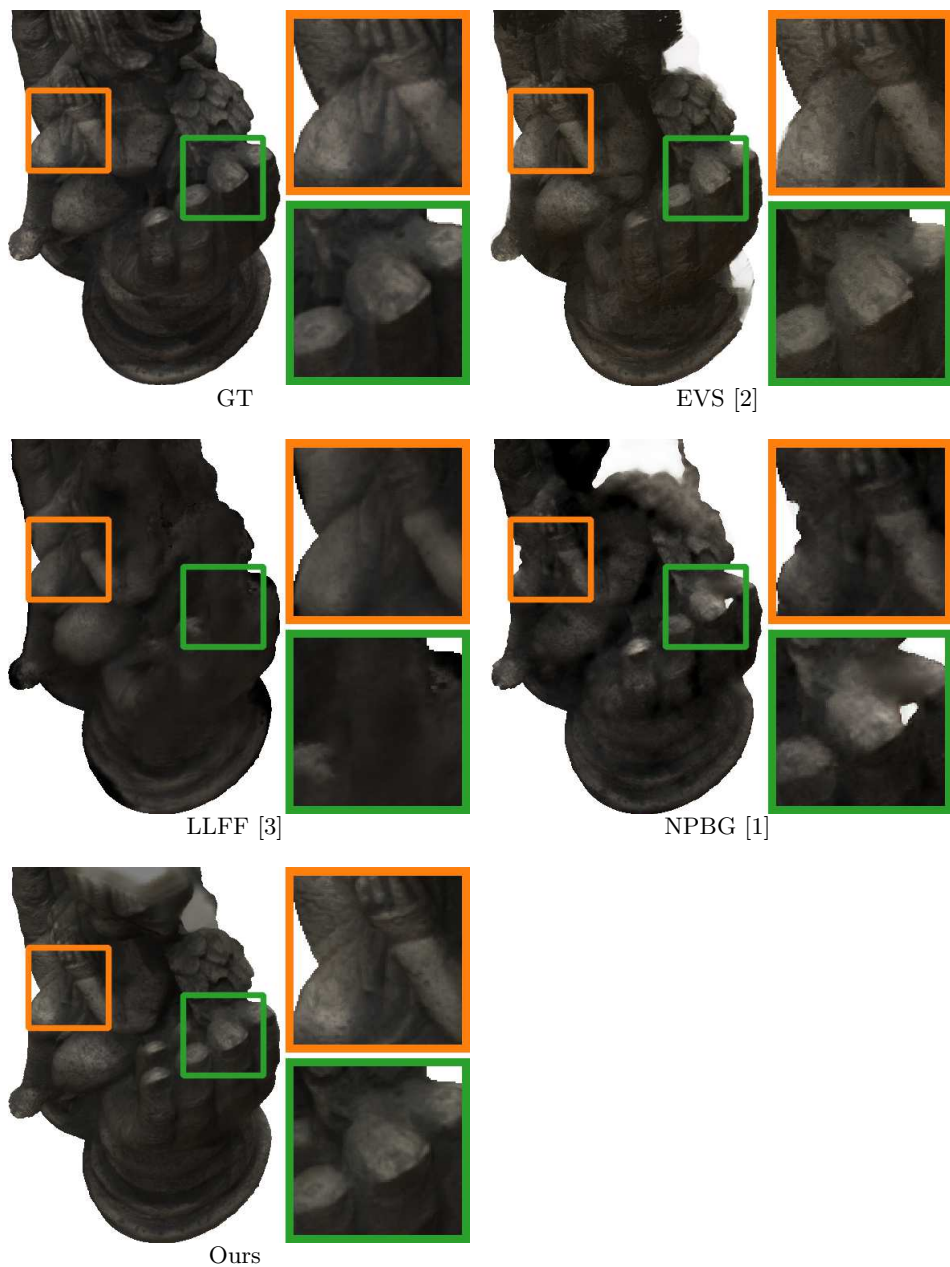


Fig. 22: DTU set 118, view 1, extrapolation





Fig. 23: DTU set 118, view 3, extrapolation



### 3 Network Details

Our recurrent mapping and blending network for novel view synthesis consists of two main components. First a VGG based encoder that extracts features from the source images, and a GRU based decoder that fuses the mapped features in the target view. Below, we present the PyTorch code of our network architecture. MappingBlendingRNN is the recurrent mapping and blending network, whereas EncoderVGGUNet is the source image encoder and DecoderGRUUnet is the GRU based decoder for blending. The code should give an overview of the number of feature maps, number of layers and layer connectivity. The complete source code of the method will be released upon acceptance.

---

```
class EncoderVGGUNet(nn.Module):
    def __init__(self, pool="average", n_encs=3, n_dec_convs=2):
        super().__init__()
        self.mean = torch.tensor([0.485, 0.456, 0.406]).view(1, 3, 1, 1)
        self.std = torch.tensor([0.229, 0.224, 0.225]).view(1, 3, 1, 1)
        vgg = torchvision.models.vgg19(pretrained=True).features

        encs = []
        enc = []
        encs_channels = []
        channels = -1
        for mod in vgg:
            if isinstance(mod, nn.Conv2d):
                channels = mod.out_channels

            if isinstance(mod, nn.MaxPool2d):
                encs.append(nn.Sequential(*enc))
                encs_channels.append(channels)
                n_encs -= 1
                if n_encs <= 0:
                    break
            if pool == "average":
                enc = [
                    nn.AvgPool2d(
                        kernel_size=2, stride=2, padding=0, ceil_mode=False
                    )
                ]
            elif pool == "max":
                enc = [
                    nn.MaxPool2d(
                        kernel_size=2, stride=2, padding=0, ceil_mode=False
                    )
                ]
            else:
                raise Exception("invalid pool")
            else:
                enc.append(mod)
        self.encs = nn.ModuleList(encs)

        cin = encs_channels[-1] + encs_channels[-2]
        decs = []
        for idx, cout in enumerate(reversed(encs_channels[:-1])):
            decs.append(self._dec(cin, cout, n_convs=n_dec_convs))
            cin = cout + encs_channels[max(-idx - 3, -len(encs_channels))]
        self.decs = nn.ModuleList(decs)

    def _dec(self, channels_in, channels_out, n_convs=2):
        mods = []
        for _ in range(n_convs):
            mods.append(
```

```

1125         nn.Conv2d(
1126             channels_in,
1127             channels_out,
1128             kernel_size=3,
1129             stride=1,
1130             padding=1,
1131             bias=False,
1132         )
1133     )
1134     mods.append(nn.ReLU())
1135     channels_in = channels_out
1136     return nn.Sequential(*mods)
1137
1138 def forward(self, x):
1139     x = (x + 1) / 2
1140     self.mean = self.mean.to(x.device)
1141     self.std = self.std.to(x.device)
1142     x = (x - self.mean) / self.std
1143
1144     feats = []
1145     for enc in self.encs:
1146         x = enc(x)
1147         feats.append(x)
1148
1149     for dec in self.decs:
1150         x0 = feats.pop()
1151         x1 = feats.pop()
1152         x0 = F.interpolate(
1153             x0, size=(x1.shape[2], x1.shape[3]), mode="nearest"
1154         )
1155         x = torch.cat((x0, x1), dim=1)
1156         x = dec(x)
1157         feats.append(x)
1158
1159     x = feats.pop()
1160     return x
1161
1162 class ConvGRU2d(nn.Module):
1163     def __init__(
1164         self,
1165         channels_x,
1166         channels_out,
1167         kernel_size=3,
1168         padding=1,
1169         bias=True,
1170     ):
1171         super().__init__()
1172         self.channels_x = channels_x
1173         self.channels_out = channels_out
1174
1175         self.conv_gates = nn.Conv2d(
1176             in_channels=channels_x + channels_out,
1177             out_channels=2 * channels_out,
1178             kernel_size=kernel_size,
1179             padding=padding,
1180             bias=bias,
1181         )
1182
1183         self.conv_can = nn.Conv2d(
1184             in_channels=channels_x + channels_out,
1185             out_channels=channels_out,
1186             kernel_size=kernel_size,
1187             padding=padding,
1188             bias=bias,
1189         )
1190
1191         self.nonlinearity = nn.ReLU()

```

```

def forward(self, x, h=None):
    if h is None:
        h = torch.zeros(
            (x.shape[0], self.channels_out, x.shape[2], x.shape[3]),
            dtype=x.dtype,
            device=x.device,
        )
    combined = torch.cat([x, h], dim=1)
    combined_conv = torch.sigmoid(self.conv_gates(combined))
    r = combined_conv[:, : self.channels_out]
    z = combined_conv[:, self.channels_out :]

    combined = torch.cat([x, r * h], dim=1)
    n = self.nonlinearity(self.conv_can(combined))

    h = z * h + (1 - z) * n
    return h

class DecoderGRUUNet(nn.Module):
    def __init__(
        self,
        channels_in,
        enc=[32, 128, 128],
        dec=[128, 32],
        n_enc_convs=2,
        n_dec_convs=2,
        gru_all=False,
        bias=False,
    ):
        super().__init__()
        self.n_rnn = 0

        stride = 1
        cin = channels_in
        encs = []
        for cout in enc:
            encs.append(
                self._enc(
                    cin,
                    cout,
                    stride=stride,
                    n_convs=n_enc_convs,
                    gru_all=gru_all,
                )
            )
            stride = 2
            cin = cout
        self.encs = nn.ModuleList(encs)

        cin = enc[-1] + enc[-2]
        decs = []
        for idx, cout in enumerate(dec):
            decs.append(
                self._dec(
                    cin,
                    cout,
                    n_convs=n_dec_convs,
                    gru_all=gru_all,
                )
            )
            cin = cout + enc[max(-idx - 3, -len(enc))]
        self.decs = nn.ModuleList(decs)

    def _enc(
        self,
        channels_in,
        channels_out,

```

```

1215         stride=2,
1216         n_convs=2,
1217         gru_all=False,
1218     ):
1219         mods = []
1220         if stride > 1:
1221             mods.append(nn.AvgPool2d(kernel_size=2))
1222         for idx in range(n_convs):
1223             if gru_all or idx == n_convs - 1:
1224                 self.n_rnn += 1
1225                 mods.append(
1226                     ConvGRU2d(
1227                         channels_in,
1228                         channels_out,
1229                         kernel_size=3,
1230                         padding=1,
1231                         bias=False,
1232                     )
1233                 )
1234             else:
1235                 mods.append(
1236                     nn.Conv2d(
1237                         channels_in,
1238                         channels_out,
1239                         kernel_size=3,
1240                         padding=1,
1241                         bias=False,
1242                     )
1243                 )
1244             mods.append(nn.ReLU())
1245             channels_in = channels_out
1246             stride = 1
1247         return nn.Sequential(*mods)
1248
1249 def _dec(
1250     self,
1251     channels_in,
1252     channels_out,
1253     n_convs=2,
1254     gru_all=False,
1255 ):
1256     mods = []
1257     for idx in range(n_convs):
1258         if gru_all or idx == n_convs - 1:
1259             self.n_rnn += 1
1260             mods.append(
1261                 ConvGRU2d(
1262                     channels_in,
1263                     channels_out,
1264                     kernel_size=3,
1265                     padding=1,
1266                     bias=False,
1267                 )
1268             )
1269         else:
1270             mods.append(
1271                 nn.Conv2d(
1272                     channels_in,
1273                     channels_out,
1274                     kernel_size=3,
1275                     padding=1,
1276                     bias=False,
1277                 )
1278             )
1279         mods.append(nn.ReLU())
1280         channels_in = channels_out
1281     return nn.Sequential(*mods)

```

```

def forward(self, x, hs=None):
    if hs is None:
        hs = [None for _ in range(self.n_rnn)]

    hidx = 0
    feats = []
    for enc in self.encs:
        for mod in enc:
            if isinstance(mod, ConvGRU2d):
                x = mod(x, hs[hidx])
                hs[hidx] = x
                hidx += 1
            else:
                x = mod(x)
        feats.append(x)

    for dec in self.decs:
        x0 = feats.pop()
        x1 = feats.pop()
        x0 = F.interpolate(
            x0, size=(x1.shape[2], x1.shape[3]), mode="nearest"
        )
        x = torch.cat((x0, x1), dim=1)
        for mod in dec:
            if isinstance(mod, ConvGRU2d):
                x = mod(x, hs[hidx])
                hs[hidx] = x
                hidx += 1
            else:
                x = mod(x)
        feats.append(x)

    x = feats.pop()

class MappingBlendingRNN(nn.Module):
    def __init__(self, enc, mapper, dec, merge_channels_out):
        super().__init__()
        self.enc = enc
        self.mapper = mapper
        self.dec = dec

        self.rgb_conv = nn.Conv2d(
            merge_channels_out,
            3,
            kernel_size=1,
            stride=1,
            padding=0,
            bias=False,
        )
        self.alpha_conv = nn.Conv2d(
            merge_channels_out,
            1,
            kernel_size=1,
            stride=1,
            padding=0,
            bias=False,
        )

    def forward(self, x, K0, R0, t0, K, R, t, dm0, nb_dms):
        bs, nv, channels, height, width = x.shape

        x = x.view(bs * nv, *x.shape[-3:])
        x = self.enc(x)
        x = x.view(bs, nv, *x.shape[-3:])

        with torch.no_grad():
            K0i = torch.inverse(K0)

```



```
1305         P0i = torch.bmm(1305
1306             R0.transpose(2, 1), torch.cat((K0i, -t0.view(bs, 3, 1)), dim=2)1306
1307         )1307
1308         P = torch.matmul(K, torch.cat((R, t.view(bs, nv, 3, 1)), dim=3))1308
1309         x = self.mapper(x, dm0, nb_dms, P, P0i)1309
1310
1311         hs = None1310
1312         rgbs = []1310
1313         alphas = []1310
1314         for vidx in range(nv):1311
1315             y, hs = self.dec(x[:, vidx], hs)1311
1316             rgbs.append(self.rgb_conv(y))1312
1317             alphas.append(self.alpha_conv(y))1313
1318
1319         rgbs = torch.stack(rgbs)1314
1320         alphas = torch.stack(alphas)1315
1321         alphas = torch.softmax(alphas, dim=0)1315
1322         x = (alphas * rgbs).sum(dim=0)1316
1323         del rgbs, alphas1317
1324         return {"rgb": x}1317
```

---

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

## References

1. Aliev, K.A., Ulyanov, D., Lempitsky, V.: Neural Point-Based Graphics. arXiv preprint arXiv:1906.08240 (2019)
2. Choi, I., Gallo, O., Troccoli, A., Kim, M.H., Kautz, J.: Extreme View Synthesis. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2019)
3. Mildenhall, B., Srinivasan, P.P., Ortiz-Cayon, R., Kalantari, N.K., Ramamoorthi, R., Ng, R., Kar, A.: Local Light Field Fusion: Practical View Synthesis with Prescriptive Sampling Guidelines. ACM Transactions on Graphics (SIGGRAPH) (2019)