

Supplementary Materials

3DFA: Aligning the Features Between Point Cloud and Query Image for Scene-Specific Visual Localization

Anonymous CVPR submission

Paper ID ****

001 **1. Implementation Details**

002 **1.1. Trainable Parameters**

003 The trainable parameters of 3DFA include

- 004 • The feature head of R2D2 network [2].
- 005 • ϕ_1 , the 2D patch encoder.
- 006 • ϕ_2 , the 3D patch encoder.
- 007 • ϕ_3 , the 2D-3D fusion layer.
- 008 • ϕ_4 , the 2D patch decoder.
- 009 • ϕ_5 , the 3D patch decoder.
- 010 • The cross attention layer in the final pose estimation mod-
011 ule.
- 012 • ϕ_6 , the reliability score layer.

013 Among all these trainable parameters, ϕ_1 to ϕ_6 are imple-
014 mented as MLPs, where ϕ_1 to ϕ_5 uses ReLU activation and
015 ϕ_6 uses Sigmoid activation. The cross attention layer as
016 well as the 3D positional encoding is borrowed from the
017 widely used point transformer [8]. Note that the keypoint
018 extraction head of R2D2 network is frozen, as we only
019 want to finetune the feature extraction part.

020 **1.2. Hyperparameters**

021 In summary, the hyperparameters of 3DFA include

- 022 • α, β , the loss factor.
- 023 • M , the final pose estimation scope.
- 024 • S , the reliability score threshold.

025 Our experimental results in the paper (Table 1~3 in main
026 paper) are obtained with $\alpha = \beta = 1, M = 5$, and $S = 0.7$.
027 The experiment on α and β is included in later part.

028 **1.3. Training and Inference**

029 3DFA is trained on RTX4090 GPU. In DL3DV-10K dataset,
030 training takes around 22 hours for one scene. This period
031 does not include the point cloud reconstruction and 2D key-
032 point extraction, which we consider as data preprocessing.
033 During inference, 3DFA is able to process 0.54 image per
034 second. One more advantage of 3DFA is that, since it re-
035 places support images by the point cloud for feature align-

ment, our approach does not require support images during
inference. This makes it easier to deploy 3DFA to real ap-
plications, as we only need the trained weights and the point
cloud.

036 **2. Data Preparation Details**

037 3DFA employs PNeRF [9] as the 3D reconstruction ap-
038 proach which also requires depth maps as point cloud
039 initialization. *7Scenes* [4] already provides ground truth
depth maps. For *DL3DV-10K* and *ScanNet++*, we use
COLMAP [3] to estimate the depth maps. If a scene is
named by hashes, we assign it a casual name. We follow
previous work [10] for the real world scale information and
apply a rough scale when the ground truth is missing.

040 **3. More Experimental Results**

041 **3.1. Comparison Schemes Setup**

042 To have a comprehensive evaluation of model performance,
043 we include both scene-agnostic and scene-specific schemes
044 in our experiments. However, since we focus on scene-
045 specific visual localization in this paper, we ensure a fair
046 comparison by feeding all networks the same set of training
047 data. All results for one scene is obtained on that specific
048 scene using the same set of support and query images. It is
049 infeasible to train those foundation models again [1, 5-7],
050 so for them we only run inference. We find that 16 images
051 are enough to yield their maximum capabilities roughly, so
052 for each query image, we feed the 16 most similar support
053 images to the foundation models.

054 **3.2. Positions of Patch-Matching Loss**

055 As mentioned in the paper, an ascending number of patches
056 result in the best localization performance. Our final 3DFA
057 uses a “3, 5, 7, 9” number of patches. Here we also explore
058 the positions of patch-matching loss inside these different
059 3DFA layers. The results in Table 1 indicate that apply-
060 ing the patch-matching loss at every layer is essential, as
061
062
063
064
065
066
067
068
069

Table 1. Median translation error (centimeters) on the ScanNet++ dataset of 3DFA when patch-matching loss is removed in a certain layer. For example, “w/o \mathcal{L}_p @3” means to remove patch-matching loss in the first 3DFA layer.

3DFA	Work	Studio	Bedroom	Living
w/o \mathcal{L}_p @3	5.0	5.3	6.0	4.8
w/o \mathcal{L}_p @5	5.1	5.3	5.9	4.8
w/o \mathcal{L}_p @7	5.1	5.4	5.9	4.8
w/o \mathcal{L}_p @9	5.2	5.5	6.0	4.7
w/ all \mathcal{L}_p	4.8	5.2	5.8	4.7

Table 2. Median translation error (centimeters) on the ScanNet++ dataset of 3DFA with different α (factor of patch-matching loss).

α	Work	Studio	Bedroom	Living
0	5.4	5.5	6.0	4.8
0.25	5.3	5.5	5.9	4.8
0.5	5.0	5.3	5.9	4.8
0.75	4.8	5.2	5.8	4.8
1	4.8	5.2	5.8	4.7
1.25	4.9	5.2	5.8	4.8
1.5	5.1	5.3	6.0	5.1
2	5.3	5.5	6.1	5.2

070 removing it from any layer leads to a measurable drop in
071 performance. This observation aligns with our design prin-
072 ciple of multiscale hierarchical feature alignment. Relying
073 on a single scale fails to yield optimal localization accuracy.

074 3.3. Loss Factor α and β

075 The two hyperparameters, α and β , are to balance the lo-
076 calization loss and patch-matching loss in 3DFA. To iden-
077 tify an effective configuration, we conduct the experiments
078 where β is fixed to be 1 and α has different values. As
079 shown in Table 2, the best choice of α lies between 0.75 to
080 1.25 and we use $\alpha = 1$ in our finalized model. These re-
081 sults suggest that maintaining a balance between the patch-
082 matching loss and the localization loss leads to the best
083 overall performance. This outcome is expected: while patch
084 matching is designed to inject spatial correspondence into
085 the 2D and 3D feature representations, the primary objec-
086 tive of the framework is still accurate localization. When the
087 patch-matching term is weighted too heavily, it can over-
088 whelm and ultimately degrade the original point cloud fea-
089 tures, resulting in diminished localization accuracy.

References

- [1] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Ground-
ing image matching in 3d with mast3r. In *European Confer-
ence on Computer Vision*, pages 71–91. Springer, 2024. 1
- [2] Jerome Revaud, Cesar De Souza, Martin Humenberger, and
Philippe Weinzaepfel. R2d2: Reliable and repeatable detec-
tor and descriptor. *Advances in neural information process-
ing systems*, 32, 2019. 1
- [3] Johannes Lutz Schönberger and Jan-Michael Frahm. Struc-
ture-from-motion revisited. In *Conference on Com-
puter Vision and Pattern Recognition (CVPR)*, 2016. 1
- [4] Jamie Shotton, Ben Glocker, Christopher Zach, Shahram
Izadi, Antonio Criminisi, and Andrew Fitzgibbon. Scene co-
ordinate regression forests for camera relocalization in rgb-d
images. In *Proceedings of the IEEE conference on computer
vision and pattern recognition*, pages 2930–2937, 2013. 1
- [5] Jianyuan Wang, Nikita Karaev, Christian Rupprecht, and
David Novotny. Vggsfm: Visual geometry grounded deep
structure from motion. In *Proceedings of the IEEE/CVF con-
ference on computer vision and pattern recognition*, pages
21686–21697, 2024. 1
- [6] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea
Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Vi-
sual geometry grounded transformer. In *Proceedings of the
Computer Vision and Pattern Recognition Conference*, pages
5294–5306, 2025. 1
- [7] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris
Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vi-
sion made easy. In *Proceedings of the IEEE/CVF Conference
on Computer Vision and Pattern Recognition*, pages 20697–
20709, 2024. 1
- [8] Xiaoyang Wu, Li Jiang, Peng-Shuai Wang, Zhijian Liu, Xi-
hui Liu, Yu Qiao, Wanli Ouyang, Tong He, and Hengshuang
Zhao. Point transformer v3: Simpler faster stronger. In *Pro-
ceedings of the IEEE/CVF conference on computer vision
and pattern recognition*, pages 4840–4851, 2024. 1
- [9] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin
Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-
nerf: Point-based neural radiance fields. In *Proceedings of
the IEEE/CVF conference on computer vision and pattern
recognition*, pages 5438–5448, 2022. 1
- [10] Boming Zhao, Luwei Yang, Mao Mao, Hujun Bao, and
Zhaopeng Cui. Pnerfloc: Visual localization with point-
based neural radiance fields. In *Proceedings of the AAAI
Conference on Artificial Intelligence*, pages 7450–7459,
2024. 1

090

091

092

093

094

095

096

097

098

099

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135