Large-displacement 3D Object Tracking with Hybrid Non-local Optimization

Supplementary Material

1 The Out-of-plane Rotation

1.1 The Decomposition of Rotation

As described in the paper, a rotation matrix \mathbf{R} can be decomposed as in-plane rotation \mathbf{R}^{in} and out-of-plane rotation \mathbf{R}^{out} , so that $\mathbf{R} = \mathbf{R}^{in}\mathbf{R}^{out}$. \mathbf{R}^{out} will change the view direction, and \mathbf{R}^{in} is a rotation around the view axis. However, note that this decomposition is not unique, because there is no standard way to change the view direction without introducing in-plane rotation.

We thus adopt a similar approach as the *gluLookAt* function, which can keep an up direction while changing the view direction. To retrieve the view direction \mathbf{v} from \mathbf{R} , we should note that in the camera coordinate space, the view axis is always the Z axis. Therefore, \mathbf{v} should be aligned with the Z axis after the rotation, i.e. $\mathbf{Rv} = [0, 0, 1]^{\mathsf{T}}$. Obviously, \mathbf{v} should be the last row of \mathbf{R} .

Given the view direction \mathbf{v} , we can compute \mathbf{R}^{out} similarly as gluLookAt, by specifying an up vector ([0, 1, 0] in our implementation). We found that this approach can better decompose in-plane and out-of-plane rotations than using Euler angles.

1.2 The Distribution of Failures

As mentioned in the paper, our basic observation is that most failures of previous 3D tracking methods are caused by the out-of-plane rotations. Figure 1 shows the statistics of tracking failures caused by different components (translation, inplane and out-of-plane rotation) for the representative methods RBOT [7] and SRT3D [5]. As can be seen, about 70% of RBOT failures and 90% of SRT3D failures are due to the out-of-plane rotation. We found that the translation and in-plane rotation introduce more failures for RBOT mainly because the optimization process is not sufficiently converged in order to achieve real-time speed, which results in larger error for the translation and in-plane rotation.

To classify the failures, the rotation error $\mathbf{R}_d = \mathbf{R}\mathbf{R}_{dt}^{-1}$ is first decomposed as in-plane part \mathbf{R}_d^{in} and out-of-plane part \mathbf{R}_d^{out} , using the method as described above. The standard 5°-5cm criteria then is applied for the classification. If the translation error is greater than 5*cm*, the translation failure is added by one. If the error of \mathbf{R}_d is greater than 5°, the one in \mathbf{R}_d^{in} and \mathbf{R}_d^{out} with larger error is counted as fail.



Fig. 1. The distribution of failures with respect to translation, in-plane and out-ofplane rotations on the *regular* variant of RBOT. *left*: The results of RBOT [7]; *right*: The results of SRT3D [5].

2 Pose Optimization

To optimize the pose, the cost function is firstly rewritten as

$$E(\xi) = \sum_{i=1}^{N} \omega_i \parallel F(\xi, i) \parallel^{\alpha}, \quad \text{where } F(\xi, i) = \mathbf{n}_i^{\top} (\mathbf{x}_i^{\xi} - \mathbf{o}_i) - d_i$$
(1)

which does not satisfy the general form of least square problem. Therefore, To solve it with *iterative reweighted least square*(IRLS) , we can further rewrite it as

$$E(\xi) = \sum_{i=1}^{N} \omega_i \psi_i F(\xi, i)^2, \quad \text{with } \psi_i = \frac{1}{\|F(\xi, i)\|^{2-\alpha}}$$
(2)

with ψ_i fixed weights computed with the current ξ , which is used to penalize the correspondences that are with large matching residuals. We solve the optimization problem using Gauss-Newton method. By linearizing $F(\xi, i)$ with its first-order Taylor approximation $F(\xi + \Delta \xi) = F(\xi) + \mathbf{J}\Delta \xi$, we have

$$E(\xi) \approx \sum_{i=1}^{N} \omega_i \psi_i \left(F(\xi, i) + \mathbf{J} \Delta \xi \right)^2 \tag{3}$$

which is quadratic with respect to $\Delta \xi$, so the optimal solution should be

$$\Delta \xi = -(\sum_{i=1}^{N} \omega_i \psi_i \mathbf{J}^{\top} \mathbf{J})^{-1} \sum_{i=1}^{N} \omega_i \psi_i \mathbf{J} F(\xi, i)$$
(4)

with **J** the Jacobian of $F(\xi, i)$

$$\mathbf{J} = \frac{\partial F(\xi, i)}{\partial \xi} = n_i^\top \frac{\partial \mathbf{x}_i^\xi}{\partial \xi}$$
(5)



Fig. 2. Results of all variants for different frame steps (S = 1, 2, 3, 4). *Ours(local)* is *Ours*⁻. The compared methods are RBOT [7], RBGT [4], SRT3D [5].

where $\mathbf{x}_{i}^{\xi} = \pi(\mathbf{K}(\mathbf{T}\widetilde{\mathbf{X}})_{3\times 1}) = \pi(\mathbf{K}(\mathbf{R}\mathbf{X} + \mathbf{t}))$, therefore

$$\frac{\partial \mathbf{x}_{i}^{\xi}}{\partial \xi} = \frac{\partial \pi (\mathbf{K}(\mathbf{RX} + \mathbf{t}))}{\partial (\mathbf{RX} + \mathbf{t})} \frac{\partial (\mathbf{RX} + \mathbf{t})}{\partial \mathbf{X}} \frac{\partial \mathbf{X}}{\partial \xi}
= \begin{bmatrix} \frac{f_{x}}{Z'} & 0 & -\frac{f_{x}X'}{(Z')^{2}} \\ 0 & \frac{f_{y}}{Z'} & -\frac{f_{y}Y'}{(Z')^{2}} \end{bmatrix} \mathbf{R} \begin{bmatrix} 1 & 0 & 0 & 0 & Z & -Y \\ 0 & 1 & 0 & -Z & 0 & X \\ 0 & 0 & 1 & Y & -X & 0 \end{bmatrix}$$
(6)

in which $\mathbf{X}' = [X', Y', Z'] = \mathbf{R}\mathbf{X} + \mathbf{t}$ is the coordinate in the camera space. After $\Delta \xi$ is computed, the object pose can be updated as

$$\mathbf{T} \leftarrow \mathbf{T} \exp(\Delta \xi) \tag{7}$$

Note that $\exp(\Delta\xi)$ should appear at the right of **T**, because in Eq. (6) the disturbance is applied in the model space.

3 Results of Other Variants

Due to the space limitation, only the results of the *regular* variant of the RBOT dataset are included in the paper. Here we present results of all variants, including *regular*, *dynamic light*, *noisy*, *unmodeled occlusion*.

Figure 2 shows the results for different frame steps. All variants have similar trends, our local and non-local methods both significantly outperform previous methods. For large displacements, great improvement is achieved with our non-local method. The accuracy of *noisy* is much lower than the accuracy of other three variants, because the severe noise would introduce significant error to the probability map, resulting in errors in the correspondences. However, when compared with previous methods, our method still performs much better for large displacements.

Table 1 contains more details and comparisons for S = 1, which is the standard benchmark of the RBOT dataset. Results of the compared methods are from the original paper. As can be seen, for S = 1 our method still achieves the best accuracy, slightly better than SRT3D [5] and significantly better than others. Since the displacement is small, using non-local search only slightly improve the accuracy. For the *noisy* variant, using non-local search for S = 1 even 4 X. Tian et al.

decrease the accuracy by 0.2%. As we have analyzed in the paper, this is mainly due to the error function, which may fail to predict the true object pose in complicated cases.

2																		~	
Arecto	Apo	Soda		South	Cetheres	0 ⁴	Solo Color	Contra Co	Cube Ballo	Driller.	Duck duck	47.00 th	Glue Glue	4.00	C. C	Lenno	Proposto	Squirre	sé. V
	Regular																		
[7] [8] [2] [6] [1] [4] [5] Ours Ours	85.0 88.8 92.8 91.9 93.0 94.6 96.4 98.8 99.8 99.8	39.0 41.3 42.6 44.8 55.2 49.4 53.2 65.1 65.6 67.1	98.9 94.0 96.8 99.7 99.3 99.5 98.8 99.6 99.5 100	82.4 85.9 87.5 89.1 85.4 91.0 93.9 96.0 95.0 97.8	79.7 86.9 90.7 89.3 96.1 93.7 93.0 98.0 96.6 97.3	87.6 89.0 86.2 90.6 93.9 96.0 92.7 96.5 92.6 93.7	95.9 98.5 99.0 97.4 98.0 97.8 99.7 100 100 100	93.3 93.7 96.9 95.9 95.6 96.6 97.1 98.4 98.7 99.4	78.1 83.1 86.8 83.9 79.5 90.2 92.5 94.1 95.0 97.4	93.0 87.3 94.6 97.6 98.2 98.2 92.5 96.9 97.1 97.6	86.8 86.2 90.4 91.8 89.7 93.4 93.7 98.0 97.4 99.3	74.6 78.5 87.0 84.4 89.1 90.3 88.5 95.3 96.1 96.9	38.9 58.6 57.6 59.0 66.5 64.4 70.0 79.3 83.3 84.7	81.0 86.3 92.5 91.3 94.0 92.1 96.0 96.9 97.7	46.8 57.9 59.9 74.3 60.6 79.0 78.8 90.3 91.5 93.4	97.5 91.7 96.5 97.4 98.6 98.8 95.5 97.4 95.8 96.7	80.7 85.0 90.6 86.4 95.6 92.9 92.5 96.2 95.2 95.2	99.4 96.2 99.5 99.7 99.6 99.8 99.6 99.8 99.7 100	79.9 82.7 85.8 86.9 88.1 89.9 90.0 94.2 94.2 94.2 95.2
Dynamic Light																			
[7] [8] [3] [2] [6] [1] [4] [5] Ours ⁻ Ours	84.9 89.7 93.5 91.8 93.8 94.3 96.5 98.2 99.7 100	$\begin{array}{c} 42.0\\ 40.2\\ 43.1\\ 42.3\\ 55.9\\ 48.3\\ 54.6\\ \textbf{65.2}\\ 64.7\\ 64.5\\ \end{array}$	99.0 92.7 96.6 98.9 99.6 99.5 99.1 99.2 99.7 99.8	81.3 86.5 88.5 89.9 85.6 90.1 93.9 95.6 95.2 97.9	84.3 86.6 92.8 91.3 97.7 94.6 93.1 97.5 97.2 97.9	88.9 89.2 86.0 87.8 93.7 96.1 94.7 98.1 93.0 94.0	95.6 98.3 99.6 97.6 97.7 97.9 99.5 100 99.8 100	92.5 93.9 95.5 94.5 96.5 97.3 97.0 98.5 98.8 99.5	77.5 81.8 85.7 84.5 78.3 90.9 93.0 94.2 94.2 94.2 97.0	94.6 88.4 96.8 98.1 98.6 99.1 93.4 97.5 98.4 98.8	86.4 83.9 91.1 91.9 91.0 92.9 93.3 97.9 97.1 99.3	77.3 76.8 90.2 86.7 91.6 91.5 92.6 96.9 97.2 97.6	52.9 55.3 68.4 66.2 72.1 72.6 74.9 86.1 86.9 87.5	77.9 79.3 86.8 90.9 90.7 94.7 91.0 95.2 95.4 97.4	47.9 54.7 59.7 73.2 63.0 80.0 79.2 89.3 91.4 92.4	96.9 88.7 96.1 97.1 98.9 98.3 95.6 97.0 96.2 97.1	81.7 81.0 91.5 89.2 94.4 95.2 89.8 95.9 95.6 96.4	99.3 95.8 99.2 99.6 100 99.8 99.5 99.9 99.9 100	81.2 81.3 86.7 87.3 88.8 90.7 90.6 94.6 94.5 95.4
Noiev																			
[7] [8] [2] [6] [1] [4] [5] <i>Ours</i> <i>Ours</i>	77.5 79.3 89.1 89.0 92.5 91.0 91.9 96.9 98.4 99.3	44.5 35.2 44.0 45.0 56.2 49.1 53.3 61.9 62.6 62.0	91.5 82.6 91.6 89.5 98.0 95.6 90.2 95.4 95.5 95.8	82.9 86.2 89.4 90.2 85.1 91.0 92.6 95.7 94.0 97.7	$51.7 \\ 65.1 \\ 75.2 \\ 68.9 \\ 91.7 \\ 76.3 \\ 67.9 \\ 84.5 \\ 90.6 \\ 90.4$	38.4 56.9 62.3 38.3 79.0 54.1 59.3 73.9 69.6 68.6	95.1 96.9 98.6 95.9 97.7 97.1 98.4 99.9 99.8 99.9	69.2 67.0 77.3 72.8 86.2 73.7 80.6 90.3 89.1 91.3	24.4 37.5 41.2 20.1 40.1 27.3 43.5 62.2 57.1 54.2	64.3 75.2 81.5 96.6 92.8 78.1 87.8 95.7 95.4	88.5 85.4 91.6 92.2 90.8 95.3 92.5 97.6 97.2 99.0	$\begin{array}{c} 11.2\\ 35.2\\ 54.5\\ 26.8\\ \textbf{70.2}\\ 30.2\\ 44.0\\ 62.2\\ 68.7\\ 64.8 \end{array}$	$\begin{array}{c} 2.9\\ 18.9\\ 31.8\\ 15.8\\ 50.9\\ 7.8\\ 31.3\\ 43.4\\ \textbf{53.2}\\ 51.6\end{array}$	46.7 63.7 65.0 66.2 84.3 73.9 72.3 84.3 87.9 89.2	32.7 35.4 46.0 52.2 49.9 56.8 62.0 78.2 76.1 75.2	57.3 64.6 78.5 58.3 91.2 71.4 59.9 73.3 76.9 74.7	44.1 66.3 69.6 65.1 89.4 70.8 71.7 83.1 89.0 87.6	96.6 93.2 97.6 98.4 99.4 98.7 98.3 99.7 99.2 100	56.6 63.6 71.4 65.0 80.5 69.6 71.5 81.7 83.4 83.2
Unmodeled Occlusion																			
[7] [8] [3] [2] [6] [1] [4] [5] Ours	80.0 83.9 89.3 86.2 91.3 92.5 90.8 96.5 98.0	$\begin{array}{r} 42.7\\ 38.1\\ 43.3\\ 46.3\\ 56.7\\ 51.5\\ 51.7\\ 66.8\\ 65.8\end{array}$	91.8 92.4 92.2 97.8 97.8 99.2 95.9 99.0 98.7	73.5 81.5 83.1 87.5 82.0 90.7 88.5 95.8 95.2	76.1 81.3 84.1 86.5 92.8 92.1 88.0 95.0 97.0	 81.7 85.5 79.0 86.3 89.9 92.2 90.5 95.9 92.0 	89.8 97.5 94.5 95.7 96.6 97.7 96.9 100 99.8	82.6 88.9 88.6 90.7 92.2 94.2 91.6 97.6 98.8	68.7 76.1 76.2 78.8 71.8 89.8 87.1 92.2 92.6	86.7 87.5 90.4 96.5 97.0 98.4 90.3 96.6 97.3		67.0 72.7 80.7 80.6 84.6 90.7 85.6 94.4 95.7	46.6 52.5 61.6 59.9 66.9 66.3 65.8 79.0 83.8	64.0 77.2 75.3 86.8 87.7 91.7 87.0 94.7 95.2	$\begin{array}{r} 43.6\\ 53.9\\ 53.1\\ 69.6\\ 56.1\\ 75.3\\ 72.7\\ 89.8\\ 89.3 \end{array}$	88.8 88.5 91.1 93.3 95.1 95.9 91.2 95.7 96.0	68.6 79.3 81.9 81.8 92.1 84.0 93.6 94.3	86.2 92.5 93.4 95.8 98.2 99.0 97.0 99.6 99.6	73.3 78.4 80.3 83.6 85.1 88.9 85.6 93.2 93.7

Table 1. Comparisons on the standard RBOT benchmark, with frame step S = 1.

4 More Results and Analysis

Figure 4 provides more visual results with very large frame step (S = 8), in which case it is nearly impossible to be successfully tracked with local methods such as SRT3D and *Ours*⁻. Our non-local method performs much better. Note that in order to achieve real-time speed, we adopted a very simple color-based segmentation method. As can be seen from Figure 4, the probability map is very inaccurate for some objects, introducing many resistant local minimum. Without a doubt, if more computing power (e.g. GPU) is available, better segmentation

					£			~		á.		â			~~		e	je,	
	$^{4}D_{0}$	500 600	:3°	South	Can	Car	ð	200	20°	Drille Drille	200	2987	Cino Cino	2007	Caro	The second	2000	South	¥.
	Ŷ	~		Ý	-	-	-	-	-	,	*	,	-	,	-	,	,	Ý	,
$\alpha=0.125$	99.8	65.6	99.5	95.0	96.6	92.6	100	98.7	95.0	97.1	97.4	96.1	83.3	96.9	91.5	95.8	95.2	99.7	94.2
$\alpha = 0.25$	99.7	65.2	99.6	95.1	96.3	92.1	100	98.7	94.8	97.1	97.2	96.1	83.5	97.1	91.7	95.9	94.8	99.9	94.2
$\alpha = 0.5$	99.3	65.6	99.4	95.2	95.5	90.1	100	98.4	94.6	96.5	97.3	94.2	82.8	97.0	90.7	95.2	92.7	99.8	93.6
$\alpha = 0.75$	99.3	65.1	99.9	95.6	92.8	86.9	100	98.3	94.3	94.8	97.5	91.1	76.9	96.1	90.1	91.6	88.3	99.9	92.1
$\alpha = 1.0$	99.3	63.8	98.9	95.5	86.1	79.7	99.9	98.3	93.9	87.6	97.2	85.1	67.9	94.7	87.3	83.4	80.0	99.5	88.8
$\alpha = 1.5$	98.9	61.2	84.7	95.1	54.2	42.2	99.3	96.5	87.7	56.2	97.4	48.7	38.4	85.6	61.8	46.0	48.4	99.1	72.3

Table 2. Tracking success rates with different α on the RBOT dataset [7].



Fig. 3. The interactive initialization process in real scenes.

and larger search range can be achieved in real time, then the accuracy of our method can be further improved.

The ablation study in the paper shows that using smaller α is helpful for achieving higher accuracy. Table 2 contains details of each object. We can find that as the increase of α , the decrease of accuracy for each object is very different. For the objects that are distinctive from the background (*ape,cat,squirrel,duck* etc.), the accuracy is only slightly decreased. On the other hand, for the objects that are indistinctive (*camera,can,glue,lamp*, etc.), the accuracy would decrease dramatically. Therefore, the robust estimation is especially important for handling errors in segmentation.

5 Real-scene Examples

We tested our method in real scenes for further evaluation. The input video is captured by an ordinary web camera (Logitech C270), with 640×480 image resolution and 25 fps frame rate. The initial pose is set interactively by aligning the current frame with a standard pose. As shown in Figure 3, the user needs only to specify a rough initial pose, then in the second frame our method can converge to the correct pose. The test scenes contain typical challenges such as fast motion, background clutter, occlusion and out of view, etc. Thanks to the advantage in handling large displacements, our method can perform significant better than previous methods for the case of fast motion. Please see the accompany video for the results.



Fig. 4. Some visual examples of large displacements (S = 8).

References

- Huang, H., Zhong, F., Qin, X.: Pixel-Wise Weighted Region-Based 3D Object Tracking using Contour Constraints. IEEE Transactions on Visualization and Computer Graphics pp. 1–1 (2021). https://doi.org/10.1109/TVCG.2021.3085197
- Huang, H., Zhong, F., Sun, Y., Qin, X.: An Occlusion-aware Edge-Based Method for Monocular 3D Object Tracking using Edge Confidence. Computer Graphics Forum 39(7), 399–409 (Oct 2020). https://doi.org/10.1111/cgf.14154
- Jiachen Li, Fan Zhong, S.X.X.Q.: 3d object tracking with adaptively weighted local bundles. Journal of Computer Science and Technology 36(3), 555 (2021). https://doi.org/10.1007/s11390-021-1272-5, https://jcst.ict.ac.cn/EN/abstract/article_2756.shtml
- 4. Stoiber, M., Pfanne, M., Strobl, K.H., Triebel, R., Albu-Schaeffer, A.: A sparse gaussian approach to region-based 6dof object tracking. In: Proceedings of the Asian Conference on Computer Vision (2020)
- 5. Stoiber, M., Pfanne, M., Strobl, K.H., Triebel, R., Albu-Schäffer, A.: Srt3d: A sparse region-based 3d object tracking approach for the real world (2021)
- Sun, X., Zhou, J., Zhang, W., Wang, Z., Yu, Q.: Robust monocular pose tracking of less-distinct objects based on contour-part model. IEEE Transactions on Circuits and Systems for Video Technology **31**(11), 4409–4421 (2021). https://doi.org/10.1109/TCSVT.2021.3053696
- Tjaden, H., Schwanecke, U., Schomer, E., Cremers, D.: A Region-Based Gauss-Newton Approach to Real-Time Monocular Multiple Object Tracking. IEEE Transactions on Pattern Analysis and Machine Intelligence 41(8), 1797–1812 (Aug 2019). https://doi.org/10.1109/TPAMI.2018.2884990
- Zhong, L., Zhao, X., Zhang, Y., Zhang, S., Zhang, L.: Occlusion-Aware Region-Based 3D Pose Tracking of Objects With Temporally Consistent Polar-Based Local Partitioning. IEEE Transactions on Image Processing 29, 5065–5078 (2020). https://doi.org/10.1109/TIP.2020.2973512